

Digital Signal Processing: A User's Guide

Collection Editor:

Douglas L. Jones

Digital Signal Processing: A User's Guide

Collection Editor:

Douglas L. Jones

Authors:

Richard Baraniuk
Benjamin Fite
Anders Gjendemsjø
Michael Haag
Don Johnson
Douglas L. Jones

Robert Nowak
Ricardo Radaelli-Sanchez
Justin Romberg
Clayton Scott
Ivan Selesnick
Melissa Selik

Online:

<<http://cnx.org/content/col10372/1.2/>>

C O N N E X I O N S

Rice University, Houston, Texas

©2008 Douglas L. Jones

This selection and arrangement of content is licensed under the Creative Commons Attribution License:
<http://creativecommons.org/licenses/by/2.0/>

Table of Contents

Preface for Digital Signal Processing: A User's Guide	1
1 Background, Review, and Reference	
1.1 Discrete-Time Signals and Systems	3
1.2 Systems in the Time-Domain	5
1.3 Discrete-Time Convolution	6
1.4 Review of Linear Algebra	12
1.5 Orthonormal Basis Expansions	22
1.6 Fourier Analysis	26
1.7 Continuous-Time Fourier Transform (CTFT)	27
1.8 Discrete-Time Fourier Transform (DTFT)	29
1.9 DFT as a Matrix Operation	34
1.10 Sampling theory	36
1.11 Z-Transform	64
1.12 Random Signals and Processes	78
Solutions	96
2 The DFT, FFT, and Practical Spectral Analysis	
2.1 The Discrete Fourier Transform	99
2.2 Spectrum Analysis	102
2.3 Fast Fourier Transform Algorithms	141
2.4 Fast Convolution	176
2.5 Chirp-z Transform	181
2.6 FFTs of prime length and Rader's conversion	183
2.7 Choosing the Best FFT Algorithm	187
Solutions	189
3 Digital Filter Design	
3.1 Overview of Digital Filter Design	191
3.2 FIR Filter Design	192
3.3 IIR Filter Design	207
Solutions	223
4 Digital Filter Structures and Quantization Error Analysis	
4.1 Filter Structures	225
4.2 Fixed-Point Numbers	239
4.3 Quantization Error Analysis	243
4.4 Overflow Problems and Solutions	255
Solutions	259
5 Adaptive Filters and Applications	
5.1 Introduction to Adaptive Filters	261
5.2 Wiener Filter Algorithm	261
5.3 The LMS Adaptive Filter Algorithm	267
5.4 Applications of Adaptive Filters	275
5.5 Other Adaptive Filter Algorithms	283
5.6 Summary of Adaptive Filtering Methods	287
Solutions	288
6 Multirate Signal Processing	
6.1 Overview of Multirate Signal Processing	289

6.2	Interpolation, Decimation, and Rate Changing by Integer Fractions	291
6.3	Efficient Multirate Filter Structures	295
6.4	Filter Design for Multirate Systems	299
6.5	Multistage Multirate Systems	302
6.6	DFT-Based Filterbanks	305
6.7	Quadrature Mirror Filterbanks (QMF)	306
6.8	M-Channel Filter Banks	310
	Solutions	312
	Glossary	313
	Bibliography	315
	Index	317
	Attributions	322

Preface for Digital Signal Processing: A User's Guide¹

Digital signal processing (DSP) has matured in the past few decades from an obscure research discipline to a large body of practical methods with very broad application. Both practicing engineers and students specializing in signal processing need a clear exposition of the ideas and methods comprising the core signal processing "toolkit" so widely used today.

This text reflects my belief that the skilled practitioner must understand the key ideas underlying the algorithms to select, apply, debug, extend, and innovate most effectively; only with real insight can the engineer make novel use of these methods in the seemingly infinite range of new problems and applications. It also reflects my belief that the needs of the typical student and the practicing engineer have converged in recent years; as the discipline of signal processing has matured, these core topics have become less a subject of active research and more a set of tools applied in the course of other research. The modern student thus has less need for exhaustive coverage of the research literature and detailed derivations and proofs as preparation for their own research on these topics, but greater need for intuition and practical guidance in their most effective use. The majority of students eventually become practicing engineers themselves and benefit from the best preparation for their future careers.

This text both explains the principles of classical signal processing methods and describes how they are used in engineering practice. It is thus much more than a recipe book; it describes the ideas behind the algorithms, gives analyses when they enhance that understanding, and includes derivations that the practitioner may need to extend when applying these methods to new situations. Analyses or derivations that are only of research interest or that do not increase intuitive understanding are left to the references. It is also much more than a theory book; it contains more description of common applications, discussion of actual implementation issues, comments on what really works in the real world, and practical "know-how" than found in the typical academic textbook. The choice of material emphasizes those methods that have found widespread practical use; techniques that have been the subject of intense research but which are rarely used in practice (for example, RLS adaptive filter algorithms) often receive only limited coverage.

The text assumes a familiarity with basic signal processing concepts such as ideal sampling theory, continuous and discrete Fourier transforms, convolution and filtering. It evolved from a set of notes for a second signal processing course, ECE 451: Digital Signal Processing II, in Electrical and Computer Engineering at the University of Illinois at Urbana-Champaign, aimed at second-semester seniors or first-semester graduate students in signal processing. Over the years, it has been enhanced substantially to include descriptions of common applications, sometimes hard-won knowledge about what actually works and what doesn't, useful tricks, important extensions known to experienced engineers but rarely discussed in academic texts, and other relevant "know-how" to aid the real-world user. This is necessarily an ongoing process, and I continue to expand and refine this component as my own practical knowledge and experience grows. The topics are the core signal processing methods that are used in the majority of signal processing applications; discrete Fourier analysis and FFTs, digital filter design, adaptive filtering, multirate signal processing, and efficient algorithm implementation and finite-precision issues. While many of these topics are covered at an intro-

¹This content is available online at <http://cnx.org/content/m13782/1.1/>.

ductory level in a first course, this text aspires to cover all of the methods, both basic and advanced, in these areas which see widespread use in practice. I have also attempted to make the individual modules and sections somewhat self-sufficient, so that those who seek specific information on a single topic can quickly find what they need. Hopefully these aspirations will eventually be achieved; in the meantime, I welcome your comments, corrections, and feedback so that I can continue to improve this text.

As of August 2006, the majority of modules are unedited transcriptions of handwritten notes and may contain typographical errors and insufficient descriptive text for documents unaccompanied by an oral lecture; I hope to have all of the modules in at least presentable shape by the end of the year.

Publication of this text in Connexions would have been impossible without the help of many people. A huge thanks to the various permanent and temporary staff at Connexions is due, in particular to those who converted the text and equations from my original handwritten notes into CNXML and MathML. My former and current faculty colleagues at the University of Illinois who have taught the second DSP course over the years have had a substantial influence on the evolution of the content, as have the students who have inspired this work and given me feedback. I am very grateful to my teachers, mentors, colleagues, collaborators, and fellow engineers who have taught me the art and practice of signal processing; this work is dedicated to you.

Chapter 1

Background, Review, and Reference

1.1 Discrete-Time Signals and Systems¹

Mathematically, analog signals are functions having as their independent variables continuous quantities, such as space and time. Discrete-time signals are functions defined on the integers; they are sequences. As with analog signals, we seek ways of decomposing discrete-time signals into simpler components. Because this approach leading to a better understanding of signal structure, we can exploit that structure to represent information (create ways of representing information with signals) and to extract information (retrieve the information thus represented). For symbolic-valued signals, the approach is different: We develop a common representation of all symbolic-valued signals so that we can embody the information they contain in a unified way. From an information representation perspective, the most important issue becomes, for both real-valued and symbolic-valued signals, efficiency: what is the most parsimonious and compact way to represent information so that it can be extracted later.

1.1.1 Real- and Complex-valued Signals

A discrete-time signal is represented symbolically as $s(n)$, where $n = \{\dots, -1, 0, 1, \dots\}$.

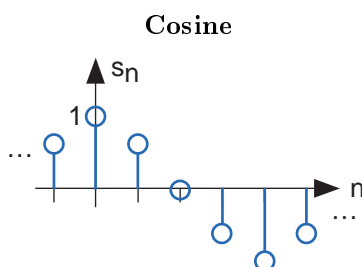


Figure 1.1: The discrete-time cosine signal is plotted as a stem plot. Can you find the formula for this signal?

We usually draw discrete-time signals as stem plots to emphasize the fact they are functions defined only on the integers. We can delay a discrete-time signal by an integer just as with analog ones. A signal delayed by m samples has the expression $s(n - m)$.

¹This content is available online at <http://cnx.org/content/m10342/2.13/>.

1.1.2 Complex Exponentials

The most important signal is, of course, the **complex exponential sequence**.

$$s(n) = e^{j2\pi fn} \quad (1.1)$$

Note that the frequency variable f is dimensionless and that adding an integer to the frequency of the discrete-time complex exponential has no effect on the signal's value.

$$\begin{aligned} e^{j2\pi(f+m)n} &= e^{j2\pi fn} e^{j2\pi mn} \\ &= e^{j2\pi fn} \end{aligned} \quad (1.2)$$

This derivation follows because the complex exponential evaluated at an integer multiple of 2π equals one. Thus, the period of a discrete-time complex exponential equals one.

1.1.3 Sinusoids

Discrete-time sinusoids have the obvious form $s(n) = A \cos(2\pi fn + \phi)$. As opposed to analog complex exponentials and sinusoids that can have their frequencies be any real value, frequencies of their discrete-time counterparts yield unique waveforms *only* when f lies in the interval $(-\frac{1}{2}, \frac{1}{2}]$. From the properties of the complex exponential, the sinusoid's period is always one; this choice of frequency interval will become evident later.

1.1.4 Unit Sample

The second-most important discrete-time signal is the **unit sample**, which is defined to be

$$\delta(n) = \begin{cases} 1 & \text{if } n = 0 \\ 0 & \text{otherwise} \end{cases} \quad (1.3)$$

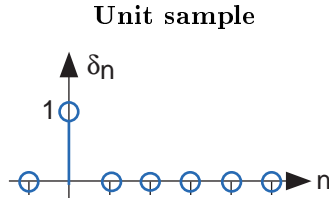


Figure 1.2: The unit sample.

Examination of a discrete-time signal's plot, like that of the cosine signal shown in Figure 1.1 (Cosine), reveals that all signals consist of a sequence of delayed and scaled unit samples. Because the value of a sequence at each integer m is denoted by $s(m)$ and the unit sample delayed to occur at m is written $\delta(n - m)$, we can decompose *any* signal as a sum of unit samples delayed to the appropriate location and scaled by the signal value.

$$s(n) = \sum_{m=-\infty}^{\infty} (s(m) \delta(n - m)) \quad (1.4)$$

This kind of decomposition is unique to discrete-time signals, and will prove useful subsequently.

1.1.5 Unit Step

The **unit sample** in discrete-time is well-defined at the origin, as opposed to the situation with analog signals.

$$u(n) = \begin{cases} 1 & \text{if } n \geq 0 \\ 0 & \text{if } n < 0 \end{cases} \quad (1.5)$$

1.1.6 Symbolic Signals

An interesting aspect of discrete-time signals is that their values do not need to be real numbers. We do have real-valued discrete-time signals like the sinusoid, but we also have signals that denote the sequence of characters typed on the keyboard. Such characters certainly aren't real numbers, and as a collection of possible signal values, they have little mathematical structure other than that they are members of a set. More formally, each element of the *symbolic-valued* signal $s(n)$ takes on one of the values $\{a_1, \dots, a_K\}$ which comprise the **alphabet** A . This technical terminology does not mean we restrict symbols to being members of the English or Greek alphabet. They could represent keyboard characters, bytes (8-bit quantities), integers that convey daily temperature. Whether controlled by software or not, discrete-time systems are ultimately constructed from digital circuits, which consist *entirely* of analog circuit elements. Furthermore, the transmission and reception of discrete-time signals, like e-mail, is accomplished with analog signals and systems. Understanding how discrete-time and analog signals and systems intertwine is perhaps the main goal of this course.

1.1.7 Discrete-Time Systems

Discrete-time systems can act on discrete-time signals in ways similar to those found in analog signals and systems. Because of the role of software in discrete-time systems, many more different systems can be envisioned and "constructed" with programs than can be with analog signals. In fact, a special class of analog signals can be converted into discrete-time signals, processed with software, and converted back into an analog signal, all without the incursion of error. For such signals, systems can be easily produced in software, with equivalent analog realizations difficult, if not impossible, to design.

1.2 Systems in the Time-Domain²

A discrete-time signal $s(n)$ is **delayed** by n_0 samples when we write $s(n - n_0)$, with $n_0 > 0$. Choosing n_0 to be negative advances the signal along the integers. As opposed to analog delays³, discrete-time delays can *only* be integer valued. In the frequency domain, delaying a signal corresponds to a linear phase shift of the signal's discrete-time Fourier transform: $(s(n - n_0) \leftrightarrow e^{-j2\pi f n_0})S(e^{j2\pi f})$.

Linear discrete-time systems have the superposition property.

Superposition

$$S(a_1 x_1(n) + a_2 x_2(n)) = a_1 S(x_1(n)) + a_2 S(x_2(n)) \quad (1.6)$$

A discrete-time system is called **shift-invariant** (analogous to time-invariant analog systems) if delaying the input delays the corresponding output.

Shift-Invariant

$$\text{If } S(x(n)) = y(n), \text{ Then } S(x(n - n_0)) = y(n - n_0) \quad (1.7)$$

We use the term shift-invariant to emphasize that delays can only have integer values in discrete-time, while in analog signals, delays can be arbitrarily valued.

²This content is available online at <http://cnx.org/content/m0508/2.7/>.

³"Simple Systems": Section Delay <http://cnx.org/content/m0006/latest/#delay>

We want to concentrate on systems that are both linear and shift-invariant. It will be these that allow us the full power of frequency-domain analysis and implementations. Because we have no physical constraints in "constructing" such systems, we need only a mathematical specification. In analog systems, the differential equation specifies the input-output relationship in the time-domain. The corresponding discrete-time specification is the **difference equation**.

The Difference Equation

$$y(n) = a_1 y(n-1) + \cdots + a_p y(n-p) + b_0 x(n) + b_1 x(n-1) + \cdots + b_q x(n-q) \quad (1.8)$$

Here, the output signal $y(n)$ is related to its *past* values $y(n-l)$, $l = \{1, \dots, p\}$, and to the current and past values of the input signal $x(n)$. The system's characteristics are determined by the choices for the number of coefficients p and q and the coefficients' values $\{a_1, \dots, a_p\}$ and $\{b_0, b_1, \dots, b_q\}$.

ASIDE: There is an asymmetry in the coefficients: where is a_0 ? This coefficient would multiply the $y(n)$ term in the difference equation (1.8: The Difference Equation). We have essentially divided the equation by it, which does not change the input-output relationship. We have thus created the convention that a_0 is always one.

As opposed to differential equations, which only provide an *implicit* description of a system (we must somehow solve the differential equation), difference equations provide an *explicit* way of computing the output for any input. We simply express the difference equation by a program that calculates each output from the previous output values, and the current and previous inputs.

1.3 Discrete-Time Convolution⁴

1.3.1 Overview

Convolution is a concept that extends to all systems that are both **linear and time-invariant**⁵ (**LTI**). The idea of **discrete-time convolution** is exactly the same as that of continuous-time convolution⁶. For this reason, it may be useful to look at both versions to help your understanding of this extremely important concept. Recall that convolution is a very powerful tool in determining a system's output from knowledge of an arbitrary input and the system's impulse response. It will also be helpful to see convolution graphically with your own eyes and to play around with it some, so experiment with the applets⁷ available on the internet. These resources will offer different approaches to this crucial concept.

1.3.2 Convolution Sum

As mentioned above, the convolution sum provides a concise, mathematical way to express the output of an LTI system based on an arbitrary discrete-time input signal and the system's response. The **convolution sum** is expressed as

$$y[n] = \sum_{k=-\infty}^{\infty} (x[k] h[n-k]) \quad (1.9)$$

As with continuous-time, convolution is represented by the symbol $*$, and can be written as

$$y[n] = x[n] * h[n] \quad (1.10)$$

By making a simple change of variables into the convolution sum, $k = n - k$, we can easily show that convolution is **commutative**:

$$x[n] * h[n] = h[n] * x[n] \quad (1.11)$$

⁴This content is available online at <<http://cnx.org/content/m10087/2.18/>>.

⁵"System Classifications and Properties" <<http://cnx.org/content/m10084/latest/>>

⁶"Continuous-Time Convolution" <<http://cnx.org/content/m10085/latest/>>

⁷<http://www.jhu.edu/~signals>

For more information on the characteristics of convolution, read about the Properties of Convolution⁸.

1.3.3 Derivation

We know that any discrete-time signal can be represented by a summation of scaled and shifted discrete-time impulses. Since we are assuming the system to be linear and time-invariant, it would seem to reason that an input signal comprised of the sum of scaled and shifted impulses would give rise to an output comprised of a sum of scaled and shifted impulse responses. This is exactly what occurs in **convolution**. Below we present a more rigorous and mathematical look at the derivation:

Letting \mathcal{H} be a DT LTI system, we start with the following equation and work our way down the convolution sum!

$$\begin{aligned}
 y[n] &= \mathcal{H}[x[n]] \\
 &= \mathcal{H}\left[\sum_{k=-\infty}^{\infty} (x[k] \delta[n-k])\right] \\
 &= \sum_{k=-\infty}^{\infty} (\mathcal{H}[x[k] \delta[n-k]]) \\
 &= \sum_{k=-\infty}^{\infty} (x[k] \mathcal{H}[\delta[n-k]]) \\
 &= \sum_{k=-\infty}^{\infty} (x[k] h[n-k])
 \end{aligned} \tag{1.12}$$

Let us take a quick look at the steps taken in the above derivation. After our initial equation, we use the DT sifting property⁹ to rewrite the function, $x[n]$, as a sum of the function times the unit impulse. Next, we can move around the \mathcal{H} operator and the summation because $\mathcal{H}[\cdot]$ is a linear, DT system. Because of this linearity and the fact that $x[k]$ is a constant, we can pull the previous mentioned constant out and simply multiply it by $\mathcal{H}[\cdot]$. Finally, we use the fact that $\mathcal{H}[\cdot]$ is time invariant in order to reach our final state - the convolution sum!

A quick graphical example may help in demonstrating why convolution works.

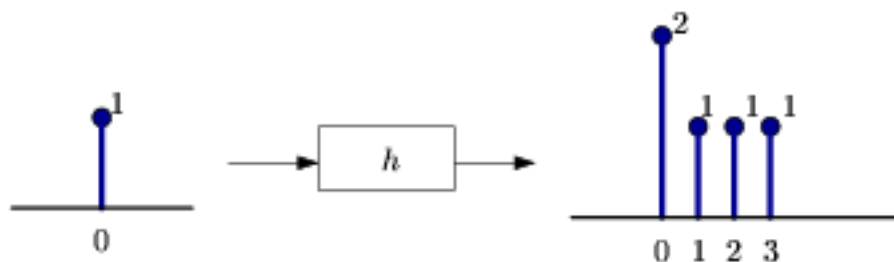


Figure 1.3: A single impulse input yields the system's impulse response.

⁸"Properties of Convolution" <<http://cnx.org/content/m10088/latest/>>

⁹"The Impulse Function": Section The Sifting Property of the Impulse <<http://cnx.org/content/m10059/latest/#sifting>>

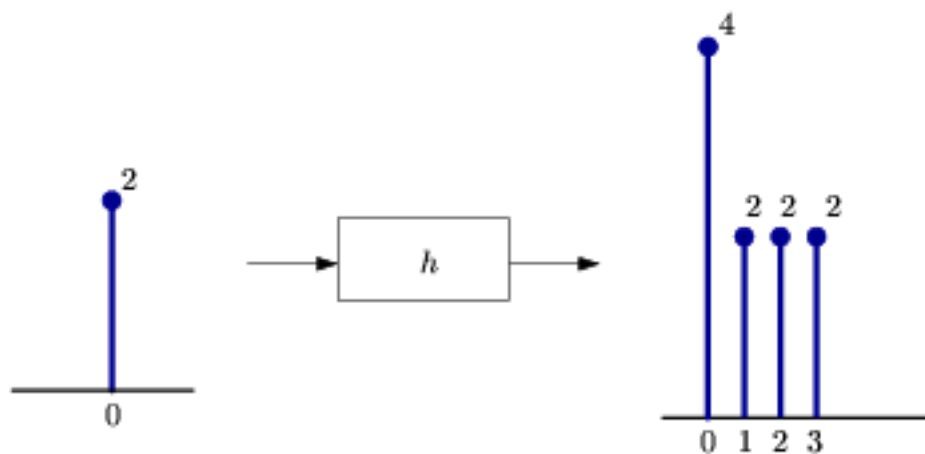


Figure 1.4: A scaled impulse input yields a scaled response, due to the scaling property of the system's linearity.

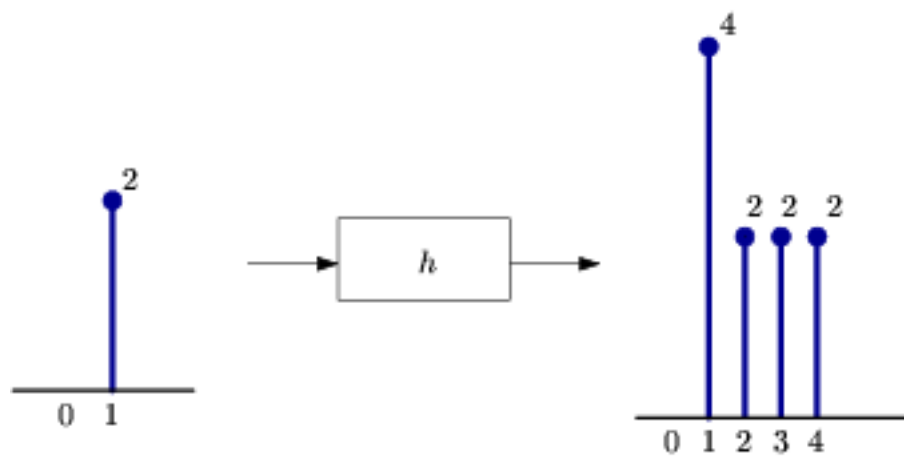


Figure 1.5: We now use the time-invariance property of the system to show that a delayed input results in an output of the same shape, only delayed by the same amount as the input.

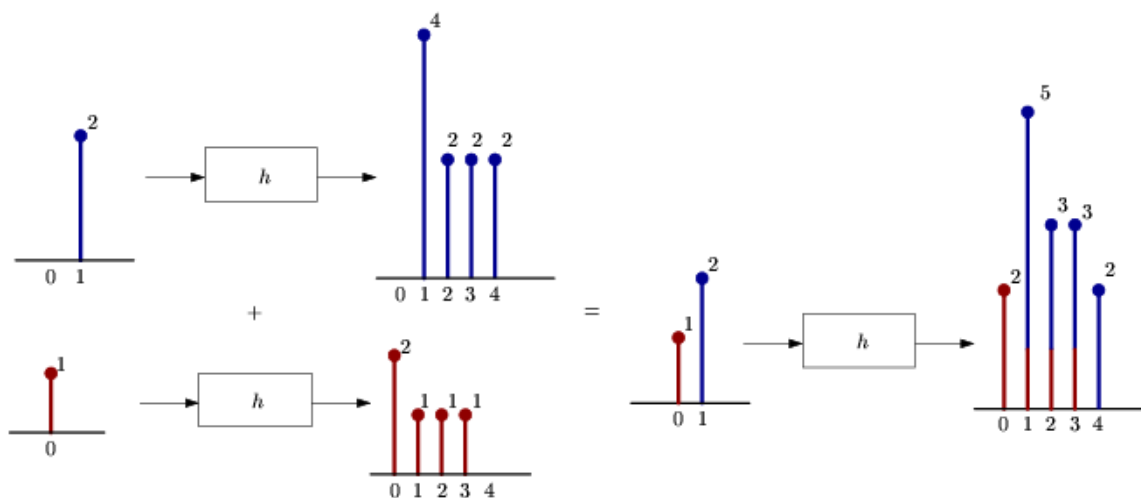


Figure 1.6: We now use the additivity portion of the linearity property of the system to complete the picture. Since any discrete-time signal is just a sum of scaled and shifted discrete-time impulses, we can find the output from knowing the input and the impulse response.

1.3.4 Convolution Through Time (A Graphical Approach)

In this section we will develop a second graphical interpretation of discrete-time convolution. We will begin this by writing the convolution sum allowing x to be a causal, length- m signal and h to be a causal, length- k , LTI system. This gives us the finite summation,

$$y[n] = \sum_{l=0}^{m-1} (x[l] h[n-l]) \quad (1.13)$$

Notice that for any given n we have a sum of the products of x_l and a time-delayed h_{-l} . This is to say that we multiply the terms of x by the terms of a time-reversed h and add them up.

Going back to the previous example:

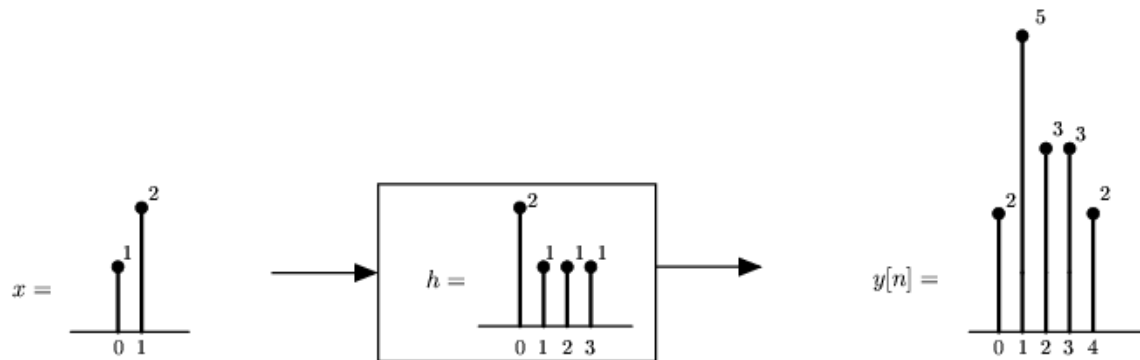


Figure 1.7: This is the end result that we are looking to find.

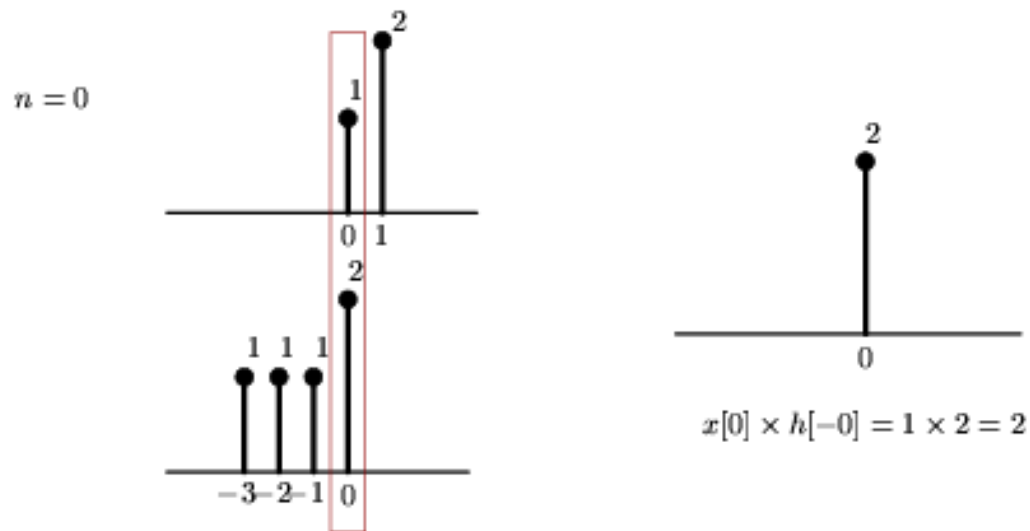


Figure 1.8: Here we reverse the impulse response, h , and begin its traverse at time 0.

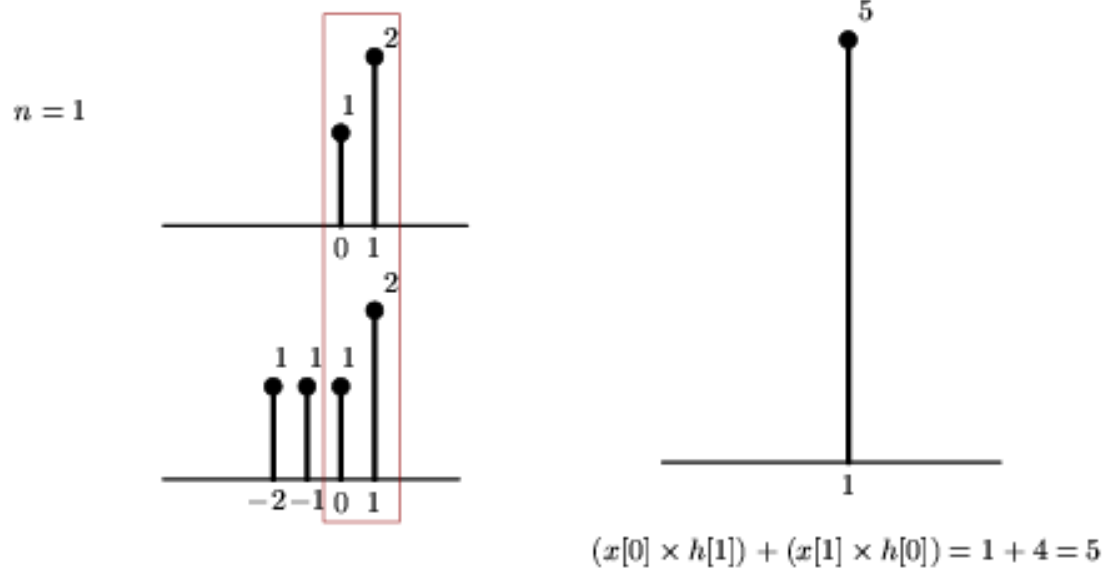


Figure 1.9: We continue the traverse. See that at time 1, we are multiplying two elements of the input signal by two elements of the impulse response.

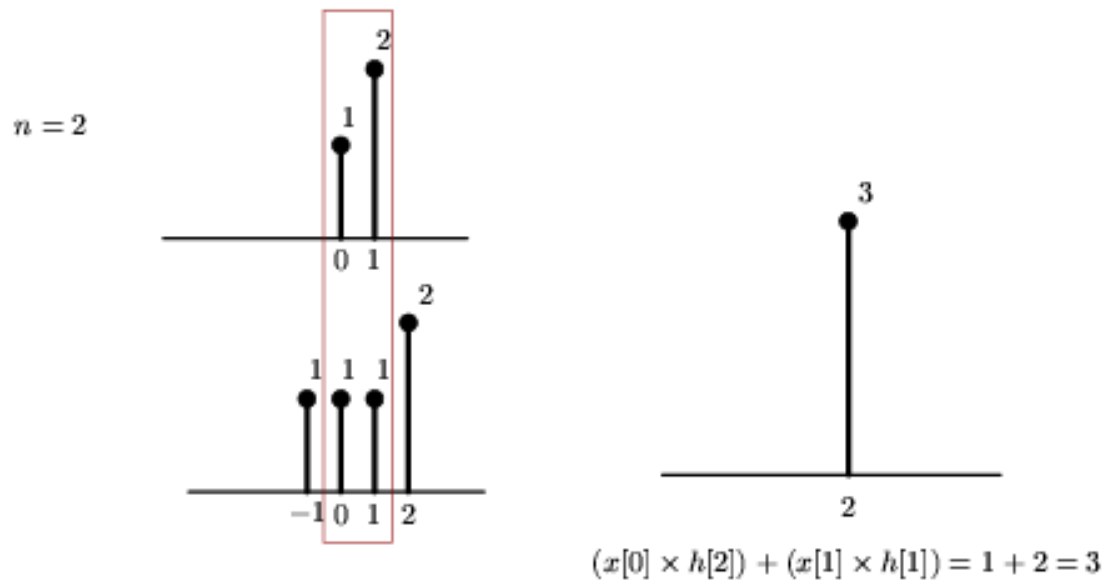


Figure 1.10

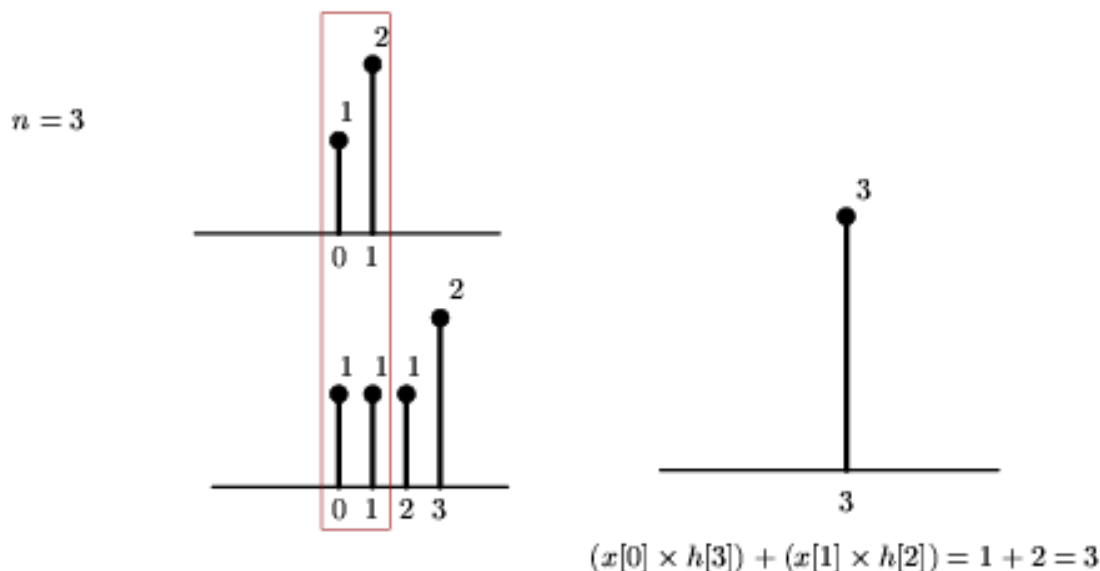


Figure 1.11: If we follow this through to one more step, $n = 4$, then we can see that we produce the same output as we saw in the initial example.

What we are doing in the above demonstration is reversing the impulse response in time and "walking it across" the input signal. Clearly, this yields the same result as scaling, shifting and summing impulse responses.

This approach of time-reversing, and sliding across is a common approach to presenting convolution, since it demonstrates how convolution builds up an output through time.

1.4 Review of Linear Algebra¹⁰

Vector spaces are the principal object of study in linear algebra. A vector space is always defined with respect to a field of scalars.

1.4.1 Fields

A field is a set F equipped with two operations, addition and multiplication, and containing two special members 0 and 1 ($0 \neq 1$), such that for all $\{a, b, c\} \in F$

1. (a) $a + b \in F$
 (b) $a + b = b + a$
 (c) $(a + b) + c = a + (b + c)$
 (d) $a + 0 = a$
 (e) there exists $-a$ such that $a + (-a) = 0$
2. (a) $ab \in F$
 (b) $ab = ba$
 (c) $(ab)c = a(bc)$

¹⁰This content is available online at <<http://cnx.org/content/m11948/1.2/>>.

- (d) $a \cdot 1 = a$
- (e) there exists a^{-1} such that $aa^{-1} = 1$
- 3. $a(b + c) = ab + ac$

More concisely

1. F is an abelian group under addition
2. F is an abelian group under multiplication
3. multiplication distributes over addition

1.4.1.1 Examples

Q, R, \mathbb{C}

1.4.2 Vector Spaces

Let F be a field, and V a set. We say V is a *vector space over F* if there exist two operations, defined for all $a \in F$, $\mathbf{u} \in V$ and $\mathbf{v} \in V$:

- vector addition: $(\mathbf{u}, \mathbf{v}) \rightarrow \mathbf{u} + \mathbf{v} \in V$
- scalar multiplication: $(a, \mathbf{v}) \rightarrow a\mathbf{v} \in V$

and if there exists an element denoted $\mathbf{0} \in V$, such that the following hold for all $a \in F$, $b \in F$, and $\mathbf{u} \in V$, $\mathbf{v} \in V$, and $\mathbf{w} \in V$

1. (a) $\mathbf{u} + (\mathbf{v} + \mathbf{w}) = (\mathbf{u} + \mathbf{v}) + \mathbf{w}$
 (b) $\mathbf{u} + \mathbf{v} = \mathbf{v} + \mathbf{u}$
 (c) $\mathbf{u} + \mathbf{0} = \mathbf{u}$
 (d) there exists $-\mathbf{u}$ such that $\mathbf{u} + (-\mathbf{u}) = \mathbf{0}$
2. (a) $a(\mathbf{u} + \mathbf{v}) = a\mathbf{u} + a\mathbf{v}$
 (b) $(a + b)\mathbf{u} = a\mathbf{u} + b\mathbf{u}$
 (c) $(ab)\mathbf{u} = a(b\mathbf{u})$
 (d) $1 \cdot \mathbf{u} = \mathbf{u}$

More concisely,

1. V is an abelian group under plus
2. Natural properties of scalar multiplication

1.4.2.1 Examples

- \mathbb{R}^N is a vector space over R
- \mathbb{C}^N is a vector space over \mathbb{C}
- \mathbb{C}^N is a vector space over R
- \mathbb{R}^N is *not* a vector space over \mathbb{C}

The elements of V are called **vectors**.

1.4.3 Euclidean Space

Throughout this course we will think of a signal as a vector

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{pmatrix} = \begin{pmatrix} x_1 & x_2 & \dots & x_N \end{pmatrix}^T$$

The samples $\{x_i\}$ could be samples from a finite duration, continuous time signal, for example.

A signal will belong to one of two vector spaces:

1.4.3.1 Real Euclidean space

$\mathbf{x} \in \mathbb{R}^N$ (over \mathbb{R})

1.4.3.2 Complex Euclidean space

$\mathbf{x} \in \mathbb{C}^N$ (over \mathbb{C})

1.4.4 Subspaces

Let V be a vector space over F .

A subset $S \subseteq V$ is called a **subspace** of V if S is a vector space over F in its own right.

Example 1.1

$V = \mathbb{R}^2$, $F = \mathbb{R}$, $S =$ any line through the origin.

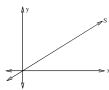


Figure 1.12: S is any line through the origin.

Are there other subspaces?

Theorem 1.1:

$S \subseteq V$ is a subspace if and only if for all $a \in F$ and $b \in F$ and for all $\mathbf{s} \in S$ and $\mathbf{t} \in S$, $a\mathbf{s} + b\mathbf{t} \in S$

1.4.5 Linear Independence

Let $\mathbf{u}_1, \dots, \mathbf{u}_k \in V$.

We say that these vectors are **linearly dependent** if there exist scalars $a_1, \dots, a_k \in F$ such that

$$\sum_{i=1}^k (a_i \mathbf{u}_i) = \mathbf{0} \quad (1.14)$$

and at least one $a_i \neq 0$.

If (1.14) only holds for the case $a_1 = \dots = a_k = 0$, we say that the vectors are **linearly independent**.

Example 1.2

$$1 \begin{pmatrix} 1 \\ -1 \\ 2 \end{pmatrix} - 2 \begin{pmatrix} -2 \\ 3 \\ 0 \end{pmatrix} + 1 \begin{pmatrix} -5 \\ 7 \\ -2 \end{pmatrix} = \mathbf{0}$$

so these vectors are linearly dependent in \mathbb{R}^3 .

1.4.6 Spanning Sets

Consider the subset $S = \{v_1, v_2, \dots, v_k\}$. Define the **span** of S

$$\langle S \rangle \equiv \text{span}(S) \equiv \left\{ \sum_{i=1}^k (a_i v_i) \mid a_i \in F \right\}$$

Fact: $\langle S \rangle$ is a subspace of V .

Example 1.3

$$V = \mathbb{R}^3, F = \mathbb{R}, S = \{v_1, v_2\}, v_1 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, v_2 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \Rightarrow \langle S \rangle = \text{xy-plane}.$$

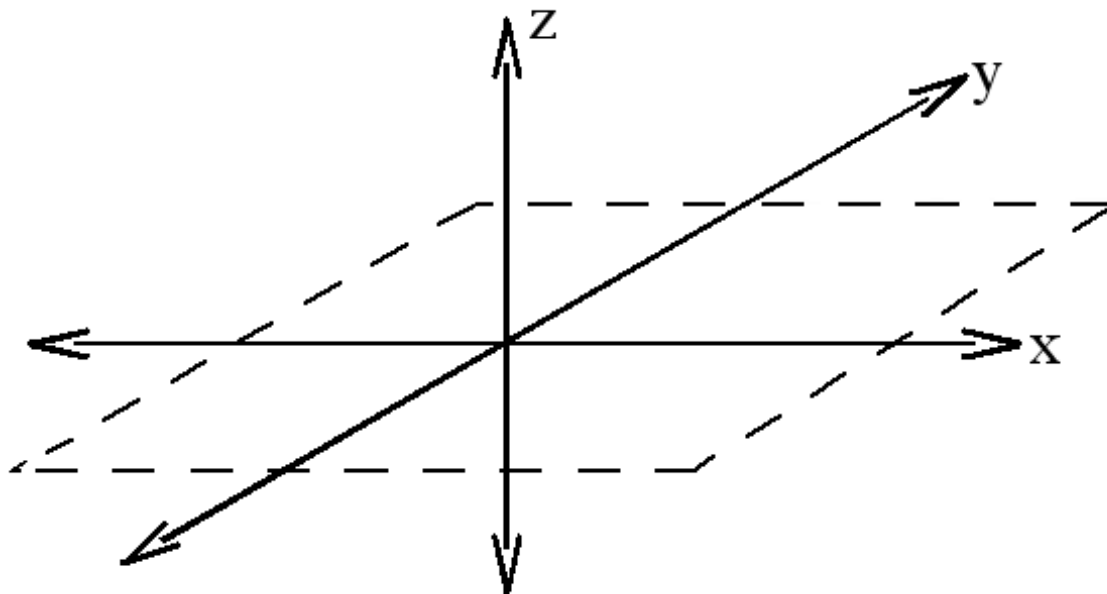


Figure 1.13: $\langle S \rangle$ is the xy-plane.

1.4.6.1 Aside

If S is infinite, the notions of linear independence and span are easily generalized:

We say S is linearly independent if, for every finite collection $u_1, \dots, u_k \in S$, (k arbitrary) we have

$$\sum_{i=1}^k (a_i u_i) = \mathbf{0} \Rightarrow \forall i : (a_i = 0)$$

The span of S is

$$\langle S \rangle = \left\{ \sum_{i=1}^k (a_i u_i) \mid a_i \in F \wedge u_i \in S \wedge k < \infty \right\}$$

NOTE: In both definitions, we only consider **finite** sums.

1.4.7 Bases

A set $B \subseteq V$ is called a **basis** for V over F if and only if

1. B is linearly independent
2. $\langle B \rangle = V$

Bases are of fundamental importance in signal processing. They allow us to decompose a signal into building blocks (basis vectors) that are often more easily understood.

Example 1.4

$V =$ (real or complex) Euclidean space, \mathbb{R}^N or \mathbb{C}^N .

$$B = \{e_1, \dots, e_N\} \equiv \text{standard basis}$$

$$e_i = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix}$$

where the 1 is in the i^{th} position.

Example 1.5

$V = \mathbb{C}^N$ over \mathbb{C} .

$$B = \{u_1, \dots, u_N\}$$

which is the DFT basis.

$$u_k = \begin{pmatrix} 1 \\ e^{-j2\pi \frac{k}{N}} \\ \vdots \\ e^{-j2\pi \frac{k}{N}(N-1)} \end{pmatrix}$$

where $j = \sqrt{-1}$.

1.4.7.1 Key Fact

If B is a basis for V , then every $\mathbf{v} \in V$ can be written uniquely (up to order of terms) in the form

$$\mathbf{v} = \sum_{i=1}^N (a_i v_i)$$

where $a_i \in F$ and $v_i \in B$.

1.4.7.2 Other Facts

- If S is a linearly independent set, then S can be extended to a basis.
- If $\langle S \rangle = V$, then S contains a basis.

1.4.8 Dimension

Let V be a vector space with basis B . The dimension of V , denoted $\dim(V)$, is the cardinality of B .

Theorem 1.2:

Every vector space has a basis.

Theorem 1.3:

Every basis for a vector space has the same cardinality.

$\Rightarrow \dim(V)$ is **well-defined**.

If $\dim(V) < \infty$, we say V is **finite dimensional**.

1.4.8.1 Examples

vector space	field of scalars	dimension
\mathbb{R}^N	\mathbb{R}	
\mathbb{C}^N	\mathbb{C}	
\mathbb{C}^N	\mathbb{R}	

Every subspace is a vector space, and therefore has its own dimension.

Example 1.6

Suppose $S = \{u_1, \dots, u_k\} \subseteq V$ is a linearly independent set. Then

$$\dim(\langle S \rangle) =$$

Facts

- If S is a subspace of V , then $\dim(S) \leq \dim(V)$.
- If $\dim(S) = \dim(V) < \infty$, then $S = V$.

1.4.9 Direct Sums

Let V be a vector space, and let $S \subseteq V$ and $T \subseteq V$ be subspaces.

We say V is the **direct sum** of S and T , written $V = (S \oplus T)$, if and only if for every $\mathbf{v} \in V$, there exist unique $\mathbf{s} \in S$ and $\mathbf{t} \in T$ such that $\mathbf{v} = \mathbf{s} + \mathbf{t}$.

If $V = (S \oplus T)$, then T is called a **complement** of S .

Example 1.7

$$V = C' = \{f : \mathbb{R} \rightarrow \mathbb{R} | f \text{ is continuous}\}$$

$$S = \text{even functions in } C'$$

$$T = \text{odd functions in } C'$$

$$f(t) = \frac{1}{2}(f(t) + f(-t)) + \frac{1}{2}(f(t) - f(-t))$$

If $f = g + h = g' + h'$, $g \in S$ and $g' \in S$, $h \in T$ and $h' \in T$, then $g - g' = h' - h$ is odd and even, which implies $g = g'$ and $h = h'$.

1.4.9.1 Facts

1. Every subspace has a complement
2. $V = (S \oplus T)$ if and only if
 - (a) $S \cap T = \{\mathbf{0}\}$
 - (b) $\langle S, T \rangle = V$
3. If $V = (S \oplus T)$, and $\dim(V) < \infty$, then $\dim(V) = \dim(S) + \dim(T)$

1.4.9.2 Proofs

Invoke a basis.

1.4.10 Norms

Let V be a vector space over F . A norm is a mapping $(V \rightarrow F)$, denoted by $\|\cdot\|$, such that for all $\mathbf{u} \in V$, $\mathbf{v} \in V$, and $\lambda \in F$

1. $\|\mathbf{u}\| > 0$ if $\mathbf{u} \neq \mathbf{0}$
2. $\|\lambda\mathbf{u}\| = |\lambda| \|\mathbf{u}\|$
3. $\|\mathbf{u} + \mathbf{v}\| \leq \|\mathbf{u}\| + \|\mathbf{v}\|$

1.4.10.1 Examples

Euclidean norms:

$$\mathbf{x} \in \mathbb{R}^N:$$

$$\|\mathbf{x}\| = \left(\sum_{i=1}^N (x_i^2) \right)^{\frac{1}{2}}$$

$$\mathbf{x} \in \mathbb{C}^N:$$

$$\|\mathbf{x}\| = \left(\sum_{i=1}^N (|x_i|^2) \right)^{\frac{1}{2}}$$

1.4.10.2 Induced Metric

Every norm induces a metric on V

$$d(\mathbf{u}, \mathbf{v}) \equiv \|\mathbf{u} - \mathbf{v}\|$$

which leads to a notion of "distance" between vectors.

1.4.11 Inner products

Let V be a vector space over F , $F = \mathbb{R}$ or \mathbb{C} . An inner product is a mapping $V \times V \rightarrow F$, denoted $\langle \cdot, \cdot \rangle$, such that

1. $\langle \mathbf{v}, \mathbf{v} \rangle \geq 0$, and $(\langle \mathbf{v}, \mathbf{v} \rangle = 0 \Leftrightarrow \mathbf{v} = \mathbf{0})$
2. $\langle \mathbf{u}, \mathbf{v} \rangle = \overline{\langle \mathbf{v}, \mathbf{u} \rangle}$
3. $\langle a\mathbf{u} + b\mathbf{v}, \mathbf{w} \rangle = a \langle \mathbf{u}, \mathbf{w} \rangle + b \langle \mathbf{v}, \mathbf{w} \rangle$

1.4.11.1 Examples

\mathbb{R}^N over \mathbb{R} :

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \mathbf{y} = \sum_{i=1}^N (x_i y_i)$$

\mathbb{C}^N over \mathbb{C} :

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^H \mathbf{y} = \sum_{i=1}^N (\overline{x_i} y_i)$$

If $\mathbf{x} = (x_1, \dots, x_N)^T \in \mathbb{C}$, then

$$\mathbf{x}^H \equiv \begin{pmatrix} \overline{x_1} \\ \vdots \\ \overline{x_N} \end{pmatrix}^T$$

is called the "Hermitian," or "conjugate transpose" of \mathbf{x} .

1.4.12 Triangle Inequality

If we define $\|\mathbf{u}\| = \sqrt{\langle \mathbf{u}, \mathbf{u} \rangle}$, then

$$\|\mathbf{u} + \mathbf{v}\| \leq \|\mathbf{u}\| + \|\mathbf{v}\|$$

Hence, every inner product induces a norm.

1.4.13 Cauchy-Schwarz Inequality

For all $\mathbf{u} \in V$, $\mathbf{v} \in V$,

$$|\langle \mathbf{u}, \mathbf{v} \rangle| \leq \|\mathbf{u}\| \|\mathbf{v}\|$$

In inner product spaces, we have a notion of the angle between two vectors:

$$\angle(\mathbf{u}, \mathbf{v}) = \arccos\left(\frac{\langle \mathbf{u}, \mathbf{v} \rangle}{\|\mathbf{u}\| \|\mathbf{v}\|}\right) \in [0, 2\pi)$$

1.4.14 Orthogonality

\mathbf{u} and \mathbf{v} are **orthogonal** if

$$\langle \mathbf{u}, \mathbf{v} \rangle = 0$$

Notation: $(\mathbf{u} \perp \mathbf{v})$.

If in addition $\|\mathbf{u}\| = \|\mathbf{v}\| = 1$, we say \mathbf{u} and \mathbf{v} are **orthonormal**.

In an orthogonal (orthonormal) *set*, each pair of vectors is orthogonal (orthonormal).



Figure 1.14: Orthogonal vectors in \mathbb{R}^2 .

1.4.15 Orthonormal Bases

An Orthonormal basis is a basis $\{v_i\}$ such that

$$\langle v_i, v_j \rangle = \delta_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}$$

Example 1.8

The standard basis for \mathbb{R}^N or \mathbb{C}^N

Example 1.9

The normalized DFT basis

$$u_k = \frac{1}{\sqrt{N}} \begin{pmatrix} 1 \\ e^{-j2\pi \frac{k}{N}} \\ \vdots \\ e^{-j2\pi \frac{k}{N} (N-1)} \end{pmatrix}$$

1.4.16 Expansion Coefficients

If the representation of \mathbf{v} with respect to $\{v_i\}$ is

$$\mathbf{v} = \sum (a_i v_i)$$

then

$$a_i = \langle v_i, \mathbf{v} \rangle$$

1.4.17 Gram-Schmidt

Every inner product space has an orthonormal basis. Any (countable) basis can be made orthogonal by the Gram-Schmidt orthogonalization process.

1.4.18 Orthogonal Compliments

Let $S \subseteq V$ be a subspace. The **orthogonal compliment** S^\perp is

$$S^\perp = \{ \mathbf{u} \mid \mathbf{u} \in V \wedge \langle \mathbf{u}, \mathbf{v} \rangle = 0 \wedge \forall \mathbf{v} : (\mathbf{v} \in S) \}$$

S^\perp is easily seen to be a subspace.

If $\dim(V) < \infty$, then $V = (S \oplus S^\perp)$.

ASIDE: If $\dim(V) = \infty$, then in order to have $V = (S \oplus S^\perp)$ we require V to be a **Hilbert Space**.

1.4.19 Linear Transformations

Loosely speaking, a linear transformation is a mapping from one vector space to another that preserves vector space operations.

More precisely, let V, W be vector spaces over the same field F . A **linear transformation** is a mapping $T : V \rightarrow W$ such that

$$T(a\mathbf{u} + b\mathbf{v}) = aT(\mathbf{u}) + bT(\mathbf{v})$$

for all $a \in F, b \in F$ and $\mathbf{u} \in V, \mathbf{v} \in V$.

In this class we will be concerned with linear transformations between (real or complex) *Euclidean spaces*, or subspaces thereof.

1.4.20 Image

$$\text{image}(T) = \{\mathbf{w} \mid \mathbf{w} \in W \wedge T(\mathbf{v}) = \mathbf{w} \text{ for some } \mathbf{v}\}$$

1.4.21 Nullspace

Also known as the kernel:

$$\ker(T) = \{\mathbf{v} \mid \mathbf{v} \in V \wedge T(\mathbf{v}) = \mathbf{0}\}$$

Both the image and the nullspace are easily seen to be subspaces.

1.4.22 Rank

$$\text{rank}(T) = \dim(\text{image}(T))$$

1.4.23 Nullity

$$\text{null}(T) = \dim(\ker(T))$$

1.4.24 Rank plus nullity theorem

$$\text{rank}(T) + \text{null}(T) = \dim(V)$$

1.4.25 Matrices

Every linear transformation T has a **matrix representation**. If $T : E^N \rightarrow E^M$, $E = \mathbb{R}$ or \mathbb{C} , then T is represented by an $M \times N$ matrix

$$A = \begin{pmatrix} a_{11} & \dots & a_{1N} \\ \vdots & \ddots & \vdots \\ a_{M1} & \dots & a_{MN} \end{pmatrix}$$

where $(a_{1i}, \dots, a_{Mi})^T = T(e_i)$ and $e_i = (0, \dots, 1, \dots, 0)^T$ is the i^{th} **standard basis** vector.

ASIDE: A linear transformation can be represented with respect to any bases of E^N and E^M , leading to a different A . We will always represent a linear transformation using the standard bases.

1.4.26 Column span

$$\text{colspan}(A) = \langle A \rangle = \text{image}(A)$$

1.4.27 Duality

If $A : \mathbb{R}^N \rightarrow \mathbb{R}^M$, then

$$\ker^\perp(A) = \text{image}(A^T)$$

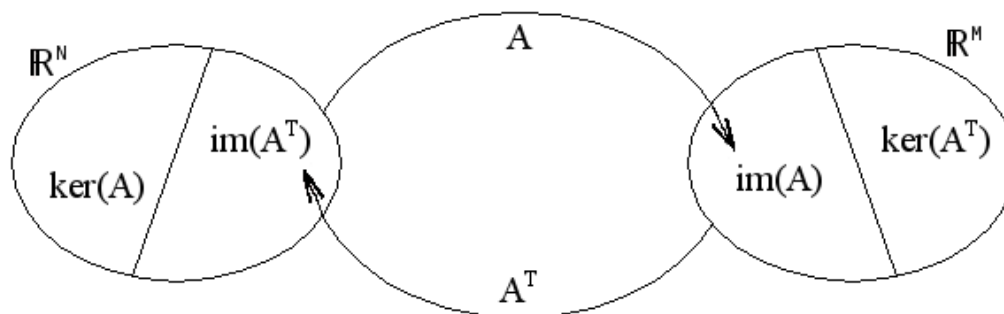


Figure 1.15

If $A : \mathbb{C}^N \rightarrow \mathbb{C}^M$, then

$$\ker^\perp(A) = \text{image}(A^H)$$

1.4.28 Inverses

The linear transformation/matrix A is **invertible** if and only if there exists a matrix B such that $AB = BA = I$ (identity).

Only *square* matrices can be invertible.

Theorem 1.4:

Let $A : F^N \rightarrow F^N$ be linear, $F = \mathbb{R}$ or \mathbb{C} . The following are equivalent:

1. A is invertible (nonsingular)
2. $\text{rank}(A) = N$
3. $\text{null}(A) = 0$
4. $\det A \neq 0$
5. The columns of A form a basis.

If $A^{-1} = A^T$ (or A^H in the complex case), we say A is **orthogonal** (or **unitary**).

1.5 Orthonormal Basis Expansions¹¹

1.5.1 Main Idea

When working with signals many times it is helpful to break up a signal into smaller, more manageable parts. Hopefully by now you have been exposed to the concept of eigenvectors¹² and there use in decomposing a

¹¹This content is available online at <<http://cnx.org/content/m10760/2.4/>>.

¹²"Eigenvectors and Eigenvalues" <<http://cnx.org/content/m10736/latest/>>

signal into one of its possible basis. By doing this we are able to simplify our calculations of signals and systems through eigenfunctions of LTI systems¹³.

Now we would like to look at an alternative way to represent signals, through the use of **orthonormal basis**. We can think of orthonormal basis as a set of building blocks we use to construct functions. We will build up the signal/vector as a weighted sum of basis elements.

Example 1.10

The complex sinusoids $\frac{1}{\sqrt{T}}e^{j\omega_0 nt}$ for all $-\infty < n < \infty$ form an orthonormal basis for $L^2([0, T])$.

In our Fourier series¹⁴ equation, $f(t) = \sum_{n=-\infty}^{\infty} (c_n e^{j\omega_0 nt})$, the $\{c_n\}$ are just another representation of $f(t)$.

NOTE: For signals/vectors in a Hilbert Space, the expansion coefficients are easy to find.

1.5.2 Alternate Representation

Recall our definition of a **basis**: A set of vectors $\{b_i\}$ in a vector space S is a basis if

1. The b_i are linearly independent.
2. The b_i span¹⁵ S . That is, we can find $\{\alpha_i\}$, where $\alpha_i \in \mathbb{C}$ (scalars) such that

$$\forall x, x \in S : \left(\mathbf{x} = \sum_i (\alpha_i b_i) \right) \quad (1.15)$$

where \mathbf{x} is a vector in S , α is a scalar in \mathbb{C} , and \mathbf{b} is a vector in S .

Condition 2 in the above definition says we can **decompose** any vector in terms of the $\{b_i\}$. Condition 1 ensures that the decomposition is **unique** (think about this at home).

NOTE: The $\{\alpha_i\}$ provide an alternate representation of \mathbf{x} .

Example 1.11

Let us look at simple example in \mathbb{R}^2 , where we have the following vector:

$$\mathbf{x} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

$$\text{Standard Basis: } \{e_0, e_1\} = \left\{ (1, 0)^T, (0, 1)^T \right\}$$

$$\mathbf{x} = e_0 + 2e_1$$

$$\text{Alternate Basis: } \{h_0, h_1\} = \left\{ (1, 1)^T, (1, -1)^T \right\}$$

$$\mathbf{x} = \frac{3}{2}h_0 + \frac{-1}{2}h_1$$

In general, given a basis $\{b_0, b_1\}$ and a vector $\mathbf{x} \in \mathbb{R}^2$, how do we find the α_0 and α_1 such that

$$\mathbf{x} = \alpha_0 b_0 + \alpha_1 b_1 \quad (1.16)$$

¹³"Eigenfunctions of LTI Systems" <<http://cnx.org/content/m10500/latest/>>

¹⁴"Fourier Series: Eigenfunction Approach" <<http://cnx.org/content/m10496/latest/>>

¹⁵"Linear Algebra: The Basics": Section Span <http://cnx.org/content/m10734/latest/#span_sec>

1.5.3 Finding the Alphas

Now let us address the question posed above about finding α_i 's in general for \mathbb{R}^2 . We start by rewriting (1.16) so that we can stack our b_i 's as columns in a 2×2 matrix.

$$\begin{pmatrix} \mathbf{x} \end{pmatrix} = \alpha_0 \begin{pmatrix} b_0 \end{pmatrix} + \alpha_1 \begin{pmatrix} b_1 \end{pmatrix} \quad (1.17)$$

$$\begin{pmatrix} \mathbf{x} \end{pmatrix} = \begin{pmatrix} \vdots & \vdots \\ b_0 & b_1 \\ \vdots & \vdots \end{pmatrix} \begin{pmatrix} \alpha_0 \\ \alpha_1 \end{pmatrix} \quad (1.18)$$

Example 1.12

Here is a simple example, which shows a little more detail about the above equations.

$$\begin{aligned} \begin{pmatrix} x[0] \\ x[1] \end{pmatrix} &= \alpha_0 \begin{pmatrix} b_0[0] \\ b_0[1] \end{pmatrix} + \alpha_1 \begin{pmatrix} b_1[0] \\ b_1[1] \end{pmatrix} \\ &= \begin{pmatrix} \alpha_0 b_0[0] + \alpha_1 b_1[0] \\ \alpha_0 b_0[1] + \alpha_1 b_1[1] \end{pmatrix} \end{aligned} \quad (1.19)$$

$$\begin{pmatrix} x[0] \\ x[1] \end{pmatrix} = \begin{pmatrix} b_0[0] & b_1[0] \\ b_0[1] & b_1[1] \end{pmatrix} \begin{pmatrix} \alpha_0 \\ \alpha_1 \end{pmatrix} \quad (1.20)$$

1.5.3.1 Simplifying our Equation

To make notation simpler, we define the following two items from the above equations:

- **Basis Matrix:**

$$B = \begin{pmatrix} \vdots & \vdots \\ b_0 & b_1 \\ \vdots & \vdots \end{pmatrix}$$

- **Coefficient Vector:**

$$\alpha = \begin{pmatrix} \alpha_0 \\ \alpha_1 \end{pmatrix}$$

This gives us the following, concise equation:

$$\mathbf{x} = B\alpha \quad (1.21)$$

which is equivalent to $\mathbf{x} = \sum_{i=0}^1 (\alpha_i b_i)$.

Example 1.13

Given a standard basis, $\left\{ \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix} \right\}$, then we have the following basis matrix:

$$B = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

To get the α_i 's, we solve for the coefficient vector in (1.21)

$$\alpha = B^{-1}\mathbf{x} \quad (1.22)$$

Where B^{-1} is the inverse matrix¹⁶ of B .

1.5.3.2 Examples

Example 1.14

Let us look at the standard basis first and try to calculate α from it.

$$B = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = I$$

Where I is the **identity matrix**. In order to solve for α let us find the inverse of B first (which is obviously very trivial in this case):

$$B^{-1} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

Therefore we get,

$$\alpha = B^{-1}\mathbf{x} = \mathbf{x}$$

Example 1.15

Let us look at a ever-so-slightly more complicated basis of $\left\{ \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ -1 \end{pmatrix} \right\} = \{h_0, h_1\}$ Then our basis matrix and inverse basis matrix becomes:

$$B = \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$$

$$B^{-1} = \begin{pmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{-1}{2} \end{pmatrix}$$

and for this example it is given that

$$\mathbf{x} = \begin{pmatrix} 3 \\ 2 \end{pmatrix}$$

Now we solve for α

$$\alpha = B^{-1}\mathbf{x} = \begin{pmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{-1}{2} \end{pmatrix} \begin{pmatrix} 3 \\ 2 \end{pmatrix} = \begin{pmatrix} 2.5 \\ 0.5 \end{pmatrix}$$

and we get

$$\mathbf{x} = 2.5h_0 + 0.5h_1$$

Exercise 1.1

(Solution on p. 96.)

Now we are given the following basis matrix and \mathbf{x} :

$$\{b_0, b_1\} = \left\{ \begin{pmatrix} 1 \\ 2 \end{pmatrix}, \begin{pmatrix} 3 \\ 0 \end{pmatrix} \right\}$$

¹⁶"Matrix Inversion" <<http://cnx.org/content/m2113/latest/>>

$$\mathbf{x} = \begin{pmatrix} 3 \\ 2 \end{pmatrix}$$

For this problem, make a sketch of the bases and then represent \mathbf{x} in terms of b_0 and b_1 .

NOTE: A change of basis simply looks at \mathbf{x} from a "different perspective." B^{-1} **transforms** \mathbf{x} from the standard basis to our new basis, $\{b_0, b_1\}$. Notice that this is a totally mechanical procedure.

1.5.4 Extending the Dimension and Space

We can also extend all these ideas past just \mathbb{R}^2 and look at them in \mathbb{R}^n and \mathbb{C}^n . This procedure extends naturally to higher (> 2) dimensions. Given a basis $\{b_0, b_1, \dots, b_{n-1}\}$ for \mathbb{R}^n , we want to find $\{\alpha_0, \alpha_1, \dots, \alpha_{n-1}\}$ such that

$$\mathbf{x} = \alpha_0 b_0 + \alpha_1 b_1 + \dots + \alpha_{n-1} b_{n-1} \quad (1.23)$$

Again, we will set up a basis matrix

$$B = \begin{pmatrix} b_0 & b_1 & b_2 & \dots & b_{n-1} \end{pmatrix}$$

where the columns equal the basis vectors and it will always be an $n \times n$ matrix (although the above matrix does not appear to be square since we left terms in vector notation). We can then proceed to rewrite (1.21)

$$\mathbf{x} = \begin{pmatrix} b_0 & b_1 & \dots & b_{n-1} \end{pmatrix} \begin{pmatrix} \alpha_0 \\ \vdots \\ \alpha_{n-1} \end{pmatrix} = B\alpha$$

and

$$\alpha = B^{-1}\mathbf{x}$$

1.6 Fourier Analysis¹⁷

Fourier analysis is fundamental to understanding the behavior of signals and systems. This is a result of the fact that sinusoids are Eigenfunctions¹⁸ of linear, time-invariant (LTI)¹⁹ systems. This is to say that if we pass any particular sinusoid through a LTI system, we get a scaled version of that same sinusoid on the output. Then, since Fourier analysis allows us to redefine the signals in terms of sinusoids, all we need to do is determine how any given system effects all possible sinusoids (its transfer function²⁰) and we have a complete understanding of the system. Furthermore, since we are able to define the passage of sinusoids through a system as multiplication of that sinusoid by the transfer function at the same frequency, we can convert the passage of any signal through a system from convolution²¹ (in time) to multiplication (in frequency). These ideas are what give Fourier analysis its power.

Now, after hopefully having sold you on the value of this method of analysis, we must examine exactly what we mean by Fourier analysis. The four Fourier transforms that comprise this analysis are the

¹⁷This content is available online at <http://cnx.org/content/m10096/2.10/>.

¹⁸"Eigenfunctions of LTI Systems" <http://cnx.org/content/m10500/latest/>

¹⁹"System Classifications and Properties" <http://cnx.org/content/m10084/latest/>

²⁰"Transfer Functions" <http://cnx.org/content/m0028/latest/>

²¹"Properties of Convolution" <http://cnx.org/content/m10088/latest/>

Fourier Series²², Continuous-Time Fourier Transform (Section 1.7), Discrete-Time Fourier Transform²³ and Discrete Fourier Transform²⁴. For this document, we will view the Laplace Transform²⁵ and Z-Transform (Section 1.11.3) as simply extensions of the CTFT and DTFT respectively. All of these transforms act essentially the same way, by converting a signal in time to an equivalent signal in frequency (sinusoids). However, depending on the nature of a specific signal i.e. whether it is finite- or infinite-length and whether it is discrete- or continuous-time) there is an appropriate transform to convert the signal into the frequency domain. Below is a table of the four Fourier transforms and when each is appropriate. It also includes the relevant convolution for the specified space.

Table of Fourier Representations

Transform	Time Domain	Frequency Domain	Convolution
Continuous-Time Fourier Series	$L^2([0, T))$	$l^2(\mathbb{Z})$	Continuous-Time Circular
Continuous-Time Fourier Transform	$L^2(\mathbb{R})$	$L^2(\mathbb{R})$	Continuous-Time Linear
Discrete-Time Fourier Transform	$l^2(\mathbb{Z})$	$L^2([0, 2\pi))$	Discrete-Time Linear
Discrete Fourier Transform	$l^2([0, N - 1])$	$l^2([0, N - 1])$	Discrete-Time Circular

1.7 Continuous-Time Fourier Transform (CTFT)²⁶

1.7.1 Introduction

Due to the large number of continuous-time signals that are present, the Fourier series²⁷ provided us the first glimpse of how we may represent some of these signals in a general manner: as a superposition of a number of sinusoids. Now, we can look at a way to represent continuous-time nonperiodic signals using the same idea of superposition. Below we will present the **Continuous-Time Fourier Transform (CTFT)**, also referred to as just the Fourier Transform (FT). Because the CTFT now deals with nonperiodic signals, we must now find a way to include *all* frequencies in the general equations.

1.7.1.1 Equations

Continuous-Time Fourier Transform

$$\mathcal{F}(\Omega) = \int_{-\infty}^{\infty} f(t) e^{-j\Omega t} dt \quad (1.24)$$

Inverse CTFT

$$f(t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} \mathcal{F}(\Omega) e^{j\Omega t} d\Omega \quad (1.25)$$

²²"Continuous-Time Fourier Series (CTFS)" <<http://cnx.org/content/m10097/latest/>>

²³"Discrete-Time Fourier Transform (DTFT)" <<http://cnx.org/content/m10108/latest/>>

²⁴"Discrete Fourier Transform(DTFT)" <<http://cnx.org/content/m0502/latest/>>

²⁵"The Laplace Transforms" <<http://cnx.org/content/m10110/latest/>>

²⁶This content is available online at <<http://cnx.org/content/m10098/2.9/>>.

²⁷"Fourier Series" <<http://cnx.org/content/m0039/latest/>>

WARNING: Do not be confused by notation - it is not uncommon to see the above formula written slightly different. One of the most common differences among many professors is the way that the exponential is written. Above we used the radial frequency variable Ω in the exponential, where $\Omega = 2\pi f$, but one will often see professors include the more explicit expression, $j2\pi ft$, in the exponential. Click here²⁸ for an overview of the notation used in Connexion's DSP modules.

The above equations for the CTFT and its inverse come directly from the Fourier series and our understanding of its coefficients. For the CTFT we simply utilize integration rather than summation to be able to express the aperiodic signals. This should make sense since for the CTFT we are simply extending the ideas of the Fourier series to include nonperiodic signals, and thus the entire frequency spectrum. Look at the Derivation of the Fourier Transform²⁹ for a more in depth look at this.

1.7.2 Relevant Spaces

The Continuous-Time Fourier Transform maps infinite-length, continuous-time signals in L^2 to infinite-length, continuous-frequency signals in L^2 . Review the Fourier Analysis (Section 1.6) for an overview of all the spaces used in Fourier analysis.

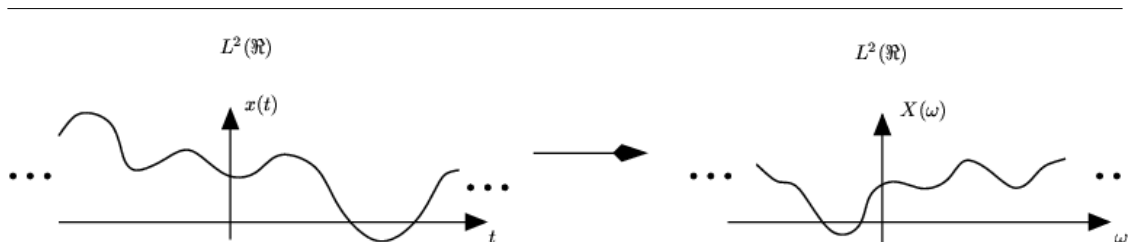


Figure 1.16: Mapping $L^2(\mathbb{R})$ in the time domain to $L^2(\mathbb{R})$ in the frequency domain.

For more information on the characteristics of the CTFT, please look at the module on Properties of the Fourier Transform³⁰.

1.7.3 Example Problems

Exercise 1.2

Find the Fourier Transform (CTFT) of the function

(Solution on p. 96.)

$$f(t) = \begin{cases} e^{-\alpha t} & \text{if } t \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad (1.26)$$

Exercise 1.3

Find the inverse Fourier transform of the square wave defined as

(Solution on p. 96.)

$$X(\Omega) = \begin{cases} 1 & \text{if } |\Omega| \leq M \\ 0 & \text{otherwise} \end{cases} \quad (1.27)$$

²⁸"DSP Notation" <<http://cnx.org/content/m10161/latest/>>

²⁹"Derivation of the Fourier Transform" <<http://cnx.org/content/m0046/latest/>>

³⁰"Properties of the Fourier Transform" <<http://cnx.org/content/m10100/latest/>>

1.8 Discrete-Time Fourier Transform (DTFT)³¹

The Fourier transform of the discrete-time signal $s(n)$ is defined to be

$$S(e^{j2\pi f}) = \sum_{n=-\infty}^{\infty} \left(s(n) e^{-(j2\pi f n)} \right) \quad (1.28)$$

Frequency here has no units. As should be expected, this definition is linear, with the transform of a sum of signals equaling the sum of their transforms. Real-valued signals have conjugate-symmetric spectra: $S(e^{-j2\pi f}) = S^*(e^{j2\pi f})$.

Exercise 1.4

(Solution on p. 96.)

A special property of the discrete-time Fourier transform is that it is periodic with period one: $S(e^{j2\pi(f+1)}) = S(e^{j2\pi f})$. Derive this property from the definition of the DTFT.

Because of this periodicity, we need only plot the spectrum over one period to understand completely the spectrum's structure; typically, we plot the spectrum over the frequency range $[-\frac{1}{2}, \frac{1}{2}]$. When the signal is real-valued, we can further simplify our plotting chores by showing the spectrum only over $[0, \frac{1}{2}]$; the spectrum at negative frequencies can be derived from positive-frequency spectral values.

When we obtain the discrete-time signal via sampling an analog signal, the Nyquist frequency³² corresponds to the discrete-time frequency $\frac{1}{2}$. To show this, note that a sinusoid having a frequency equal to the Nyquist frequency $\frac{1}{2T_s}$ has a sampled waveform that equals

$$\cos\left(2\pi \frac{1}{2T_s} n T_s\right) = \cos(\pi n) = (-1)^n$$

The exponential in the DTFT at frequency $\frac{1}{2}$ equals $e^{-j\frac{1}{2}2\pi n} = e^{-j\pi n} = (-1)^n$, meaning that discrete-time frequency equals analog frequency multiplied by the sampling interval

$$f_D = f_A T_s \quad (1.29)$$

f_D and f_A represent discrete-time and analog frequency variables, respectively. The aliasing figure³³ provides another way of deriving this result. As the duration of each pulse in the periodic sampling signal $p_{T_s}(t)$ narrows, the amplitudes of the signal's spectral repetitions, which are governed by the Fourier series coefficients³⁴ of $p_{T_s}(t)$, become increasingly equal. Examination of the periodic pulse signal³⁵ reveals that as Δ decreases, the value of c_0 , the largest Fourier coefficient, decreases to zero: $|c_0| = \frac{A\Delta}{T_s}$. Thus, to maintain a mathematically viable Sampling Theorem, the amplitude A must increase as $\frac{1}{\Delta}$, becoming infinitely large as the pulse duration decreases. Practical systems use a small value of Δ , say $0.1 \cdot T_s$ and use amplifiers to rescale the signal. Thus, the sampled signal's spectrum becomes periodic with period $\frac{1}{T_s}$. Thus, the Nyquist frequency $\frac{1}{2T_s}$ corresponds to the frequency $\frac{1}{2}$.

³¹This content is available online at <<http://cnx.org/content/m10247/2.28/>>.

³²"The Sampling Theorem" <<http://cnx.org/content/m0050/latest/#para1>>

³³"The Sampling Theorem", Figure 2: aliasing <<http://cnx.org/content/m0050/latest/#alias>>

³⁴"Definition of the Complex Fourier Series", (9) <<http://cnx.org/content/m0042/latest/#eqn2>>

³⁵"Definition of the Complex Fourier Series", Figure 1 <<http://cnx.org/content/m0042/latest/#pps>>

Example 1.16

Let's compute the discrete-time Fourier transform of the exponentially decaying sequence $s(n) = a^n u(n)$, where $u(n)$ is the unit-step sequence. Simply plugging the signal's expression into the Fourier transform formula,

$$\begin{aligned} S(e^{j2\pi f}) &= \sum_{n=-\infty}^{\infty} (a^n u(n) e^{-j2\pi f n}) \\ &= \sum_{n=0}^{\infty} \left((ae^{-j2\pi f})^n \right) \end{aligned} \quad (1.30)$$

This sum is a special case of the **geometric series**.

$$\sum_{n=0}^{\infty} (\alpha^n) = \forall \alpha, |\alpha| < 1 : \left(\frac{1}{1 - \alpha} \right) \quad (1.31)$$

Thus, as long as $|a| < 1$, we have our Fourier transform.

$$S(e^{j2\pi f}) = \frac{1}{1 - ae^{-j2\pi f}} \quad (1.32)$$

Using Euler's relation, we can express the magnitude and phase of this spectrum.

$$|S(e^{j2\pi f})| = \frac{1}{\sqrt{(1 - a \cos(2\pi f))^2 + a^2 \sin^2(2\pi f)}} \quad (1.33)$$

$$\angle(S(e^{j2\pi f})) = - \left(\tan^{-1} \left(\frac{a \sin(2\pi f)}{1 - a \cos(2\pi f)} \right) \right) \quad (1.34)$$

No matter what value of a we choose, the above formulae clearly demonstrate the periodic nature of the spectra of discrete-time signals. Figure 1.17 (Spectrum of exponential signal) shows indeed that the spectrum is a periodic function. We need only consider the spectrum between $-\left(\frac{1}{2}\right)$ and $\frac{1}{2}$ to unambiguously define it. When $a > 0$, we have a lowpass spectrum—the spectrum diminishes as frequency increases from 0 to $\frac{1}{2}$ —with increasing a leading to a greater low frequency content; for $a < 0$, we have a highpass spectrum (Figure 1.18 (Spectra of exponential signals)).

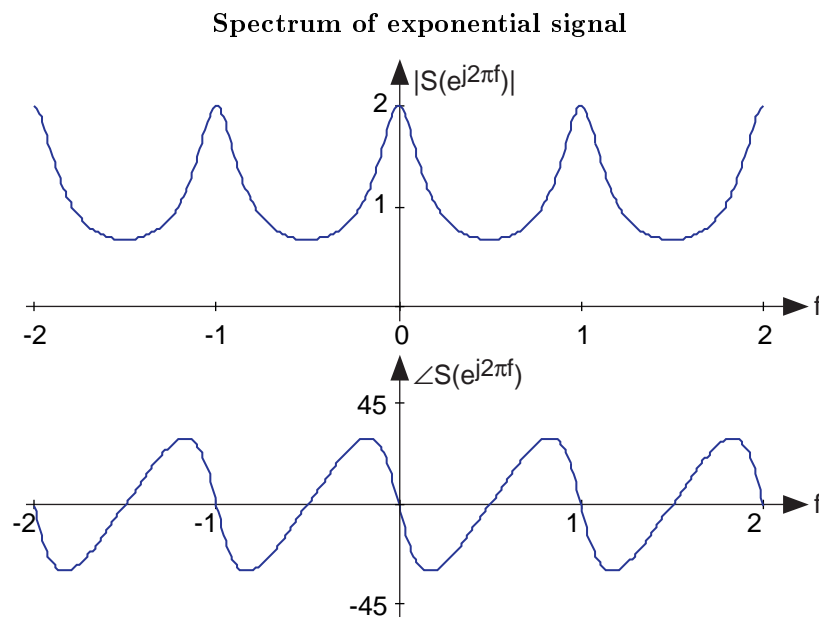


Figure 1.17: The spectrum of the exponential signal ($a = 0.5$) is shown over the frequency range $[-2, 2]$, clearly demonstrating the periodicity of all discrete-time spectra. The angle has units of degrees.

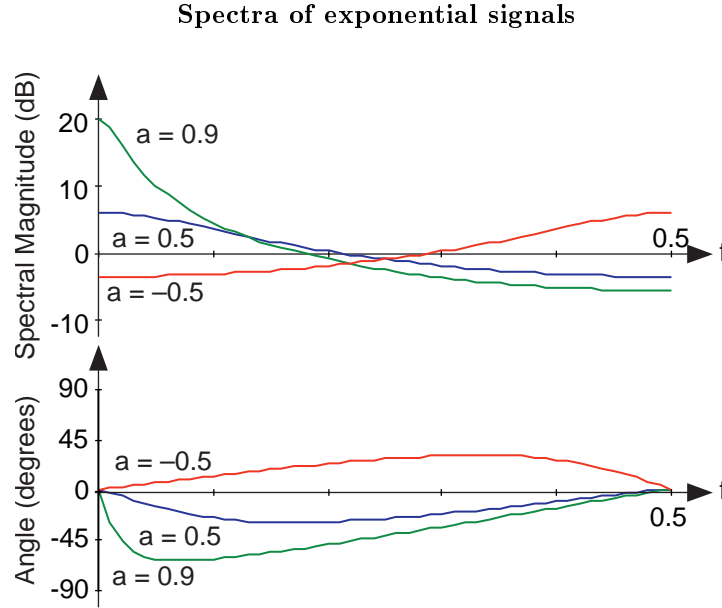


Figure 1.18: The spectra of several exponential signals are shown. What is the apparent relationship between the spectra for $a = 0.5$ and $a = -0.5$?

Example 1.17

Analogous to the analog pulse signal, let's find the spectrum of the length- N pulse sequence.

$$s(n) = \begin{cases} 1 & \text{if } 0 \leq n \leq N-1 \\ 0 & \text{otherwise} \end{cases} \quad (1.35)$$

The Fourier transform of this sequence has the form of a truncated geometric series.

$$S(e^{j2\pi f}) = \sum_{n=0}^{N-1} (e^{(j2\pi f)n}) \quad (1.36)$$

For the so-called finite geometric series, we know that

$$\sum_{n=n_0}^{N+n_0-1} (\alpha^n) = \alpha^{n_0} \frac{1 - \alpha^N}{1 - \alpha} \quad (1.37)$$

for all values of α .

Exercise 1.5

(Solution on p. 97.)

Derive this formula for the finite geometric series sum. The "trick" is to consider the difference between the series' sum and the sum of the series multiplied by α .

Applying this result yields (Figure 1.19 (Spectrum of length-ten pulse).)

$$\begin{aligned} S(e^{j2\pi f}) &= \frac{1 - e^{-(j2\pi f)N}}{1 - e^{-(j2\pi f)}} \\ &= e^{-(j\pi f(N-1))} \frac{\sin(\pi f N)}{\sin(\pi f)} \end{aligned} \quad (1.38)$$

The ratio of sine functions has the generic form of $\frac{\sin(Nx)}{\sin(x)}$, which is known as the **discrete-time sinc function** $dsinc(x)$. Thus, our transform can be concisely expressed as $S(e^{j2\pi f}) = e^{-(j\pi f(N-1))} dsinc(\pi f)$. The discrete-time pulse's spectrum contains many ripples, the number of which increase with N , the pulse's duration.

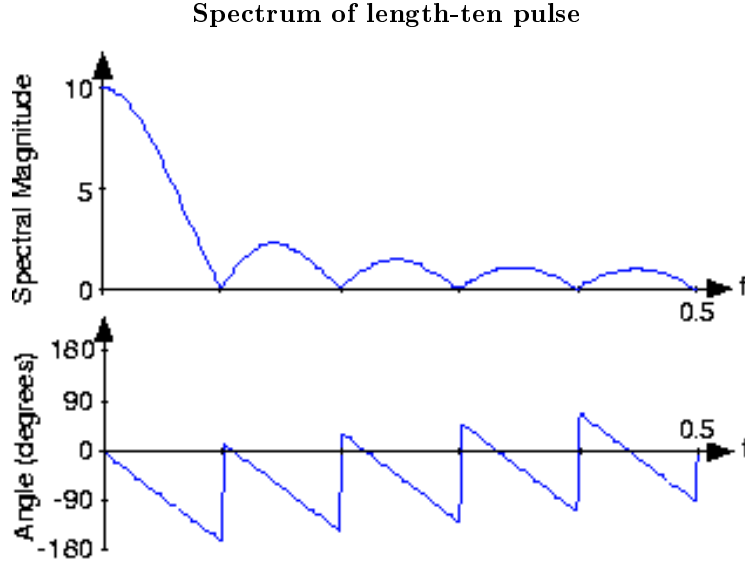


Figure 1.19: The spectrum of a length-ten pulse is shown. Can you explain the rather complicated appearance of the phase?

The inverse discrete-time Fourier transform is easily derived from the following relationship:

$$\int_{-(\frac{1}{2})}^{\frac{1}{2}} e^{-(j2\pi fm)} e^{j2\pi fn} df = \begin{cases} 1 & \text{if } m = n \\ 0 & \text{if } m \neq n \end{cases} \quad (1.39)$$

Therefore, we find that

$$\begin{aligned} \int_{-(\frac{1}{2})}^{\frac{1}{2}} S(e^{j2\pi f}) e^{j2\pi fn} df &= \int_{-(\frac{1}{2})}^{\frac{1}{2}} \sum_m (s(m) e^{-(j2\pi fm)} e^{j2\pi fn}) df \\ &= \sum_m \left(s(m) \int_{-(\frac{1}{2})}^{\frac{1}{2}} e^{-(j2\pi f)(m-n)} df \right) \\ &= s(n) \end{aligned} \quad (1.40)$$

The Fourier transform pairs in discrete-time are

$$S(e^{j2\pi f}) = \sum_{n=-\infty}^{\infty} (s(n) e^{-(j2\pi fn)}) \quad (1.41)$$

$$s(n) = \int_{-(\frac{1}{2})}^{\frac{1}{2}} S(e^{j2\pi f}) e^{j2\pi fn} df \quad (1.42)$$

The properties of the discrete-time Fourier transform mirror those of the analog Fourier transform. The DTFT properties table ³⁶ shows similarities and differences. One important common property is Parseval's Theorem.

$$\sum_{n=-\infty}^{\infty} (|s(n)|^2) = \int_{-\frac{1}{2}}^{\frac{1}{2}} (|S(e^{j2\pi f})|^2) df \quad (1.43)$$

To show this important property, we simply substitute the Fourier transform expression into the frequency-domain expression for power.

$$\begin{aligned} \int_{-\frac{1}{2}}^{\frac{1}{2}} (|S(e^{j2\pi f})|^2) df &= \int_{-\frac{1}{2}}^{\frac{1}{2}} (\sum_n (s(n) e^{-j2\pi f n})) \sum_m (\overline{s(n)} e^{j2\pi f m}) df \\ &= \sum_{(n,m)} \left(s(n) \overline{s(n)} \int_{-\frac{1}{2}}^{\frac{1}{2}} e^{j2\pi f(m-n)} df \right) \end{aligned} \quad (1.44)$$

Using the orthogonality relation (1.39), the integral equals $\delta(m-n)$, where $\delta(n)$ is the unit sample (Figure 1.2: Unit sample). Thus, the double sum collapses into a single sum because nonzero values occur only when $n = m$, giving Parseval's Theorem as a result. We term $\sum_n (s^2(n))$ the energy in the discrete-time signal $s(n)$ in spite of the fact that discrete-time signals don't consume (or produce for that matter) energy. This terminology is a carry-over from the analog world.

Exercise 1.6

(Solution on p. 97.)

Suppose we obtained our discrete-time signal from values of the product $s(t)p_{T_s}(t)$, where the duration of the component pulses in $p_{T_s}(t)$ is Δ . How is the discrete-time signal energy related to the total energy contained in $s(t)$? Assume the signal is bandlimited and that the sampling rate was chosen appropriate to the Sampling Theorem's conditions.

1.9 DFT as a Matrix Operation³⁷

1.9.1 Matrix Review

Recall:

- Vectors in \mathbb{R}^N :

$$\forall x_i, x_i \in \mathbb{R} : \mathbf{x} = \begin{pmatrix} x_0 \\ x_1 \\ \dots \\ x_{N-1} \end{pmatrix}$$

- Vectors in \mathbb{C}^N :

$$\forall x_i, x_i \in \mathbb{C} : \mathbf{x} = \begin{pmatrix} x_0 \\ x_1 \\ \dots \\ x_{N-1} \end{pmatrix}$$

- Transposition:

1. transpose:

$$\mathbf{x}^T = \begin{pmatrix} x_0 & x_1 & \dots & x_{N-1} \end{pmatrix}$$

³⁶"Discrete-Time Fourier Transform Properties" <<http://cnx.org/content/m0506/latest/>>

³⁷This content is available online at <<http://cnx.org/content/m10962/2.5/>>.

2. conjugate:

$$\mathbf{x}^H = \begin{pmatrix} \overline{x_0} & \overline{x_1} & \dots & \overline{x_{N-1}} \end{pmatrix}$$

- Inner product³⁸:

1. real:

$$\mathbf{x}^T \mathbf{y} = \sum_{i=0}^{N-1} (x_i y_i)$$

2. complex:

$$\mathbf{x}^H \mathbf{y} = \sum_{i=0}^{N-1} (\overline{x_i} y_i)$$

- Matrix Multiplication:

$$\mathbf{A}\mathbf{x} = \begin{pmatrix} a_{00} & a_{01} & \dots & a_{0,N-1} \\ a_{10} & a_{11} & \dots & a_{1,N-1} \\ \vdots & \vdots & \dots & \vdots \\ a_{N-1,0} & a_{N-1,1} & \dots & a_{N-1,N-1} \end{pmatrix} \begin{pmatrix} x_0 \\ x_1 \\ \dots \\ x_{N-1} \end{pmatrix} = \begin{pmatrix} y_0 \\ y_1 \\ \dots \\ y_{N-1} \end{pmatrix}$$

$$y_k = \sum_{n=0}^{N-1} (a_{kn} x_n)$$

- Matrix Transposition:

$$\mathbf{A}^T = \begin{pmatrix} a_{00} & a_{10} & \dots & a_{N-1,0} \\ a_{01} & a_{11} & \dots & a_{N-1,1} \\ \vdots & \vdots & \dots & \vdots \\ a_{0,N-1} & a_{1,N-1} & \dots & a_{N-1,N-1} \end{pmatrix}$$

Matrix transposition involved simply swapping the rows with columns.

$$\mathbf{A}^H = \overline{\mathbf{A}^T}$$

The above equation is Hermitian transpose.

$$[\mathbf{A}^T]_{kn} = \mathbf{A}_{nk}$$

$$[\mathbf{A}^H]_{kn} = \overline{[\mathbf{A}]_{nk}}$$

1.9.2 Representing DFT as Matrix Operation

Now let's represent the DFT³⁹ in vector-matrix notation.

$$\mathbf{x} = \begin{pmatrix} x[0] \\ x[1] \\ \dots \\ x[N-1] \end{pmatrix}$$

³⁸"Inner Products" <<http://cnx.org/content/m10755/latest/>>

³⁹"Discrete Fourier Transform (DFT)" <<http://cnx.org/content/m10249/latest/>>

$$\mathbf{X} = \begin{pmatrix} X[0] \\ X[1] \\ \dots \\ X[N-1] \end{pmatrix} \in \mathbb{C}^N$$

Here \mathbf{x} is the vector of time samples and \mathbf{X} is the vector of DFT coefficients. How are \mathbf{x} and \mathbf{X} related:

$$X[k] = \sum_{n=0}^{N-1} \left(x[n] e^{-j\frac{2\pi}{N}kn} \right)$$

where

$$a_{kn} = \left(e^{-j\frac{2\pi}{N}} \right)^{kn} = W_N^{kn}$$

so

$$\mathbf{X} = W\mathbf{x}$$

where \mathbf{X} is the DFT vector, W is the matrix and \mathbf{x} the time domain vector.

$$W_{kn} = \left(e^{-j\frac{2\pi}{N}} \right)^{kn}$$

$$\mathbf{X} = W \begin{pmatrix} x[0] \\ x[1] \\ \dots \\ x[N-1] \end{pmatrix}$$

IDFT:

$$x[n] = \frac{1}{N} \sum_{k=0}^{N-1} \left(X[k] \left(e^{j\frac{2\pi}{N}} \right)^{nk} \right)$$

where

$$\left(e^{j\frac{2\pi}{N}} \right)^{nk} = \overline{W_N^{nk}}$$

$\overline{W_N^{nk}}$ is the matrix Hermitian transpose. So,

$$\mathbf{x} = \frac{1}{N} W^H \mathbf{X}$$

where \mathbf{x} is the time vector, $\frac{1}{N} W^H$ is the inverse DFT matrix, and \mathbf{X} is the DFT vector.

1.10 Sampling theory

1.10.1 Introduction⁴⁰

Contents of Sampling chapter

- Introduction(Current module)
- Proof (Section 1.10.2)
- Illustrations (Section 1.10.3)

⁴⁰This content is available online at <<http://cnx.org/content/m11419/1.29/>>.

- Matlab Example⁴¹
- Hold operation⁴²
- System view (Section 1.10.4)
- Aliasing applet⁴³
- Exercises⁴⁴
- Table of formulas⁴⁵

1.10.1.1 Why sample?

This section introduces sampling. Sampling is the necessary fundament for all digital signal processing and communication. Sampling can be defined as the process of measuring an analog signal at distinct points.

Digital representation of analog signals offers advantages in terms of

- robustness towards noise, meaning we can send more bits/s
- use of flexible processing equipment, in particular the computer
- more reliable processing equipment
- easier to adapt complex algorithms

1.10.1.2 Claude E. Shannon



Figure 1.20: Claude Elwood Shannon (1916-2001)

Claude Shannon⁴⁶ has been called the father of information theory, mainly due to his landmark papers on the "Mathematical theory of communication"⁴⁷. Harry Nyquist⁴⁸ was the first to state the sampling theorem

⁴¹"Sampling and reconstruction with Matlab" <<http://cnx.org/content/m11549/latest/>>

⁴²"Hold operation" <<http://cnx.org/content/m11458/latest/>>

⁴³"Aliasing Applet" <<http://cnx.org/content/m11448/latest/>>

⁴⁴"Exercises" <<http://cnx.org/content/m11442/latest/>>

⁴⁵"Table of Formulas" <<http://cnx.org/content/m11450/latest/>>

⁴⁶<http://www.research.att.com/~njas/doc/ces5.html>

⁴⁷<http://cm.bell-labs.com/cm/ms/what/shannonday/shannon1948.pdf>

⁴⁸http://www.wikipedia.org/wiki/Harry_Nyquist

in 1928, but it was not proven until Shannon proved it 21 years later in the paper "Communications in the presence of noise"⁴⁹.

1.10.1.3 Notation

In this chapter we will be using the following notation

- Original analog signal $x(t)$
- Sampling frequency F_s
- Sampling interval T_s (Note that: $F_s = \frac{1}{T_s}$)
- Sampled signal $x_s(n)$. (Note that $x_s(n) = x(nT_s)$)
- Real angular frequency Ω
- Digital angular frequency ω . (Note that: $\omega = \Omega T_s$)

1.10.1.4 The Sampling Theorem

THE SAMPLING THEOREM: When sampling an analog signal the sampling frequency must be greater than twice the highest frequency component of the analog signal to be able to reconstruct the original signal from the sampled version.

Finished? Have a look at: Proof (Section 1.10.2); Illustrations (Section 1.10.3); Matlab Example⁵⁰; Aliasing applet⁵¹; Hold operation⁵²; System view (Section 1.10.4); Exercises⁵³

1.10.2 Proof⁵⁴

SAMPLING THEOREM: In order to recover the signal $x(t)$ from its samples exactly, it is necessary to sample $x(t)$ at a rate greater than twice its highest frequency component.

1.10.2.1 Introduction

As mentioned earlier (p. 37), sampling is the necessary fundament when we want to apply digital signal processing on analog signals.

Here we present the proof of the sampling theorem. The proof is divided in two. First we find an expression for the spectrum of the signal resulting from sampling the original signal $x(t)$. Next we show that the signal $x(t)$ can be recovered from the samples. Often it is easier using the frequency domain when carrying out a proof, and this is also the case here.

Key points in the proof

- We find an equation (1.52) for the spectrum of the sampled signal
- We find a simple method to reconstruct (1.58) the original signal
- The sampled signal has a periodic spectrum...
- ...and the period is $2\pi F_s$

⁴⁹<http://www.stanford.edu/class/ee104/shannonpaper.pdf>

⁵⁰"Sampling and reconstruction with Matlab" <<http://cnx.org/content/m11549/latest/>>

⁵¹"Aliasing Applet" <<http://cnx.org/content/m11448/latest/>>

⁵²"Hold operation" <<http://cnx.org/content/m11458/latest/>>

⁵³"Exercises" <<http://cnx.org/content/m11442/latest/>>

⁵⁴This content is available online at <<http://cnx.org/content/m11423/1.27/>>.

1.10.2.2 Proof part 1 - Spectral considerations

By sampling $x(t)$ every T_s second we obtain $x_s(n)$. The inverse fourier transform of this time discrete signal⁵⁵ is

$$x_s(n) = \frac{1}{2\pi} \int_{-\pi}^{\pi} X_s(e^{j\omega}) e^{j\omega n} d\omega \quad (1.45)$$

For convenience we express the equation in terms of the real angular frequency Ω using $\omega = \Omega T_s$. We then obtain

$$x_s(n) = \frac{T_s}{2\pi} \int_{-\frac{\pi}{T_s}}^{\frac{\pi}{T_s}} X_s(e^{j\Omega T_s}) e^{j\Omega T_s n} d\Omega \quad (1.46)$$

The inverse fourier transform of a continuous signal is

$$x(t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} X(j\Omega) e^{j\Omega t} d\Omega \quad (1.47)$$

From this equation we find an expression for $x(nT_s)$

$$x(nT_s) = \frac{1}{2\pi} \int_{-\infty}^{\infty} X(j\Omega) e^{j\Omega n T_s} d\Omega \quad (1.48)$$

To account for the difference in region of integration we split the integration in (1.48) into subintervals of length $\frac{2\pi}{T_s}$ and then take the sum over the resulting integrals to obtain the complete area.

$$x(nT_s) = \frac{1}{2\pi} \sum_{k=-\infty}^{\infty} \left(\int_{\frac{(2k-1)\pi}{T_s}}^{\frac{(2k+1)\pi}{T_s}} X(j\Omega) e^{j\Omega n T_s} d\Omega \right) \quad (1.49)$$

Then we change the integration variable, setting $\Omega = \eta + \frac{2\pi k}{T_s}$

$$x(nT_s) = \frac{1}{2\pi} \sum_{k=-\infty}^{\infty} \left(\int_{-\frac{\pi}{T_s}}^{\frac{\pi}{T_s}} X\left(j\left(\eta + \frac{2\pi k}{T_s}\right)\right) e^{j\left(\eta + \frac{2\pi k}{T_s}\right)n T_s} d\eta \right) \quad (1.50)$$

We obtain the final form by observing that $e^{j2\pi k n} = 1$, reinserting $\eta = \Omega$ and multiplying by $\frac{T_s}{T_s}$

$$x(nT_s) = \frac{T_s}{2\pi} \int_{-\frac{\pi}{T_s}}^{\frac{\pi}{T_s}} \sum_{k=-\infty}^{\infty} \left(\frac{1}{T_s} X\left(j\left(\Omega + \frac{2\pi k}{T_s}\right)\right) e^{j\Omega n T_s} \right) d\Omega \quad (1.51)$$

To make $x_s(n) = x(nT_s)$ for all values of n , the integrands in (1.46) and (1.51) have to agree, that is

$$X_s(e^{j\Omega T_s}) = \frac{1}{T_s} \sum_{k=-\infty}^{\infty} \left(X\left(j\left(\Omega + \frac{2\pi k}{T_s}\right)\right) \right) \quad (1.52)$$

This is a central result. We see that the digital spectrum consists of a sum of shifted versions of the original, analog spectrum. Observe the periodicity!

We can also express this relation in terms of the digital angular frequency $\omega = \Omega T_s$

$$X_s(e^{j\omega}) = \frac{1}{T_s} \sum_{k=-\infty}^{\infty} \left(X\left(j\frac{\omega + 2\pi k}{T_s}\right) \right) \quad (1.53)$$

This concludes the first part of the proof. Now we want to find a reconstruction formula, so that we can recover $x(t)$ from $x_s(n)$.

⁵⁵"Discrete time signals" <<http://cnx.org/content/m11476/latest/>>

1.10.2.3 Proof part II - Signal reconstruction

For a bandlimited (Figure 1.22) signal the inverse fourier transform is

$$x(t) = \frac{1}{2\pi} \int_{-\frac{\pi}{T_s}}^{\frac{\pi}{T_s}} X(j\Omega) e^{j\Omega t} d\Omega \quad (1.54)$$

In the interval we are integrating we have: $X_s(e^{j\Omega T_s}) = \frac{X(j\Omega)}{T_s}$. Substituting this relation into (1.54) we get

$$x(t) = \frac{T_s}{2\pi} \int_{-\frac{\pi}{T_s}}^{\frac{\pi}{T_s}} X_s(e^{j\Omega T_s}) e^{j\Omega t} d\Omega \quad (1.55)$$

Using the DTFT⁵⁶ relation for $X_s(e^{j\Omega T_s})$ we have

$$x(t) = \frac{T_s}{2\pi} \int_{-\frac{\pi}{T_s}}^{\frac{\pi}{T_s}} \sum_{n=-\infty}^{\infty} \left(x_s(n) e^{-(j\Omega n T_s)} e^{j\Omega t} \right) d\Omega \quad (1.56)$$

Interchanging integration and summation (under the assumption of convergence) leads to

$$x(t) = \frac{T_s}{2\pi} \sum_{n=-\infty}^{\infty} \left(x_s(n) \int_{-\frac{\pi}{T_s}}^{\frac{\pi}{T_s}} e^{j\Omega(t-nT_s)} d\Omega \right) \quad (1.57)$$

Finally we perform the integration and arrive at the important reconstruction formula

$$x(t) = \sum_{n=-\infty}^{\infty} \left(x_s(n) \frac{\sin\left(\frac{\pi}{T_s}(t-nT_s)\right)}{\frac{\pi}{T_s}(t-nT_s)} \right) \quad (1.58)$$

(Thanks to R.Loos for pointing out an error in the proof.)

1.10.2.4 Summary

$$\text{SPECTRUM SAMPLED SIGNAL: } X_s(e^{j\Omega T_s}) = \frac{1}{T_s} \sum_{k=-\infty}^{\infty} \left(X\left(j\left(\Omega + \frac{2\pi k}{T_s}\right)\right) \right)$$

$$\text{RECONSTRUCTION FORMULA: } x(t) = \sum_{n=-\infty}^{\infty} \left(x_s(n) \frac{\sin\left(\frac{\pi}{T_s}(t-nT_s)\right)}{\frac{\pi}{T_s}(t-nT_s)} \right)$$

Go to Introduction (Section 1.10.1); Illustrations (Section 1.10.3); Matlab Example⁵⁷; Hold operation⁵⁸; Aliasing applet⁵⁹; System view (Section 1.10.4); Exercises⁶⁰ ?

1.10.3 Illustrations⁶¹

In this module we illustrate the processes involved in sampling and reconstruction. To see how all these processes work together as a whole, take a look at the system view (Section 1.10.4). In Sampling and reconstruction with Matlab⁶² we provide a Matlab script for download. The matlab script shows the process of sampling and reconstruction **live**.

⁵⁶"Table of Formulas" <<http://cnx.org/content/m11450/latest/>>

⁵⁷"Sampling and reconstruction with Matlab" <<http://cnx.org/content/m11549/latest/>>

⁵⁸"Hold operation" <<http://cnx.org/content/m11458/latest/>>

⁵⁹"Aliasing Applet" <<http://cnx.org/content/m11448/latest/>>

⁶⁰"Exercises" <<http://cnx.org/content/m11442/latest/>>

⁶¹This content is available online at <<http://cnx.org/content/m11443/1.33/>>.

⁶²"Sampling and reconstruction with Matlab" <<http://cnx.org/content/m11549/latest/>>

1.10.3.1 Basic examples

Example 1.18

To sample an analog signal with 3000 Hz as the highest frequency component requires sampling at 6000 Hz or above.

Example 1.19

The sampling theorem can also be applied in two dimensions, i.e. for image analysis. A 2D sampling theorem has a simple physical interpretation in image analysis: Choose the sampling interval such that it is less than or equal to half of the smallest interesting detail in the image.

1.10.3.2 The process of sampling

We start off with an analog signal. This can for example be the sound coming from your stereo at home or your friend talking.

The signal is then sampled uniformly. Uniform sampling implies that we sample every T_s seconds. In Figure 1.21 we see an analog signal. The analog signal has been sampled at times $t = nT_s$.

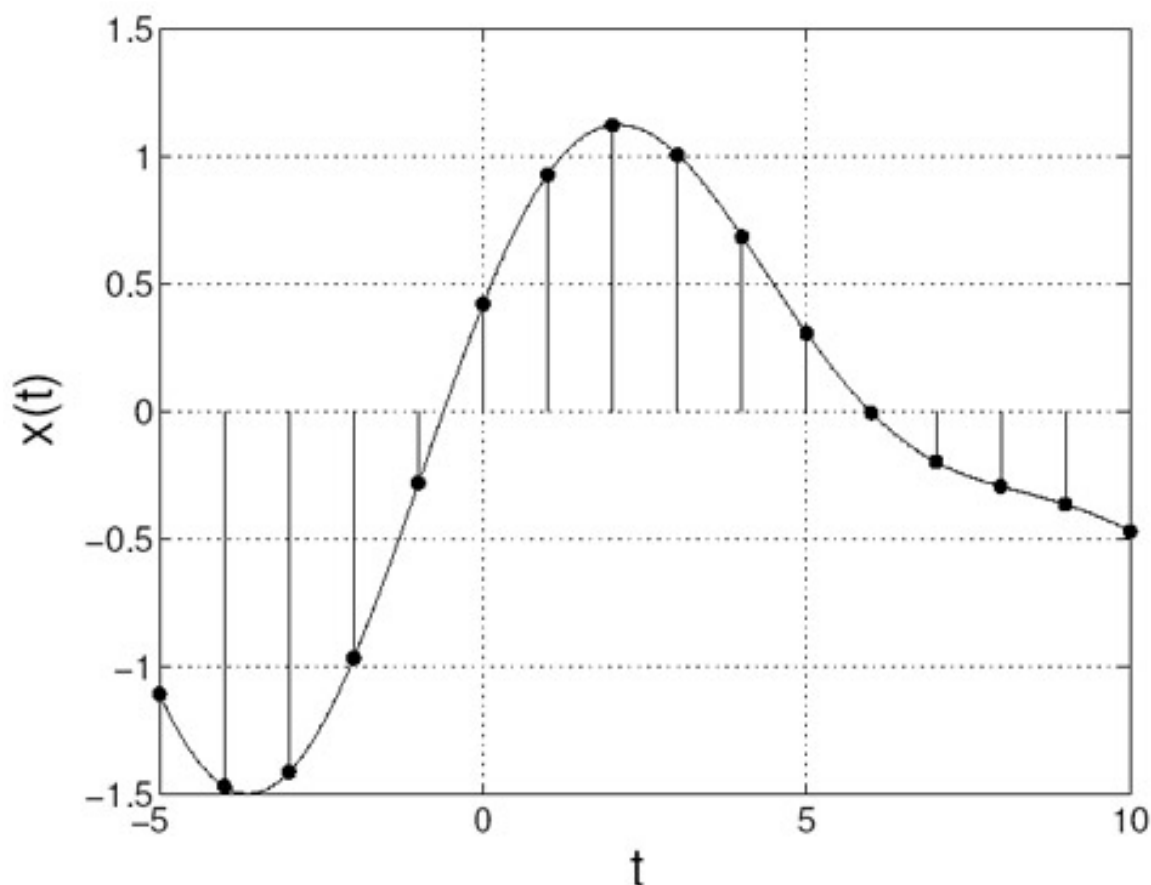


Figure 1.21: Analog signal, samples are marked with dots.

In signal processing it is often more convenient and easier to work in the frequency domain. So let's look at the signal in frequency domain, Figure 1.22. For illustration purposes we take the frequency content of the signal as a triangle. (If you Fourier transform the signal in Figure 1.21 you will not get such a nice triangle.)

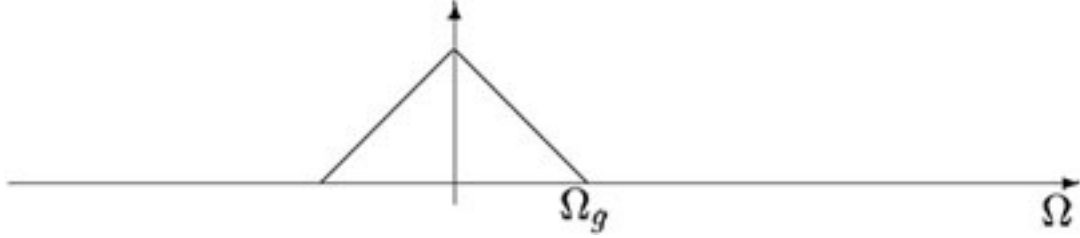


Figure 1.22: The spectrum $X(j\Omega)$.

Notice that the signal in Figure 1.22 is bandlimited. We can see that the signal is bandlimited because $X(j\Omega)$ is zero outside the interval $[-\Omega_g, \Omega_g]$. Equivalently we can state that the signal has no angular frequencies above Ω_g , corresponding to no frequencies above $F_g = \frac{\Omega_g}{2\pi}$.

Now let's take a look at the sampled signal in the frequency domain. While proving (Section 1.10.2) the sampling theorem we found the the spectrum of the sampled signal consists of a sum of shifted versions of the analog spectrum. Mathematically this is described by the following equation:

$$X_s(e^{j\Omega T_s}) = \frac{1}{T_s} \sum_{k=-\infty}^{\infty} \left(X \left(j \left(\Omega + \frac{2\pi k}{T_s} \right) \right) \right) \quad (1.59)$$

1.10.3.2.1 Sampling fast enough

In Figure 1.23 we show the result of sampling $x(t)$ according to the sampling theorem (Section 1.10.1.4: The Sampling Theorem). This means that when sampling the signal in Figure 1.21 / Figure 1.22 we use $F_s \geq 2F_g$. Observe in Figure 1.23 that we have the same spectrum as in Figure 1.22 for $\Omega \in [-\Omega_g, \Omega_g]$, except for the scaling factor $\frac{1}{T_s}$. This is a consequence of the sampling frequency. As mentioned in the proof (Key points in the proof, p. 38) the spectrum of the sampled signal is periodic with period $2\pi F_s = \frac{2\pi}{T_s}$.

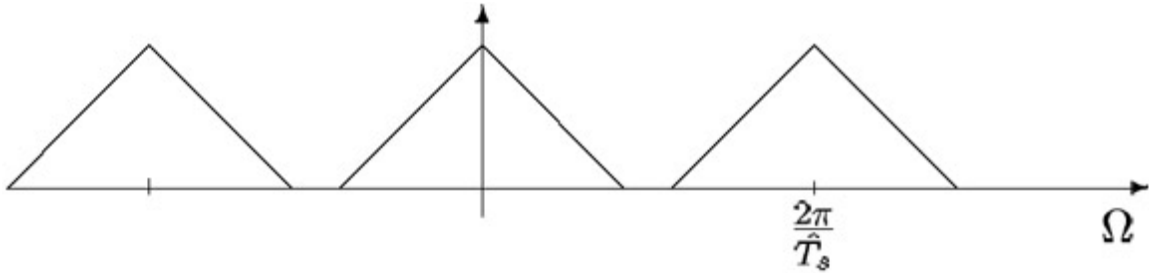


Figure 1.23: The spectrum X_s . Sampling frequency is OK.

So now we are, according to the sample theorem (Section 1.10.1.4: The Sampling Theorem), able to reconstruct the original signal *exactly*. How we can do this will be explored further down under reconstruction (Section 1.10.3.3: Reconstruction). But first we will take a look at what happens when we sample too slowly.

1.10.3.2.2 Sampling too slowly

If we sample $x(t)$ too slowly, that is $F_s < 2F_g$, we will get overlap between the repeated spectra, see Figure 1.24. According to (1.59) the resulting spectra is the sum of these. This overlap gives rise to the concept of aliasing.

ALIASING: If the sampling frequency is less than twice the highest frequency component, then frequencies in the original signal that are above half the sampling rate will be "aliased" and will appear in the resulting signal as lower frequencies.

The consequence of aliasing is that we cannot recover the original signal, so aliasing has to be avoided. Sampling too slowly will produce a sequence $x_s(n)$ that could have originated from a number of signals. So there is *no* chance of recovering the original signal. To learn more about aliasing, take a look at this module⁶³. (Includes an applet for demonstration!)

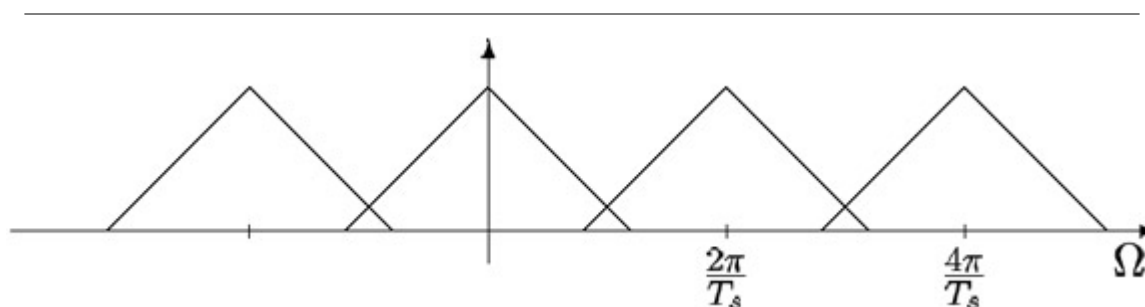


Figure 1.24: The spectrum X_s . Sampling frequency is too low.

To avoid aliasing we have to sample fast enough. But if we can't sample fast enough (possibly due to costs) we can include an Anti-Aliasing filter. This will not able us to get an exact reconstruction but can still be a good solution.

ANTI-ALIASING FILTER: Typically a low-pass filter that is applied before sampling to ensure that no components with frequencies greater than half the sample frequency remain.

Example 1.20

The stagecoach effect

In older western movies you can observe aliasing on a stagecoach when it starts to roll. At first the spokes appear to turn forward, but as the stagecoach increase its speed the spokes appear to turn backward. This comes from the fact that the sampling rate, here the number of frames per second, is too low. We can view each frame as a sample of an image that is changing continuously in time. (Applet illustrating the stagecoach effect⁶⁴)

⁶³"Aliasing Applet" <<http://cnx.org/content/m11448/latest/>>

⁶⁴<http://flowers.ofthenight.org/wagonWheel/wagonWheel.html>

1.10.3.3 Reconstruction

Given the signal in Figure 1.23 we want to recover the original signal, but the question is how?

When there is no overlapping in the spectrum, the spectral component given by $k = 0$ (see (1.59)), is equal to the spectrum of the analog signal. This offers an opportunity to use a simple reconstruction process. Remember what you have learned about filtering. What we want is to change signal in Figure 1.23 into that of Figure 1.22. To achieve this we have to remove all the extra components generated in the sampling process. To remove the extra components we apply an ideal analog low-pass filter as shown in Figure 1.25. As we see the ideal filter is rectangular in the frequency domain. A rectangle in the frequency domain corresponds to a sinc⁶⁵ function in time domain (and vice versa).

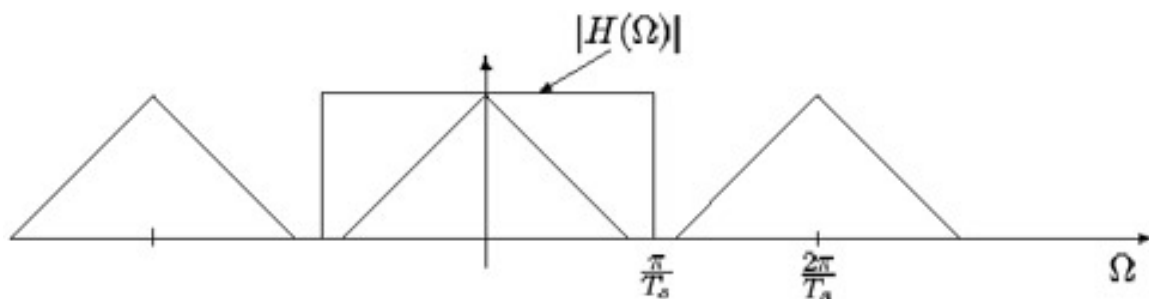


Figure 1.25: $H(j\Omega)$ The ideal reconstruction filter.

Then we have reconstructed the original spectrum, and as we know *if two signals are identical in the frequency domain, they are also identical in the time domain*. End of reconstruction.

1.10.3.4 Conclusions

The Shannon sampling theorem requires that the input signal prior to sampling is band-limited to at most half the sampling frequency. Under this condition the samples give an exact signal representation. It is truly remarkable that such a broad and useful class signals can be represented that easily!

We also looked into the problem of reconstructing the signals from its samples. Again the simplicity of the *principle* is striking: linear filtering by an ideal low-pass filter will do the job. However, the ideal filter is impossible to create, but that is another story..

Go to? Introduction (Section 1.10.1); Proof (Section 1.10.2); Illustrations (Section 1.10.3); Matlab Example⁶⁶; Aliasing applet⁶⁷; Hold operation⁶⁸; System view (Section 1.10.4); Exercises⁶⁹

1.10.4 Systems view of sampling and reconstruction⁷⁰

1.10.4.1 Ideal reconstruction system

Figure 1.26 shows the ideal reconstruction system based on the results of the Sampling theorem proof (Section 1.10.2).

⁶⁵http://ccrma-www.stanford.edu/~jos/Interpolation/sinc_function.html

⁶⁶"Sampling and reconstruction with Matlab" <<http://cnx.org/content/m11549/latest/>>

⁶⁷"Aliasing Applet" <<http://cnx.org/content/m11448/latest/>>

⁶⁸"Hold operation" <<http://cnx.org/content/m11458/latest/>>

⁶⁹"Exercises" <<http://cnx.org/content/m11442/latest/>>

⁷⁰This content is available online at <<http://cnx.org/content/m11465/1.20/>>.

Figure 1.26 consists of a sampling device which produces a time-discrete sequence $x_s(n)$. The reconstruction filter, $h(t)$, is an ideal analog sinc⁷¹ filter, with $h(t) = \text{sinc}\left(\frac{t}{T_s}\right)$. We can't apply the time-discrete sequence $x_s(n)$ directly to the analog filter $h(t)$. To solve this problem we turn the sequence into an analog signal using delta functions⁷². Thus we write $x_s(t) = \sum_{n=-\infty}^{\infty} (x_s(n) \delta(t - nT))$.

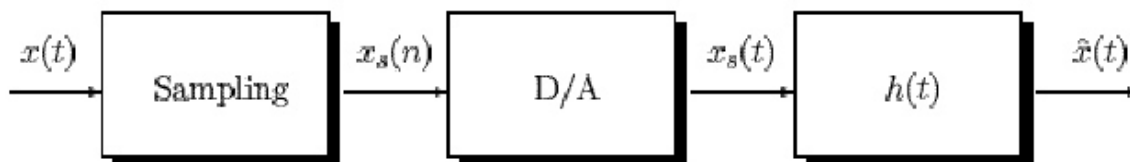


Figure 1.26: Ideal reconstruction system

But when will the system produce an output $\hat{x}(t) = x(t)$? According to the sampling theorem (Section 1.10.1.4: The Sampling Theorem) we have $\hat{x}(t) = x(t)$ when the sampling frequency, F_s , is at least twice the highest frequency component of $x(t)$.

1.10.4.2 Ideal system including anti-aliasing

To be sure that the reconstructed signal is free of aliasing it is customary to apply a lowpass filter, an anti-aliasing filter (p. 43), before sampling as shown in Figure 1.27.

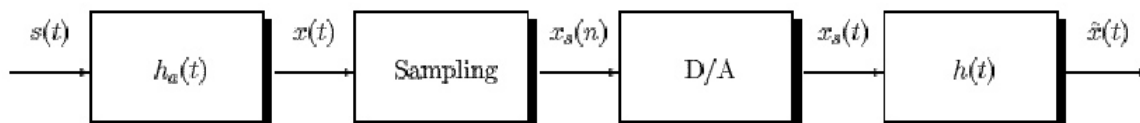


Figure 1.27: Ideal reconstruction system with anti-aliasing filter (p. 43)

Again we ask the question of when the system will produce an output $\hat{x}(t) = s(t)$? If the signal is entirely confined within the passband of the lowpass filter we will get perfect reconstruction if F_s is high enough.

But if the anti-aliasing filter removes the "higher" frequencies, (which in fact is the job of the anti-aliasing filter), we will *never* be able to *exactly* reconstruct the original signal, $s(t)$. If we sample fast enough we can reconstruct $x(t)$, which in most cases is satisfying.

The reconstructed signal, $\hat{x}(t)$, will not have aliased frequencies. This is essential for further use of the signal.

1.10.4.3 Reconstruction with hold operation

To make our reconstruction system realizable there are many things to look into. Among them are the fact that any practical reconstruction system must input finite length pulses into the reconstruction filter. This can be accomplished by the hold operation⁷³. To alleviate the distortion caused by the hold operator we apply the output from the hold device to a compensator. The compensation can be as accurate as we wish, this is cost and application consideration.

⁷¹http://ccrma-www.stanford.edu/~jos/Interpolation/sinc_function.html

⁷²"Table of Formulas" <<http://cnx.org/content/m11450/latest/>>

⁷³"Hold operation" <<http://cnx.org/content/m11458/latest/>>

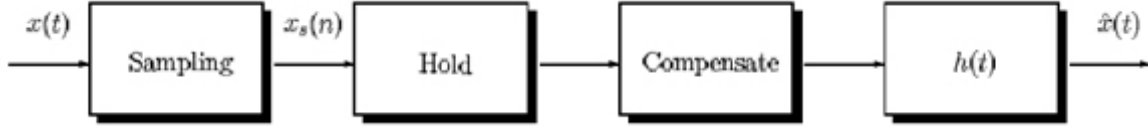


Figure 1.28: More practical reconstruction system with a hold component⁷⁴

By the use of the hold component the reconstruction will not be exact, but as mentioned above we can get as close as we want.

Introduction (Section 1.10.1); Proof (Section 1.10.2); Illustrations (Section 1.10.3); Matlab example⁷⁵; Hold operation⁷⁶; Aliasing applet⁷⁷; Exercises⁷⁸

1.10.5 Sampling CT Signals: A Frequency Domain Perspective⁷⁹

1.10.5.1 Understanding Sampling in the Frequency Domain

We want to relate $x_c(t)$ directly to $x[n]$. Compute the CTFT of

$$\begin{aligned}
 x_s(t) &= \sum_{n=-\infty}^{\infty} (x_c(nT) \delta(t - nT)) \\
 X_s(\Omega) &= \int_{-\infty}^{\infty} \left(\sum_{n=-\infty}^{\infty} (x_c(nT) \delta(t - nT)) \right) e^{(-j)\Omega t} dt \\
 &= \sum_{n=-\infty}^{\infty} \left(x_c(nT) \int_{-\infty}^{\infty} \delta(t - nT) e^{(-j)\Omega t} dt \right) \\
 &= \sum_{n=-\infty}^{\infty} (x[n] e^{(-j)\Omega nT}) \\
 &= \sum_{n=-\infty}^{\infty} (x[n] e^{(-j)\omega n}) \\
 &= X(\omega)
 \end{aligned} \tag{1.60}$$

where $\omega \equiv \Omega T$ and $X(\omega)$ is the DTFT of $x[n]$.

RECALL:

$$\begin{aligned}
 X_s(\Omega) &= \frac{1}{T} \sum_{k=-\infty}^{\infty} (X_c(\Omega - k\Omega_s)) \\
 X(\omega) &= \frac{1}{T} \sum_{k=-\infty}^{\infty} (X_c(\Omega - k\Omega_s)) \\
 &= \frac{1}{T} \sum_{k=-\infty}^{\infty} \left(X_c\left(\frac{\omega - 2\pi k}{T}\right) \right)
 \end{aligned} \tag{1.61}$$

where this last part is 2π -periodic.

⁷⁴"Hold operation" <<http://cnx.org/content/m11458/latest/>>

⁷⁵"Sampling and reconstruction with Matlab" <<http://cnx.org/content/m11549/latest/>>

⁷⁶"Hold operation" <<http://cnx.org/content/m11458/latest/>>

⁷⁷"Aliasing Applet" <<http://cnx.org/content/m11448/latest/>>

⁷⁸"Exercises" <<http://cnx.org/content/m11442/latest/>>

⁷⁹This content is available online at <<http://cnx.org/content/m10994/2.2/>>.

1.10.5.1.1 Sampling

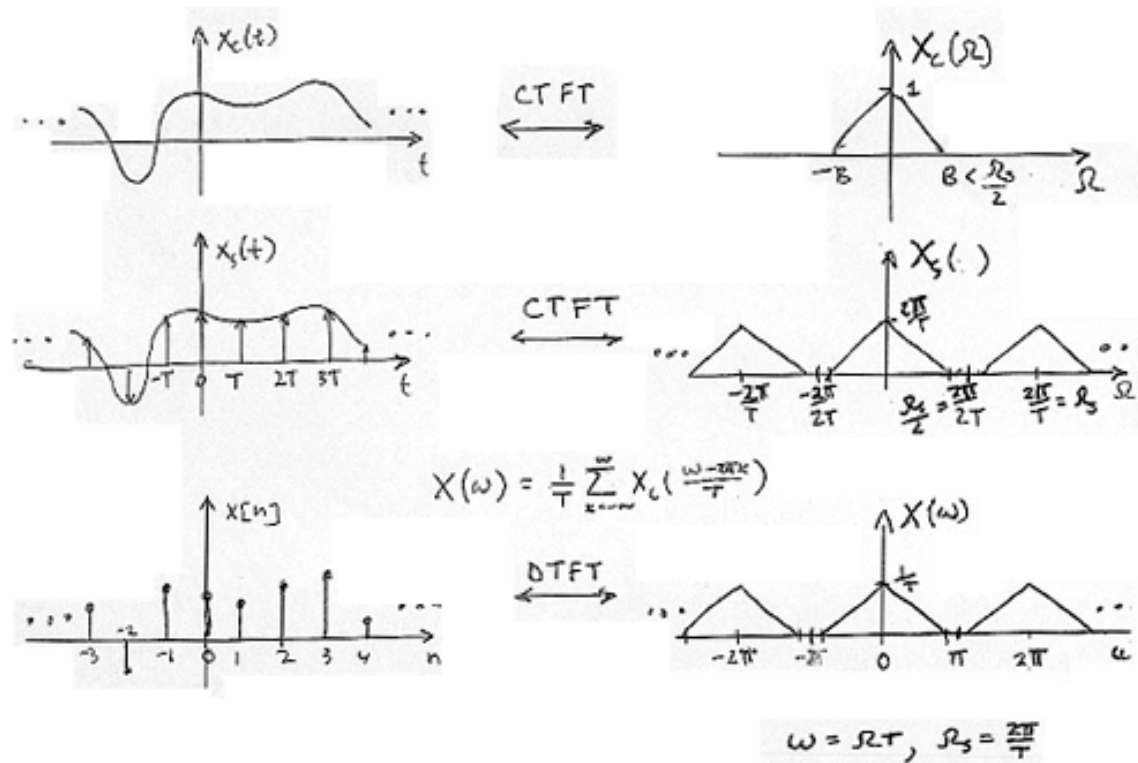


Figure 1.29

Example 1.21: Speech

Speech is intelligible if bandlimited by a CT lowpass filter to the band ± 4 kHz. We can sample speech as slowly as _____?

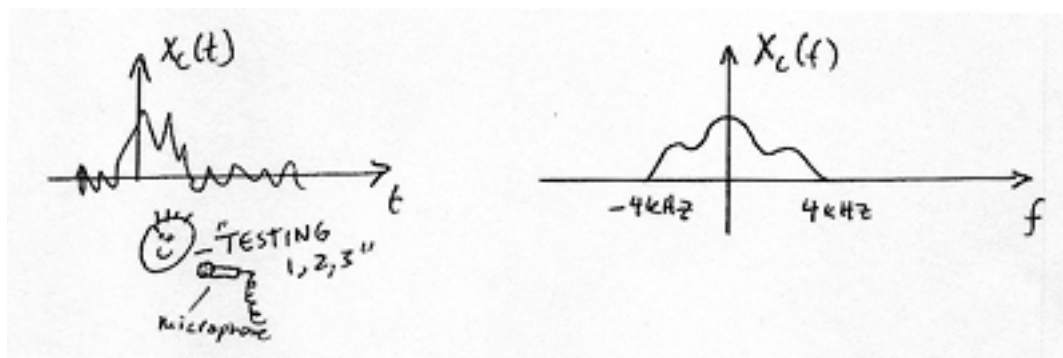
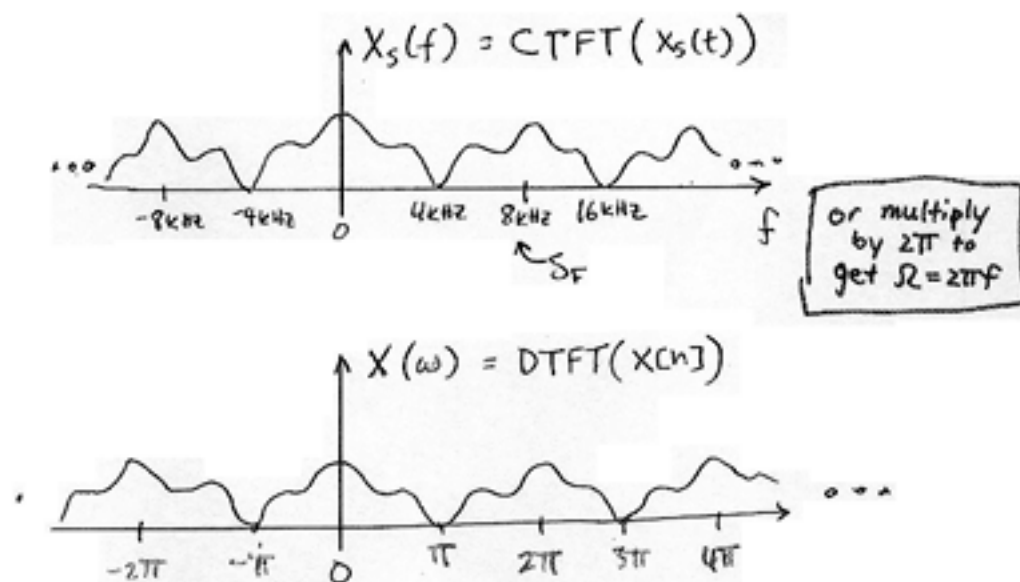


Figure 1.30

Figure 1.31: Note that there is no mention of T or Ω_s !

1.10.5.2 Relating $x[n]$ to sampled $x(t)$

Recall the following equality:

$$x_s(t) = \sum_n (x(nT) \delta(t - nT))$$

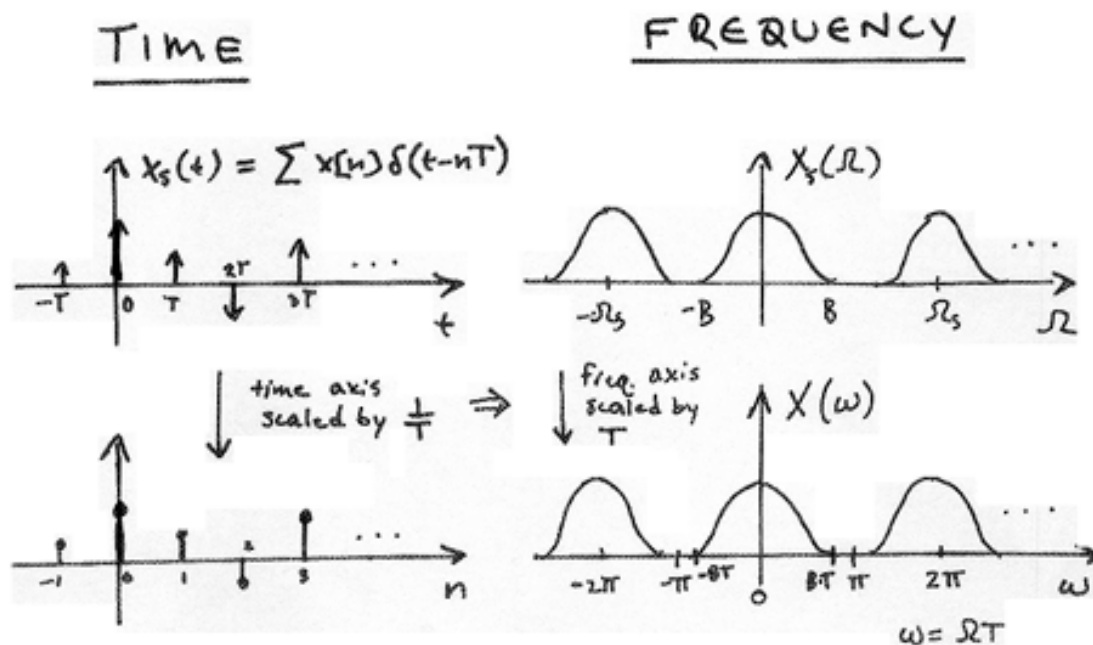


Figure 1.32

Recall the CTFT relation:

$$\left(x(\alpha t) \leftrightarrow \frac{1}{\alpha} X\left(\frac{\Omega}{\alpha}\right) \right) \quad (1.62)$$

where α is a scaling of time and $\frac{1}{\alpha}$ is a scaling in frequency.

$$X_s(\Omega) \equiv X(\Omega T) \quad (1.63)$$

1.10.6 The DFT: Frequency Domain with a Computer Analysis⁸⁰

1.10.6.1 Introduction

We just covered ideal (and non-ideal) (time) sampling of CT signals (Section 1.10.5). This enabled DT signal processing solutions for CT applications (Figure 1.33):

⁸⁰This content is available online at <http://cnx.org/content/m10992/2.3/>.



Figure 1.33

Much of the theoretical analysis of such systems relied on frequency domain representations. How do we carry out these frequency domain analysis on the computer? Recall the following relationships:

$$x[n] \xleftrightarrow{DTFT} X(\omega)$$

$$x(t) \xleftrightarrow{CTFT} X(\Omega)$$

where ω and Ω are continuous frequency variables.

1.10.6.1.1 Sampling DTFT

Consider the DTFT of a discrete-time (DT) signal $x[n]$. Assume $x[n]$ is of finite duration N (i.e., an N -point signal).

$$X(\omega) = \sum_{n=0}^{N-1} \left(x[n] e^{(-j)\omega n} \right) \quad (1.64)$$

where $X(\omega)$ is the continuous function that is indexed by the real-valued parameter $-\pi \leq \omega \leq \pi$. The other function, $x[n]$, is a discrete function that is indexed by integers.

We want to work with $X(\omega)$ on a computer. Why not just *sample* $X(\omega)$?

$$\begin{aligned} X[k] &= X\left(\frac{2\pi}{N}k\right) \\ &= \sum_{n=0}^{N-1} \left(x[n] e^{(-j)2\pi \frac{k}{N}n} \right) \end{aligned} \quad (1.65)$$

In (1.65) we sampled at $\omega = \frac{2\pi}{N}k$ where $k = \{0, 1, \dots, N-1\}$ and $X[k]$ for $k = \{0, \dots, N-1\}$ is called the **Discrete Fourier Transform (DFT)** of $x[n]$.

Example 1.22

Finite Duration DT Signal

Image not finished

Figure 1.34

The DTFT of the image in Figure 1.34 (Finite Duration DT Signal) is written as follows:

$$X(\omega) = \sum_{n=0}^{N-1} \left(x[n] e^{(-j)\omega n} \right) \quad (1.66)$$

where ω is any 2π -interval, for example $-\pi \leq \omega \leq \pi$.

Sample $X(\omega)$

Image not finished

Figure 1.35

where again we sampled at $\omega = \frac{2\pi}{N}k$ where $k = \{0, 1, \dots, M-1\}$. For example, we take

$$M = 10$$

. In the following section (Section 1.10.6.1.1.1: Choosing M) we will discuss in more detail how we should choose M , the number of samples in the 2π interval.

(This is precisely how we would plot $X(\omega)$ in Matlab.)

1.10.6.1.1.1 Choosing M

1.10.6.1.1.1.1 Case 1

Given N (length of $x[n]$), choose ($M \gg N$) to obtain a dense sampling of the DTFT (Figure 1.36):

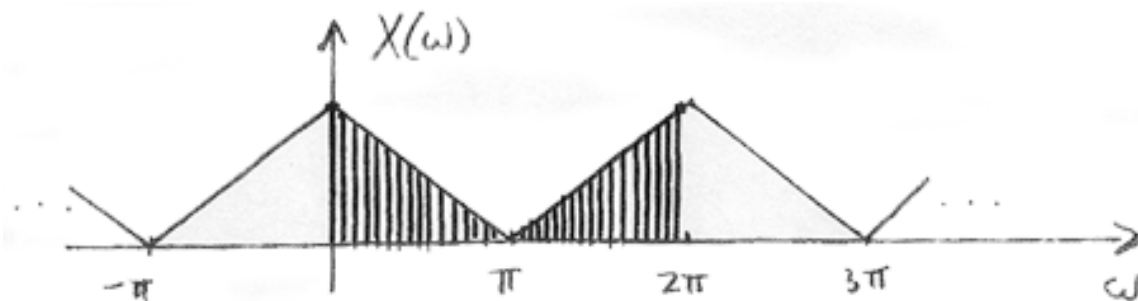


Figure 1.36

1.10.6.1.1.1.2 Case 2

Choose M as small as possible (to minimize the amount of computation).

In general, we require $M \geq N$ in order to represent all information in

$$\forall n, n = \{0, \dots, N-1\} : (x[n])$$

Let's concentrate on $M = N$:

$$x[n] \stackrel{DFT}{\leftrightarrow} X[k]$$

for $n = \{0, \dots, N-1\}$ and $k = \{0, \dots, N-1\}$

$$numbers \leftrightarrow N \text{ numbers}$$

1.10.6.2 Discrete Fourier Transform (DFT)

Define

$$X[k] \equiv X\left(\frac{2\pi k}{N}\right) \quad (1.67)$$

where $N = \text{length}(x[n])$ and $k = \{0, \dots, N-1\}$. In this case, $M = N$.

DFT

$$X[k] = \sum_{n=0}^{N-1} \left(x[n] e^{(-j)2\pi \frac{k}{N} n} \right) \quad (1.68)$$

Inverse DFT (IDFT)

$$x[n] = \frac{1}{N} \sum_{k=0}^{N-1} \left(X[k] e^{j2\pi \frac{k}{N} n} \right) \quad (1.69)$$

1.10.6.2.1 Interpretation

Represent $x[n]$ in terms of a sum of N complex sinusoids⁸¹ of amplitudes $X[k]$ and frequencies

$$\forall k, k \in \{0, \dots, N-1\} : \left(\omega_k = \frac{2\pi k}{N} \right)$$

THINK: Fourier Series with fundamental frequency $\frac{2\pi}{N}$

1.10.6.2.1.1 Remark 1

IDFT treats $x[n]$ as though it were N -periodic.

$$x[n] = \frac{1}{N} \sum_{k=0}^{N-1} \left(X[k] e^{j2\pi \frac{k}{N} n} \right) \quad (1.70)$$

where $n \in \{0, \dots, N-1\}$

Exercise 1.7

What about other values of n ?

(Solution on p. 97.)

1.10.6.2.1.2 Remark 2

Proof that the IDFT inverts the DFT for $n \in \{0, \dots, N-1\}$

$$\begin{aligned} \frac{1}{N} \sum_{k=0}^{N-1} \left(X[k] e^{j2\pi \frac{k}{N} n} \right) &= \frac{1}{N} \sum_{k=0}^{N-1} \left(\sum_{m=0}^{N-1} \left(x[m] e^{(-j)2\pi \frac{k}{N} m} e^{j2\pi \frac{k}{N} n} \right) \right) \\ &= ??? \end{aligned} \quad (1.71)$$

Example 1.23: Computing DFT

Given the following discrete-time signal (Figure 1.37) with $N = 4$, we will compute the DFT using two different methods (the DFT Formula and Sample DTFT):

⁸¹"The Complex Exponential" <<http://cnx.org/content/m10060/latest/>>

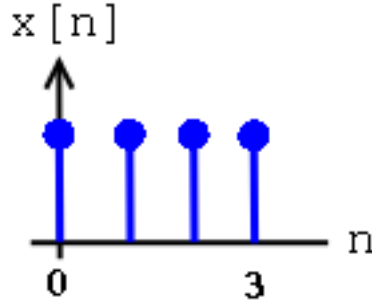


Figure 1.37

1. DFT Formula

$$\begin{aligned}
 X[k] &= \sum_{n=0}^{N-1} \left(x[n] e^{(-j)2\pi \frac{k}{N}n} \right) \\
 &= 1 + e^{(-j)2\pi \frac{k}{4}} + e^{(-j)2\pi \frac{k}{4}2} + e^{(-j)2\pi \frac{k}{4}3} \\
 &= 1 + e^{(-j)\frac{\pi}{2}k} + e^{(-j)\pi k} + e^{(-j)\frac{3}{2}\pi k}
 \end{aligned} \tag{1.72}$$

Using the above equation, we can solve and get the following results:

$$x[0] = 4$$

$$x[1] = 0$$

$$x[2] = 0$$

$$x[3] = 0$$

2. Sample DTFT. Using the same figure, Figure 1.37, we will take the DTFT of the signal and get the following equations:

$$\begin{aligned}
 X(\omega) &= \sum_{n=0}^3 \left(e^{(-j)\omega n} \right) \\
 &= \frac{1 - e^{(-j)4\omega}}{1 - e^{(-j)\omega}} \\
 &= ???
 \end{aligned} \tag{1.73}$$

Our sample points will be:

$$\omega_k = \frac{2\pi k}{4} = \frac{\pi}{2}k$$

where $k = \{0, 1, 2, 3\}$ (Figure 1.38).

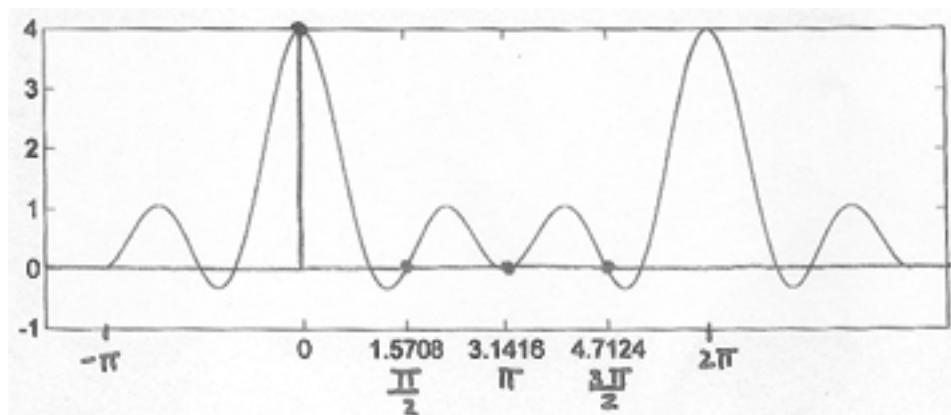


Figure 1.38

1.10.6.3 Periodicity of the DFT

DFT $X[k]$ consists of *samples* of DTFT, so $X(\omega)$, a 2π -periodic DTFT signal, can be converted to $X[k]$, an N -periodic DFT.

$$X[k] = \sum_{n=0}^{N-1} \left(x[n] e^{(-j)2\pi \frac{k}{N}n} \right) \quad (1.74)$$

where $e^{(-j)2\pi \frac{k}{N}n}$ is an N -periodic basis function (See Figure 1.39).

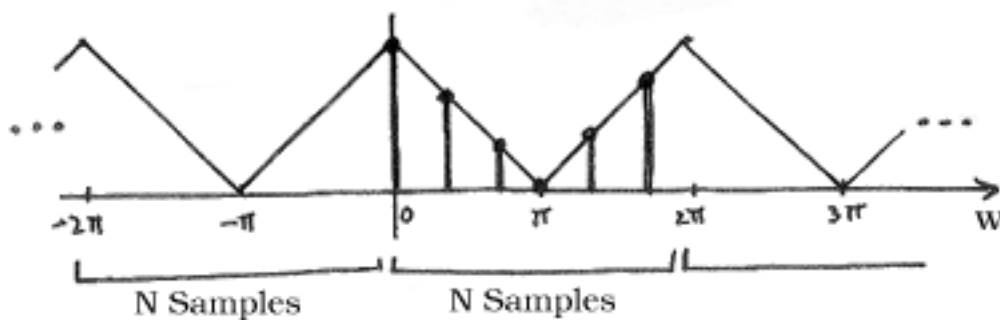
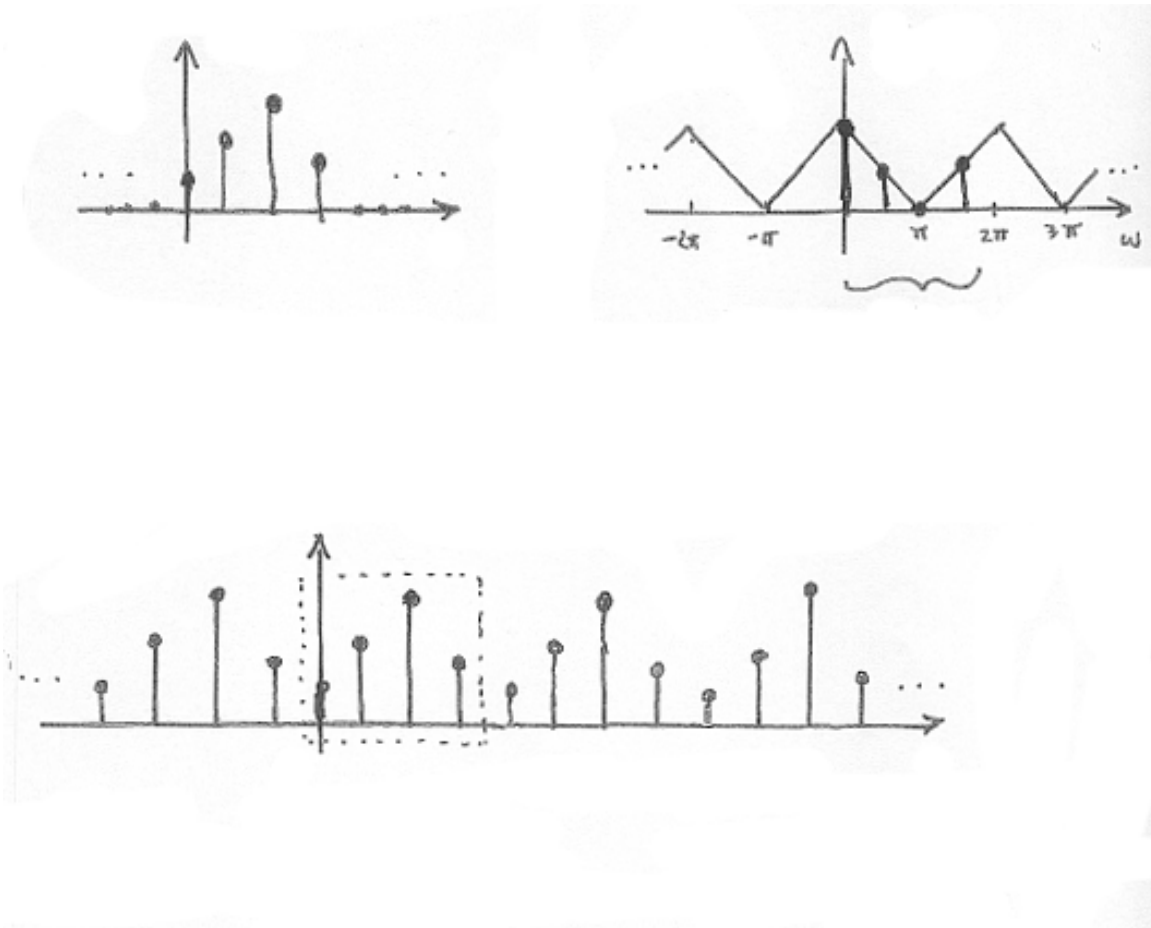


Figure 1.39

Also, recall,

$$\begin{aligned} x[n] &= \frac{1}{N} \sum_{k=0}^{N-1} \left(X[k] e^{j2\pi \frac{k}{N}n} \right) \\ &= \frac{1}{N} \sum_{k=0}^{N-1} \left(X[k] e^{j2\pi \frac{k}{N}(n+mN)} \right) \\ &= ??? \end{aligned} \quad (1.75)$$

Example 1.24: Illustration**Figure 1.40**

NOTE: When we deal with the DFT, we need to remember that, in effect, this treats the signal as an N -periodic sequence.

1.10.6.4 A Sampling Perspective

Think of sampling the continuous function $X(\omega)$, as depicted in Figure 1.41. $S(\omega)$ will represent the sampling function applied to $X(\omega)$ and is illustrated in Figure 1.41 as well. This will result in our discrete-time sequence, $X[k]$.

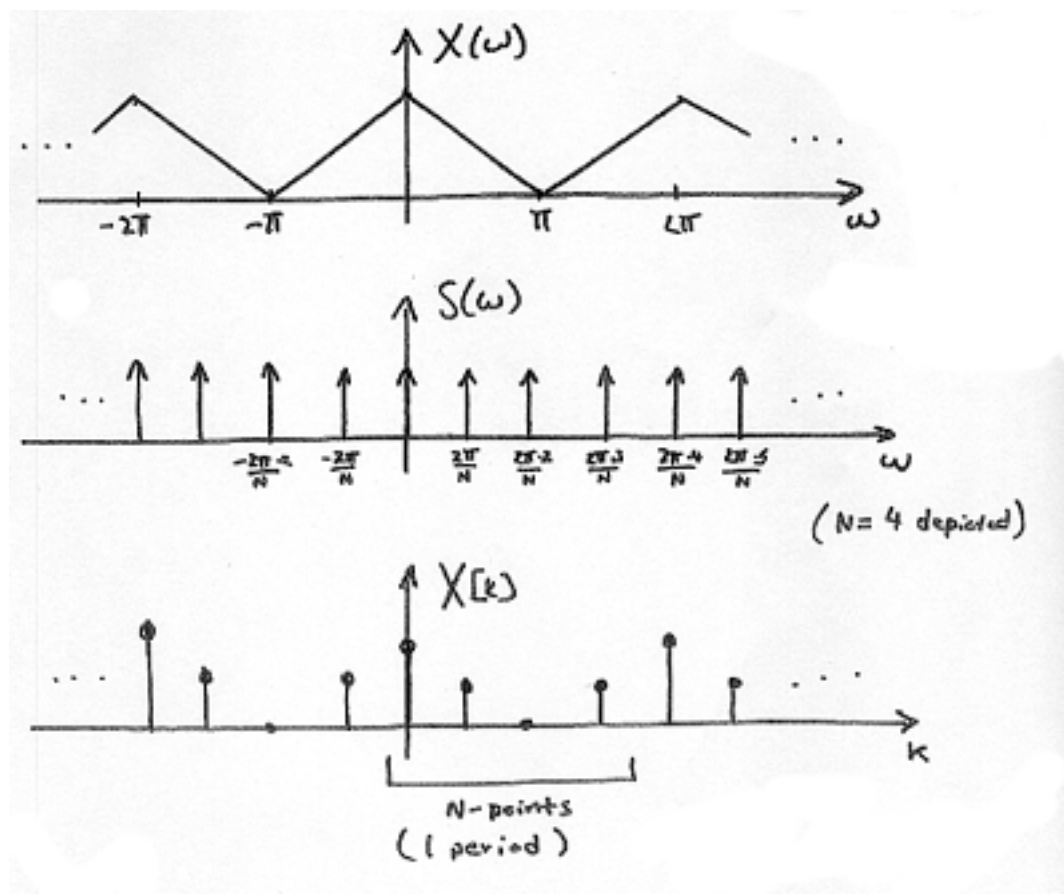


Figure 1.41

RECALL: Remember the multiplication in the frequency domain is equal to convolution in the time domain!

1.10.6.4.1 Inverse DTFT of $S(\omega)$

$$\sum_{k=-\infty}^{\infty} \left(\delta \left(\omega - \frac{2\pi k}{N} \right) \right) \quad (1.76)$$

Given the above equation, we can take the DTFT and get the following equation:

$$N \sum_{m=-\infty}^{\infty} (\delta[n - mN]) \equiv S[n] \quad (1.77)$$

Exercise 1.8

Why does (1.77) equal $S[n]$?

(Solution on p. 97.)

So, in the time-domain we have (Figure 1.42):

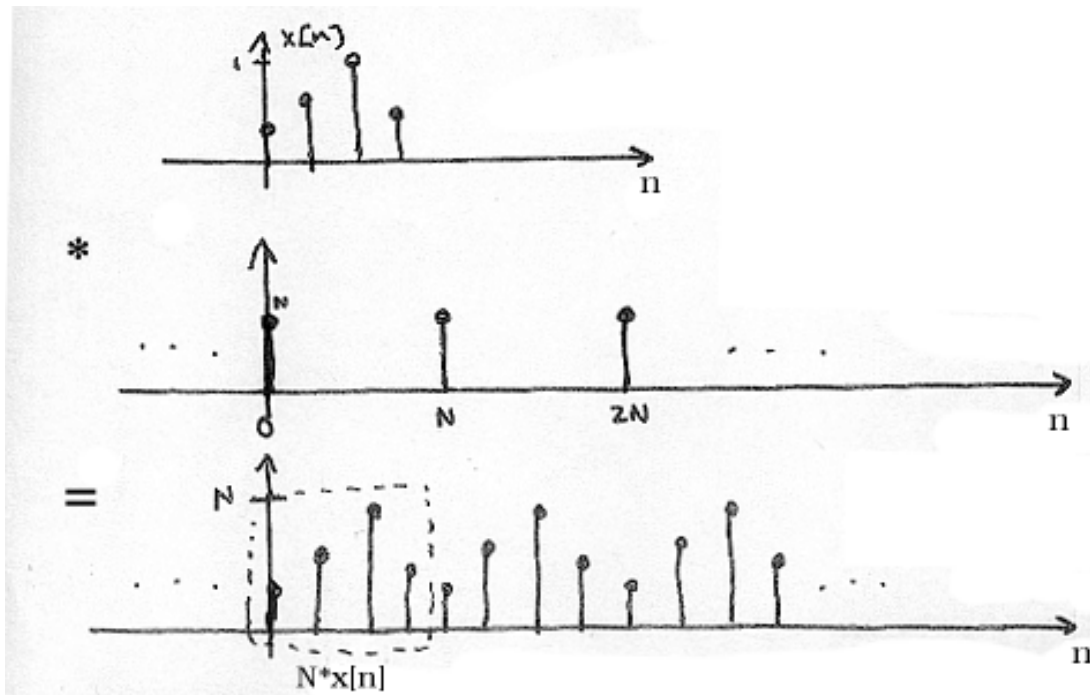


Figure 1.42

1.10.6.5 Connections

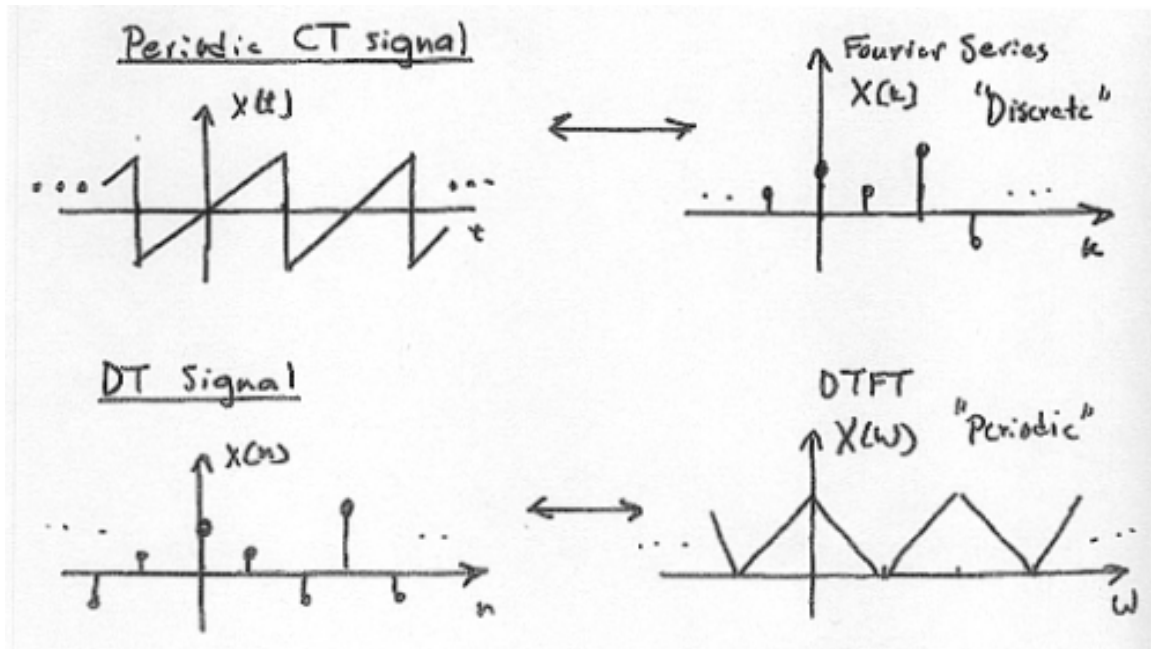


Figure 1.43

Combine signals in Figure 1.43 to get signals in Figure 1.44.

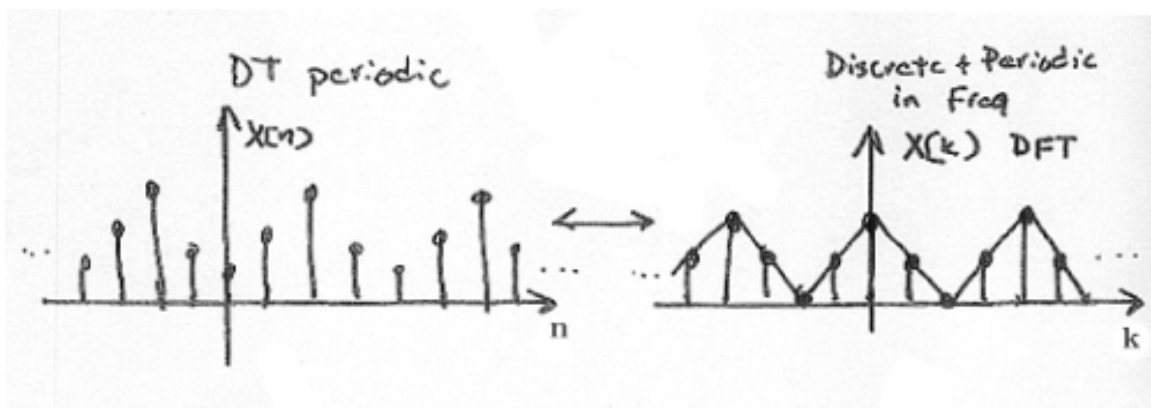


Figure 1.44

1.10.7 Discrete-Time Processing of CT Signals⁸²

1.10.7.1 DT Processing of CT Signals

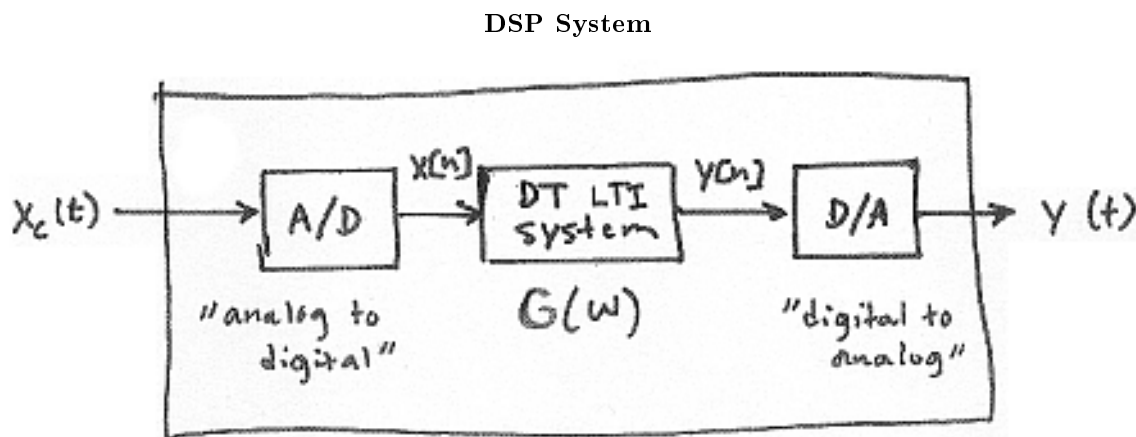


Figure 1.45

1.10.7.1.1 Analysis

$$Y_c(\Omega) = H_{LP}(\Omega) Y(\Omega T) \quad (1.78)$$

where we know that $Y(\omega) = X(\omega) G(\omega)$ and $G(\omega)$ is the frequency response of the DT LTI system. Also, remember that

$$\omega \equiv \Omega T$$

So,

$$Y_c(\Omega) = H_{LP}(\Omega) G(\Omega T) X(\Omega T) \quad (1.79)$$

where $Y_c(\Omega)$ and $H_{LP}(\Omega)$ are CTFTs and $G(\Omega T)$ and $X(\Omega T)$ are DTFTs.

RECALL:

$$X(\omega) = \frac{2\pi}{T} \sum_{k=-\infty}^{\infty} \left(X_c \left(\frac{\omega - 2\pi k}{T} \right) \right)$$

OR

$$X(\Omega T) = \frac{2\pi}{T} \sum_{k=-\infty}^{\infty} (X_c(\Omega - k\Omega_s))$$

Therefore our final output signal, $Y_c(\Omega)$, will be:

$$Y_c(\Omega) = H_{LP}(\Omega) G(\Omega T) \left(\frac{2\pi}{T} \sum_{k=-\infty}^{\infty} (X_c(\Omega - k\Omega_s)) \right) \quad (1.80)$$

⁸²This content is available online at <<http://cnx.org/content/m10993/2.2/>>.

Now, if $X_c(\Omega)$ is bandlimited to $[-(\frac{\Omega_s}{2}), \frac{\Omega_s}{2}]$ and we use the usual lowpass reconstruction filter in the D/A, Figure 1.46:

Image not finished

Figure 1.46

Then,

$$Y_c(\Omega) = \begin{cases} G(\Omega T) X_c(\Omega) & \text{if } |\Omega| < \frac{\Omega_s}{2} \\ 0 & \text{otherwise} \end{cases} \quad (1.81)$$

1.10.7.1.2 Summary

For bandlimited signals sampled at or above the Nyquist rate, we can relate the input and output of the DSP system by:

$$Y_c(\Omega) = G_{eff}(\Omega) X_c(\Omega) \quad (1.82)$$

where

$$G_{eff}(\Omega) = \begin{cases} G(\Omega T) & \text{if } |\Omega| < \frac{\Omega_s}{2} \\ 0 & \text{otherwise} \end{cases}$$

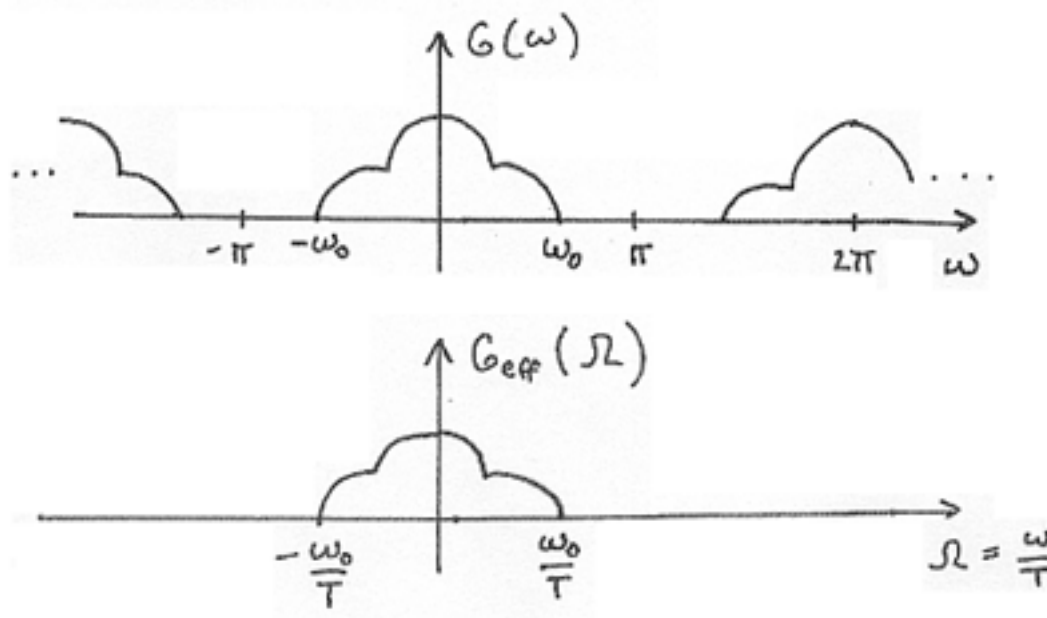


Figure 1.47

1.10.7.1.2.1 Note

$G_{eff}(\Omega)$ is LTI if and only if the following two conditions are satisfied:

1. $G(\omega)$ is LTI (in DT).
2. $X_c(T)$ is bandlimited and sampling rate equal to or greater than Nyquist. For example, if we had a simple pulse described by

$$X_c(t) = u(t - T_0) - u(t - T_1)$$

where $T_1 > T_0$. If the sampling period $T > T_1 - T_0$, then some samples might "miss" the pulse while others might not be "missed." This is what we term **time-varying behavior**.

Example 1.25

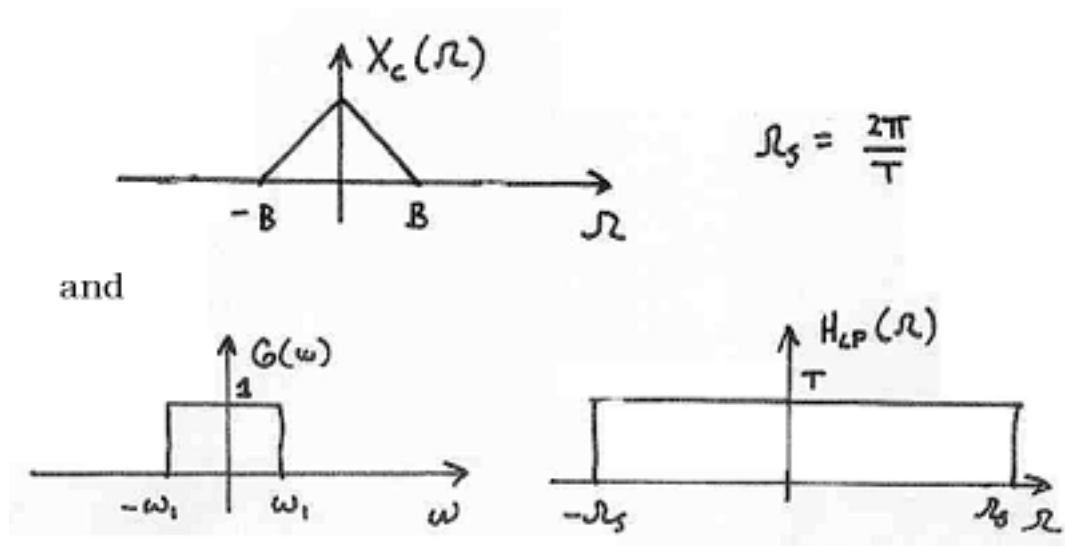


Figure 1.48

If $\frac{2\pi}{T} > 2B$ and $\omega_1 < BT$, determine and sketch $Y_c(\Omega)$ using Figure 1.48.

1.10.7.2 Application: 60Hz Noise Removal

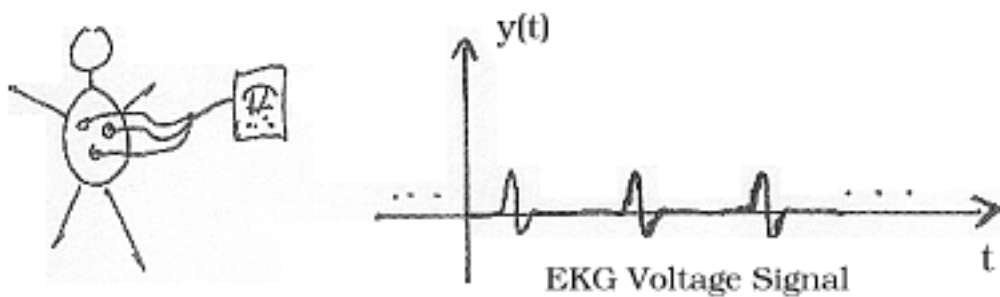


Figure 1.49

Unfortunately, in real-world situations electrodes also pick up ambient 60 Hz signals from lights, computers, *etc.*. In fact, usually this "60 Hz noise" is much greater in amplitude than the EKG signal shown in Figure 1.49. Figure 1.50 shows the EKG signal; it is barely noticeable as it has become overwhelmed by noise.

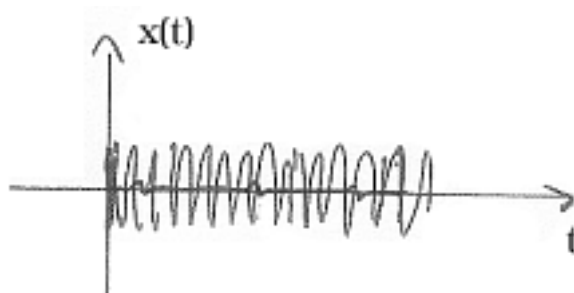


Figure 1.50: Our EKG signal, $y(t)$, is overwhelmed by noise.

1.10.7.2.1 DSP Solution

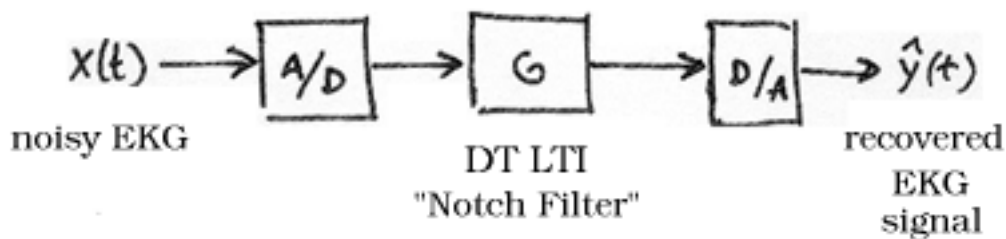


Figure 1.51

Image not finished

Figure 1.52

1.10.7.2.2 Sampling Period/Rate

First we must note that $|Y(\Omega)|$ is **bandlimited** to ± 60 Hz. Therefore, the minimum rate should be 120 Hz. In order to get the best results we should set

$$f_s = 240 \text{ Hz}$$

$$\Omega_s = 2\pi \left(240 \frac{\text{rad}}{\text{s}} \right)$$

Image not finished

Figure 1.53

1.10.7.2.3 Digital Filter

Therefore, we want to design a digital filter that will remove the 60Hz component and preserve the rest.

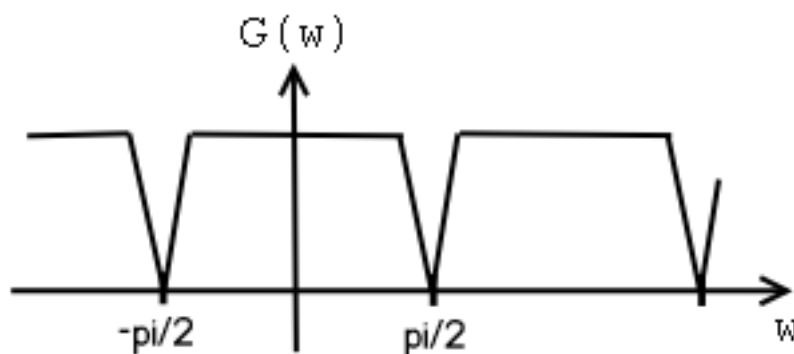


Figure 1.54

1.11 Z-Transform

1.11.1 Difference Equation⁸³

1.11.1.1 Introduction

One of the most important concepts of DSP is to be able to properly represent the input/output relationship to a given LTI system. A linear constant-coefficient **difference equation** (LCCDE) serves as a way to express just this relationship in a discrete-time system. Writing the sequence of inputs and outputs, which represent the characteristics of the LTI system, as a difference equation help in understanding and manipulating a system.

⁸³This content is available online at <http://cnx.org/content/m10595/2.5/>.

Definition 1: difference equation

An equation that shows the relationship between consecutive values of a sequence and the differences among them. They are often rearranged as a recursive formula so that a systems output can be computed from the input signal and past outputs.

Example

$$y[n] + 7y[n-1] + 2y[n-2] = x[n] - 4x[n-1] \quad (1.83)$$

1.11.1.2 General Formulas from the Difference Equation

As stated briefly in the definition above, a difference equation is a very useful tool in describing and calculating the output of the system described by the formula for a given sample n . The key property of the difference equation is its ability to help easily find the transform, $H(z)$, of a system. In the following two subsections, we will look at the general form of the difference equation and the general conversion to a z-transform directly from the difference equation.

1.11.1.2.1 Difference Equation

The general form of a linear, constant-coefficient difference equation (LCCDE), is shown below:

$$\sum_{k=0}^N (a_k y[n-k]) = \sum_{k=0}^M (b_k x[n-k]) \quad (1.84)$$

We can also write the general form to easily express a recursive output, which looks like this:

$$y[n] = - \left(\sum_{k=1}^N (a_k y[n-k]) \right) + \sum_{k=0}^M (b_k x[n-k]) \quad (1.85)$$

From this equation, note that $y[n-k]$ represents the outputs and $x[n-k]$ represents the inputs. The value of N represents the **order** of the difference equation and corresponds to the memory of the system being represented. Because this equation relies on past values of the output, in order to compute a numerical solution, certain past outputs, referred to as the **initial conditions**, must be known.

1.11.1.2.2 Conversion to Z-Transform

Using the above formula, (1.84), we can easily generalize the **transfer function**, $H(z)$, for any difference equation. Below are the steps taken to convert any difference equation into its transfer function, i.e. z-transform. The first step involves taking the Fourier Transform⁸⁴ of all the terms in (1.84). Then we use the linearity property to pull the transform inside the summation and the time-shifting property of the z-transform to change the time-shifting terms to exponentials. Once this is done, we arrive at the following equation: $a_0 = 1$.

$$Y(z) = - \left(\sum_{k=1}^N (a_k Y(z) z^{-k}) \right) + \sum_{k=0}^M (b_k X(z) z^{-k}) \quad (1.86)$$

$$\begin{aligned} H(z) &= \frac{Y(z)}{X(z)} \\ &= \frac{\sum_{k=0}^M (b_k z^{-k})}{1 + \sum_{k=1}^N (a_k z^{-k})} \end{aligned} \quad (1.87)$$

⁸⁴"Derivation of the Fourier Transform" <<http://cnx.org/content/m0046/latest/>>

1.11.1.2.3 Conversion to Frequency Response

Once the z-transform has been calculated from the difference equation, we can go one step further to define the frequency response of the system, or filter, that is being represented by the difference equation.

NOTE: Remember that the reason we are dealing with these formulas is to be able to aid us in filter design. A LCCDE is one of the easiest ways to represent FIR filters. By being able to find the frequency response, we will be able to look at the basic properties of any filter represented by a simple LCCDE.

Below is the general formula for the frequency response of a z-transform. The conversion is simple a matter of taking the z-transform formula, $H(z)$, and replacing every instance of z with e^{jw} .

$$\begin{aligned} H(w) &= H(z) \big|_{z=e^{jw}} \\ &= \frac{\sum_{k=0}^M (b_k e^{-(jwk)})}{\sum_{k=0}^N (a_k e^{-(jwk)})} \end{aligned} \quad (1.88)$$

Once you understand the derivation of this formula, look at the module concerning Filter Design from the Z-Transform⁸⁵ for a look into how all of these ideas of the Z-transform (Section 1.11.2), Difference Equation, and Pole/Zero Plots (Section 1.11.4) play a role in filter design.

1.11.1.3 Example

Example 1.26: Finding Difference Equation

Below is a basic example showing the opposite of the steps above: given a transfer function one can easily calculate the systems difference equation.

$$H(z) = \frac{(z+1)^2}{(z-\frac{1}{2})(z+\frac{3}{4})} \quad (1.89)$$

Given this transfer function of a time-domain filter, we want to find the difference equation. To begin with, expand both polynomials and divide them by the highest order z .

$$\begin{aligned} H(z) &= \frac{(z+1)(z+1)}{(z-\frac{1}{2})(z+\frac{3}{4})} \\ &= \frac{z^2+2z+1}{z^2+2z+1-\frac{3}{8}} \\ &= \frac{1+2z^{-1}+z^{-2}}{1+\frac{1}{4}z^{-1}-\frac{3}{8}z^{-2}} \end{aligned} \quad (1.90)$$

From this transfer function, the coefficients of the two polynomials will be our a_k and b_k values found in the general difference equation formula, (1.84). Using these coefficients and the above form of the transfer function, we can easily write the difference equation:

$$x[n] + 2x[n-1] + x[n-2] = y[n] + \frac{1}{4}y[n-1] - \frac{3}{8}y[n-2] \quad (1.91)$$

In our final step, we can rewrite the difference equation in its more common form showing the recursive nature of the system.

$$y[n] = x[n] + 2x[n-1] + x[n-2] + \frac{-1}{4}y[n-1] + \frac{3}{8}y[n-2] \quad (1.92)$$

⁸⁵"Filter Design using the Pole/Zero Plot of a Z-Transform" <<http://cnx.org/content/m10548/latest/>>

1.11.1.4 Solving a LCCDE

In order for a linear constant-coefficient difference equation to be useful in analyzing a LTI system, we must be able to find the systems output based upon a known input, $x(n)$, and a set of initial conditions. Two common methods exist for solving a LCCDE: the **direct method** and the **indirect method**, the later being based on the z-transform. Below we will briefly discuss the formulas for solving a LCCDE using each of these methods.

1.11.1.4.1 Direct Method

The final solution to the output based on the direct method is the sum of two parts, expressed in the following equation:

$$y(n) = y_h(n) + y_p(n) \quad (1.93)$$

The first part, $y_h(n)$, is referred to as the **homogeneous solution** and the second part, $y_p(n)$, is referred to as **particular solution**. The following method is very similar to that used to solve many differential equations, so if you have taken a differential calculus course or used differential equations before then this should seem very familiar.

1.11.1.4.1.1 Homogeneous Solution

We begin by assuming that the input is zero, $x(n) = 0$. Now we simply need to solve the homogeneous difference equation:

$$\sum_{k=0}^N (a_k y[n-k]) = 0 \quad (1.94)$$

In order to solve this, we will make the assumption that the solution is in the form of an exponential. We will use lambda, λ , to represent our exponential terms. We now have to solve the following equation:

$$\sum_{k=0}^N (a_k \lambda^{n-k}) = 0 \quad (1.95)$$

We can expand this equation out and factor out all of the lambda terms. This will give us a large polynomial in parenthesis, which is referred to as the **characteristic polynomial**. The roots of this polynomial will be the key to solving the homogeneous equation. If there are all distinct roots, then the general solution to the equation will be as follows:

$$y_h(n) = C_1(\lambda_1)^n + C_2(\lambda_2)^n + \cdots + C_N(\lambda_N)^n \quad (1.96)$$

However, if the characteristic equation contains multiple roots then the above general solution will be slightly different. Below we have the modified version for an equation where λ_1 has K multiple roots:

$$y_h(n) = C_1(\lambda_1)^n + C_1 n(\lambda_1)^n + C_1 n^2(\lambda_1)^n + \cdots + C_1 n^{K-1}(\lambda_1)^n + C_2(\lambda_2)^n + \cdots + C_N(\lambda_N)^n \quad (1.97)$$

1.11.1.4.1.2 Particular Solution

The particular solution, $y_p(n)$, will be any solution that will solve the general difference equation:

$$\sum_{k=0}^N (a_k y_p(n-k)) = \sum_{k=0}^M (b_k x(n-k)) \quad (1.98)$$

In order to solve, our guess for the solution to $y_p(n)$ will take on the form of the input, $x(n)$. After guessing at a solution to the above equation involving the particular solution, one only needs to plug the solution into the difference equation and solve it out.

1.11.1.4.2 Indirect Method

The indirect method utilizes the relationship between the difference equation and z-transform, discussed earlier (Section 1.11.1.2: General Formulas from the Difference Equation), to find a solution. The basic idea is to convert the difference equation into a z-transform, as described above (Section 1.11.1.2.2: Conversion to Z-Transform), to get the resulting output, $Y(z)$. Then by inverse transforming this and using partial-fraction expansion, we can arrive at the solution.

1.11.2 The Z Transform: Definition⁸⁶

1.11.2.1 Basic Definition of the Z-Transform

The **z-transform** of a sequence is defined as

$$X(z) = \sum_{n=-\infty}^{\infty} (x[n] z^{-n}) \quad (1.99)$$

Sometimes this equation is referred to as the **bilateral z-transform**. At times the z-transform is defined as

$$X(z) = \sum_{n=0}^{\infty} (x[n] z^{-n}) \quad (1.100)$$

which is known as the **unilateral z-transform**.

There is a close relationship between the z-transform and the **Fourier transform** of a discrete time signal, which is defined as

$$X(e^{j\omega}) = \sum_{n=-\infty}^{\infty} (x[n] e^{-j\omega n}) \quad (1.101)$$

Notice that that when the z^{-n} is replaced with $e^{-j\omega n}$ the z-transform reduces to the Fourier Transform. When the Fourier Transform exists, $z = e^{j\omega}$, which is to have the magnitude of z equal to unity.

1.11.2.2 The Complex Plane

In order to get further insight into the relationship between the Fourier Transform and the Z-Transform it is useful to look at the complex plane or **z-plane**. Take a look at the complex plane:

⁸⁶This content is available online at <<http://cnx.org/content/m10549/2.9/>>.

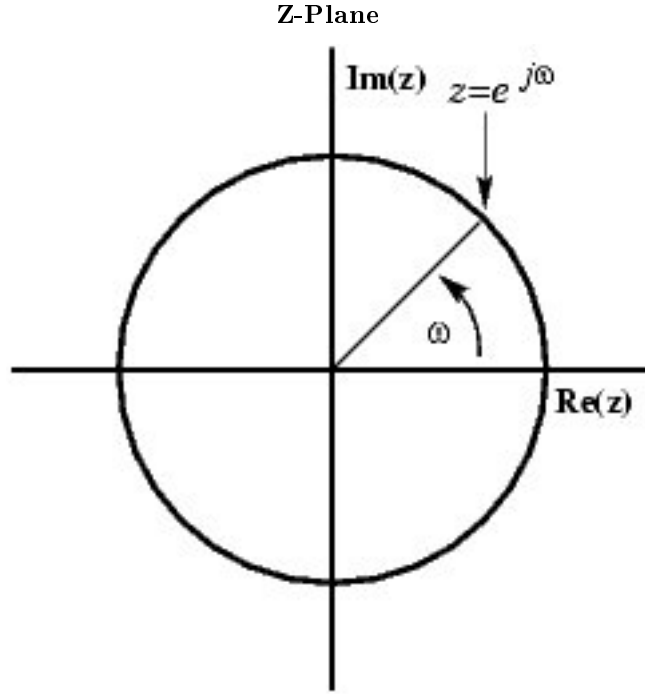


Figure 1.55

The Z-plane is a complex plane with an imaginary and real axis referring to the complex-valued variable z . The position on the complex plane is given by $re^{j\omega}$, and the angle from the positive, real axis around the plane is denoted by ω . $X(z)$ is defined everywhere on this plane. $X(e^{j\omega})$ on the other hand is defined only where $|z| = 1$, which is referred to as the unit circle. So for example, $\omega = 0$ at $z = 1$ and $\omega = \pi$ at $z = -1$. This is useful because, by representing the Fourier transform as the z-transform on the unit circle, the periodicity of Fourier transform is easily seen.

1.11.2.3 Region of Convergence

The region of convergence, known as the **ROC**, is important to understand because it defines the region where the z-transform exists. The ROC for a given $x[n]$, is defined as the range of z for which the z-transform converges. Since the z-transform is a **power series**, it converges when $x[n]z^{-n}$ is absolutely summable. Stated differently,

$$\sum_{n=-\infty}^{\infty} (|x[n]z^{-n}|) < \infty \quad (1.102)$$

must be satisfied for convergence. This is best illustrated by looking at the different ROC's of the z-transforms of $\alpha^n u[n]$ and $\alpha^n u[n-1]$.

Example 1.27

For

$$x[n] = \alpha^n u[n] \quad (1.103)$$

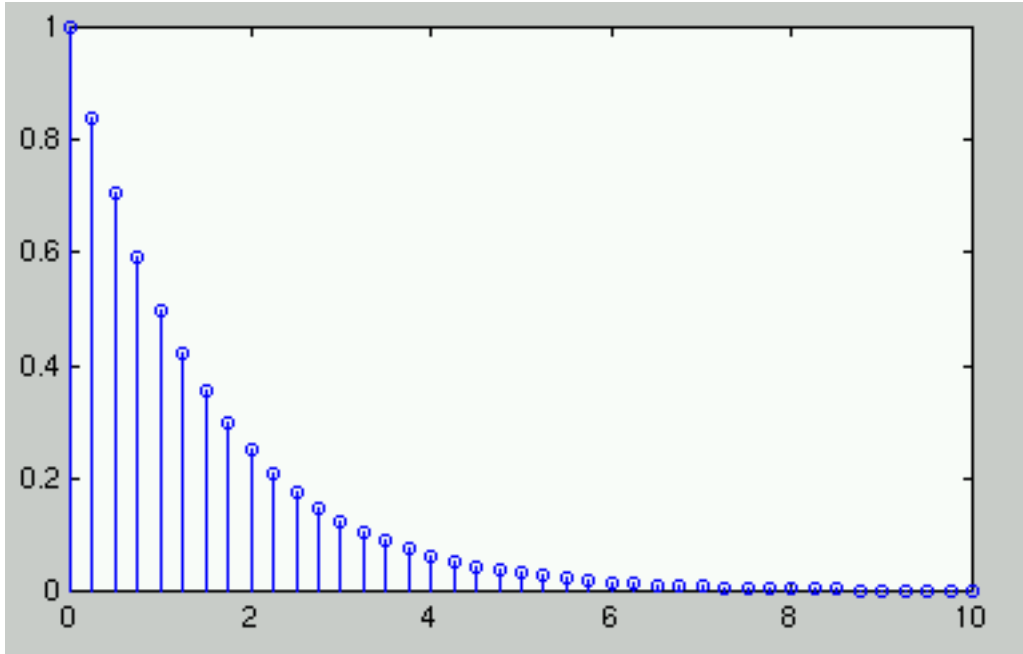


Figure 1.56: $x[n] = \alpha^n u[n]$ where $\alpha = 0.5$.

$$\begin{aligned}
 X(z) &= \sum_{n=-\infty}^{\infty} (x[n] z^{-n}) \\
 &= \sum_{n=-\infty}^{\infty} (\alpha^n u[n] z^{-n}) \\
 &= \sum_{n=0}^{\infty} (\alpha^n z^{-n}) \\
 &= \sum_{n=0}^{\infty} ((\alpha z^{-1})^n)
 \end{aligned} \tag{1.104}$$

This sequence is an example of a right-sided exponential sequence because it is nonzero for $n \geq 0$. It only converges when $|\alpha z^{-1}| < 1$. When it converges,

$$\begin{aligned}
 X(z) &= \frac{1}{1 - \alpha z^{-1}} \\
 &= \frac{z}{z - \alpha}
 \end{aligned} \tag{1.105}$$

If $|\alpha z^{-1}| \geq 1$, then the series, $\sum_{n=0}^{\infty} ((\alpha z^{-1})^n)$ does not converge. Thus the ROC is the range of values where

$$|\alpha z^{-1}| < 1 \tag{1.106}$$

or, equivalently,

$$|z| > |\alpha| \tag{1.107}$$

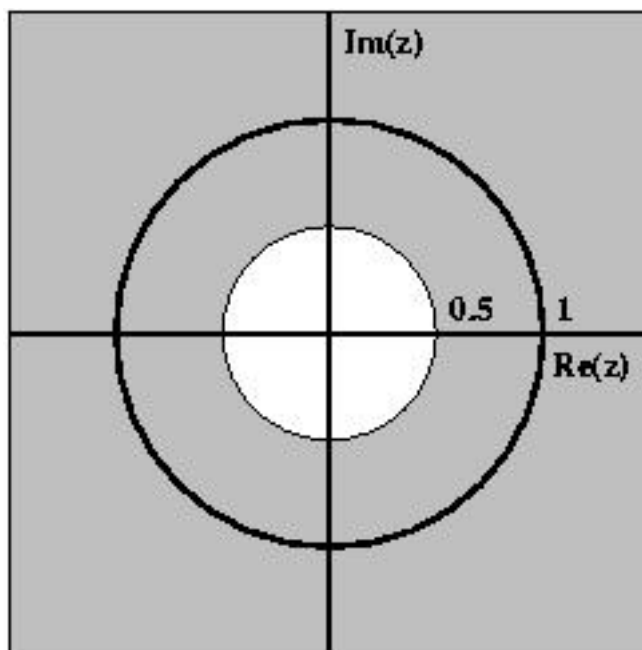


Figure 1.57: ROC for $x[n] = \alpha^n u[n]$ where $\alpha = 0.5$

Example 1.28

For

$$x[n] = -(\alpha^n) u[-n-1] \quad (1.108)$$

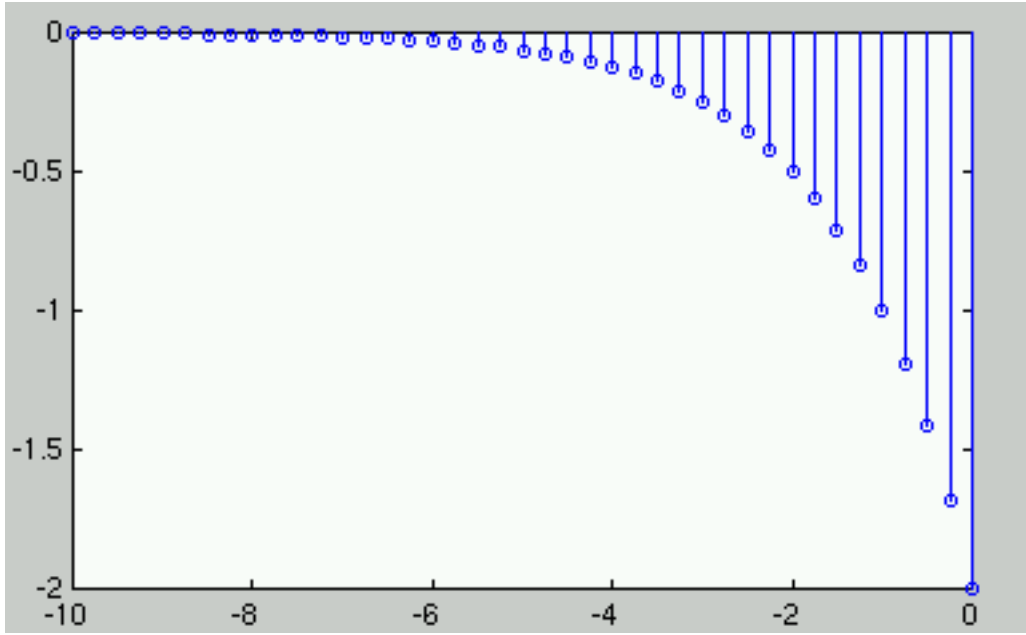


Figure 1.58: $x[n] = -(\alpha^n)u[-n-1]$ where $\alpha = 0.5$.

$$\begin{aligned}
 X(z) &= \sum_{n=-\infty}^{\infty} (x[n] z^{-n}) \\
 &= \sum_{n=-\infty}^{\infty} ((-\alpha^n) u[-n-1] z^{-n}) \\
 &= - \left(\sum_{n=-\infty}^{-1} (\alpha^n z^{-n}) \right) \\
 &= - \left(\sum_{n=-\infty}^{-1} ((\alpha^{-1} z)^{-n}) \right) \\
 &= - \left(\sum_{n=1}^{\infty} ((\alpha^{-1} z)^n) \right) \\
 &= 1 - \sum_{n=0}^{\infty} ((\alpha^{-1} z)^n)
 \end{aligned} \tag{1.109}$$

The ROC in this case is the range of values where

$$|\alpha^{-1} z| < 1 \tag{1.110}$$

or, equivalently,

$$|z| < |\alpha| \tag{1.111}$$

If the ROC is satisfied, then

$$\begin{aligned}
 X(z) &= 1 - \frac{1}{1 - \alpha^{-1} z} \\
 &= \frac{z}{z - \alpha}
 \end{aligned} \tag{1.112}$$

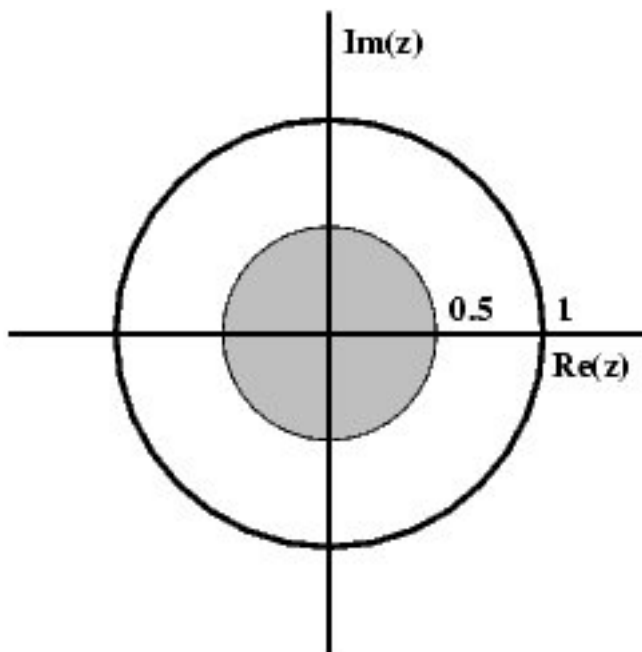


Figure 1.59: ROC for $x[n] = -(\alpha^n)u[-n-1]$

1.11.3 Table of Common **z**-Transforms⁸⁷

The table below provides a number of unilateral and bilateral **z-transforms** (**Section 1.11.2**). The table also specifies the region of convergence⁸⁸.

NOTE: The notation for z found in the table below may differ from that found in other tables. For example, the basic z -transform of $u[n]$ can be written as either of the following two expressions, which are equivalent:

$$\frac{z}{z-1} = \frac{1}{1-z^{-1}} \quad (1.113)$$

⁸⁷This content is available online at <http://cnx.org/content/m10119/2.13/>.

⁸⁸"Region of Convergence for the Z-transform" <http://cnx.org/content/m10622/latest/>

Signal	Z-Transform	ROC
$\delta[n - k]$	z^{-k}	All z
$u[n]$	$\frac{z}{z-1}$	$ z > 1$
$-(u[-n - 1])$	$\frac{z}{z-1}$	$ z < 1$
$nu[n]$	$\frac{z}{(z-1)^2}$	$ z > 1$
$n^2u[n]$	$\frac{z(z+1)}{(z-1)^3}$	$ z > 1$
$n^3u[n]$	$\frac{z(z^2+4z+1)}{(z-1)^4}$	$ z > 1$
$(-\alpha^n)u[-n - 1]$	$\frac{z}{z-\alpha}$	$ z < \alpha $
$\alpha^n u[n]$	$\frac{z}{z-\alpha}$	$ z > \alpha $
$n\alpha^n u[n]$	$\frac{\alpha z}{(z-\alpha)^2}$	$ z > \alpha $
$n^2\alpha^n u[n]$	$\frac{\alpha z(z+\alpha)}{(z-\alpha)^3}$	$ z > \alpha $
$\frac{\prod_{k=1}^m (n-k+1)}{\alpha^m m!} \alpha^n u[n]$	$\frac{z}{(z-\alpha)^{m+1}}$	
$\gamma^n \cos(\alpha n) u[n]$	$\frac{z(z-\gamma \cos(\alpha))}{z^2 - (2\gamma \cos(\alpha))z + \gamma^2}$	$ z > \gamma $
$\gamma^n \sin(\alpha n) u[n]$	$\frac{z\gamma \sin(\alpha)}{z^2 - (2\gamma \cos(\alpha))z + \gamma^2}$	$ z > \gamma $

1.11.4 Understanding Pole/Zero Plots on the Z-Plane⁸⁹

1.11.4.1 Introduction to Poles and Zeros of the Z-Transform

Once the Z-transform of a system has been determined, one can use the information contained in function's polynomials to graphically represent the function and easily observe many defining characteristics. The Z-transform will have the below structure, based on Rational Functions⁹⁰:

$$X(z) = \frac{P(z)}{Q(z)} \quad (1.114)$$

The two polynomials, $P(z)$ and $Q(z)$, allow us to find the poles and zeros⁹¹ of the Z-Transform.

Definition 2: zeros

1. The value(s) for z where $P(z) = 0$.
2. The complex frequencies that make the overall gain of the filter transfer function zero.

Definition 3: poles

1. The value(s) for z where $Q(z) = 0$.
2. The complex frequencies that make the overall gain of the filter transfer function infinite.

Example 1.29

Below is a simple transfer function with the poles and zeros shown below it.

$$H(z) = \frac{z + 1}{\left(z - \frac{1}{2}\right)\left(z + \frac{3}{4}\right)}$$

The zeros are: $\{-1\}$

The poles are: $\left\{\frac{1}{2}, -\left(\frac{3}{4}\right)\right\}$

⁸⁹This content is available online at <<http://cnx.org/content/m10556/2.8/>>.

⁹⁰"Rational Functions" <<http://cnx.org/content/m10593/latest/>>

⁹¹"Poles and Zeros" <<http://cnx.org/content/m10112/latest/>>

1.11.4.2 The Z-Plane

Once the poles and zeros have been found for a given Z-Transform, they can be plotted onto the Z-Plane. The Z-plane is a complex plane with an imaginary and real axis referring to the complex-valued variable z . The position on the complex plane is given by $re^{j\theta}$ and the angle from the positive, real axis around the plane is denoted by θ . When mapping poles and zeros onto the plane, poles are denoted by an "x" and zeros by an "o". The below figure shows the Z-Plane, and examples of plotting zeros and poles onto the plane can be found in the following section.

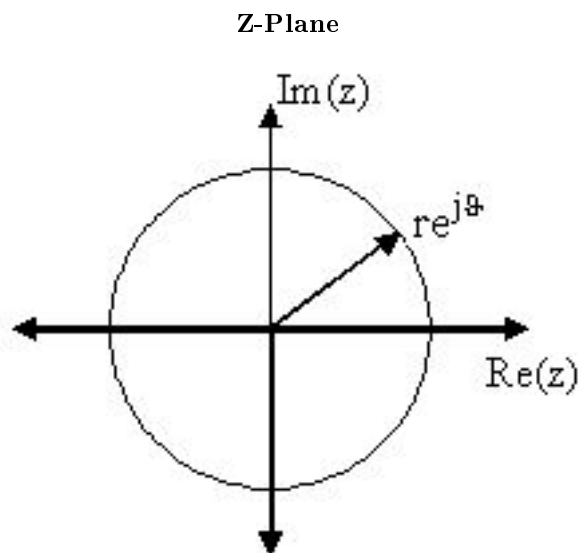


Figure 1.60

1.11.4.3 Examples of Pole/Zero Plots

This section lists several examples of finding the poles and zeros of a transfer function and then plotting them onto the Z-Plane.

Example 1.30: Simple Pole/Zero Plot

$$H(z) = \frac{z}{\left(z - \frac{1}{2}\right)\left(z + \frac{3}{4}\right)}$$

The zeros are: $\{0\}$

The poles are: $\left\{\frac{1}{2}, -\left(\frac{3}{4}\right)\right\}$

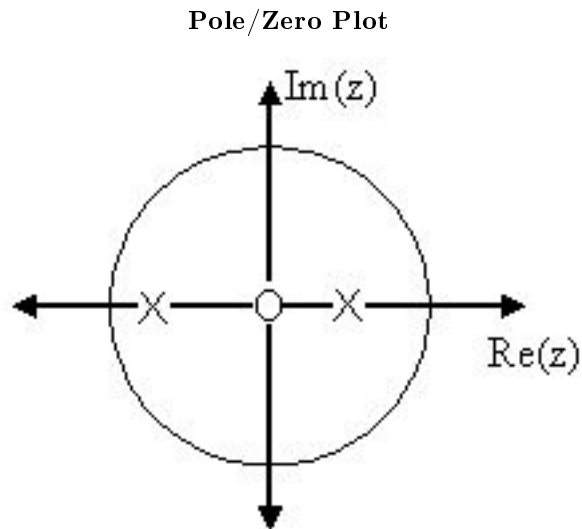


Figure 1.61: Using the zeros and poles found from the transfer function, the one zero is mapped to zero and the two poles are placed at $\frac{1}{2}$ and $-\frac{3}{4}$

Example 1.31: Complex Pole/Zero Plot

$$H(z) = \frac{(z - j)(z + j)}{(z - (\frac{1}{2} - \frac{1}{2}j))(z - (\frac{1}{2} + \frac{1}{2}j))}$$

The zeros are: $\{j, -j\}$

The poles are: $\{-1, \frac{1}{2} + \frac{1}{2}j, \frac{1}{2} - \frac{1}{2}j\}$

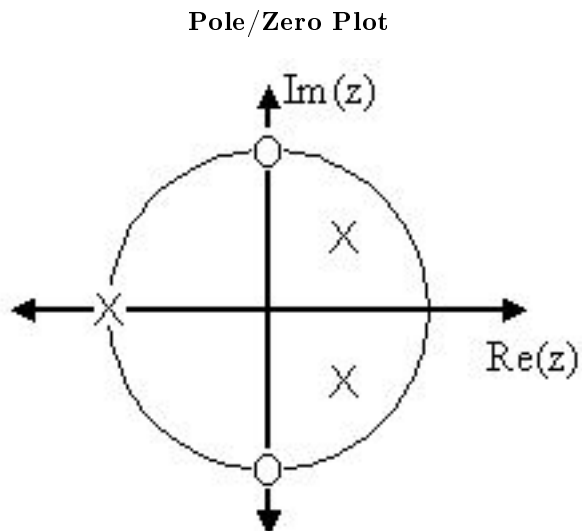


Figure 1.62: Using the zeros and poles found from the transfer function, the zeros are mapped to $\pm j$, and the poles are placed at -1 , $\frac{1}{2} + \frac{1}{2}j$ and $\frac{1}{2} - \frac{1}{2}j$

MATLAB - If access to MATLAB is readily available, then you can use its functions to easily create pole/zero plots. Below is a short program that plots the poles and zeros from the above example onto the Z-Plane.

```
% Set up vector for zeros
z = [j ; -j];

% Set up vector for poles
p = [-1 ; .5+.5j ; .5-.5j];

figure(1);
zplane(z,p);
title('Pole/Zero Plot for Complex Pole/Zero Plot Example');
```

1.11.4.4 Pole/Zero Plot and Region of Convergence

The region of convergence (ROC) for $X(z)$ in the complex Z-plane can be determined from the pole/zero plot. Although several regions of convergence may be possible, where each one corresponds to a different impulse response, there are some choices that are more practical. A ROC can be chosen to make the transfer function causal and/or stable depending on the pole/zero plot.

Filter Properties from ROC

- If the ROC extends outward from the outermost pole, then the system is **causal**.
- If the ROC includes the unit circle, then the system is **stable**.

Below is a pole/zero plot with a possible ROC of the Z-transform in the Simple Pole/Zero Plot (Example 1.30: Simple Pole/Zero Plot) discussed earlier. The shaded region indicates the ROC chosen for the filter. From this figure, we can see that the filter will be both causal and stable since the above listed conditions are both met.

Example 1.32

$$H(z) = \frac{z}{\left(z - \frac{1}{2}\right)\left(z + \frac{3}{4}\right)}$$

Region of Convergence for the Pole/Zero Plot

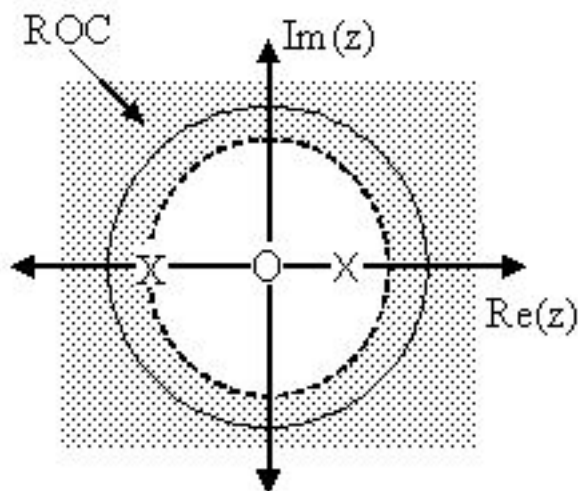


Figure 1.63: The shaded area represents the chosen ROC for the transfer function.

1.11.4.5 Frequency Response and the Z-Plane

The reason it is helpful to understand and create these pole/zero plots is due to their ability to help us easily design a filter. Based on the location of the poles and zeros, the magnitude response of the filter can be quickly understood. Also, by starting with the pole/zero plot, one can design a filter and obtain its transfer function very easily. Refer to this module⁹² for information on the relationship between the pole/zero plot and the frequency response.

1.12 Random Signals and Processes

1.12.1 Introduction to Random Signals and Processes⁹³

Before now, you have probably dealt strictly with the theory behind signals and systems, as well as look at some the basic characteristics of signals⁹⁴ and systems⁹⁵. In doing so you have developed an important foundation; however, most electrical engineers do not get to work in this type of fantasy world. In many cases the signals of interest are very complex due to the randomness of the world around them, which leaves them noisy and often corrupted. This often causes the information contained in the signal to be hidden and distorted. For this reason, it is important to understand these random signals and how to recover the necessary information.

1.12.1.1 Signals: Deterministic vs. Stochastic

For this study of signals and systems, we will divide signals into two groups: those that have a fixed behavior and those that change randomly. As most of you have probably already dealt with the first type, we will

⁹²"Filter Design using the Pole/Zero Plot of a Z-Transform" <<http://cnx.org/content/m10548/latest/>>

⁹³This content is available online at <<http://cnx.org/content/m10649/2.2/>>.

⁹⁴"Signal Classifications and Properties" <<http://cnx.org/content/m10057/latest/>>

⁹⁵"System Classifications and Properties" <<http://cnx.org/content/m10084/latest/>>

focus on introducing you to random signals. Also, note that we will be dealing strictly with discrete-time signals since they are the signals we deal with in DSP and most real-world computations, but these same ideas apply to continuous-time signals.

1.12.1.1.1 Deterministic Signals

Most introductions to signals and systems deal strictly with **deterministic signals**. Each value of these signals are fixed and can be determined by a mathematical expression, rule, or table. Because of this, future values of any deterministic signal can be calculated from past values. For this reason, these signals are relatively easy to analyze as they do not change, and we can make accurate assumptions about their past and future behavior.

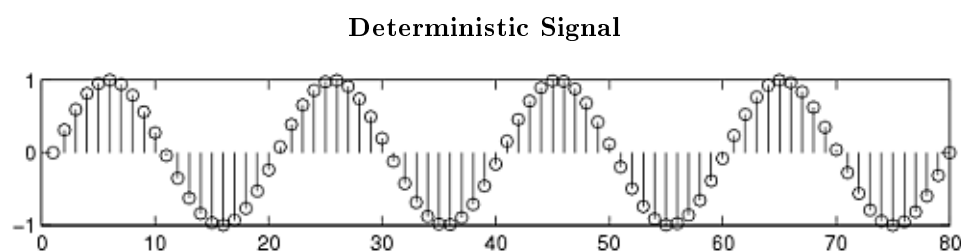


Figure 1.64: An example of a deterministic signal, the sine wave.

1.12.1.1.2 Stochastic Signals

Unlike deterministic signals, **stochastic signals**, or **random signals**, are not so nice. Random signals cannot be characterized by a simple, well-defined mathematical equation and their future values cannot be predicted. Rather, we must use probability and statistics to analyze their behavior. Also, because of their randomness, average values (Section 1.12.3) from a collection of signals are usually studied rather than analyzing one individual signal.

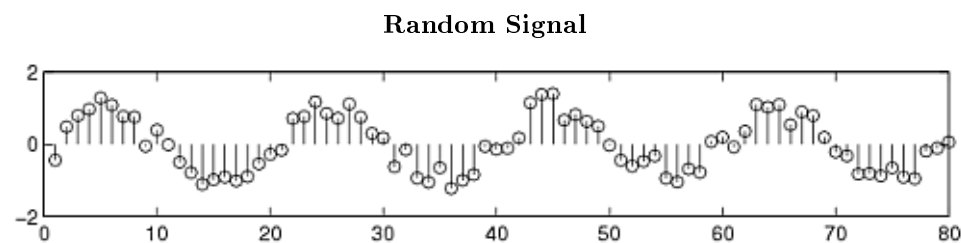


Figure 1.65: We have taken the above sine wave and added random noise to it to come up with a noisy, or random, signal. These are the types of signals that we wish to learn how to deal with so that we can recover the original sine wave.

1.12.1.2 Random Process

As mentioned above, in order to study random signals, we want to look at a collection of these signals rather than just one instance of that signal. This collection of signals is called a **random process**.

Definition 4: random process

A family or ensemble of signals that correspond to every possible outcome of a certain signal measurement. Each signal in this collection is referred to as a **realization** or **sample function** of the process.

Example

As an example of a random process, let us look at the Random Sinusoidal Process below. We use $f[n] = A \sin(\omega n + \phi)$ to represent the sinusoid with a given amplitude and phase. Note that the phase and amplitude of each sinusoid is based on a random number, thus making this a random process.

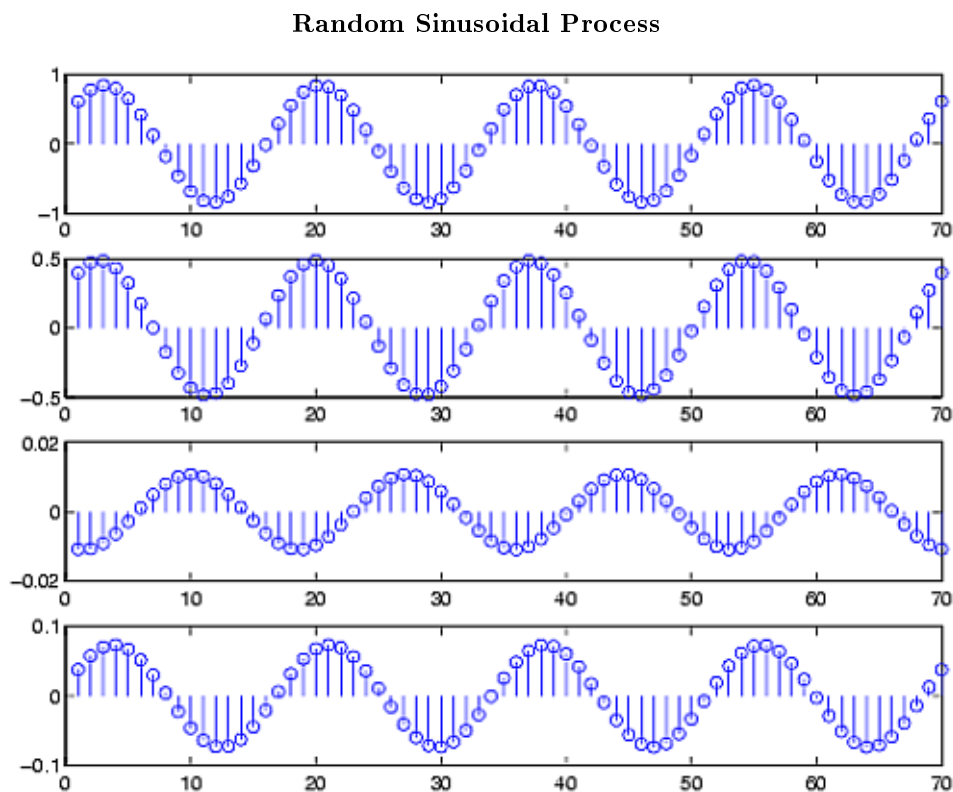


Figure 1.66: A random sinusoidal process, with the amplitude and phase being random numbers.

A random process is usually denoted by $X(t)$ or $X[n]$, with $x(t)$ or $x[n]$ used to represent an individual signal or waveform from this process.

In many notes and books, you might see the following notation and terms used to describe different types of random processes. For a **discrete random process**, sometimes just called a **random sequence**, t represents time that has a finite number of values. If t can take on any value of time, we have a **continuous**

random process. Often times discrete and continuous refer to the amplitude of the process, and process or sequence refer to the nature of the time variable. For this study, we often just use **random process** to refer to a general collection of discrete-time signals, as seen above in Figure 1.66 (Random Sinusoidal Process).

1.12.2 Stationary and Nonstationary Random Processes⁹⁶

1.12.2.1 Introduction

From the definition of a random process (Section 1.12.1), we know that all random processes are composed of random variables, each at its own unique point in time. Because of this, random processes have all the properties of random variables, such as mean, correlation, variances, etc.. When dealing with groups of signals or sequences it will be important for us to be able to show whether or not these statistical properties hold true for the entire random process. To do this, the concept of **stationary processes** has been developed. The general definition of a stationary process is:

Definition 5: stationary process

a random process where all of its statistical properties do not vary with time

Processes whose statistical properties do change are referred to as **nonstationary**.

Understanding the basic idea of stationarity will help you to be able to follow the more concrete and mathematical definition to follow. Also, we will look at various levels of stationarity used to describe the various types of stationarity characteristics a random process can have.

1.12.2.2 Distribution and Density Functions

In order to properly define what it means to be stationary from a mathematical standpoint, one needs to be somewhat familiar with the concepts of distribution and density functions. If you can remember your statistics then feel free to skip this section!

Recall that when dealing with a single random variable, the **probability distribution function** is a simply tool used to identify the probability that our observed random variable will be less than or equal to a given number. More precisely, let X be our random variable, and let x be our given value; from this we can define the distribution function as

$$F_x(x) = Pr[X \leq x] \quad (1.115)$$

This same idea can be applied to instances where we have multiple random variables as well. There may be situations where we want to look at the probability of event X and Y both occurring. For example, below is an example of a second-order **joint distribution function**.

$$F_x(x, y) = Pr[X \leq x, Y \leq y] \quad (1.116)$$

While the distribution function provides us with a full view of our variable or processes probability, it is not always the most useful for calculations. Often times we will want to look at its derivative, the **probability density function (pdf)**. We define the pdf as

$$f_x(x) = \frac{d}{dx} F_x(x) \quad (1.117)$$

$$f_x(x) dx = Pr[x < X \leq x + dx] \quad (1.118)$$

(1.118) reveals some of the physical significance of the density function. This equation tells us the probability that our random variable falls within a given interval can be approximated by $f_x(x) dx$. From the pdf, we can now use our knowledge of integrals to evaluate probabilities from the above approximation. Again we can also define a **joint density function** which will include multiple random variables just as was done

⁹⁶This content is available online at <<http://cnx.org/content/m10684/2.2/>>.

for the distribution function. The density function is used for a variety of calculations, such as finding the expected value or proving a random variable is stationary, to name a few.

NOTE: The above examples explain the distribution and density functions in terms of a single random variable, X . When we are dealing with signals and random processes, remember that we will have a set of random variables where a different random variable will occur at each time instance of the random process, $X(t_k)$. In other words, the distribution and density function will also need to take into account the choice of time.

1.12.2.3 Stationarity

Below we will now look at a more in depth and mathematical definition of a stationary process. As was mentioned previously, various levels of stationarity exist and we will look at the most common types.

1.12.2.3.1 First-Order Stationary Process

A random process is classified as **first-order stationary** if its first-order probability density function remains equal regardless of any shift in time to its time origin. If we let x_{t_1} represent a given value at time t_1 , then we define a first-order stationary as one that satisfies the following equation:

$$f_x(x_{t_1}) = f_x(x_{t_1+\tau}) \quad (1.119)$$

The physical significance of this equation is that our density function, $f_x(x_{t_1})$, is completely independent of t_1 and thus any time shift, τ .

The most important result of this statement, and the identifying characteristic of any first-order stationary process, is the fact that the mean is a constant, independent of any time shift. Below we show the results for a random process, X , that is a discrete-time signal, $x[n]$.

$$\begin{aligned} \overline{X} &= m_x[n] \\ &= E[x[n]] \\ &= \text{constant (independent of } n) \end{aligned} \quad (1.120)$$

1.12.2.3.2 Second-Order and Strict-Sense Stationary Process

A random process is classified as **second-order stationary** if its second-order probability density function does not vary over any time shift applied to both values. In other words, for values x_{t_1} and x_{t_2} then we will have the following be equal for an arbitrary time shift τ .

$$f_x(x_{t_1}, x_{t_2}) = f_x(x_{t_1+\tau}, x_{t_2+\tau}) \quad (1.121)$$

From this equation we see that the absolute time does not affect our functions, rather it only really depends on the time difference between the two variables. Looked at another way, this equation can be described as

$$Pr[X(t_1) \leq x_1, X(t_2) \leq x_2] = Pr[X(t_1 + \tau) \leq x_1, X(t_2 + \tau) \leq x_2] \quad (1.122)$$

These random processes are often referred to as **strict sense stationary (SSS)** when *all* of the distribution functions of the process are unchanged regardless of the time shift applied to them.

For a second-order stationary process, we need to look at the autocorrelation function (Section 1.12.5) to see its most important property. Since we have already stated that a second-order stationary process depends only on the time difference, then all of these types of processes have the following property:

$$\begin{aligned} R_{xx}(t, t + \tau) &= E[X(t + \tau)] \\ &= R_{xx}(\tau) \end{aligned} \quad (1.123)$$

1.12.2.3.3 Wide-Sense Stationary Process

As you begin to work with random processes, it will become evident that the strict requirements of a SSS process is more than is often necessary in order to adequately approximate our calculations on random processes. We define a final type of stationarity, referred to as **wide-sense stationary (WSS)**, to have slightly more relaxed requirements but ones that are still enough to provide us with adequate results. In order to be WSS a random process only needs to meet the following two requirements.

1. $\bar{X} = E[x[n]] = \text{constant}$
2. $E[X(t + \tau)] = R_{xx}(\tau)$

Note that a second-order (or SSS) stationary process will always be WSS; however, the reverse will not always hold true.

1.12.3 Random Processes: Mean and Variance⁹⁷

In order to study the characteristics of a random process (Section 1.12.1), let us look at some of the basic properties and operations of a random process. Below we will focus on the operations of the random signals that compose our random processes. We will denote our random process with X and a random variable from a random process or signal by x .

1.12.3.1 Mean Value

Finding the average value of a set of random signals or random variables is probably the most fundamental concepts we use in evaluating random processes through any sort of statistical method. *The mean of a random process is the average of all realizations of that process.* In order to find this average, we must look at a random signal over a range of time (possible values) and determine our average from this set of values. The **mean**, or average, of a random process, $x(t)$, is given by the following equation:

$$\begin{aligned}
 m_x(t) &= \mu_x(t) \\
 &= \bar{X} \\
 &= E[X] \\
 &= \int_{-\infty}^{\infty} x f(x) dx
 \end{aligned}
 \tag{1.124}$$

This equation may seem quite cluttered at first glance, but we want to introduce you to the various notations used to represent the mean of a random signal or process. Throughout texts and other readings, remember that these will all equal the same thing. The symbol, $\mu_x(t)$, and the X with a bar over it are often used as a short-hand to represent an average, so you might see it in certain textbooks. The other important notation used is, $E[X]$, which represents the "expected value of X " or the mathematical expectation. This notation is very common and will appear again.

If the random variables, which make up our random process, are discrete or quantized values, such as in a binary process, then the integrals become summations over all the possible values of the random variable. In this case, our expected value becomes

$$E[x[n]] = \sum_x (\alpha \text{Pr}[x[n] = \alpha]) \tag{1.125}$$

If we have two random signals or variables, their averages can reveal how the two signals interact. If the product of the two individual averages of both signals do *not* equal the average of the product of the two signals, then the two signals are said to be **linearly independent**, also referred to as **uncorrelated**.

In the case where we have a random process in which only one sample can be viewed at a time, then we will often not have all the information available to calculate the mean using the density function as shown

⁹⁷This content is available online at <<http://cnx.org/content/m10656/2.3/>>.

above. In this case we must estimate the mean through the time-average mean (Section 1.12.3.4: Time Averages), discussed later. For fields such as signal processing that deal mainly with discrete signals and values, then these are the averages most commonly used.

1.12.3.1.1 Properties of the Mean

- The expected value of a constant, α , is the constant:

$$E[\alpha] = \alpha \quad (1.126)$$

- Adding a constant, α , to each term increases the expected value by that constant:

$$E[X + \alpha] = E[X] + \alpha \quad (1.127)$$

- Multiplying the random variable by a constant, α , multiplies the expected value by that constant.

$$E[\alpha X] = \alpha E[X] \quad (1.128)$$

- The expected value of the sum of two or more random variables, is the sum of each individual expected value.

$$E[X + Y] = E[X] + E[Y] \quad (1.129)$$

1.12.3.2 Mean-Square Value

If we look at the second **moment** of the term (we now look at x^2 in the integral), then we will have the **mean-square value** of our random process. As you would expect, this is written as

$$\begin{aligned} \overline{X^2} &= E[X^2] \\ &= \int_{-\infty}^{\infty} x^2 f(x) dx \end{aligned} \quad (1.130)$$

This equation is also often referred to as the **average power** of a process or signal.

1.12.3.3 Variance

Now that we have an idea about the average value or values that a random process takes, we are often interested in seeing just how spread out the different random values might be. To do this, we look at the **variance** which is a measure of this spread. The variance, often denoted by σ^2 , is written as follows:

$$\begin{aligned} \sigma^2 &= Var(X) \\ &= E[(X - E[X])^2] \\ &= \int_{-\infty}^{\infty} (x - \overline{X})^2 f(x) dx \end{aligned} \quad (1.131)$$

Using the rules for the expected value, we can rewrite this formula as the following form, which is commonly seen:

$$\begin{aligned} \sigma^2 &= \overline{X^2} - (\overline{X})^2 \\ &= E[X^2] - (E[X])^2 \end{aligned} \quad (1.132)$$

1.12.3.3.1 Standard Deviation

Another common statistical tool is the standard deviation. Once you know how to calculate the variance, the standard deviation is simply *the square root of the variance*, or σ .

1.12.3.3.2 Properties of Variance

- The variance of a constant, α , equals zero:

$$\begin{aligned} Var(\alpha) &= \sigma(\alpha)^2 \\ &= 0 \end{aligned} \quad (1.133)$$

- Adding a constant, α , to a random variable does not affect the variance because the mean increases by the same value:

$$\begin{aligned} Var(X + \alpha) &= \sigma(X + \alpha)^2 \\ &= \sigma(X)^2 \end{aligned} \quad (1.134)$$

- Multiplying the random variable by a constant, α , increases the variance by the square of the constant:

$$\begin{aligned} Var(\alpha X) &= \sigma(\alpha X)^2 \\ &= \alpha^2 \sigma(X)^2 \end{aligned} \quad (1.135)$$

- The variance of the sum of two random variables only equals the sum of the variances if the variables are **independent**.

$$\begin{aligned} Var(X + Y) &= \sigma(X + Y)^2 \\ &= \sigma(X)^2 + \sigma(Y)^2 \end{aligned} \quad (1.136)$$

Otherwise, if the random variables are *not* independent, then we must also include the covariance of the product of the variables as follows:

$$Var(X + Y) = \sigma(X)^2 + 2Cov(X, Y) + \sigma(Y)^2 \quad (1.137)$$

1.12.3.4 Time Averages

In the case where we can not view the entire ensemble of the random process, we must use time averages to estimate the values of the mean and variance for the process. Generally, this will only give us acceptable results for independent and **ergodic** processes, meaning those processes in which each signal or member of the process seems to have the same statistical behavior as the entire process. The time averages will also only be taken over a finite interval since we will only be able to see a finite part of the sample.

1.12.3.4.1 Estimating the Mean

For the ergodic random process, $x(t)$, we will estimate the mean using the time averaging function defined as

$$\begin{aligned} \bar{X} &= E[X] \\ &= \frac{1}{T} \int_0^T X(t) dt \end{aligned} \quad (1.138)$$

However, for most real-world situations we will be dealing with discrete values in our computations and signals. We will represent this mean as

$$\begin{aligned} \bar{X} &= E[X] \\ &= \frac{1}{N} \sum_{n=1}^N X[n] \end{aligned} \quad (1.139)$$

1.12.3.4.2 Estimating the Variance

Once the mean of our random process has been estimated then we can simply use those values in the following variance equation (introduced in one of the above sections)

$$\sigma_x^2 = \overline{X^2} - (\overline{X})^2 \quad (1.140)$$

1.12.3.5 Example

Let us now look at how some of the formulas and concepts above apply to a simple example. We will just look at a single, continuous random variable for this example, but the calculations and methods are the same for a random process. For this example, we will consider a random variable having the probability density function described below and shown in Figure 1.67 (Probability Density Function).

$$f(x) = \begin{cases} \frac{1}{10} & \text{if } 10 \leq x \leq 20 \\ 0 & \text{otherwise} \end{cases} \quad (1.141)$$

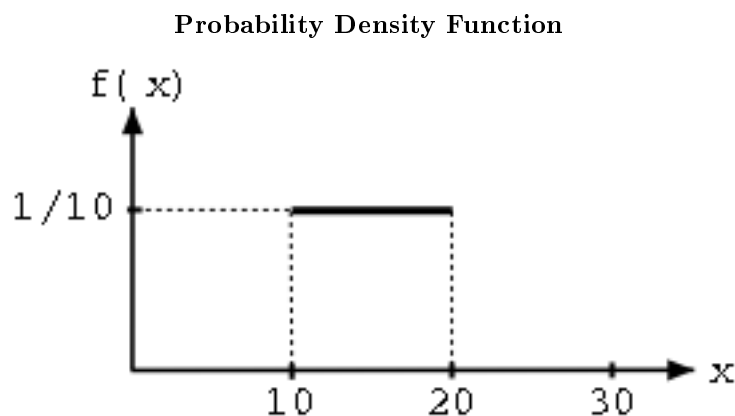


Figure 1.67: A uniform probability density function.

First, we will use (1.124) to solve for the mean value.

$$\begin{aligned} \overline{X} &= \int_{10}^{20} x \frac{1}{10} dx \\ &= \frac{1}{10} \left(\frac{x^2}{2} \right) \Big|_{x=10}^{20} \\ &= \frac{1}{10} (200 - 50) \\ &= 15 \end{aligned} \quad (1.142)$$

Using (1.130) we can obtain the mean-square value for the above density function.

$$\begin{aligned}
 \overline{X^2} &= \int_{10}^{20} x^2 \frac{1}{10} dx \\
 &= \frac{1}{10} \left(\frac{x^3}{3} \right) \Big|_{x=10}^{20} \\
 &= \frac{1}{10} \left(\frac{8000}{3} - \frac{1000}{3} \right) \\
 &= 233.33
 \end{aligned} \tag{1.143}$$

And finally, let us solve for the variance of this function.

$$\begin{aligned}
 \sigma^2 &= \overline{X^2} - (\overline{X})^2 \\
 &= 233.33 - 15^2 \\
 &= 8.33
 \end{aligned} \tag{1.144}$$

1.12.4 Correlation and Covariance of a Random Signal⁹⁸

When we take the expected value (Section 1.12.3), or average, of a random process (Section 1.12.1.2: Random Process), we measure several important characteristics about how the process behaves in general. This proves to be a very important observation. However, suppose we have several random processes measuring different aspects of a system. The relationship between these different processes will also be an important observation. The covariance and correlation are two important tools in finding these relationships. Below we will go into more details as to what these words mean and how these tools are helpful. Note that much of the following discussions refer to just random variables, but keep in mind that these variables can represent random signals or random processes.

1.12.4.1 Covariance

To begin with, when dealing with more than one random process, it should be obvious that it would be nice to be able to have a number that could quickly give us an idea of how similar the processes are. To do this, we use the **covariance**, which is analogous to the variance of a single variable.

Definition 6: Covariance

A measure of how much the deviations of two or more variables or processes match.

For two processes, X and Y , if they are *not* closely related then the covariance will be small, and if they are similar then the covariance will be large. Let us clarify this statement by describing what we mean by "related" and "similar." Two processes are "closely related" if their distribution spreads are almost equal and they are around the same, or a very slightly different, mean.

Mathematically, covariance is often written as σ_{xy} and is defined as

$$\begin{aligned}
 cov(X, Y) &= \sigma_{xy} \\
 &= E[(X - \overline{X})(Y - \overline{Y})]
 \end{aligned} \tag{1.145}$$

This can also be reduced and rewritten in the following two forms:

$$\sigma_{xy} = \overline{(xy)} - (\overline{x})(\overline{y}) \tag{1.146}$$

$$\sigma_{xy} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (X - \overline{X})(Y - \overline{Y}) f(x, y) dx dy \tag{1.147}$$

⁹⁸This content is available online at <<http://cnx.org/content/m10673/2.3/>>.

1.12.4.1.1 Useful Properties

- If X and Y are independent and uncorrelated or one of them has zero mean value, then

$$\sigma_{xy} = 0$$

- If X and Y are orthogonal, then

$$\sigma_{xy} = -(E[X] E[Y])$$

- The covariance is symmetric

$$\text{cov}(X, Y) = \text{cov}(Y, X)$$

- Basic covariance identity

$$\text{cov}(X + Y, Z) = \text{cov}(X, Z) + \text{cov}(Y, Z)$$

- Covariance of equal variables

$$\text{cov}(X, X) = \text{Var}(X)$$

1.12.4.2 Correlation

For anyone who has any kind of statistical background, you should be able to see that the idea of dependence/independence among variables and signals plays an important role when dealing with random processes. Because of this, the **correlation** of two variables provides us with a measure of how the two variables affect one another.

Definition 7: Correlation

A measure of how much one random variable depends upon the other.

This measure of association between the variables will provide us with a clue as to how well the value of one variable can be predicted from the value of the other. The correlation is equal to the average of the product of two random variables and is defined as

$$\begin{aligned} \text{cor}(X, Y) &= E[XY] \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xyf(x, y) dx dy \end{aligned} \quad (1.148)$$

1.12.4.2.1 Correlation Coefficient

It is often useful to express the correlation of random variables with a range of numbers, like a percentage. For a given set of variables, we use the **correlation coefficient** to give us the linear relationship between our variables. The correlation coefficient of two variables is defined in terms of their covariance and standard deviations (Section 1.12.3.3.1: Standard Deviation), denoted by σ_x , as seen below

$$\rho = \frac{\text{cov}(X, Y)}{\sigma_x \sigma_y} \quad (1.149)$$

where we will always have

$$-1 \leq \rho \leq 1$$

This provides us with a quick and easy way to view the correlation between our variables. If there is no relationship between the variables then the correlation coefficient will be zero and if there is a perfect positive match it will be one. If there is a perfect inverse relationship, where one set of variables increases while the other decreases, then the correlation coefficient will be negative one. This type of correlation is often referred to more specifically as the **Pearson's Correlation Coefficient**, or Pearson's Product Moment Correlation.

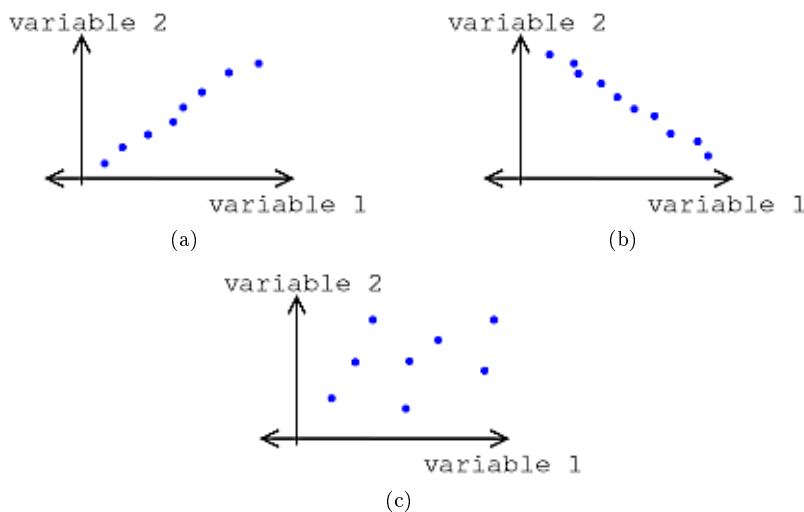


Figure 1.68: Types of Correlation (a) Positive Correlation (b) Negative Correlation (c) Uncorrelated (No Correlation)

NOTE: So far we have dealt with correlation simply as a number relating the relationship between any two variables. However, since our goal will be to relate random processes to each other, which deals with signals as a function of time, we will want to continue this study by looking at **correlation functions (Section 1.12.5)**.

1.12.4.3 Example

Now let us take just a second to look at a simple example that involves calculating the covariance and correlation of two sets of random numbers. We are given the following data sets:

$$x = \{3, 1, 6, 3, 4\}$$

$$y = \{1, 5, 3, 4, 3\}$$

To begin with, for the covariance we will need to find the expected value (Section 1.12.3), or mean, of x and y .

$$\bar{x} = \frac{1}{5} (3 + 1 + 6 + 3 + 4) = 3.4$$

$$\bar{y} = \frac{1}{5} (1 + 5 + 3 + 4 + 3) = 3.2$$

$$\overline{xy} = \frac{1}{5} (3 + 5 + 18 + 12 + 12) = 10$$

Next we will solve for the standard deviations of our two sets using the formula below (for a review click here (Section 1.12.3.3: Variance)).

$$\sigma = \sqrt{E[(X - E[X])^2]}$$

$$\sigma_x = \sqrt{\frac{1}{5} (0.16 + 5.76 + 6.76 + 0.16 + 0.36)} = 1.625$$

$$\sigma_y = \sqrt{\frac{1}{6} (4.84 + 3.24 + 0.04 + 0.64 + 0.04)} = 1.327$$

Now we can finally calculate the covariance using one of the two formulas found above. Since we calculated the three means, we will use that formula (1.146) since it will be much simpler.

$$\sigma_{xy} = 10 - 3.4 \times 3.2 = -0.88$$

And for our last calculation, we will solve for the correlation coefficient, ρ .

$$\rho = \frac{-0.88}{1.625 \times 1.327} = -0.408$$

1.12.4.3.1 Matlab Code for Example

The above example can be easily calculated using Matlab. Below I have included the code to find all of the values above.

```
x = [3 1 6 3 4];
y = [1 5 3 4 3];

mx = mean(x)
my = mean(y)
mxy = mean(x.*y)

% Standard Dev. from built-in Matlab Functions
std(x,1)
std(y,1)

% Standard Dev. from Equation Above (same result as std(?,1))
sqrt( 1/5 * sum((x-mx).^2))
sqrt( 1/5 * sum((y-my).^2))

cov(x,y,1)

corrcoef(x,y)
```

1.12.5 Autocorrelation of Random Processes⁹⁹

Before diving into a more complex statistical analysis of random signals and processes (Section 1.12.1), let us quickly review the idea of correlation (Section 1.12.4). Recall that the correlation of two signals or variables is the expected value of the product of those two variables. Since our focus will be to discover more about a random process, a collection of random signals, then imagine us dealing with two samples of

⁹⁹This content is available online at <<http://cnx.org/content/m10676/2.4/>>.

a random process, where each sample is taken at a different point in time. Also recall that the key property of these random processes is that they are now functions of time; imagine them as a collection of signals. The expected value (Section 1.12.3) of the product of these two variables (or samples) will now depend on how quickly they change in regards to *time*. For example, if the two variables are taken from almost the same time period, then we should expect them to have a high correlation. We will now look at a correlation function that relates a pair of random variables from the same process to the time separations between them, where the argument to this correlation function will be the time difference. For the correlation of signals from two different random process, look at the crosscorrelation function (Section 1.12.6).

1.12.5.1 Autocorrelation Function

The first of these correlation functions we will discuss is the **autocorrelation**, where each of the random variables we will deal with come from the same random process.

Definition 8: Autocorrelation

the expected value of the product of a random variable or signal realization with a time-shifted version of itself

With a simple calculation and analysis of the autocorrelation function, we can discover a few important characteristics about our random process. These include:

1. How quickly our random signal or processes changes with respect to the time function
2. Whether our process has a periodic component and what the expected frequency might be

As was mentioned above, the autocorrelation function is simply the expected value of a product. Assume we have a pair of random variables from the same process, $X_1 = X(t_1)$ and $X_2 = X(t_2)$, then the autocorrelation is often written as

$$\begin{aligned} R_{xx}(t_1, t_2) &= E[X_1 X_2] \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x_1 x_2 f(x_1, x_2) dx_2 dx_1 \end{aligned} \quad (1.150)$$

The above equation is valid for stationary and nonstationary random processes. For stationary processes (Section 1.12.2), we can generalize this expression a little further. Given a wide-sense stationary processes, it can be proven that the expected values from our random process will be independent of the origin of our time function. Therefore, we can say that our autocorrelation function will depend on the time difference and not some absolute time. For this discussion, we will let $\tau = t_2 - t_1$, and thus we generalize our autocorrelation expression as

$$\begin{aligned} R_{xx}(t, t + \tau) &= R_{xx}(\tau) \\ &= E[X(t) X(t + \tau)] \end{aligned} \quad (1.151)$$

for the continuous-time case. In most DSP course we will be more interested in dealing with real signal sequences, and thus we will want to look at the discrete-time case of the autocorrelation function. The formula below will prove to be more common and useful than (1.150):

$$R_{xx}[n, n + m] = \sum_{n=-\infty}^{\infty} (x[n] x[n + m]) \quad (1.152)$$

And again we can generalize the notation for our autocorrelation function as

$$\begin{aligned} R_{xx}[n, n + m] &= R_{xx}[m] \\ &= E[X[n] X[n + m]] \end{aligned} \quad (1.153)$$

1.12.5.1.1 Properties of Autocorrelation

Below we will look at several properties of the autocorrelation function that hold for *stationary* random processes.

- Autocorrelation is an even function for τ

$$R_{xx}(\tau) = R_{xx}(-\tau)$$

- The mean-square value can be found by evaluating the autocorrelation where $\tau = 0$, which gives us

$$R_{xx}(0) = \overline{X^2}$$

- The autocorrelation function will have its largest value when $\tau = 0$. This value can appear again, for example in a periodic function at the values of the equivalent periodic points, but will never be exceeded.

$$R_{xx}(0) \geq |R_{xx}(\tau)|$$

- If we take the autocorrelation of a period function, then $R_{xx}(\tau)$ will also be periodic with the same frequency.

1.12.5.1.2 Estimating the Autocorrelation with Time-Averaging

Sometimes the whole random process is not available to us. In these cases, we would still like to be able to find out some of the characteristics of the stationary random process, even if we just have part of one sample function. In order to do this we can *estimate* the autocorrelation from a given interval, 0 to T seconds, of the sample function.

$$R_{xx}(\tau) = \frac{1}{T-\tau} \int_0^{T-\tau} x(t) x(t+\tau) dt \quad (1.154)$$

However, a lot of times we will not have sufficient information to build a complete continuous-time function of one of our random signals for the above analysis. If this is the case, we can treat the information we do know about the function as a discrete signal and use the discrete-time formula for estimating the autocorrelation.

$$R_{xx}[m] = \frac{1}{N-m} \sum_{n=0}^{N-m-1} (x[n] x[n+m]) \quad (1.155)$$

1.12.5.2 Examples

Below we will look at a variety of examples that use the autocorrelation function. We will begin with a simple example dealing with Gaussian White Noise (GWN) and a few basic statistical properties that will prove very useful in these and future calculations.

Example 1.33

We will let $x[n]$ represent our GWN. For this problem, it is important to remember the following fact about the mean of a GWN function:

$$E[x[n]] = 0$$

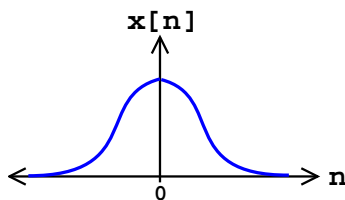


Figure 1.69: Gaussian density function. By examination, can easily see that the above statement is true - the mean equals zero.

Along with being *zero-mean*, recall that GWN is always *independent*. With these two facts, we are now ready to do the short calculations required to find the autocorrelation.

$$R_{xx}[n, n+m] = E[x[n] x[n+m]]$$

Since the function, $x[n]$, is independent, then we can take the product of the individual expected values of both functions.

$$R_{xx}[n, n+m] = E[x[n]] E[x[n+m]]$$

Now, looking at the above equation we see that we can break it up further into two conditions: one when m and n are equal and one when they are not equal. When they are equal we can combine the expected values. We are left with the following piecewise function to solve:

$$R_{xx}[n, n+m] = \begin{cases} E[x[n]] E[x[n+m]] & \text{if } m \neq 0 \\ E[x^2[n]] & \text{if } m = 0 \end{cases}$$

We can now solve the two parts of the above equation. The first equation is easy to solve as we have already stated that the expected value of $x[n]$ will be zero. For the second part, you should recall from statistics that the expected value of the square of a function is equal to the variance. Thus we get the following results for the autocorrelation:

$$R_{xx}[n, n+m] = \begin{cases} 0 & \text{if } m \neq 0 \\ \sigma^2 & \text{if } m = 0 \end{cases}$$

Or in a more concise way, we can represent the results as

$$R_{xx}[n, n+m] = \sigma^2 \delta[m]$$

1.12.6 Crosscorrelation of Random Processes¹⁰⁰

Before diving into a more complex statistical analysis of random signals and processes (Section 1.12.1), let us quickly review the idea of correlation (Section 1.12.4). Recall that the correlation of two signals or variables is the expected value of the product of those two variables. Since our main focus is to discover more about random processes, a collection of random signals, we will deal with two random processes in this discussion, where in this case we will deal with samples from two *different* random processes. We will analyze the expected value (Section 1.12.3.1: Mean Value) of the product of these two variables and how they correlate to one another, where the argument to this correlation function will be the time difference. For the correlation of signals from the same random process, look at the autocorrelation function (Section 1.12.5).

¹⁰⁰This content is available online at <<http://cnx.org/content/m10686/2.2/>>.

1.12.6.1 Crosscorrelation Function

When dealing with multiple random processes, it is also important to be able to describe the relationship, if any, between the processes. For example, this may occur if more than one random signal is applied to a system. In order to do this, we use the **crosscorrelation function**, where the variables are instances from two different wide sense stationary random processes.

Definition 9: Crosscorrelation

if two processes are wide sense stationary, the expected value of the product of a random variable from one random process with a time-shifted, random variable from a different random process

Looking at the generalized formula for the crosscorrelation, we will represent our two random processes by allowing $U = U(t)$ and $V = V(t - \tau)$. We will define the crosscorrelation function as

$$\begin{aligned} R_{uv}(t, t - \tau) &= E[UV] \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} uv f(u, v) dv du \end{aligned} \quad (1.156)$$

Just as the case with the autocorrelation function, if our input and output, denoted as $U(t)$ and $V(t)$, are at least jointly wide sense stationary, then the crosscorrelation does not depend on absolute time; it is just a function of the time difference. This means we can simplify our writing of the above function as

$$R_{uv}(\tau) = E[UV] \quad (1.157)$$

or if we deal with two real signal sequences, $x[n]$ and $y[n]$, then we arrive at a more commonly seen formula for the discrete crosscorrelation function. See the formula below and notice the similarities between it and the convolution (Section 1.3) of two signals:

$$\begin{aligned} R_{xy}(n, n - m) &= R_{xy}(m) \\ &= \sum_{n=-\infty}^{\infty} (x[n] y[n - m]) \end{aligned} \quad (1.158)$$

1.12.6.1.1 Properties of Crosscorrelation

Below we will look at several properties of the crosscorrelation function that hold for two *wide sense stationary* (WSS) random processes.

- Crosscorrelation is *not* an even function; however, it does have a unique symmetry property:

$$R_{xy}(-\tau) = R_{yx}(\tau) \quad (1.159)$$

- The maximum value of the crosscorrelation is not always when the shift equals zero; however, we can prove the following property revealing to us what value the maximum cannot exceed.

$$|R_{xy}(\tau)| \leq \sqrt{R_{xx}(0) R_{yy}(0)} \quad (1.160)$$

- When two random processes are statistically independent then we have

$$R_{xy}(\tau) = R_{yx}(\tau) \quad (1.161)$$

1.12.6.2 Examples**Exercise 1.9***(Solution on p. 97.)*

Let us begin by looking at a simple example showing the relationship between two sequences. Using (1.158), find the crosscorrelation of the sequences

$$x[n] = \{\dots, 0, 0, 2, -3, 6, 1, 3, 0, 0, \dots\}$$

$$y[n] = \{\dots, 0, 0, 1, -2, 4, 1, -3, 0, 0, \dots\}$$

for each of the following possible time shifts: $m = \{0, 3, -1\}$.

Solutions to Exercises in Chapter 1

Solution to Exercise 1.1 (p. 25)

In order to represent \mathbf{x} in terms of b_0 and b_1 we will follow the same steps we used in the above example.

$$B = \begin{pmatrix} 1 & 2 \\ 3 & 0 \end{pmatrix}$$

$$B^{-1} = \begin{pmatrix} 0 & \frac{1}{2} \\ \frac{1}{3} & \frac{-1}{6} \end{pmatrix}$$

$$\alpha = B^{-1}\mathbf{x} = \begin{pmatrix} 1 \\ \frac{2}{3} \end{pmatrix}$$

And now we can write \mathbf{x} in terms of b_0 and b_1 .

$$\mathbf{x} = b_0 + \frac{2}{3}b_1$$

And we can easily substitute in our known values of b_0 and b_1 to verify our results.

Solution to Exercise 1.2 (p. 28)

In order to calculate the Fourier transform, all we need to use is (1.24) (Continuous-Time Fourier Transform), complex exponentials¹⁰¹, and basic calculus.

$$\begin{aligned} \mathcal{F}(\Omega) &= \int_{-\infty}^{\infty} f(t) e^{-(j\Omega t)} dt \\ &= \int_0^{\infty} e^{-(\alpha t)} e^{-(j\Omega t)} dt \\ &= \int_0^{\infty} e^{(-t)(\alpha + j\Omega)} dt \\ &= 0 - \frac{-1}{\alpha + j\Omega} \end{aligned} \tag{1.162}$$

$$\mathcal{F}(\Omega) = \frac{1}{\alpha + j\Omega} \tag{1.163}$$

Solution to Exercise 1.3 (p. 28)

Here we will use (1.25) (Inverse CTFT) to find the inverse FT given that $t \neq 0$.

$$\begin{aligned} x(t) &= \frac{1}{2\pi} \int_{-M}^M e^{j\Omega t} d\Omega \\ &= \frac{1}{2\pi} e^{j\Omega t} \Big|_{\Omega, \Omega=e^{jw}} \\ &= \frac{1}{\pi t} \sin(Mt) \end{aligned} \tag{1.164}$$

$$x(t) = \frac{M}{\pi} \left(\text{sinc} \frac{Mt}{\pi} \right) \tag{1.165}$$

Solution to Exercise 1.4 (p. 29)

$$\begin{aligned} S(e^{j2\pi(f+1)}) &= \sum_{n=-\infty}^{\infty} (s(n) e^{-(j2\pi(f+1)n)}) \\ &= \sum_{n=-\infty}^{\infty} (e^{-(j2\pi n)} s(n) e^{-(j2\pi f n)}) \\ &= \sum_{n=-\infty}^{\infty} (s(n) e^{-(j2\pi f n)}) \\ &= S(e^{j2\pi f}) \end{aligned} \tag{1.166}$$

¹⁰¹"The Complex Exponential" <<http://cnx.org/content/m10060/latest/>>

Solution to Exercise 1.5 (p. 32)

$$\alpha \sum_{n=n_0}^{N+n_0-1} (\alpha^n) - \sum_{n=n_0}^{N+n_0-1} (\alpha^n) = \alpha^{N+n_0} - \alpha^{n_0}$$

which, after manipulation, yields the geometric sum formula.

Solution to Exercise 1.6 (p. 34)

If the sampling frequency exceeds the Nyquist frequency, the spectrum of the samples equals the analog spectrum, but over the normalized analog frequency fT . Thus, the energy in the sampled signal equals the original signal's energy multiplied by T .

Solution to Exercise 1.7 (p. 52)

$$x[n+N] = ???$$

Solution to Exercise 1.8 (p. 56)

$S[n]$ is N -periodic, so it has the following Fourier Series¹⁰²:

$$\begin{aligned} c_k &= \frac{1}{N} \int_{-\frac{N}{2}}^{\frac{N}{2}} \delta[n] e^{(-j)2\pi \frac{k}{N}n} dn \\ &= \frac{1}{N} \end{aligned} \tag{1.167}$$

$$S[n] = \sum_{k=-\infty}^{\infty} \left(e^{(-j)2\pi \frac{k}{N}n} \right) \tag{1.168}$$

where the DTFT of the exponential in the above equation is equal to $\delta\left(\omega - \frac{2\pi k}{N}\right)$.

Solution to Exercise 1.9 (p. 95)

1. For $m = 0$, we should begin by finding the product sequence $s[n] = x[n]y[n]$. Doing this we get the following sequence:

$$s[n] = \{\dots, 0, 0, 2, 6, 24, 1, -9, 0, 0, \dots\}$$

and so from the sum in our crosscorrelation function we arrive at the answer of

$$R_{xy}(0) = 22$$

2. For $m = 3$, we will approach it the same way as we did above; however, we will now shift $y[n]$ to the right. Then we can find the product sequence $s[n] = x[n]y[n-3]$, which yields

$$s[n] = \{\dots, 0, 0, 0, 0, 0, 1, -6, 0, 0, \dots\}$$

and from the crosscorrelation function we arrive at the answer of

$$R_{xy}(3) = -6$$

3. For $m = -1$, we will again take the same approach; however, we will now shift $y[n]$ to the left. Then we can find the product sequence $s[n] = x[n]y[n+1]$, which yields

$$s[n] = \{\dots, 0, 0, -4, -12, 6, -3, 0, 0, 0, \dots\}$$

and from the crosscorrelation function we arrive at the answer of

$$R_{xy}(-1) = -13$$

¹⁰²"Fourier Series: Eigenfunction Approach" <<http://cnx.org/content/m10496/latest/>>

Chapter 2

The DFT, FFT, and Practical Spectral Analysis

2.1 The Discrete Fourier Transform

2.1.1 DFT Definition and Properties¹

2.1.1.1 DFT

The discrete Fourier transform (DFT) (Section 1.10.6) is the primary transform used for numerical computation in digital signal processing. It is very widely used for spectrum analysis (Section 2.2.1), fast convolution (Section 2.4), and many other applications. The DFT transforms N discrete-time samples to the same number of discrete frequency samples, and is defined as

$$X(k) = \sum_{n=0}^{N-1} \left(x(n) e^{-j\frac{2\pi nk}{N}} \right) \quad (2.1)$$

The DFT is widely used in part because it can be computed very efficiently using fast Fourier transform (FFT)² algorithms.

2.1.1.2 IDFT

The inverse DFT (IDFT) transforms N discrete-frequency samples to the same number of discrete-time samples. The IDFT has a form very similar to the DFT,

$$x(n) = \frac{1}{N} \sum_{k=0}^{N-1} \left(X(k) e^{j\frac{2\pi nk}{N}} \right) \quad (2.2)$$

and can thus also be computed efficiently using FFTs³.

2.1.1.3 DFT and IDFT properties

2.1.1.3.1 Periodicity

Due to the N -sample periodicity of the complex exponential basis functions $e^{j\frac{2\pi nk}{N}}$ in the DFT and IDFT, the resulting transforms are also periodic with N samples.

¹This content is available online at <<http://cnx.org/content/m12019/1.5/>>.

²The DFT, FFT, and Practical Spectral Analysis <<http://cnx.org/content/col10281/latest/>>

³The DFT, FFT, and Practical Spectral Analysis <<http://cnx.org/content/col10281/latest/>>

$$X(k + N) = X(k)$$

$$x(n) = x(n + N)$$

2.1.1.3.2 Circular Shift

A shift in time corresponds to a phase shift that is linear in frequency. Because of the periodicity induced by the DFT and IDFT, the shift is *circular*, or modulo N samples.

$$\left(x((n - m) \bmod N) \Leftrightarrow X(k) e^{-j \frac{2\pi km}{N}} \right)$$

The modulus operator $p \bmod N$ means the *remainder* of p when divided by N . For example,

$$9 \bmod 5 = 4$$

and

$$-1 \bmod 5 = 4$$

2.1.1.3.3 Time Reversal

$$(x((-n) \bmod N) = x((N - n) \bmod N) \Leftrightarrow X((N - k) \bmod N) = X((-k) \bmod N))$$

Note: time-reversal maps $(0 \Leftrightarrow 0)$, $(1 \Leftrightarrow N - 1)$, $(2 \Leftrightarrow N - 2)$, etc. as illustrated in the figure below.

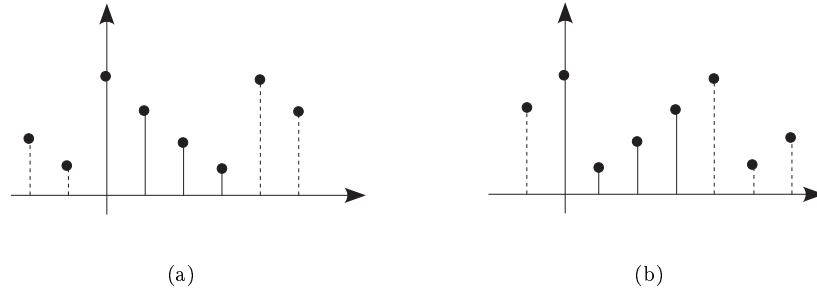


Figure 2.1: Illustration of circular time-reversal (a) Original signal (b) Time-reversed

2.1.1.3.4 Complex Conjugate

$$\left(\overline{x(n)} \Leftrightarrow \overline{X((-k) \bmod N)} \right)$$

2.1.1.3.5 Circular Convolution Property

Circular convolution is defined as

$$\left(x(n) * h(n) \doteq \sum_{m=0}^{N-1} (x(m) x((n-m) \bmod N)) \right)$$

Circular convolution of two discrete-time signals corresponds to multiplication of their DFTs:

$$(x(n) * h(n) \Leftrightarrow X(k) H(k))$$

2.1.1.3.6 Multiplication Property

A similar property relates multiplication in time to circular convolution in frequency.

$$\left(x(n) h(n) \Leftrightarrow \frac{1}{N} X(k) * H(k) \right)$$

2.1.1.3.7 Parseval's Theorem

Parseval's theorem relates the energy of a length- N discrete-time signal (or one period) to the energy of its DFT.

$$\sum_{n=0}^{N-1} (|x(n)|^2) = \frac{1}{N} \sum_{k=0}^{N-1} (|X(k)|^2)$$

2.1.1.3.8 Symmetry

The continuous-time Fourier transform (Section 1.7), the DTFT (2.3), and DFT (2.5) are all defined as transforms of complex-valued data to complex-valued spectra. However, in practice signals are often real-valued. The DFT of a real-valued discrete-time signal has a special symmetry, in which the real part of the transform values are **DFT even symmetric** and the imaginary part is **DFT odd symmetric**, as illustrated in the equation and figure below.

$$x(n) \text{ real} \Leftrightarrow X(k) = \overline{X((N-k) \bmod N)} \text{ (This implies } X(0), X(\frac{N}{2}) \text{ are real-valued.)}$$

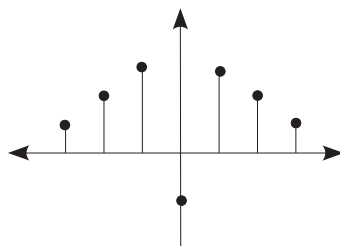
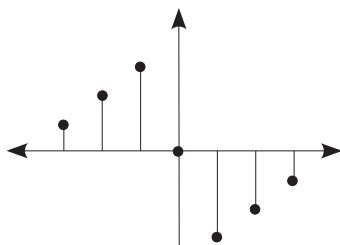
(a) Real part of $X(k)$ is even(b) Imaginary part of $X(k)$ is odd

Figure 2.2: DFT symmetry of real-valued signal (a) Even-symmetry in DFT sense (b) Odd-symmetry in DFT sense

2.2 Spectrum Analysis

2.2.1 Spectrum Analysis Using the Discrete Fourier Transform⁴

2.2.1.1 Discrete-Time Fourier Transform

The Discrete-Time Fourier Transform (DTFT) (Section 1.8) is the primary theoretical tool for understanding the frequency content of a discrete-time (sampled) signal. The DTFT (Section 1.8) is defined as

$$X(\omega) = \sum_{n=-\infty}^{\infty} (x(n) e^{-j\omega n}) \quad (2.3)$$

The inverse DTFT (IDTFT) is defined by an integral formula, because it operates on a continuous-frequency DTFT spectrum:

$$x(n) = \frac{1}{2\pi} \int_{-\pi}^{\pi} X(k) e^{j\omega n} d\omega \quad (2.4)$$

The DTFT is very useful for theory and analysis, but is not practical for numerically computing a spectrum digitally, because

⁴This content is available online at <<http://cnx.org/content/m12032/1.6/>>.

1. infinite time samples means
 - infinite computation
 - infinite delay
2. The transform is continuous in the discrete-time frequency, ω

For practical computation of the frequency content of real-world signals, the Discrete Fourier Transform (DFT) is used.

2.2.1.2 Discrete Fourier Transform

The DFT transforms N samples of a discrete-time signal to the same number of discrete frequency samples, and is defined as

$$X(k) = \sum_{n=0}^{N-1} \left(x(n) e^{-j\frac{2\pi nk}{N}} \right) \quad (2.5)$$

The DFT is invertible by the inverse discrete Fourier transform (IDFT):

$$x(n) = \frac{1}{N} \sum_{k=0}^{N-1} \left(X(k) e^{j\frac{2\pi nk}{N}} \right) \quad (2.6)$$

The DFT (2.5) and IDFT (2.6) are a self-contained, one-to-one transform pair for a length- N discrete-time signal. (That is, the DFT (2.5) is not *merely* an approximation to the DTFT (2.3) as discussed next.) However, the DFT (2.5) is very often used as a practical approximation to the DTFT (2.3).

2.2.1.3 Relationships Between DFT and DTFT

2.2.1.3.1 DFT and Discrete Fourier Series

The DFT (2.5) gives the discrete-time Fourier series coefficients of a periodic sequence ($x(n) = x(n + N)$) of period N samples, or

$$X(\omega) = \frac{2\pi}{N} \sum \left(X(k) \delta \left(\omega - \frac{2\pi k}{N} \right) \right) \quad (2.7)$$

as can easily be confirmed by computing the inverse DTFT of the corresponding line spectrum:

$$\begin{aligned} x(n) &= \frac{1}{2\pi} \int_{-\pi}^{\pi} \left(\frac{2\pi}{N} \sum \left(X(k) \delta \left(\omega - \frac{2\pi k}{N} \right) \right) \right) e^{j\omega n} d\omega \\ &= \frac{1}{N} \sum_{k=0}^{N-1} \left(X(k) e^{j\frac{2\pi nk}{N}} \right) \\ &= IDFT(X(k)) \\ &= x(n) \end{aligned} \quad (2.8)$$

The DFT can thus be used to *exactly* compute the relative values of the N line spectral components of the DTFT of any periodic discrete-time sequence with an integer-length period.

2.2.1.3.2 DFT and DTFT of finite-length data

When a discrete-time sequence happens to equal zero for all samples except for those between 0 and $N - 1$, the infinite sum in the DTFT (2.3) equation becomes the same as the finite sum from 0 to $N - 1$ in the DFT (2.5) equation. By matching the arguments in the exponential terms, we observe that the DFT values *exactly* equal the DTFT for specific DTFT frequencies $\omega_k = \frac{2\pi k}{N}$. That is, the DFT computes exact samples of the DTFT at N equally spaced frequencies $\omega_k = \frac{2\pi k}{N}$, or

$$X \left(\omega_k = \frac{2\pi k}{N} \right) = \sum_{n=-\infty}^{\infty} \left(x(n) e^{-j\omega_k n} \right) = \sum_{n=0}^{N-1} \left(x(n) e^{-j\frac{2\pi nk}{N}} \right) = X(k)$$

2.2.1.3.3 DFT as a DTFT approximation

In most cases, the signal is neither exactly periodic nor truly of finite length; in such cases, the DFT of a finite block of N consecutive discrete-time samples does *not* exactly equal samples of the DTFT at specific frequencies. Instead, the DFT (2.5) gives frequency samples of a windowed (truncated) DTFT (2.3)

$$\hat{X}\left(\omega_k = \frac{2\pi k}{N}\right) = \sum_{n=0}^{N-1} \left(x(n) e^{-j\omega_k n}\right) = \sum_{n=-\infty}^{\infty} \left(x(n) w(n) e^{-j\omega_k n}\right) = X(k)$$

where $w(n) = \begin{cases} 1 & \text{if } 0 \leq n < N \\ 0 & \text{if else} \end{cases}$ Once again, $X(k)$ *exactly* equals $X(\omega_k)$ a DTFT frequency sample only when $\forall n, n \notin [0, N-1] : (x(n) = 0)$

2.2.1.4 Relationship between continuous-time FT and DFT

The goal of spectrum analysis is often to determine the frequency content of an analog (continuous-time) signal; very often, as in most modern spectrum analyzers, this is actually accomplished by sampling the analog signal, windowing (truncating) the data, and computing and plotting the magnitude of its DFT. It is thus essential to relate the DFT frequency samples back to the original analog frequency. Assuming that the analog signal is bandlimited and the sampling frequency exceeds twice that limit so that no frequency aliasing occurs, the relationship between the continuous-time Fourier frequency Ω (in radians) and the DTFT frequency ω imposed by sampling is $\omega = \Omega T$ where T is the sampling period. Through the relationship $\omega_k = \frac{2\pi k}{N}$ between the DTFT frequency ω and the DFT frequency index k , the correspondence between the DFT frequency index and the original analog frequency can be found:

$$\Omega = \frac{2\pi k}{NT}$$

or in terms of analog frequency f in Hertz (cycles per second rather than radians)

$$f = \frac{k}{NT}$$

for k in the range k between 0 and $\frac{N}{2}$. It is important to note that $k \in [\frac{N}{2} + 1, N-1]$ correspond to *negative* frequencies due to the periodicity of the DTFT and the DFT.

Exercise 2.1

(Solution on p. 189.)

In general, will DFT frequency values $X(k)$ *exactly* equal samples of the analog Fourier transform X_a at the corresponding frequencies? That is, will $X(k) = X_a\left(\frac{2\pi k}{NT}\right)$?

2.2.1.5 Zero-Padding

If more than N equally spaced frequency samples of a length- N signal are desired, they can easily be obtained by **zero-padding** the discrete-time signal and computing a DFT of the longer length. In particular, if LN DTFT (2.3) samples are desired of a length- N sequence, one can compute the length- LN DFT (2.5) of a length- LN zero-padded sequence

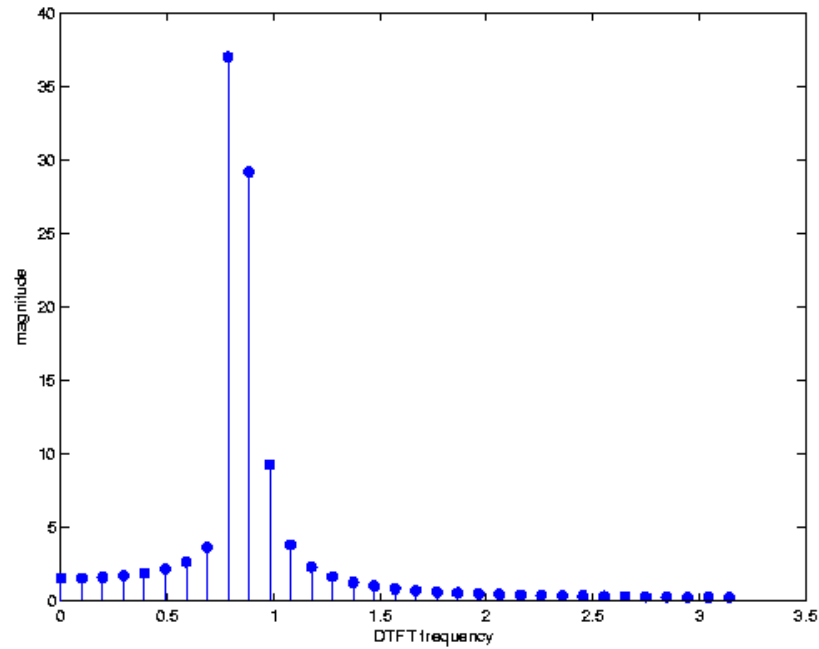
$$z(n) = \begin{cases} x(n) & \text{if } 0 \leq n \leq N-1 \\ 0 & \text{if } N \leq n \leq LN-1 \end{cases}$$

$$X\left(w_k = \frac{2\pi k}{LN}\right) = \sum_{n=0}^{N-1} \left(x(n) e^{-j\left(\frac{2\pi k n}{LN}\right)}\right) = \sum_{n=0}^{LN-1} \left(z(n) e^{-j\left(\frac{2\pi k n}{LN}\right)}\right) = DFT_{LN}[z[n]]$$

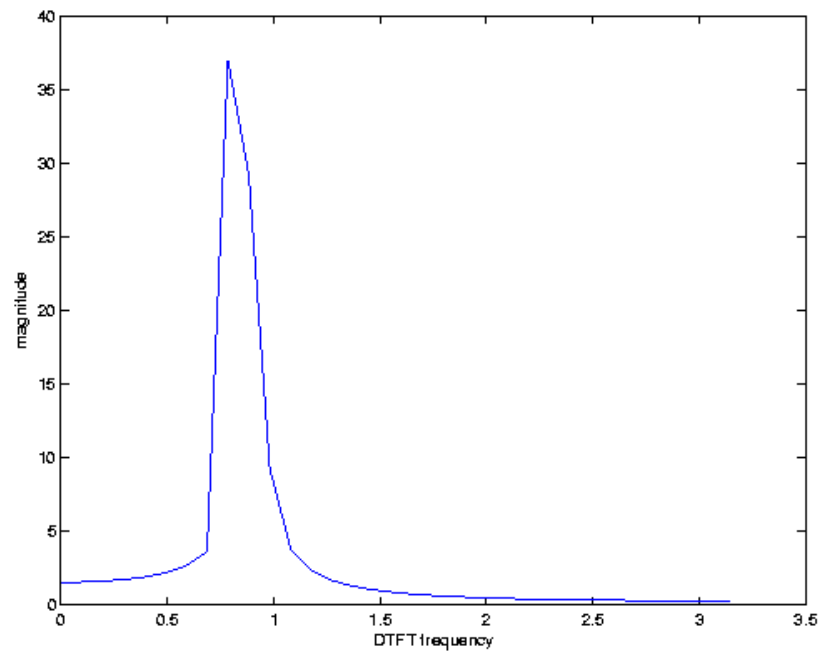
Note that zero-padding *interpolates* the spectrum. One should always zero-pad (by about at least a factor of 4) when using the DFT (2.5) to approximate the DTFT (2.3) to get a clear picture of the DTFT (2.3). While performing computations on zeros may at first seem inefficient, using FFT (Section 2.3.1) algorithms, which generally expect the same number of input and output samples, actually makes this approach very efficient.

Figure 2.3 (Spectrum without zero-padding) shows the magnitude of the DFT values corresponding to the non-negative frequencies of a real-valued length-64 DFT of a length-64 signal, both in a "stem" format to emphasize the discrete nature of the DFT frequency samples, and as a line plot to emphasize its use as an approximation to the continuous-in-frequency DTFT. From this figure, it appears that the signal has a single dominant frequency component.

Spectrum without zero-padding



(a) Stem plot

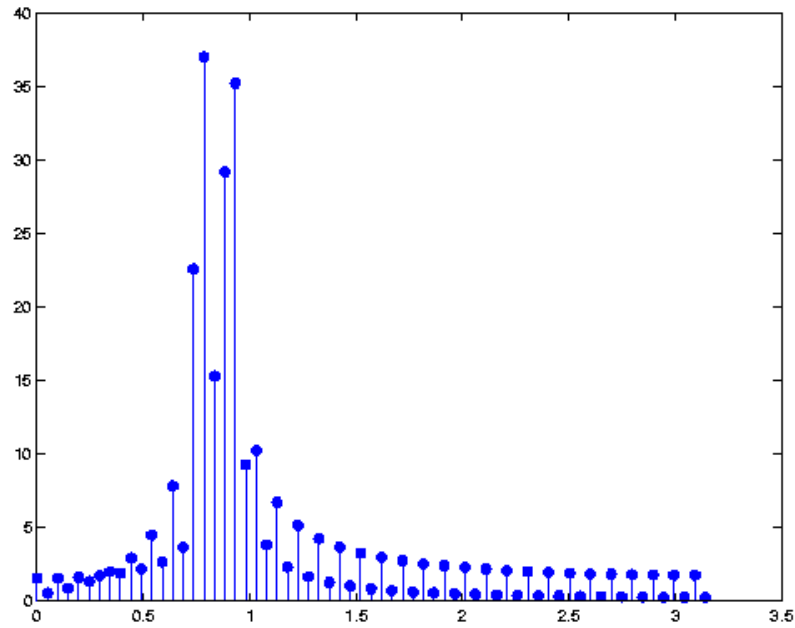


(b) Line Plot

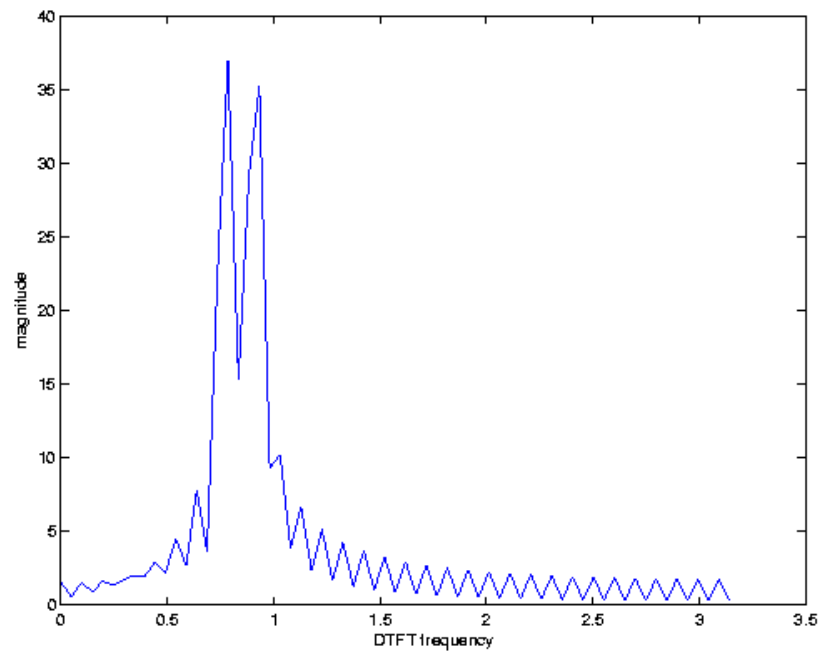
Figure 2.3: Magnitude DFT spectrum of 64 samples of a signal with a length-64 DFT (no zero padding)

Zero-padding by a factor of two by appending 64 zero values to the signal and computing a length-128 DFT yields Figure 2.4 (Spectrum with factor-of-two zero-padding). It can now be seen that the signal consists of at least two narrowband frequency components; the gap between them fell between DFT samples in Figure 2.3 (Spectrum without zero-padding), resulting in a misleading picture of the signal's spectral content. This is sometimes called the **picket-fence effect**, and is a result of insufficient sampling in frequency. While zero-padding by a factor of two has revealed more structure, it is unclear whether the peak magnitudes are reliably rendered, and the jagged linear interpolation in the line graph does not yet reflect the smooth, continuously-differentiable spectrum of the DTFT of a finite-length truncated signal. Errors in the apparent peak magnitude due to insufficient frequency sampling is sometimes referred to as **scalloping loss**.

Spectrum with factor-of-two zero-padding



(a) Stem plot

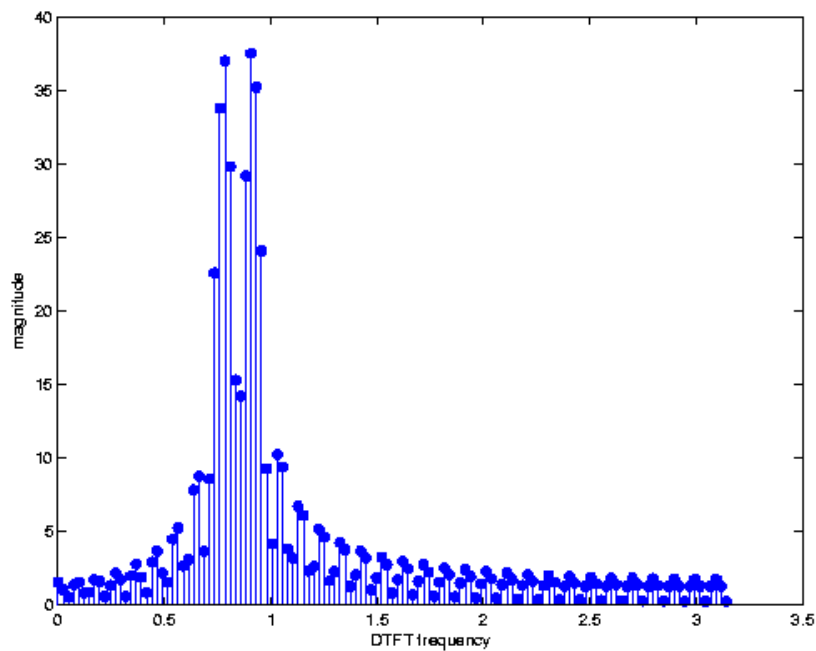


(b) Line Plot

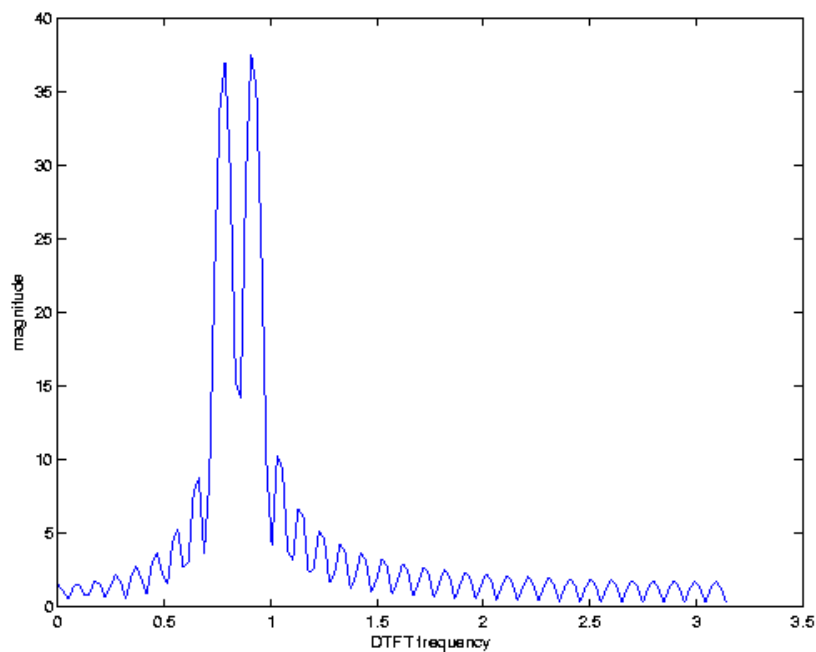
Figure 2.4: Magnitude DFT spectrum of 64 samples of a signal with a length-128 DFT (double-length zero-padding)

Zero-padding to four times the length of the signal, as shown in Figure 2.5 (Spectrum with factor-of-four zero-padding), clearly shows the spectral structure and reveals that the magnitude of the two spectral lines are nearly identical. The line graph is still a bit rough and the peak magnitudes and frequencies may not be precisely captured, but the spectral characteristics of the truncated signal are now clear.

Spectrum with factor-of-four zero-padding



(a) Stem plot

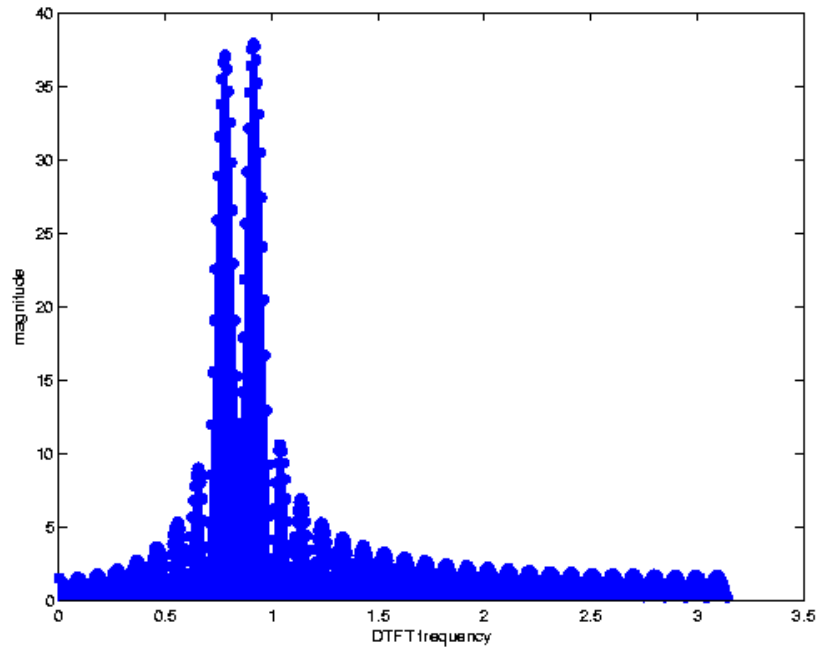


(b) Line Plot

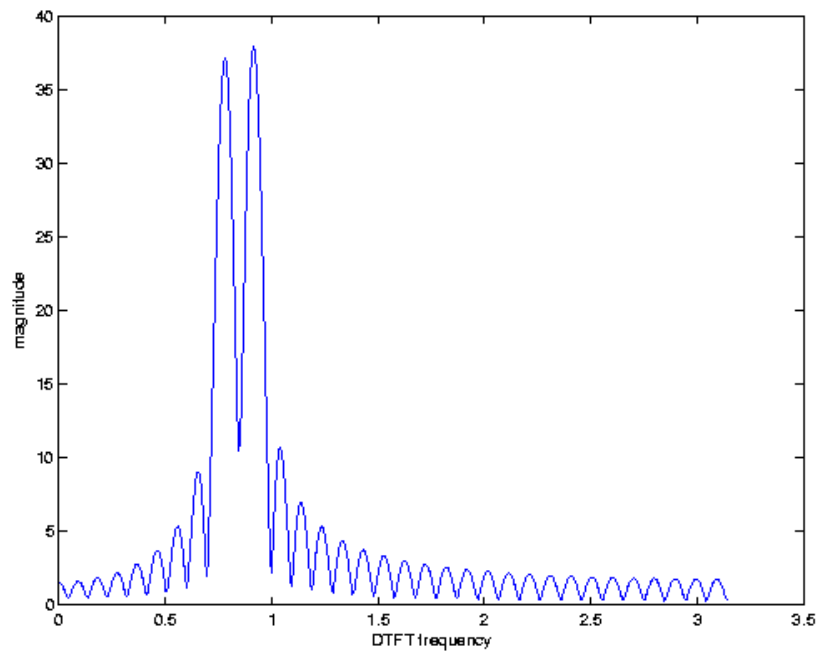
Figure 2.5: Magnitude DFT spectrum of 64 samples of a signal with a length-256 zero-padded DFT (four times zero-padding)

Zero-padding to a length of 1024, as shown in Figure 2.6 (Spectrum with factor-of-sixteen zero-padding) yields a spectrum that is smooth and continuous to the resolution of the computer screen, and produces a very accurate rendition of the DTFT of the *truncated* signal.

Spectrum with factor-of-sixteen zero-padding



(a) Stem plot



(b) Line Plot

Figure 2.6: Magnitude DFT spectrum of 64 samples of a signal with a length-1024 zero-padded DFT. The spectrum now looks smooth and continuous and reveals all the structure of the DTFT of a truncated signal.

The signal used in this example actually consisted of two pure sinusoids of equal magnitude. The slight difference in magnitude of the two dominant peaks, the breadth of the peaks, and the sinc-like lesser **side lobe** peaks throughout frequency are artifacts of the truncation, or windowing, process used to practically approximate the DFT. These problems and partial solutions to them are discussed in the following section.

2.2.1.6 Effects of Windowing

Applying the DTFT multiplication property

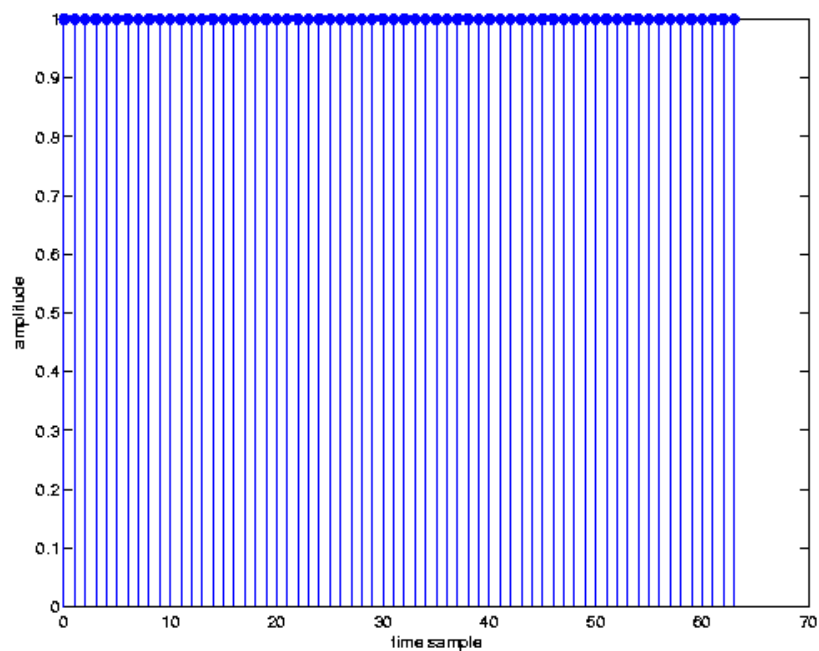
$$X(\hat{\omega}_k) = \sum_{n=-\infty}^{\infty} \left(x(n) w(n) e^{-j\omega_k n} \right) = \frac{1}{2\pi} X(\omega_k) * W(\omega_k)$$

we find that the DFT (2.5) of the windowed (truncated) signal produces samples not of the true (desired) DTFT spectrum $X(\omega)$, but of a *smoothed* version $X(\omega) * W(\omega)$. We want this to resemble $X(\omega)$ as closely as possible, so $W(\omega)$ should be as close to an impulse as possible. The **window** $w(n)$ need not be a simple **truncation** (or **rectangle**, or **boxcar**) window; other shapes can also be used as long as they limit the sequence to at most N consecutive nonzero samples. All good windows are impulse-like, and represent various tradeoffs between three criteria:

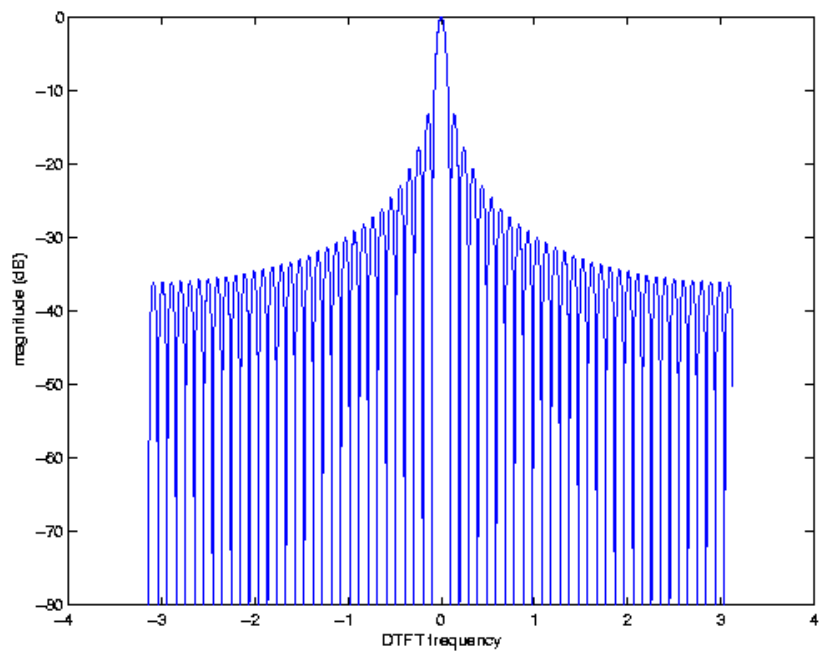
1. main lobe width: (limits resolution of closely-spaced peaks of equal height)
2. height of first sidelobe: (limits ability to see a small peak near a big peak)
3. slope of sidelobe drop-off: (limits ability to see small peaks further away from a big peak)

Many different window functions⁵ have been developed for truncating and shaping a length- N signal segment for spectral analysis. The simple **truncation** window has a periodic sinc DTFT, as shown in Figure 2.7. It has the narrowest main-lobe width, $\frac{2\pi}{N}$ at the -3 dB level and $\frac{4\pi}{N}$ between the two zeros surrounding the main lobe, of the common window functions, but also the largest side-lobe peak, at about -13 dB. The side-lobes also taper off relatively slowly.

⁵http://en.wikipedia.org/wiki/Window_function



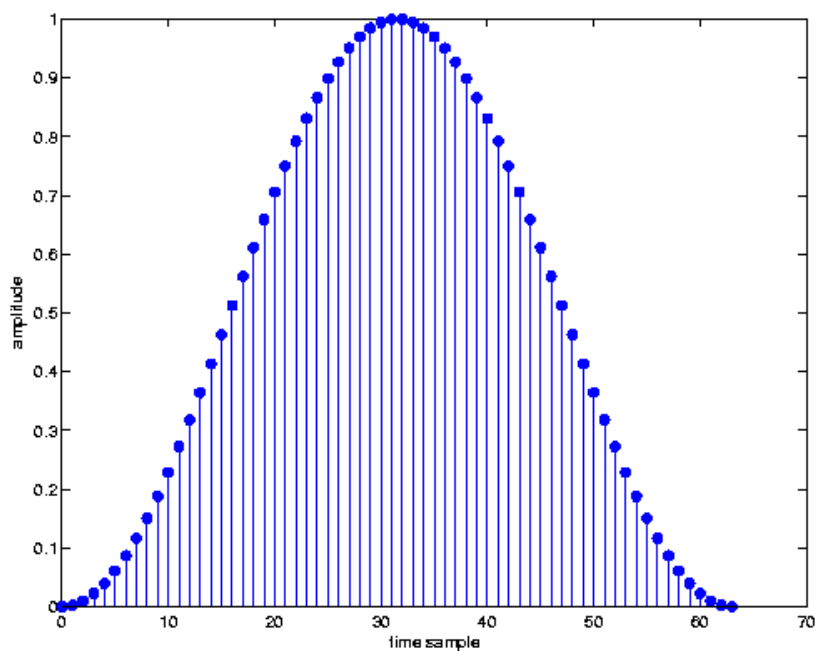
(a) Rectangular window



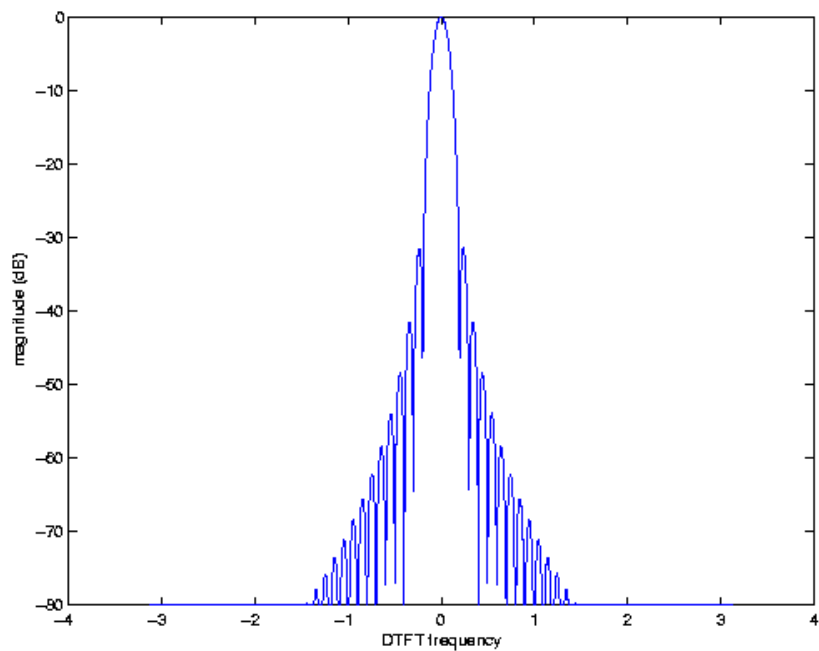
(b) Magnitude of boxcar window spectrum

Figure 2.7: Length-64 truncation (boxcar) window and its magnitude DFT spectrum

The **Hann window** (sometimes also called the **hanning** window), illustrated in Figure 2.8, takes the form $w[n] = 0.5 - 0.5\cos\left(\frac{2\pi n}{N-1}\right)$ for n between 0 and $N-1$. It has a main-lobe width (about $\frac{3\pi}{N}$ at the -3 dB level and $\frac{8\pi}{N}$ between the two zeros surrounding the main lobe) considerably larger than the rectangular window, but the largest side-lobe peak is much lower, at about -31.5 dB. The side-lobes also taper off much faster. For a given length, this window is worse than the boxcar window at separating closely-spaced spectral components of similar magnitude, but better for identifying smaller-magnitude components at a greater distance from the larger components.



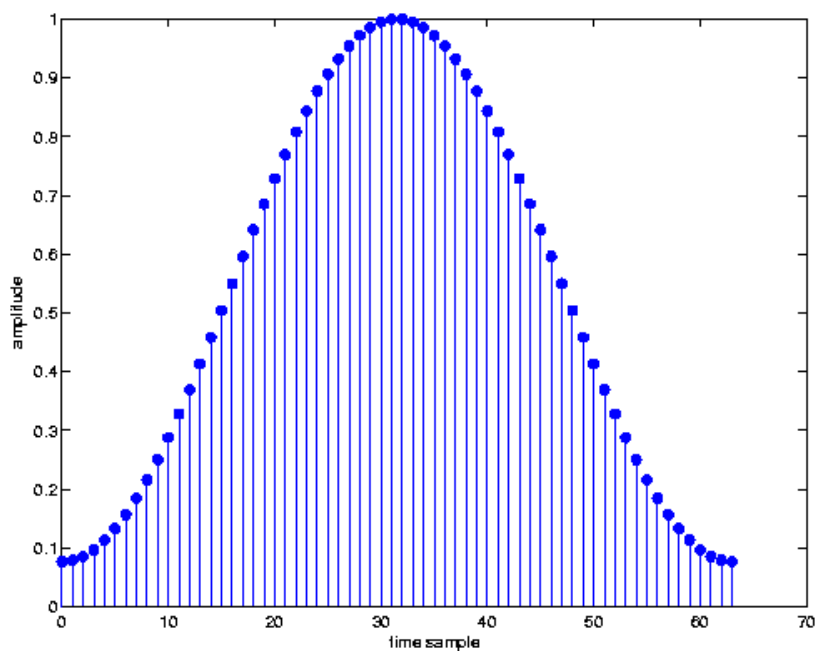
(a) Hann window



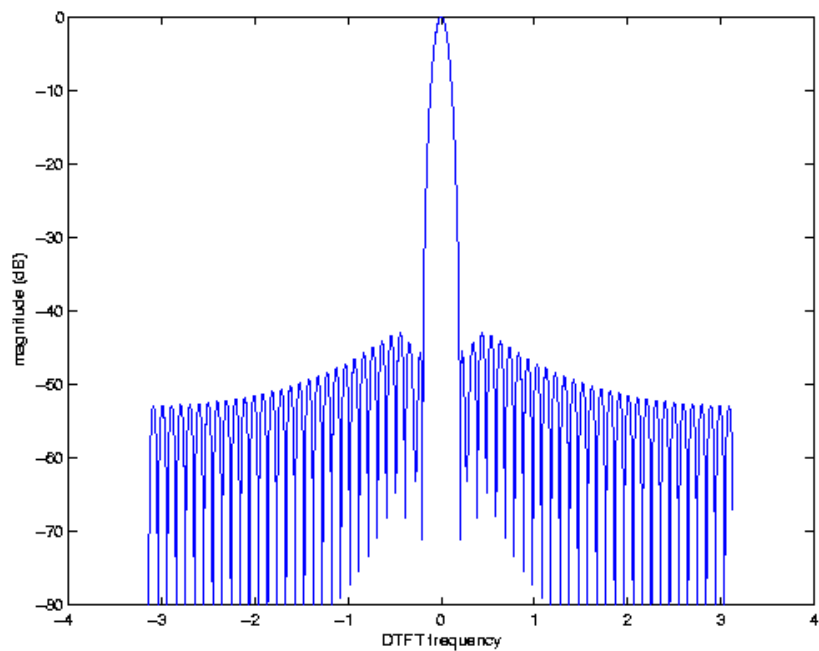
(b) Magnitude of Hann window spectrum

Figure 2.8: Length-64 Hann window and its magnitude DFT spectrum

The **Hamming window**, illustrated in Figure 2.9, has a form similar to the Hann window but with slightly different constants: $w[n] = 0.538 - 0.462 \cos\left(\frac{2\pi n}{N-1}\right)$ for n between 0 and $N-1$. Since it is composed of the same Fourier series harmonics as the Hann window, it has a similar main-lobe width (a bit less than $\frac{3\pi}{N}$ at the -3 dB level and $\frac{8\pi}{N}$ between the two zeros surrounding the main lobe), but the largest side-lobe peak is much lower, at about -42.5 dB. However, the side-lobes also taper off much more slowly than with the Hann window. For a given length, the Hamming window is better than the Hann (and of course the boxcar) windows at separating a small component relatively near to a large component, but worse than the Hann for identifying very small components at considerable frequency separation. Due to their shape and form, the Hann and Hamming windows are also known as **raised-cosine windows**.



(a) Hamming window



(b) Magnitude of Hamming window spectrum

Figure 2.9: Length-64 Hamming window and its magnitude DFT spectrum

NOTE: Standard even-length windows are symmetric around a point halfway between the window samples $\frac{N}{2} - 1$ and $\frac{N}{2}$. For some applications such as time-frequency analysis (Section 2.2.3), it may be important to align the window perfectly to a sample. In such cases, a **DFT-symmetric** window that is symmetric around the $\frac{N}{2}$ -th sample can be used. For example, the DFT-symmetric Hamming window is $w[n] = 0.538 - 0.462\cos\left(\frac{2\pi n}{N}\right)$. A DFT-symmetric window has a purely real-valued DFT and DTFT. DFT-symmetric versions of windows, such as the Hamming and Hann windows, composed of few discrete Fourier series terms of period N , have few non-zero DFT terms (only when *not* zero-padded) and can be used efficiently in running FFTs (Section 2.3.2).

The main-lobe width of a window is an inverse function of the window-length N ; for any type of window, a longer window will always provide better resolution.

Many other windows exist that make various other tradeoffs between main-lobe width, height of largest side-lobe, and side-lobe rolloff rate. The Kaiser window⁶ family, based on a modified Bessel function, has an adjustable parameter that allows the user to tune the tradeoff over a continuous range. The Kaiser window has near-optimal time-frequency resolution and is widely used. A list of many different windows can be found here⁷.

Example 2.1

Figure 2.10 shows 64 samples of a real-valued signal composed of several sinusoids of various frequencies and amplitudes.

⁶http://en.wikipedia.org/wiki/Kaiser_window

⁷http://en.wikipedia.org/wiki/Window_function

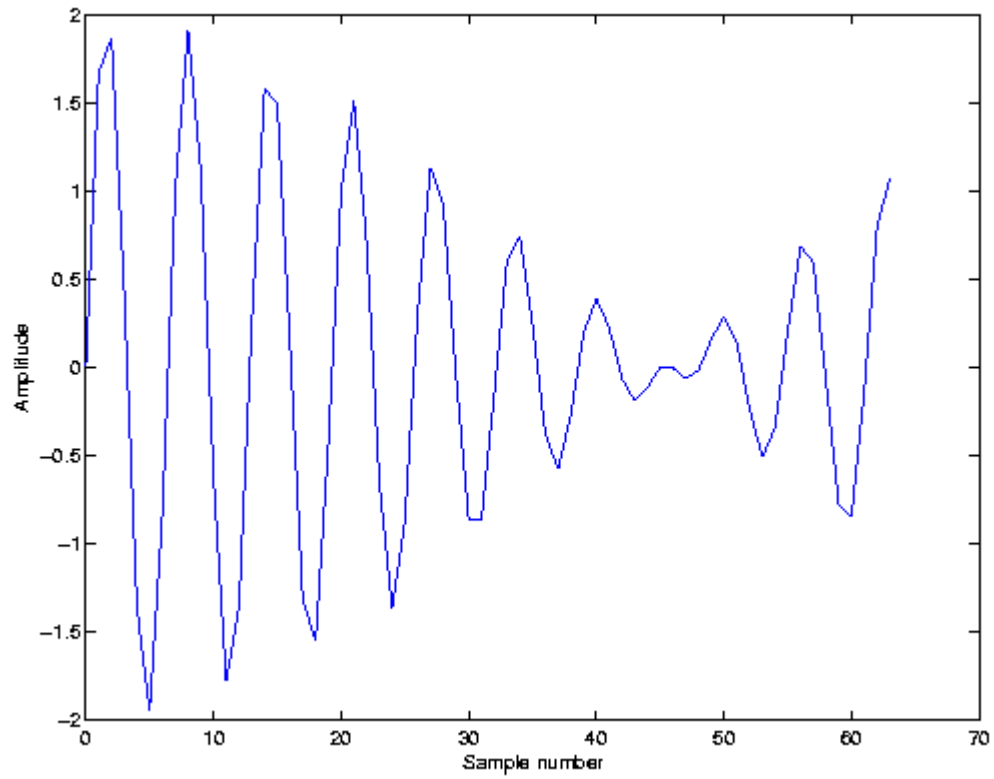


Figure 2.10: 64 samples of an unknown signal

Figure 2.11 shows the magnitude (in dB) of the positive frequencies of a length-1024 zero-padded DFT of this signal (that is, using a simple truncation, or rectangular, window).

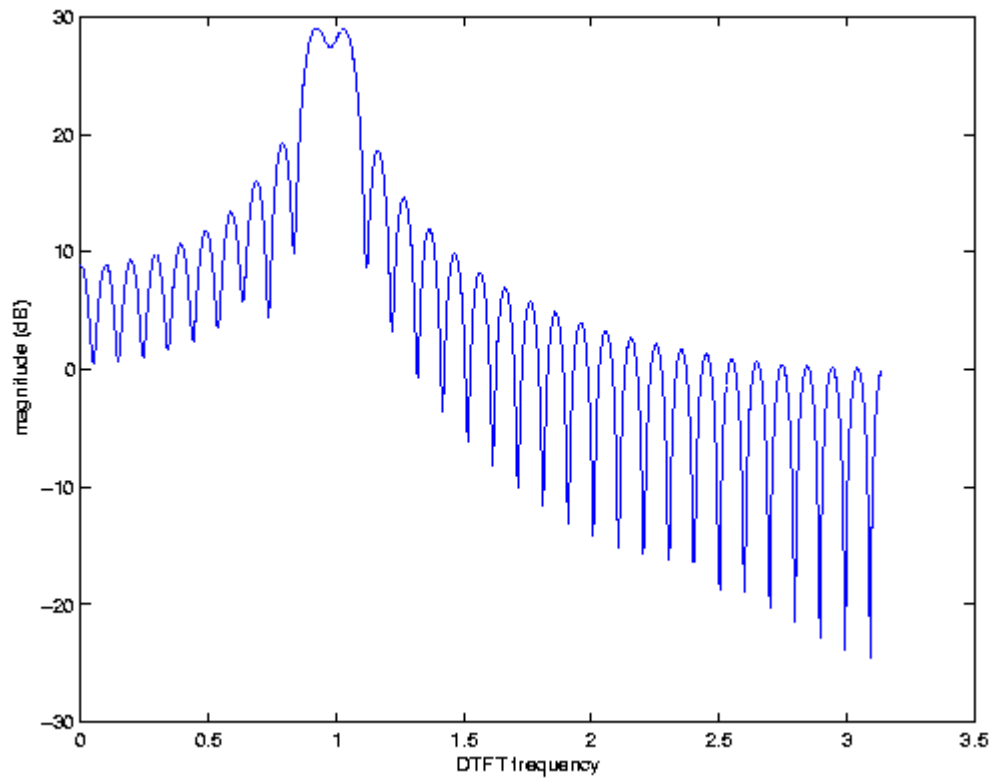


Figure 2.11: Magnitude (in dB) of the zero-padded DFT spectrum of the signal in Figure 2.10 using a simple length-64 rectangular window

From this spectrum, it is clear that the signal has two large, nearby frequency components with frequencies near 1 radian of essentially the same magnitude.

Figure 2.12 shows the spectral estimate produced using a length-64 Hamming window applied to the same signal shown in Figure 2.10.

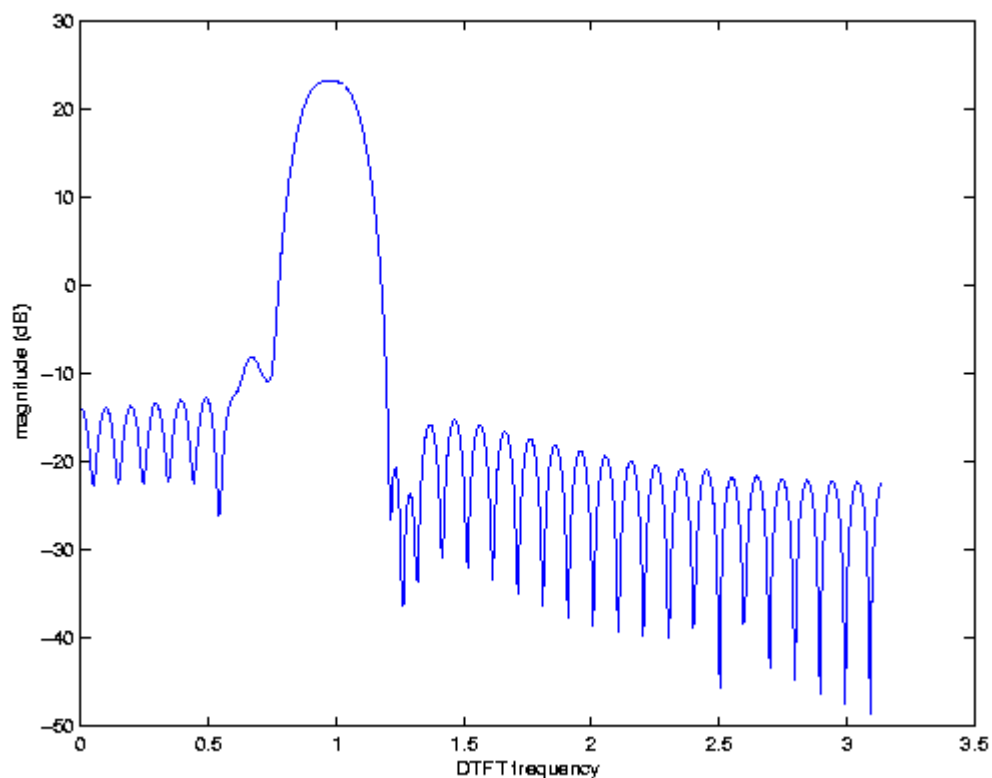


Figure 2.12: Magnitude (in dB) of the zero-padded DFT spectrum of the signal in Figure 2.10 using a length-64 Hamming window

The two large spectral peaks can no longer be resolved; they blur into a single broad peak due to the reduced spectral resolution of the broader main lobe of the Hamming window. However, the lower side-lobes reveal a third component at a frequency of about 0.7 radians at about 35 dB lower magnitude than the larger components. This component was entirely buried under the side-lobes when the rectangular window was used, but now stands out well above the much lower nearby side-lobes of the Hamming window.

Figure 2.13 shows the spectral estimate produced using a length-64 Hann window applied to the same signal shown in Figure 2.10.

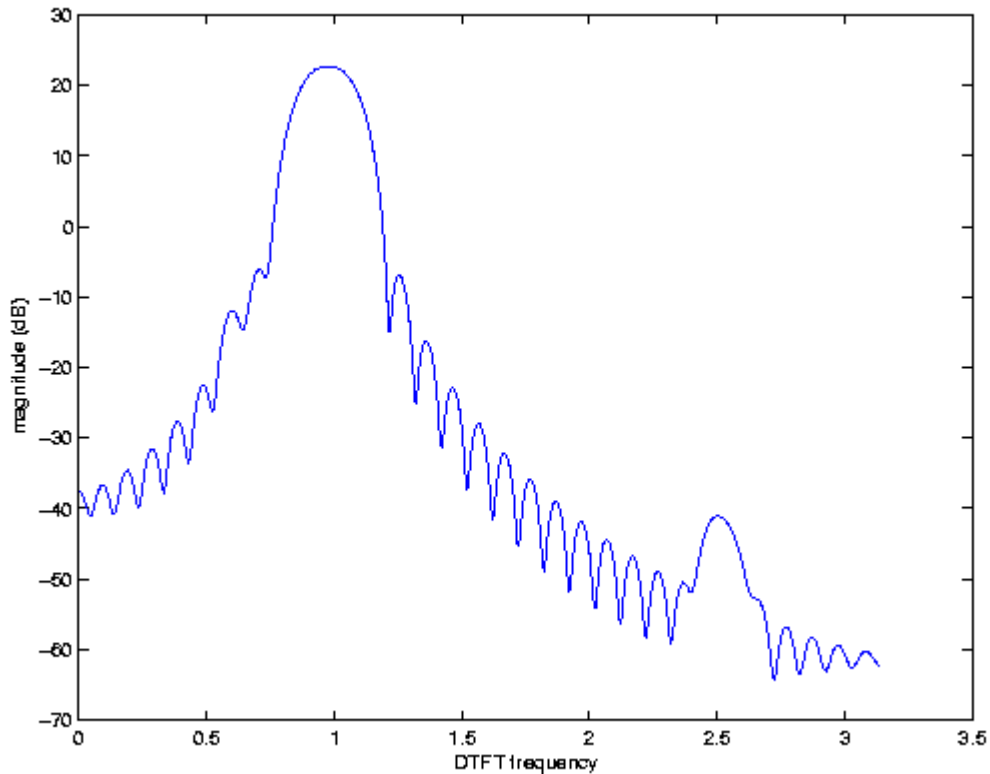


Figure 2.13: Magnitude (in dB) of the zero-padded DFT spectrum of the signal in Figure 2.10 using a length-64 Hann window

The two large components again merge into a single peak, and the smaller component observed with the Hamming window is largely lost under the higher nearby side-lobes of the Hann window. However, due to the much faster side-lobe rolloff of the Hann window's spectrum, a fourth component at a frequency of about 2.5 radians with a magnitude about 65 dB below that of the main peaks is now clearly visible.

This example illustrates that no single window is best for all spectrum analyses. The best window depends on the nature of the signal, and different windows may be better for different components of the same signal. A skilled spectrum analyst may apply several different windows to a signal to gain a fuller understanding of the data.

2.2.2 Classical Statistical Spectral Estimation⁸

Many signals are either partly or wholly stochastic, or random. Important examples include human speech, vibration in machines, and CDMA⁹ communication signals. Given the ever-present noise in electronic systems, it can be argued that almost *all* signals are at least partly stochastic. Such signals may have a

⁸This content is available online at <<http://cnx.org/content/m12014/1.3/>>.

⁹<http://en.wikipedia.org/wiki/Cdma>

distinct *average* spectral structure that reveals important information (such as for speech recognition or early detection of damage in machinery). Spectrum analysis of any single block of data using window-based deterministic spectrum analysis (Section 2.2.1), however, produces a random spectrum that may be difficult to interpret. For such situations, the classical statistical spectrum estimation methods described in this module can be used.

The goal in classical statistical spectrum analysis is to estimate $E \left[(|X(\omega)|)^2 \right]$, the **power spectral density (PSD)** across frequency of the stochastic signal. That is, the goal is to find the expected (mean, or average) energy density of the signal as a function of frequency. (For zero-mean signals, this equals the variance of each frequency sample.) Since the spectrum of each block of signal samples is itself random, we must average the squared spectral magnitudes over a number of blocks of data to find the mean. There are two main classical approaches, the periodogram (Section 2.2.2.1: Periodogram method) and auto-correlation (Section 2.2.2.2: Auto-correlation-based approach) methods.

2.2.2.1 Periodogram method

The periodogram method divides the signal into a number of shorter (and often overlapped) blocks of data, computes the squared magnitude of the windowed (Section 2.2.1.6: Effects of Windowing) (and usually zero-padded (Section 2.2.1.5: Zero-Padding)) DFT (2.5), $X_i(\omega_k)$, of each block, and averages them to estimate the power spectral density. The squared magnitudes of the DFTs of L possibly overlapped length- N windowed blocks of signal (each probably with zero-padding (Section 2.2.1.5: Zero-Padding)) are averaged to estimate the power spectral density:

$$\hat{X}(\omega_k) = \frac{1}{L} \sum_{i=1}^L \left((|X_i(\omega_k)|)^2 \right)$$

For a fixed *total* number of samples, this introduces a tradeoff: Larger individual data blocks provides better frequency resolution due to the use of a longer window, but it means there are less blocks to average, so the estimate has higher variance and appears more noisy. The best tradeoff depends on the application. Overlapping blocks by a factor of two to four increases the number of averages and reduces the variance, but since the same data is being reused, still more overlapping does not further reduce the variance. As with any window-based spectrum estimation (Section 2.2.1.6: Effects of Windowing) procedure, the window function introduces broadening and sidelobes into the power spectrum estimate. That is, the periodogram produces an estimate of the *windowed* spectrum $\hat{X}(\omega) = E \left[(|X(\omega) * W_M|)^2 \right]$, not of $E \left[(|X(\omega)|)^2 \right]$.

Example 2.2

Figure 2.14 shows the non-negative frequencies of the DFT (zero-padded to 1024 total samples) of 64 samples of a real-valued stochastic signal.

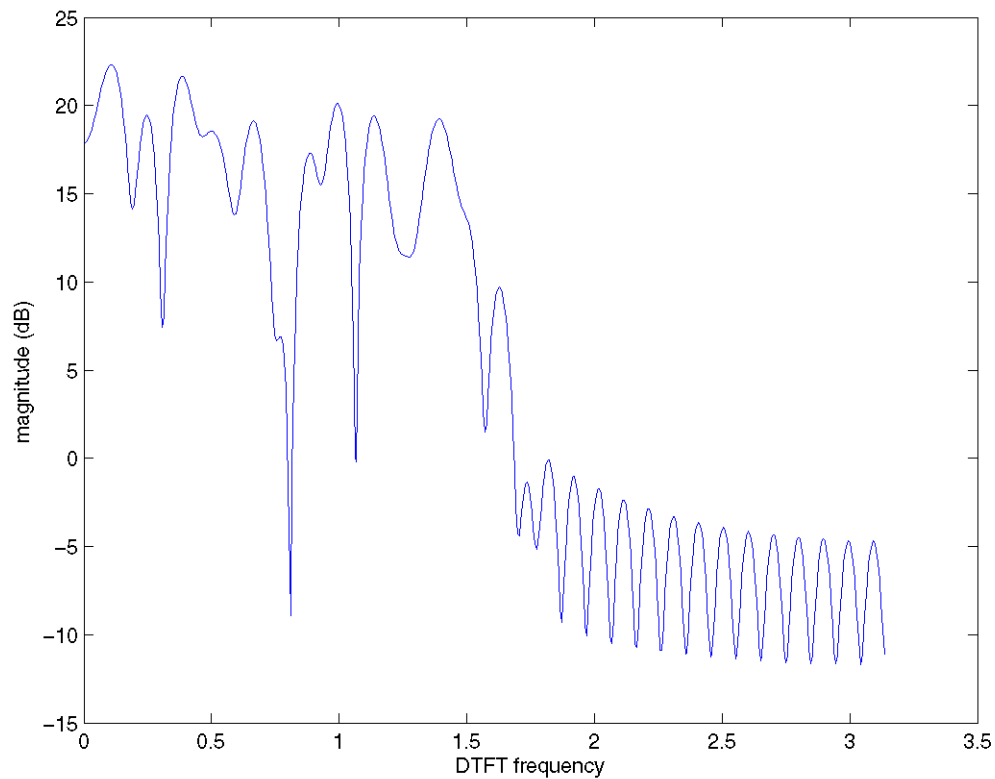


Figure 2.14: DFT magnitude (in dB) of 64 samples of a stochastic signal

With no averaging, the power spectrum is very noisy and difficult to interpret other than noting a significant reduction in spectral energy above about half the Nyquist frequency. Various peaks and valleys appear in the lower frequencies, but it is impossible to say from this figure whether they represent actual structure in the power spectral density (PSD) or simply random variation in this single realization. Figure 2.15 shows the same frequencies of a length-1024 DFT of a length-1024 signal. While the frequency resolution has improved, there is still no averaging, so it remains difficult to understand the power spectral density of this signal. Certain small peaks in frequency might represent narrowband components in the spectrum, or may just be random noise peaks.

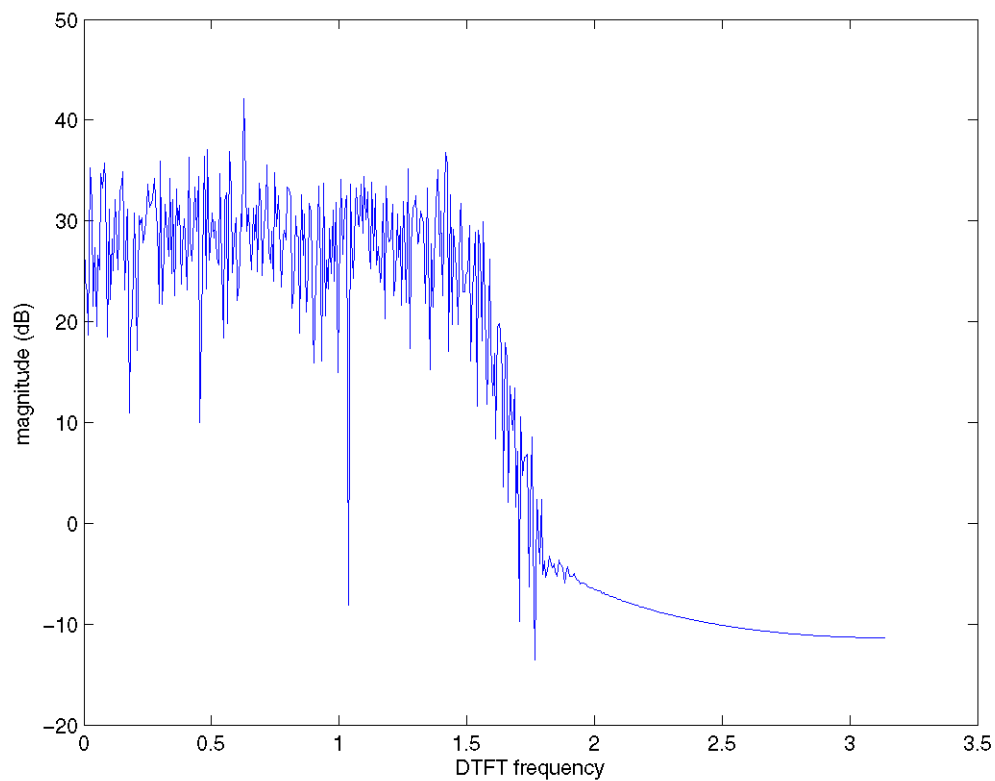


Figure 2.15: DFT magnitude (in dB) of 1024 samples of a stochastic signal

In Figure 2.16, a power spectral density computed from averaging the squared magnitudes of length-1024 zero-padded DFTs of 508 length-64 blocks of data (overlapped by a factor of four, or a 16-sample step between blocks) are shown.

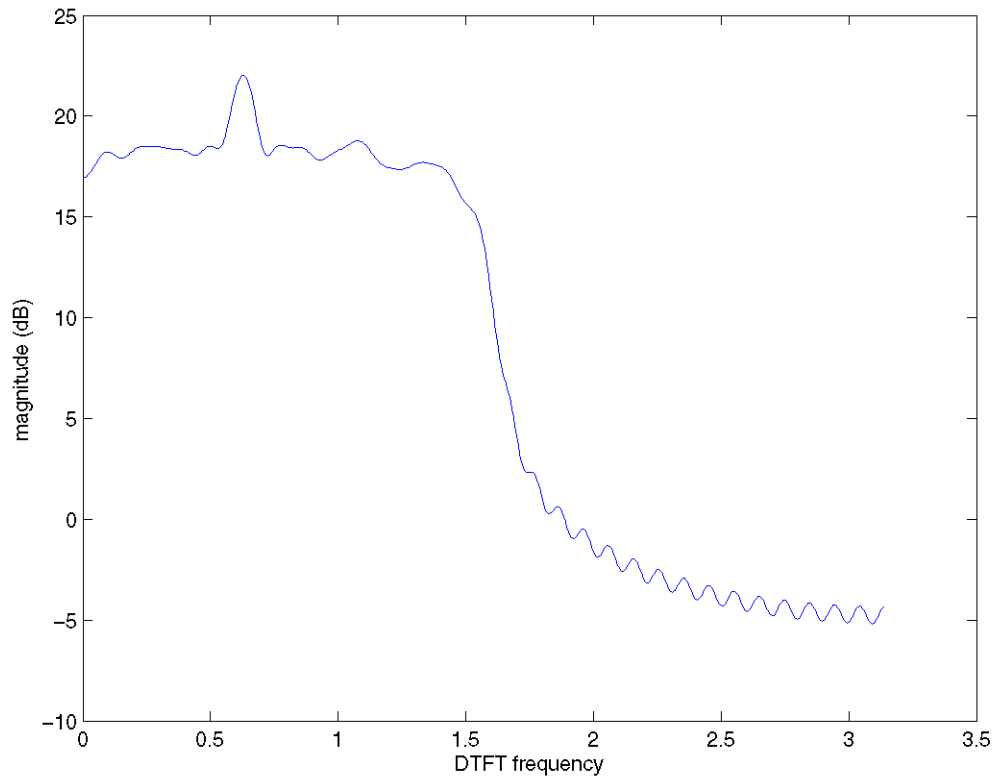


Figure 2.16: Power spectrum density estimate (in dB) of 1024 samples of a stochastic signal

While the frequency resolution corresponds to that of a length-64 truncation window, the averaging greatly reduces the variance of the spectral estimate and allows the user to reliably conclude that the signal consists of lowpass broadband noise with a flat power spectrum up to half the Nyquist frequency, with a stronger narrowband frequency component at around 0.65 radians.

2.2.2.2 Auto-correlation-based approach

The averaging necessary to estimate a power spectral density can be performed in the discrete-time domain, rather than in frequency, using the auto-correlation method. The squared magnitude of the frequency response, from the DTFT multiplication and conjugation properties, corresponds in the discrete-time domain to the signal convolved with the time-reverse of itself,

$$\left(|X(\omega)|\right)^2 = X(\omega) X^*(\omega) \leftrightarrow (x(n), x^*(-n)) = r(n)$$

or its **auto-correlation**

$$r(n) = \sum (x(k) x^*(n+k))$$

We can thus compute the squared magnitude of the spectrum of a signal by computing the DFT of its auto-correlation. For stochastic signals, the power spectral density is an expectation, or average, and by

linearity of expectation can be found by transforming the average of the auto-correlation. For a finite block of N signal samples, the average of the autocorrelation values, $r(n)$, is

$$r(n) = \frac{1}{N-n} \sum_{k=0}^{N-(1-n)} (x(k) x^*(n+k))$$

Note that with increasing **lag**, n , fewer values are averaged, so they introduce more noise into the estimated power spectrum. By windowing (Section 2.2.1.6: Effects of Windowing) the auto-correlation before transforming it to the frequency domain, a less noisy power spectrum is obtained, at the expense of less resolution. The multiplication property of the DTFT shows that the windowing smooths the resulting power spectrum via convolution with the DTFT of the window:

$$X(\omega) = \sum_{n=-M}^M \left(r(n) w(n) e^{-j\omega n} \right) = \left(E \left[|X(\omega)|^2 \right] \right) * W(\omega)$$

This yields another important interpretation of how the auto-correlation method works: it estimates the power spectral density by *averaging the power spectrum over nearby frequencies*, through convolution with the window function's transform, to reduce variance. Just as with the periodogram approach, there is always a variance vs. resolution tradeoff. The periodogram and the auto-correlation method give similar results for a similar amount of averaging; the user should simply note that in the periodogram case, the window introduces smoothing of the spectrum via frequency convolution *before* squaring the magnitude, whereas the periodogram convolves the squared magnitude with $W(\omega)$.

2.2.3 Short Time Fourier Transform¹⁰

2.2.3.1 Short Time Fourier Transform

The Fourier transforms (FT, DTFT, DFT, *etc.*) do not clearly indicate how the frequency content of a signal changes over time.

That information is hidden in the phase - it is not revealed by the plot of the magnitude of the spectrum.

NOTE: To see how the frequency content of a signal changes over time, we can cut the signal into blocks and compute the spectrum of each block.

To improve the result,

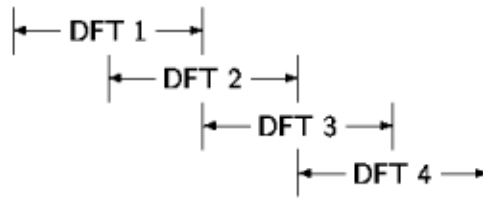
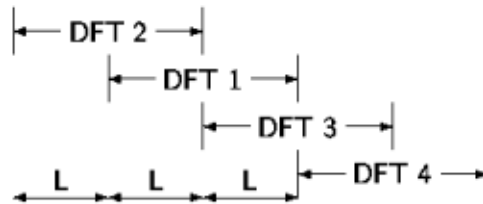
1. blocks are overlapping
2. each block is multiplied by a window that is tapered at its endpoints.

Several parameters must be chosen:

- Block length, R .
- The type of window.
- Amount of overlap between blocks. (Figure 2.17 (STFT: Overlap Parameter))
- Amount of zero padding, if any.

¹⁰This content is available online at <<http://cnx.org/content/m10570/2.4/>>.

STFT: Overlap Parameter

NO OVERLAP**R/4 OVERLAP****R/2 OVERLAP****The parameter L**

L is the number of samples between adjacent blocks.

Figure 2.17

The short-time Fourier transform is defined as

$$\begin{aligned}
 X(\omega, m) &= (STFT(x(n)) := DTFT(x(n-m)w(n))) \\
 &= \sum_{n=-\infty}^{\infty} (x(n-m)w(n)e^{-j\omega n}) \\
 &= \sum_{n=0}^{R-1} (x(n-m)w(n)e^{-j\omega n})
 \end{aligned} \tag{2.9}$$

where $w(n)$ is the window function of length R .

1. The STFT of a signal $x(n)$ is a function of two variables: time and frequency.
2. The block length is determined by the support of the window function $w(n)$.
3. A graphical display of the magnitude of the STFT, $|X(\omega, m)|$, is called the **spectrogram** of the signal. It is often used in speech processing.
4. The STFT of a signal is invertible.
5. One can choose the block length. A long block length will provide higher frequency resolution (because the main-lobe of the window function will be narrow). A short block length will provide higher time resolution because less averaging across samples is performed for each STFT value.
6. A **narrow-band spectrogram** is one computed using a relatively long block length R , (long window function).
7. A **wide-band spectrogram** is one computed using a relatively short block length R , (short window function).

2.2.3.1.1 Sampled STFT

To numerically evaluate the STFT, we sample the frequency axis ω in N equally spaced samples from $\omega = 0$ to $\omega = 2\pi$.

$$\forall k, 0 \leq k \leq N-1 : \left(\omega_k = \frac{2\pi}{N}k \right) \quad (2.10)$$

We then have the discrete STFT,

$$\begin{aligned} (X^d(k, m) := X(\frac{2\pi}{N}k, m)) &= \sum_{n=0}^{R-1} (x(n-m)w(n)e^{-j\omega n}) \\ &= \sum_{n=0}^{R-1} (x(n-m)w(n)W_N^{-(kn)}) \\ &= DFT_N(x(n-m)w(n)|_{n=0}^{R-1}, 0, \dots, 0) \end{aligned} \quad (2.11)$$

where $0, \dots, 0$ is $N-R$.

In this definition, the overlap between adjacent blocks is $R-1$. The signal is shifted along the window one sample at a time. That generates more points than is usually needed, so we also sample the STFT along the time direction. That means we usually evaluate

$$X^d(k, Lm)$$

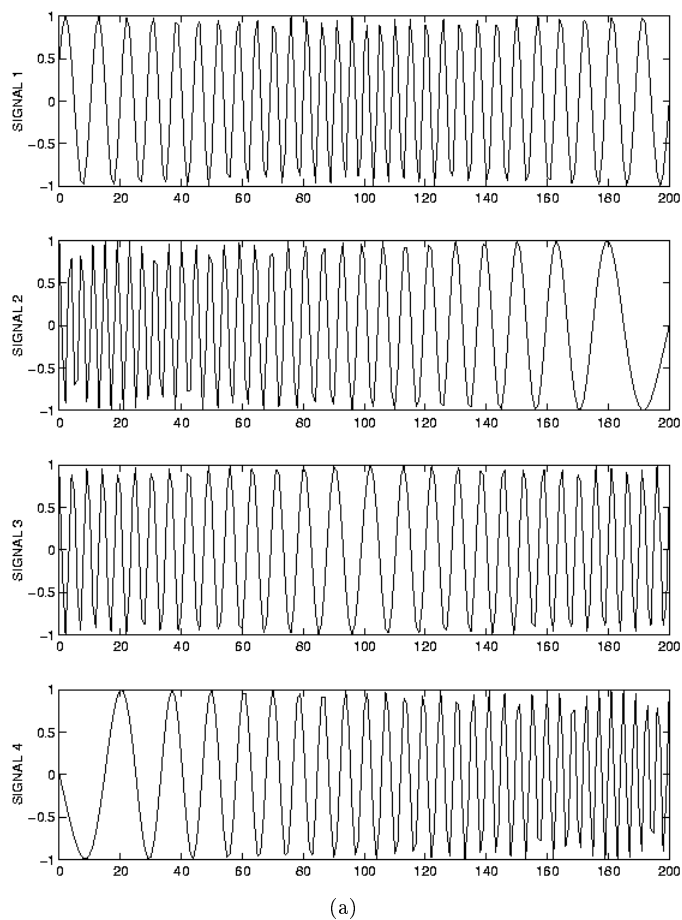
where L is the time-skip. The relation between the time-skip, the number of overlapping samples, and the block length is

$$Overlap = R - L$$

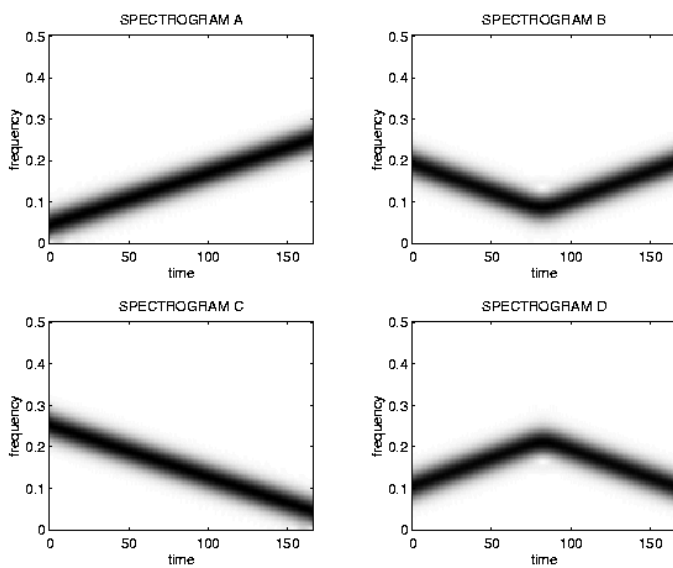
Exercise 2.2

Match each signal to its spectrogram in Figure 2.18.

(Solution on p. 189.)



(a)



(b)

Figure 2.18

2.2.3.1.2 Spectrogram Example

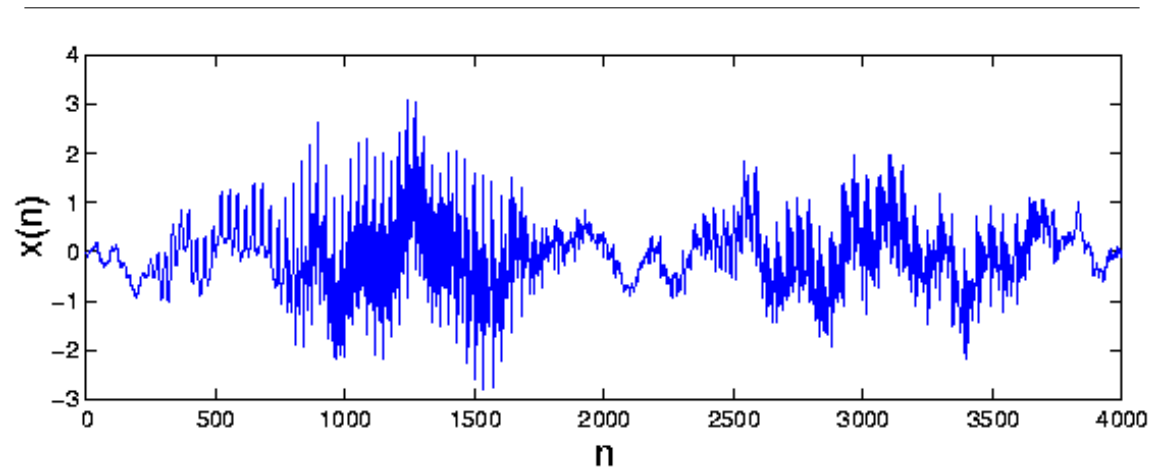


Figure 2.19

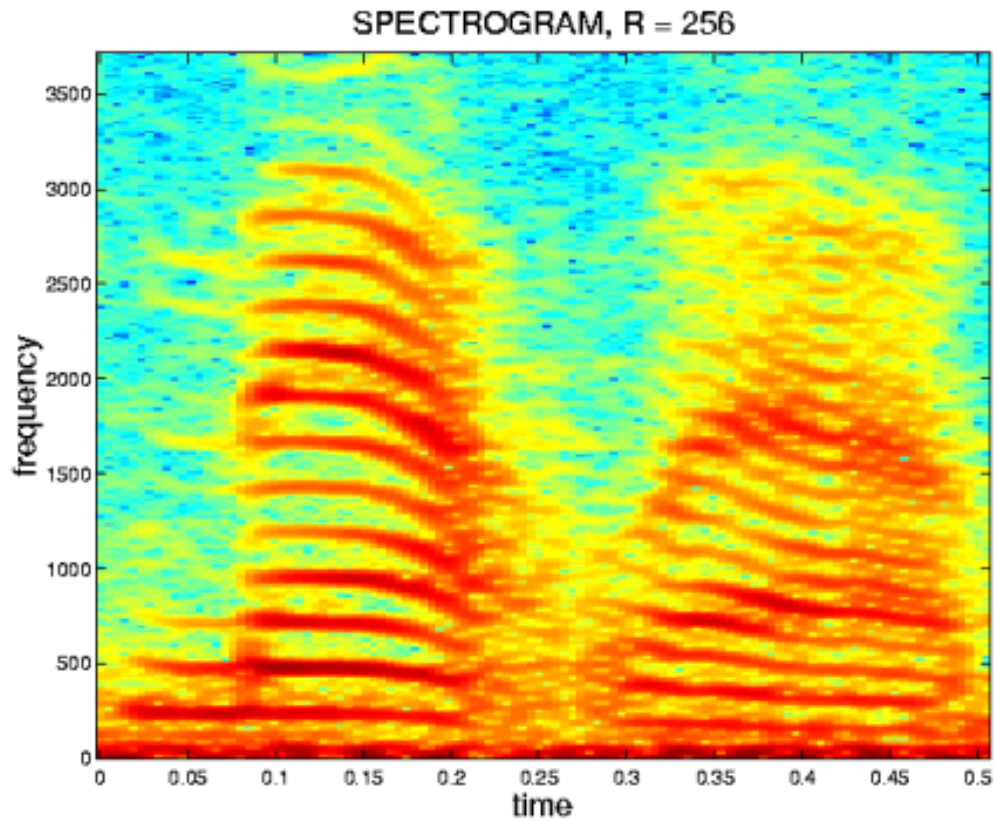


Figure 2.20

The matlab program for producing the figures above (Figure 2.19 and Figure 2.20).

```
% LOAD DATA
load mtlb;
x = mtlb;

figure(1), clf
plot(0:4000,x)
xlabel('n')
ylabel('x(n)')

% SET PARAMETERS
R = 256;           % R: block length
window = hamming(R); % window function of length R
N = 512;          % N: frequency discretization
L = 35;           % L: time lapse between blocks
fs = 7418;        % fs: sampling frequency
```

```
overlap = R - L;

% COMPUTE SPECTROGRAM
[B,f,t] = specgram(x,N,fs>window,overlap);

% MAKE PLOT
figure(2), clf
imagesc(t,f,log10(abs(B)));
colormap('jet')
axis xy
xlabel('time')
ylabel('frequency')
title('SPECTROGRAM, R = 256')
```

2.2.3.1.3 Effect of window length R

Narrow-band spectrogram: better frequency resolution

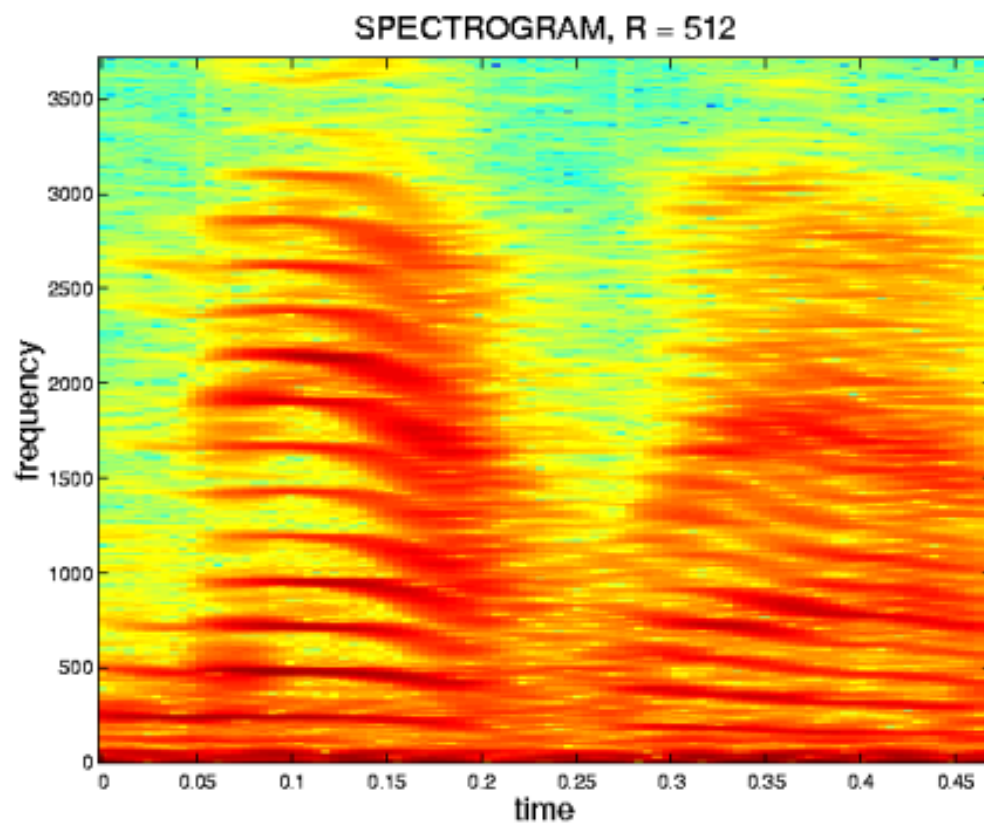


Figure 2.21

Wide-band spectrogram: better time resolution

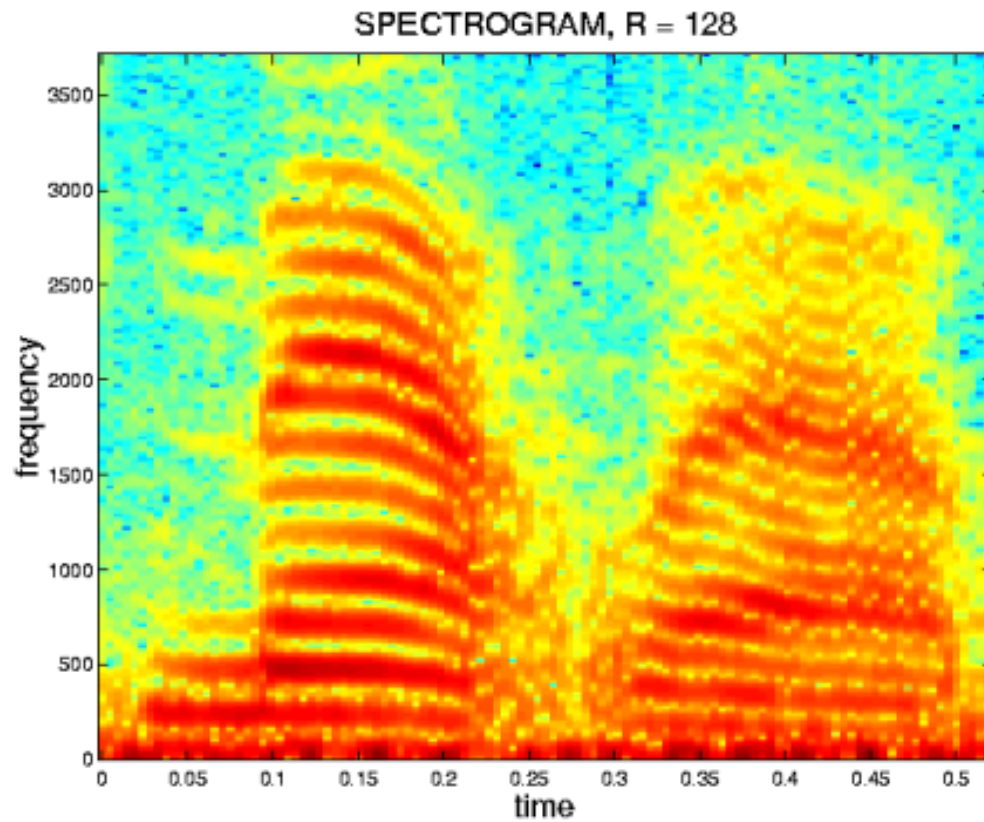
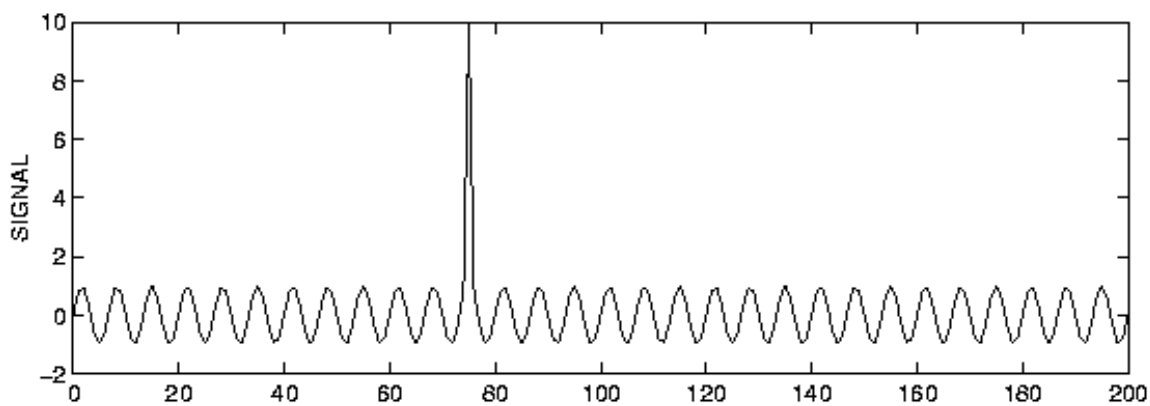


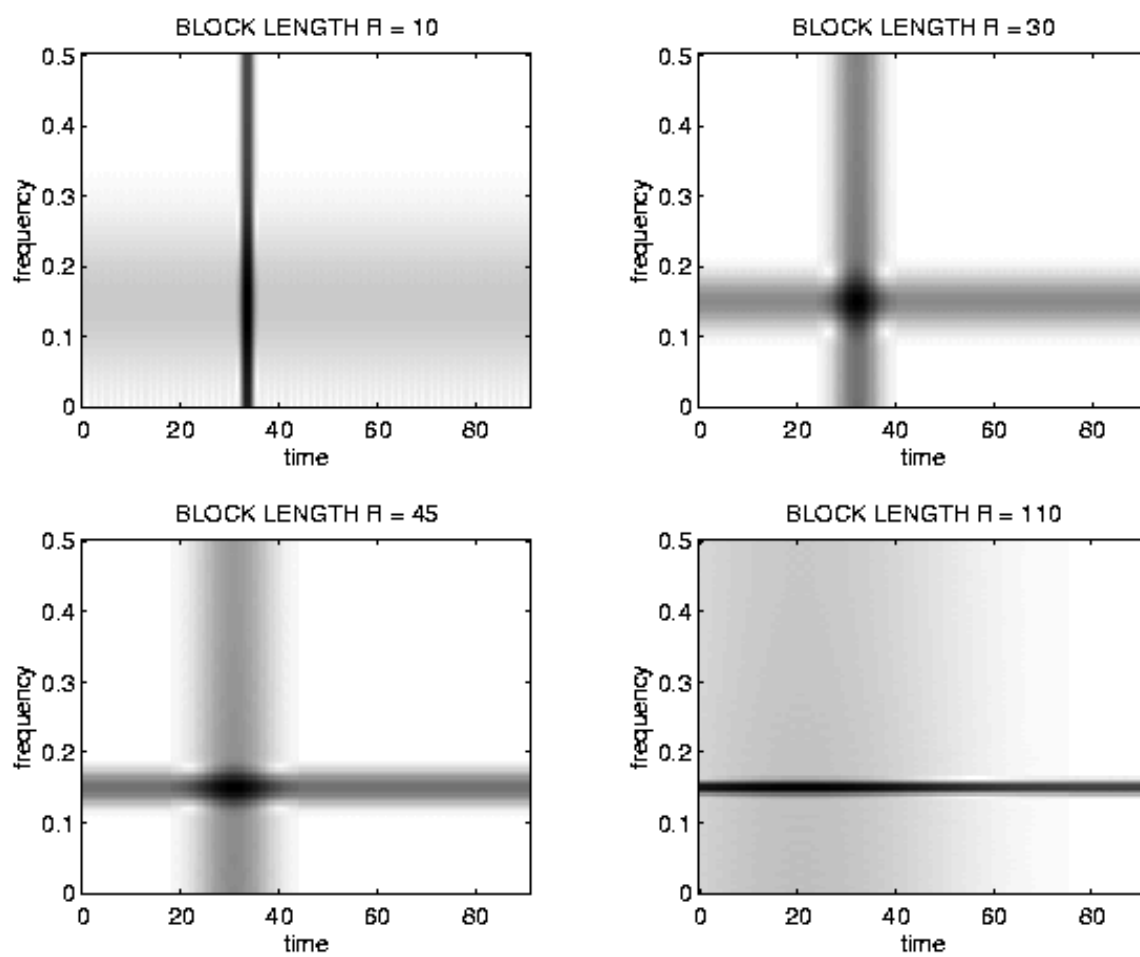
Figure 2.22

Here is another example to illustrate the frequency/time resolution trade-off (See figures - Figure 2.21 (Narrow-band spectrogram: better frequency resolution), Figure 2.22 (Wide-band spectrogram: better time resolution), and Figure 2.23 (Effect of Window Length R)).

Effect of Window Length R



(a)



(b)

Figure 2.23

2.2.3.1.4 Effect of L and N

A spectrogram is computed with different parameters:

$$L \in \{1, 10\}$$

$$N \in \{32, 256\}$$

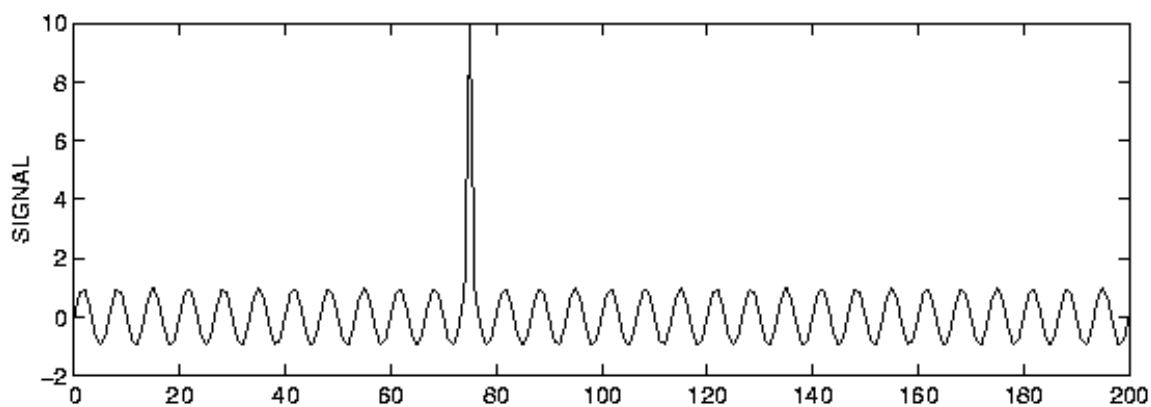
- L = time lapse between blocks.
- N = FFT length (Each block is zero-padded to length N .)

In each case, the block length is 30 samples.

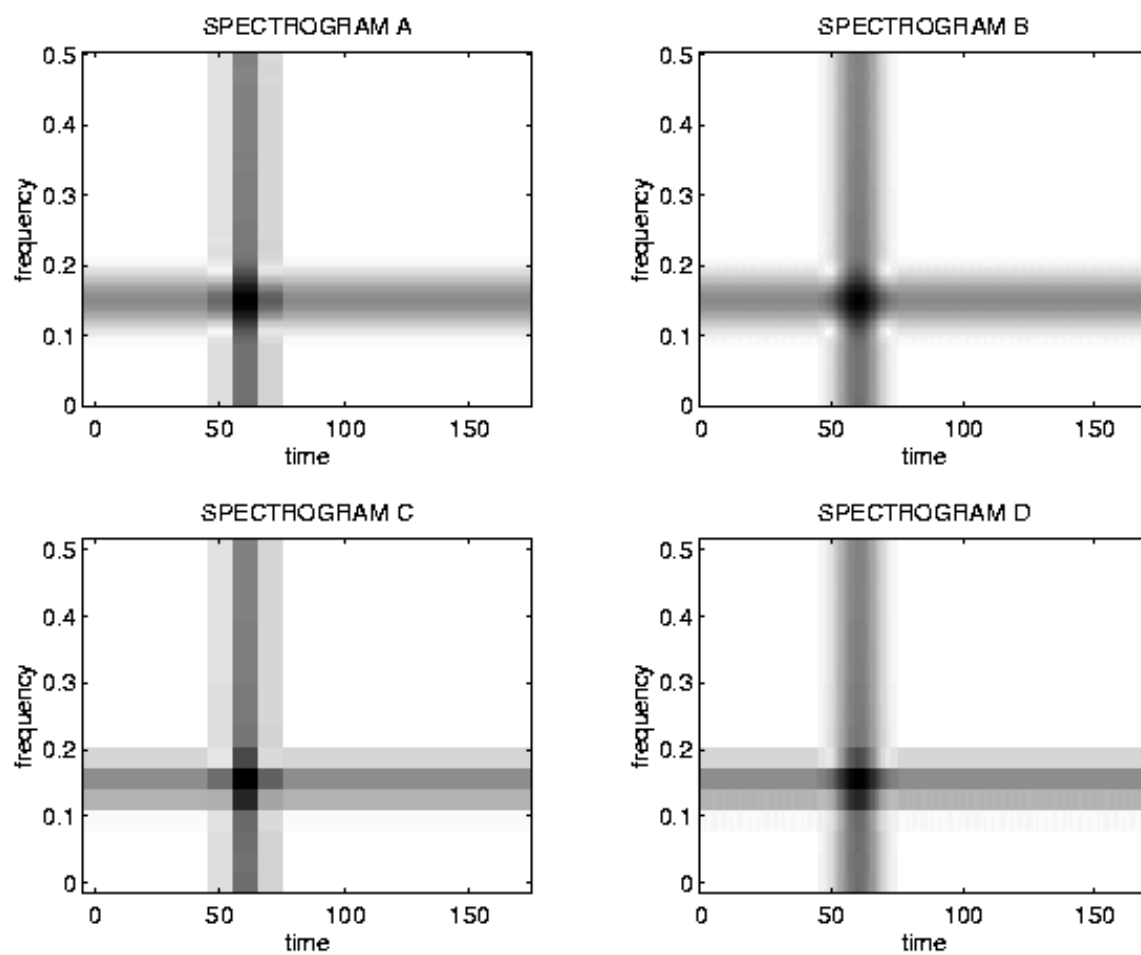
Exercise 2.3

(Solution on p. 189.)

For each of the four spectrograms in Figure 2.24 can you tell what L and N are?



(a)



(b)

Figure 2.24

L and N do not effect the time resolution or the frequency resolution. They only affect the 'pixelation'.

2.2.3.1.5 Effect of R and L

Shown below are four spectrograms of the same signal. Each spectrogram is computed using a different set of parameters.

$$R \in \{120, 256, 1024\}$$

$$L \in \{35, 250\}$$

where

- R = block length
- L = time lapse between blocks.

Exercise 2.4

(Solution on p. 189.)

For each of the four spectrograms in Figure 2.25, match the above values of L and R .

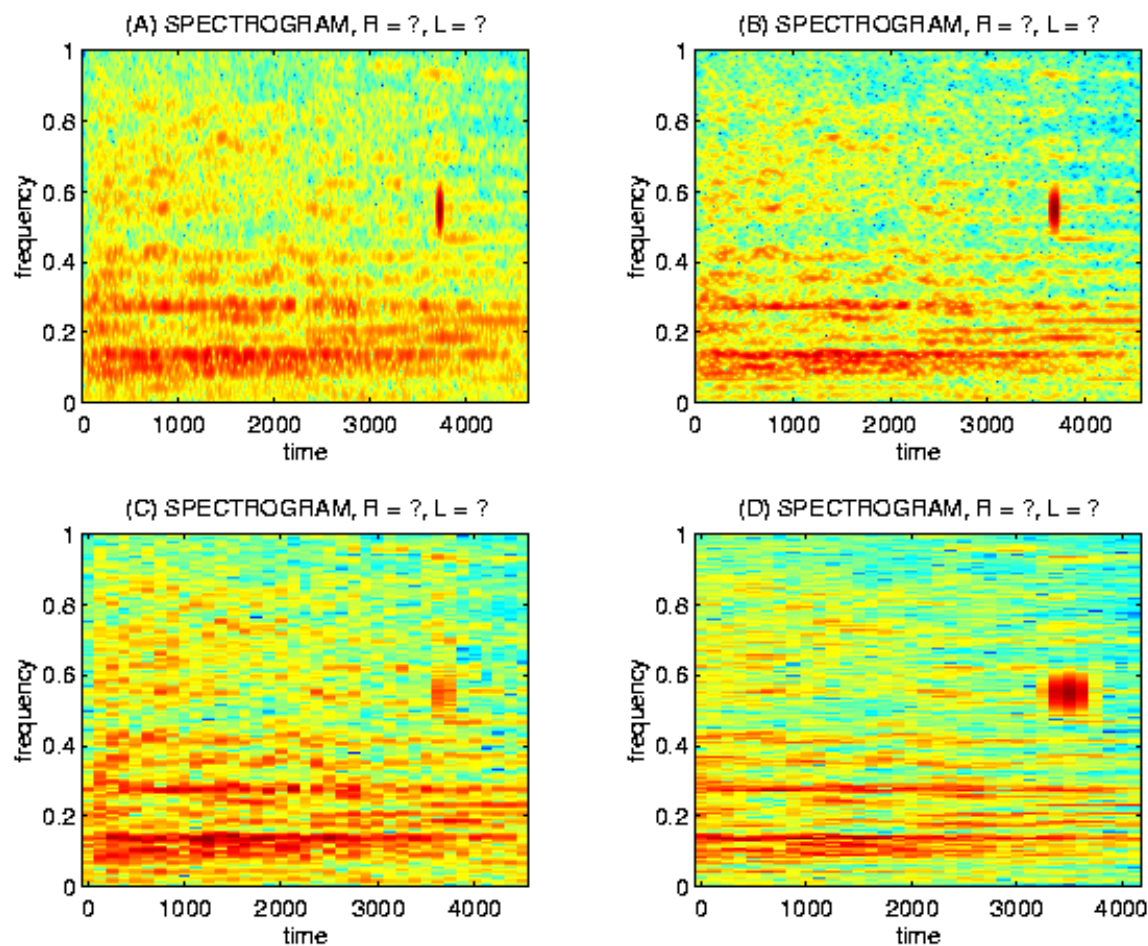


Figure 2.25

If you like, you may listen to this signal with the `soundsc` command; the data is in the file: `stft_data.m`. Here (Figure 2.26) is a figure of the signal.

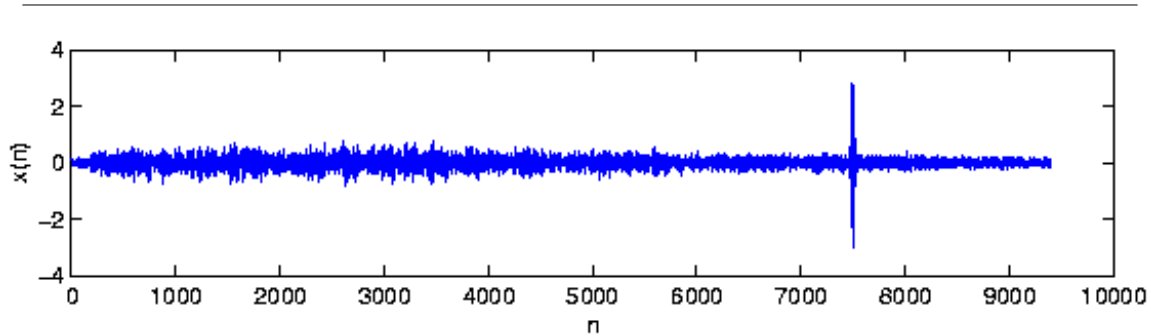


Figure 2.26

2.3 Fast Fourier Transform Algorithms

2.3.1 Overview of Fast Fourier Transform (FFT) Algorithms¹¹

A fast Fourier transform¹², or FFT¹³, is not a new transform, but is a computationally efficient algorithm for the computing the DFT (Section 2.1.1). The length- N DFT, defined as

$$X(k) = \sum_{n=0}^{N-1} \left(x(n) e^{-j \frac{2\pi nk}{N}} \right) \quad (2.12)$$

where $X(k)$ and $x(n)$ are in general complex-valued and $0 \leq k, n \leq N-1$, requires N complex multiplies to compute each $X(k)$. Direct computation of all N frequency samples thus requires N^2 complex multiplies and $N(N-1)$ complex additions. (This assumes precomputation of the DFT coefficients $(W_N^{nk} \doteq e^{-j \frac{2\pi nk}{N}})$; otherwise, the cost is even higher.) For the large DFT lengths used in many applications, N^2 operations may be prohibitive. (For example, digital terrestrial television broadcast in Europe uses $N = 2048$ or 8192 OFDM channels, and the SETI¹⁴ project uses up to length-4194304 DFTs.) DFTs are thus almost always computed in practice by an FFT algorithm¹⁵. FFTs are very widely used in signal processing, for applications such as spectrum analysis (Section 2.2.1) and digital filtering via fast convolution (Section 2.4).

2.3.1.1 History of the FFT

It is now known that C.F. Gauss¹⁶ invented an FFT in 1805 or so to assist the computation of planetary orbits via discrete Fourier series. Various FFT algorithms were independently invented over the next two centuries, but FFTs achieved widespread awareness and impact only with the Cooley and Tukey algorithm published in 1965, which came at a time of increasing use of digital computers and when the vast range of

¹¹This content is available online at <http://cnx.org/content/m12026/1.3/>.

¹²The DFT, FFT, and Practical Spectral Analysis <http://cnx.org/content/col10281/latest/>

¹³The DFT, FFT, and Practical Spectral Analysis <http://cnx.org/content/col10281/latest/>

¹⁴<http://en.wikipedia.org/wiki/SETI>

¹⁵The DFT, FFT, and Practical Spectral Analysis <http://cnx.org/content/col10281/latest/>

¹⁶http://en.wikipedia.org/wiki/Carl_Friedrich_Gauss

applications of numerical Fourier techniques was becoming apparent. Cooley and Tukey's algorithm spawned a surge of research in FFTs and was also partly responsible for the emergence of Digital Signal Processing (DSP) as a distinct, recognized discipline. Since then, many different algorithms have been rediscovered or developed, and efficient FFTs now exist for all DFT lengths.

2.3.1.2 Summary of FFT algorithms

The main strategy behind most FFT algorithms is to factor a length- N DFT into a number of shorter-length DFTs, the outputs of which are reused multiple times (usually in additional short-length DFTs!) to compute the final results. The lengths of the short DFTs correspond to integer factors of the DFT length, N , leading to different algorithms for different lengths and factors. By far the most commonly used FFTs select $N = 2^M$ to be a power of two, leading to the very efficient power-of-two FFT algorithms (Section 2.3.4.1), including the decimation-in-time radix-2 FFT (Section 2.3.4.2.1) and the decimation-in-frequency radix-2 FFT (Section 2.3.4.2.2) algorithms, the radix-4 FFT (Section 2.3.4.3) ($N = 4^M$), and the split-radix FFT (Section 2.3.4.4). Power-of-two algorithms gain their high efficiency from extensive reuse of intermediate results and from the low complexity of length-2 and length-4 DFTs, which require no multiplications. Algorithms for lengths with repeated common factors (Section 2.3.6) (such as 2 or 4 in the radix-2 and radix-4 algorithms, respectively) require extra **twiddle factor** multiplications between the short-length DFTs, which together lead to a computational complexity of $O(N \log N)$, a very considerable savings over direct computation of the DFT.

The other major class of algorithms is the Prime-Factor Algorithms (PFA) (Section 2.3.7). In PFAs, the short-length DFTs must be of relatively prime lengths. These algorithms gain efficiency by reuse of intermediate computations and by eliminating twiddle-factor multiplies, but require more operations than the power-of-two algorithms to compute the short DFTs of various prime lengths. In the end, the computational costs of the prime-factor and the power-of-two algorithms are comparable for similar lengths, as illustrated in Choosing the Best FFT Algorithm (Section 2.7). Prime-length DFTs cannot be factored into shorter DFTs, but in different ways both Rader's conversion (Section 2.6) and the chirp z-transform (Section 2.5) convert prime-length DFTs into convolutions of other lengths that can be computed efficiently using FFTs via fast convolution (Section 2.4).

Some applications require only a few DFT frequency samples, in which case Goertzel's algorithm (Section 2.3.3) halves the number of computations relative to the DFT sum. Other applications involve successive DFTs of overlapped blocks of samples, for which the running FFT (Section 2.3.2) can be more efficient than separate FFTs of each block.

2.3.2 Running FFT¹⁷

Some applications need DFT (2.5) frequencies of the most recent N samples on an ongoing basis. One example is DTMF¹⁸, or touch-tone telephone dialing, in which a detection circuit must constantly monitor the line for two simultaneous frequencies indicating that a telephone button is depressed. In such cases, most of the data in each successive block of samples is the same, and it is possible to efficiently update the DFT value from the previous sample to compute that of the current sample. Figure 2.27 illustrates successive length-4 blocks of data for which successive DFT values may be needed. The **running FFT** algorithm described here can be used to compute successive DFT values at a cost of only two complex multiplies and additions per DFT frequency.

¹⁷This content is available online at <http://cnx.org/content/m12029/1.5/>.

¹⁸<http://en.wikipedia.org/wiki/DTMF>

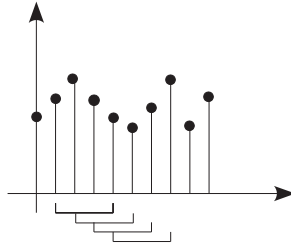


Figure 2.27: The running FFT efficiently computes DFT values for successive overlapped blocks of samples.

The running FFT algorithm is derived by expressing each DFT sample, $X_{n+1}(\omega_k)$, for the next block at time $n + 1$ in terms of the previous value, $X_n(\omega_k)$, at time n .

$$X_n(\omega_k) = \sum_{p=0}^{N-1} \left(x(n-p) e^{-j\omega_k p} \right)$$

$$X_{n+1}(\omega_k) = \sum_{p=0}^{N-1} \left(x(n+1-p) e^{-j\omega_k p} \right)$$

Let $q = p - 1$:

$$X_{n+1}(\omega_k) = \sum_{q=-1}^{N-2} \left(x(n-q) e^{-j\omega_k(q-1)} \right) = e^{j\omega_k} \sum_{q=0}^{N-2} \left(x(n-q) e^{-j\omega_k q} \right) + x(n+1)$$

Now let's add and subtract $e^{-j\omega_k(N-2)}x(n-N+1)$:

$$\begin{aligned} X_{n+1}(\omega_k) &= e^{j\omega_k} \sum_{q=0}^{N-2} \left(x(n-q) e^{-j\omega_k q} \right) + e^{j\omega_k} x(n-(N-1)) e^{-j\omega_k(N-1)} - \\ &e^{-j\omega_k(N-2)}x(n-N+1) + x(n+1) = e^{j\omega_k} \sum_{q=0}^{N-1} \left(x(n-q) e^{-j\omega_k q} \right) + x(n+1) - \\ &e^{-j\omega_k}x(n-N+1) = e^{j\omega_k} X_n(\omega_k) + x(n+1) - e^{-j\omega_k(N-2)}x(n-N+1) \end{aligned} \quad (2.13)$$

This running FFT algorithm requires only two complex multiplies and adds per update, rather than N if each DFT value were recomputed according to the DFT equation. Another advantage of this algorithm is that it works for *any* ω_k , rather than just the standard DFT frequencies. This can make it advantageous for applications, such as DTMF detection, where only a few arbitrary frequencies are needed.

Successive computation of a specific DFT frequency for overlapped blocks can also be thought of as a length- N FIR filter¹⁹. The running FFT is an efficient recursive implementation of this filter for this special case. Figure 2.28 shows a block diagram of the running FFT algorithm. The running FFT is one way to compute DFT filterbanks (Section 6.6). If a window other than rectangular is desired, a running FFT requires either a fast recursive implementation of the corresponding windowed, modulated impulse response, or it must have few non-zero coefficients so that it can be applied after the running FFT update via frequency-domain convolution. DFT-symmetric raised-cosine windows (Section 2.2.1.6: Effects of Windowing) are an example.

¹⁹Digital Filter Design <<http://cnx.org/content/col10285/latest/>>

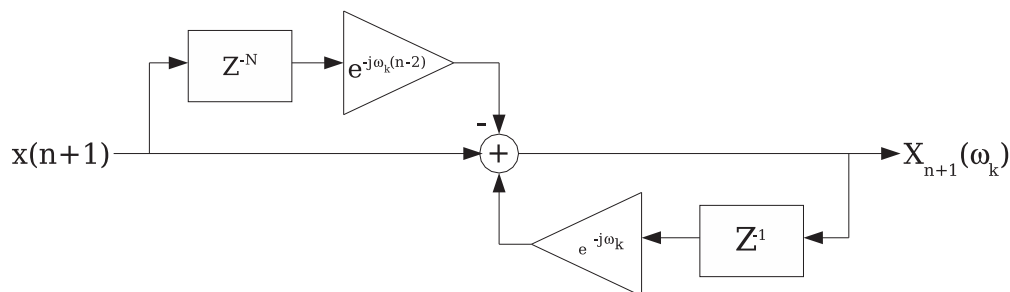


Figure 2.28: Block diagram of the running FFT computation, implemented as a recursive filter

2.3.3 Goertzel's Algorithm²⁰

Some applications require only a few DFT frequencies. One example is frequency-shift keying (FSK)²¹ demodulation, in which typically two frequencies are used to transmit binary data; another example is DTMF²², or touch-tone telephone dialing, in which a detection circuit must constantly monitor the line for two simultaneous frequencies indicating that a telephone button is depressed. *Goertzel's algorithm*[15] reduces the number of real-valued multiplications by almost a factor of two relative to direct computation via the DFT equation (2.5). Goertzel's algorithm is thus useful for computing a *few* frequency values; if many or most DFT values are needed, FFT algorithms (Section 2.3.1) that compute all DFT samples in $O(N \log N)$ operations are faster. Goertzel's algorithm can be derived by converting the DFT equation (Section 2.1.1) into an equivalent form as a convolution, which can be efficiently implemented as a digital filter. For increased clarity, in the equations below the complex exponential is denoted as $e^{-j\frac{2\pi k}{N}} = W_N^k$. Note that because W_N^{-Nk} always equals 1, the DFT equation (Section 2.1.1) can be rewritten as a convolution, or filtering operation:

$$\begin{aligned}
 X(k) &= \sum_{n=0}^{N-1} (x(n) 1 W_N^{nk}) \\
 &= \sum_{n=0}^{N-1} (x(n) W_N^{-Nk} W_N^{nk}) \\
 &= \sum_{n=0}^{N-1} (x(n) W_N^{(N-n)(-k)}) \\
 &= (((W_N^{-k} x(0) + x(1)) W_N^{-k} + x(2)) W_N^{-k} + \cdots + x(N-1)) W_N^{-k}
 \end{aligned} \tag{2.14}$$

Note that this last expression can be written in terms of a recursive difference equation (Section 1.11.1)

$$y(n) = W_N^{-k} y(n-1) + x(n)$$

where $y(-1) = 0$. The DFT coefficient equals the output of the difference equation at time $n = N$:

$$X(k) = y(N)$$

²⁰This content is available online at <<http://cnx.org/content/m12024/1.5/>>.

²¹http://en.wikipedia.org/wiki/Frequency-shift_keying

²²<http://en.wikipedia.org/wiki/DTMF>

Expressing the difference equation as a z-transform (Section 1.11.1) and multiplying both numerator and denominator by $1 - W_N^k z^{-1}$ gives the transfer function

$$\frac{Y(z)}{X(z)} = H(z) = \frac{1}{1 - W_N^{-k} z^{-1}} = \frac{1 - W_N^k z^{-1}}{1 - ((W_N^k + W_N^{-k}) z^{-1} - z^{-2})} = \frac{1 - W_N^k z^{-1}}{1 - (2\cos(\frac{2\pi k}{N}) z^{-1} - z^{-2})}$$

This system can be realized by the structure in Figure 2.29

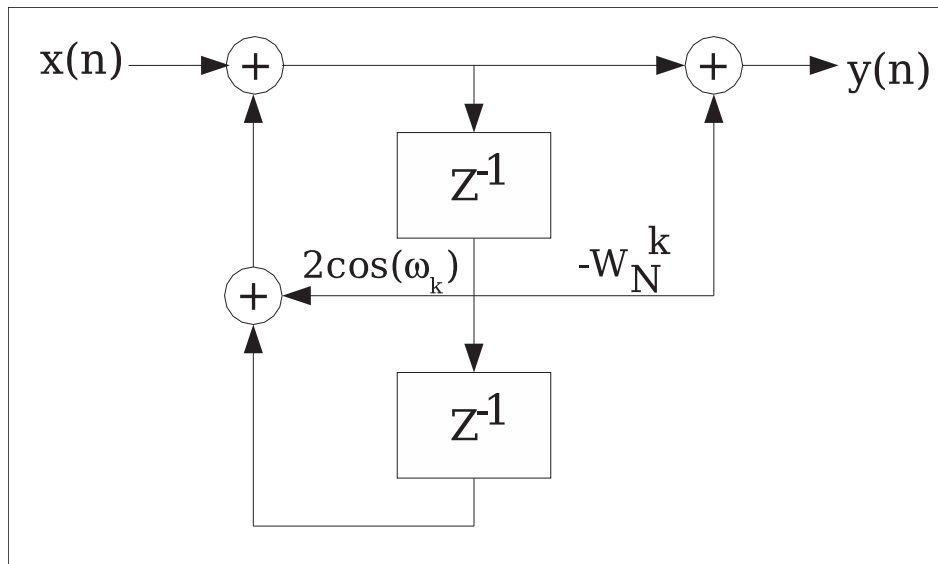


Figure 2.29

We want $y(n)$ not for all n , but only for $n = N$. We can thus compute only the *recursive* part, or just the left side of the flow graph in Figure 2.29, for $n = [0, 1, \dots, N]$, which involves only a *real/complex* product rather than a complex/complex product as in a direct DFT (2.5), plus one complex multiply to get $y(N) = X(k)$.

NOTE: The input $x(N)$ at time $n = N$ must equal 0! A slightly more efficient alternate implementation²³ that computes the full recursion only through $n = N - 1$ and combines the nonzero operations of the final recursion with the final complex multiply can be found here²⁴, complete with pseudocode (for real-valued data).

If the data are real-valued, only real/real multiplications and real additions are needed until the final multiply.

COST: The computational cost of Goertzel's algorithm is thus $2N + 2$ real multiplies and $4N - 2$ real adds, a reduction of almost a factor of two in the number of real multiplies relative to direct computation via the DFT equation. If the data are real-valued, this cost is almost halved again.

For certain frequencies, additional simplifications requiring even fewer multiplications are possible. (For example, for the DC ($k = 0$) frequency, all the multipliers equal 1 and only additions are needed.) A correspondence by *C.G. Boncelet, Jr.*[7] describes some of these additional simplifications. Once again, Goertzel's and Boncelet's algorithms are efficient for a few DFT frequency samples; if more than $\log N$ frequencies are needed, $O(N \log N)$ FFT algorithms (Section 2.3.1) that compute all frequencies simultaneously will be more efficient.

²³<http://www.mstarlabs.com/dsp/goertzel/goertzel.html>

²⁴<http://www.mstarlabs.com/dsp/goertzel/goertzel.html>

2.3.4 Power-of-Two FFTs

2.3.4.1 Power-of-two FFTs²⁵

FFTs of length $N = 2^M$ equal to a power of two are, by far, the most commonly used. These algorithms are very efficient, relatively simple, and a single program can compute power-of-two FFTs of different lengths. As with most FFT algorithms, they gain their efficiency by computing *all* DFT (Section 2.1.1) points simultaneously through extensive reuse of intermediate computations; they are thus efficient when many DFT frequency samples are needed. The simplest power-of-two FFTs are the decimation-in-time radix-2 FFT (Section 2.3.4.2.1) and the decimation-in-frequency radix-2 FFT (Section 2.3.4.2.2); they reduce the length- $N = 2^M$ DFT to a series of length-2 DFT computations with **twiddle-factor** complex multiplications between them. The radix-4 FFT algorithm (Section 2.3.4.3) similarly reduces a length- $N = 4^M$ DFT to a series of length-4 DFT computations with twiddle-factor multiplies in between. Radix-4 FFTs require only 75% as many complex multiplications as the radix-2 algorithms, although the number of complex additions remains the same. Radix-8 and higher-radix FFT algorithms can be derived using multi-dimensional index maps (Section 2.3.6) to reduce the computational complexity a bit more. However, the split-radix algorithm (Section 2.3.4.4) and its recent extensions combine the best elements of the radix-2 and radix-4 algorithms to obtain lower complexity than either or than any higher radix, requiring only two-thirds as many complex multiplies as the radix-2 algorithms. All of these algorithms obtain huge savings over direct computation of the DFT, reducing the complexity from $O(N^2)$ to $O(N \log N)$.

The efficiency of an FFT implementation depends on more than just the number of computations. Efficient FFT programming tricks (Section 2.3.5) can make up to a several-fold difference in the run-time of FFT programs. Alternate FFT structures (Section 2.3.4.2.3) can lead to a more convenient data flow for certain hardware. As discussed in choosing the best FFT algorithm (Section 2.7), certain hardware is designed for, and thus most efficient for, FFTs of specific lengths or radices.

2.3.4.2 Radix-2 Algorithms

2.3.4.2.1 Decimation-in-time (DIT) Radix-2 FFT²⁶

The radix-2 decimation-in-time and decimation-in-frequency (Section 2.3.4.2.2) fast Fourier transforms (FFTs) are the simplest FFT algorithms (Section 2.3.1). Like all FFTs, they gain their speed by reusing the results of smaller, intermediate computations to compute multiple DFT frequency outputs.

2.3.4.2.1.1 Decimation in time

The radix-2 decimation-in-time algorithm rearranges the discrete Fourier transform (DFT) equation (Section 2.1.1) into two parts: a sum over the even-numbered discrete-time indices $n = [0, 2, 4, \dots, N-2]$ and a sum over the odd-numbered indices $n = [1, 3, 5, \dots, N-1]$ as in (2.15):

$$\begin{aligned}
 X(k) &= \sum_{n=0}^{N-1} \left(x(n) e^{-j \frac{2\pi nk}{N}} \right) \\
 &= \sum_{n=0}^{\frac{N}{2}-1} \left(x(2n) e^{-j \frac{2\pi (2n)k}{N}} \right) + \sum_{n=0}^{\frac{N}{2}-1} \left(x(2n+1) e^{-j \frac{2\pi (2n+1)k}{N}} \right) \\
 &= \sum_{n=0}^{\frac{N}{2}-1} \left(x(2n) e^{-j \frac{2\pi nk}{\frac{N}{2}}} \right) + e^{-j \frac{2\pi k}{N}} \sum_{n=0}^{\frac{N}{2}-1} \left(x(2n+1) e^{-j \frac{2\pi nk}{\frac{N}{2}}} \right) \\
 &= DFT_{\frac{N}{2}} [[x(0), x(2), \dots, x(N-2)]] + W_N^k DFT_{\frac{N}{2}} [[x(1), x(3), \dots, x(N-1)]]
 \end{aligned} \tag{2.15}$$

The mathematical simplifications in (2.15) reveal that all DFT frequency outputs $X(k)$ can be computed as the sum of the outputs of two length- $\frac{N}{2}$ DFTs, of the even-indexed and odd-indexed discrete-time samples, respectively, where the odd-indexed short DFT is multiplied by a so-called **twiddle factor** term $W_N^k = e^{-j \frac{2\pi k}{N}}$. This is called a **decimation in time** because the time samples are rearranged in alternating

²⁵This content is available online at <<http://cnx.org/content/m12059/1.2/>>.

²⁶This content is available online at <<http://cnx.org/content/m12016/1.7/>>.

groups, and a **radix-2** algorithm because there are two groups. Figure 2.30 graphically illustrates this form of the DFT computation, where for convenience the frequency outputs of the length- $\frac{N}{2}$ DFT of the even-indexed time samples are denoted $G(k)$ and those of the odd-indexed samples as $H(k)$. Because of the periodicity with $\frac{N}{2}$ frequency samples of these length- $\frac{N}{2}$ DFTs, $G(k)$ and $H(k)$ can be used to compute *two* of the length- N DFT frequencies, namely $X(k)$ and $X(k + \frac{N}{2})$, but with a different twiddle factor. This reuse of these short-length DFT outputs gives the FFT its computational savings.

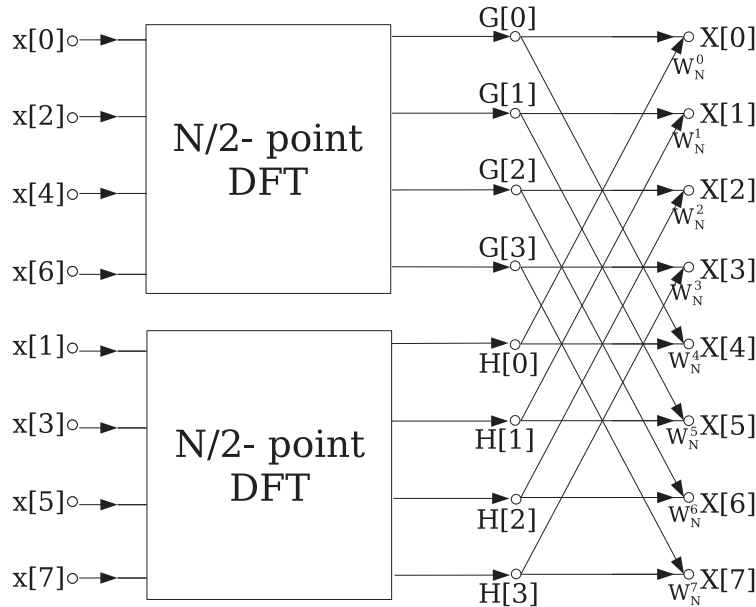


Figure 2.30: Decimation in time of a length- N DFT into two length- $\frac{N}{2}$ DFTs followed by a combining stage.

Whereas direct computation of all N DFT frequencies according to the DFT equation (Section 2.1.1) would require N^2 complex multiplies and $N^2 - N$ complex additions (for complex-valued data), by reusing the results of the two short-length DFTs as illustrated in Figure 2.30, the computational cost is now

New Operation Counts

- $2\left(\frac{N}{2}\right)^2 + N = \frac{N^2}{2} + N$ complex multiplies
- $2\frac{N}{2}\left(\frac{N}{2} - 1\right) + N = \frac{N^2}{2}$ complex additions

This simple reorganization and reuse has reduced the total computation by almost a factor of two over direct DFT (Section 2.1.1) computation!

2.3.4.2.1.2 Additional Simplification

A basic **butterfly** operation is shown in Figure 2.31, which requires only $\frac{N}{2}$ **twiddle-factor** multiplies per **stage**. It is worthwhile to note that, after merging the twiddle factors to a single term on the lower branch, the remaining butterfly is actually a length-2 DFT! The theory of multi-dimensional index maps

(Section 2.3.6) shows that this must be the case, and that FFTs of any factorable length may consist of successive stages of shorter-length FFTs with twiddle-factor multiplications in between.

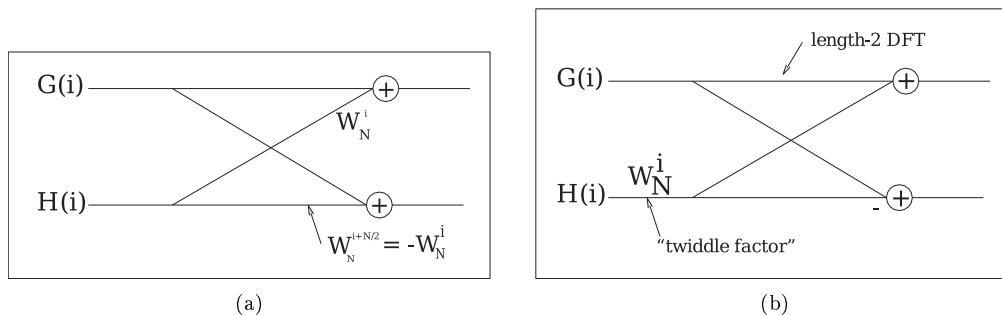


Figure 2.31: Radix-2 DIT butterfly simplification: both operations produce the same outputs

2.3.4.2.1.3 Radix-2 decimation-in-time FFT

The same radix-2 decimation in time can be applied recursively to the two length $\frac{N}{2}$ DFT (Section 2.1.1)s to save computation. When successively applied until the shorter and shorter DFTs reach length-2, the result is the radix-2 DIT FFT algorithm (Figure 2.32).

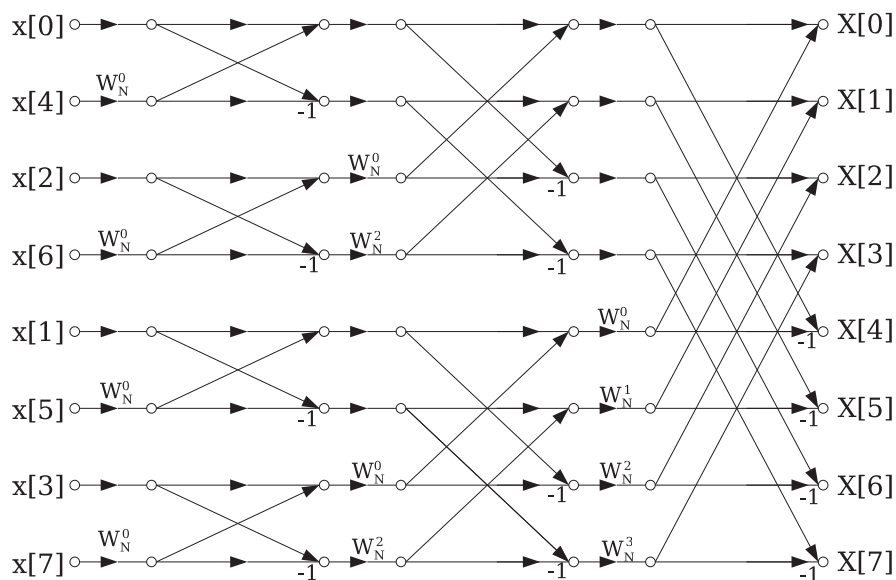


Figure 2.32: Radix-2 Decimation-in-Time FFT algorithm for a length-8 signal

The full radix-2 decimation-in-time decomposition illustrated in Figure 2.32 using the simplified butterflies (Figure 2.31) involves $M = \log_2 N$ stages, each with $\frac{N}{2}$ butterflies per stage. Each butterfly requires 1 complex multiply and 2 adds per butterfly. The total cost of the algorithm is thus

Computational cost of radix-2 DIT FFT

- $\frac{N}{2} \log_2 N$ complex multiplies
- $N \log_2 N$ complex adds

This is a remarkable savings over direct computation of the DFT. For example, a length-1024 DFT would require 1048576 complex multiplications and 1047552 complex additions with direct computation, but only 5120 complex multiplications and 10240 complex additions using the radix-2 FFT, a savings by a factor of 100 or more. The relative savings increase with longer FFT lengths, and are less for shorter lengths.

Modest additional reductions in computation can be achieved by noting that certain twiddle factors, namely Using special butterflies for W_N^0 , $W_N^{\frac{N}{2}}$, $W_N^{\frac{N}{4}}$, $W_N^{\frac{N}{8}}$, $W_N^{\frac{3N}{8}}$, require no multiplications, or fewer real multiplies than other ones. By implementing special butterflies for these twiddle factors as discussed in FFT algorithm and programming tricks, the computational cost of the radix-2 decimation-in-time FFT can be reduced to

- $2N \log_2 N - 7N + 12$ real multiplies
- $3N \log_2 N - 3N + 4$ real additions

NOTE: In a decimation-in-time radix-2 FFT as illustrated in Figure 2.32, the input is in **bit-reversed** order (hence "decimation-in-time"). That is, if the time-sample index n is written as a binary number, the order is that binary number reversed. The bit-reversal process is illustrated for a length- $N = 8$ example below.

Example 2.3: N=8

In-order index	In-order index in binary	Bit-reversed binary	Bit-reversed index
0	000	000	0
1	001	100	4
2	010	010	2
3	011	110	6
4	100	001	1
5	101	101	5
6	110	011	3
7	111	111	7

It is important to note that, if the input signal data are placed in bit-reversed order before beginning the FFT computations, the outputs of each butterfly throughout the computation can be placed in the same memory locations from which the inputs were fetched, resulting in an **in-place algorithm** that requires no extra memory to perform the FFT. Most FFT implementations are in-place, and overwrite the input data with the intermediate values and finally the output.

2.3.4.2.1.4 Example FFT Code

The following function, written in the C programming language, implements a radix-2 decimation-in-time FFT. It is designed for computing the DFT of complex-valued inputs to produce complex-valued outputs,

with the real and imaginary parts of each number stored in separate double-precision floating-point arrays. It is an in-place algorithm, so the intermediate and final output values are stored in the same array as the input data, which is overwritten. After initializations, the program first bit-reverses the discrete-time samples, as is typical with a decimation-in-time algorithm (but see alternate FFT structures (Section 2.3.4.2.3) for DIT algorithms with other input orders), then computes the FFT in stages according to the above description.

This FFT program (p. 150) uses a standard three-loop structure for the main FFT computation. The outer loop steps through the stages (each column in Figure 2.32); the middle loop steps through "flights" (butterflies with the same twiddle factor from each short-length DFT at each stage), and the inner loop steps through the individual butterflies. This ordering minimizes the number of fetches or computations of the twiddle-factor values. Since the bit-reverse of a bit-reversed index is the original index, bit-reversal can be performed fairly simply by swapping pairs of data.

NOTE: While of $O(N \log N)$ complexity and thus much faster than a direct DFT, this simple program is optimized for clarity, not for speed. A speed-optimized program making use of additional efficient FFT algorithm and programming tricks (Section 2.3.5) will compute a DFT several times faster on most machines.

```

/*****/
/* fft.c */
/* (c) Douglas L. Jones */
/* University of Illinois at Urbana-Champaign */
/* January 19, 1992 */
/* */
/* fft: in-place radix-2 DIT DFT of a complex input */
/* */
/* input: */
/* n: length of FFT: must be a power of two */
/* m: n = 2**m */
/* input/output */
/* x: double array of length n with real part of data */
/* y: double array of length n with imag part of data */
/* */
/* Permission to copy and use this program is granted */
/* under a Creative Commons "Attribution" license */
/* http://creativecommons.org/licenses/by/1.0/ */
/*****/
fft(n,m,x,y)
int n,m;
double x[],y[];
{
    int i,j,k,n1,n2;
    double c,s,e,a,t1,t2;

    j = 0; /* bit-reverse */
    n2 = n/2;
    for (i=1; i < n - 1; i++)
    {
        n1 = n2;
        while ( j >= n1 )

```

```

    {
        j = j - n1;
        n1 = n1/2;
    }
    j = j + n1;

    if (i < j)
    {
        t1 = x[i];
        x[i] = x[j];
        x[j] = t1;
        t1 = y[i];
        y[i] = y[j];
        y[j] = t1;
    }
}

n1 = 0; /* FFT */
n2 = 1;

for (i=0; i < m; i++)
{
    n1 = n2;
    n2 = n2 + n2;
    e = -6.283185307179586/n2;
    a = 0.0;

    for (j=0; j < n1; j++)
    {
        c = cos(a);
        s = sin(a);
        a = a + e;

        for (k=j; k < n; k=k+n2)
        {
            t1 = c*x[k+n1] - s*y[k+n1];
            t2 = s*x[k+n1] + c*y[k+n1];
            x[k+n1] = x[k] - t1;
            y[k+n1] = y[k] - t2;
            x[k] = x[k] + t1;
            y[k] = y[k] + t2;
        }
    }
}

return;
}

```

2.3.4.2.2 Decimation-in-Frequency (DIF) Radix-2 FFT²⁷

The radix-2 decimation-in-frequency and decimation-in-time (Section 2.3.4.2.1) fast Fourier transforms (FFTs) are the simplest FFT algorithms (Section 2.3.1). Like all FFTs, they compute the discrete Fourier transform (DFT) (Section 2.1.1)

$$\begin{aligned} X(k) &= \sum_{n=0}^{N-1} \left(x(n) e^{-j \frac{2\pi nk}{N}} \right) \\ &= \sum_{n=0}^{N-1} x(n) W_N^{nk} \end{aligned} \quad (2.16)$$

where for notational convenience $W_N^k = e^{-j \frac{2\pi k}{N}}$. FFT algorithms gain their speed by reusing the results of smaller, intermediate computations to compute multiple DFT frequency outputs.

2.3.4.2.2.1 Decimation in frequency

The radix-2 decimation-in-frequency algorithm rearranges the discrete Fourier transform (DFT) equation (2.16) into two parts: computation of the even-numbered discrete-frequency indices $X(k)$ for $k = [0, 2, 4, \dots, N-2]$ (or $X(2r)$ as in (2.17)) and computation of the odd-numbered indices $k = [1, 3, 5, \dots, N-1]$ (or $X(2r+1)$ as in (2.18))

$$\begin{aligned} X(2r) &= \sum_{n=0}^{N-1} x(n) W_N^{2rn} \\ &= \sum_{n=0}^{\frac{N}{2}-1} x(n) W_N^{2rn} + \sum_{n=0}^{\frac{N}{2}-1} \left(x\left(n + \frac{N}{2}\right) W_N^{2r\left(n + \frac{N}{2}\right)} \right) \\ &= \sum_{n=0}^{\frac{N}{2}-1} x(n) W_N^{2rn} + \sum_{n=0}^{\frac{N}{2}-1} x\left(n + \frac{N}{2}\right) W_N^{2rn} \\ &= \sum_{n=0}^{\frac{N}{2}-1} \left(x(n) + x\left(n + \frac{N}{2}\right) \right) W_N^{rn} \\ &= DFT_{\frac{N}{2}} \left[x(n) + x\left(n + \frac{N}{2}\right) \right] \end{aligned} \quad (2.17)$$

$$\begin{aligned} X(2r+1) &= \sum_{n=0}^{N-1} x(n) W_N^{(2r+1)n} \\ &= \sum_{n=0}^{\frac{N}{2}-1} \left(x(n) + W_N^{\frac{N}{2}} x\left(n + \frac{N}{2}\right) \right) W_N^{(2r+1)n} \\ &= \sum_{n=0}^{\frac{N}{2}-1} \left((x(n) - x\left(n + \frac{N}{2}\right)) W_N^n \right) W_N^{\frac{N}{2}} \\ &= DFT_{\frac{N}{2}} \left[(x(n) - x\left(n + \frac{N}{2}\right)) W_N^n \right] \end{aligned} \quad (2.18)$$

The mathematical simplifications in (2.17) and (2.18) reveal that both the even-indexed and odd-indexed frequency outputs $X(k)$ can each be computed by a length- $\frac{N}{2}$ DFT. The inputs to these DFTs are sums or differences of the first and second halves of the input signal, respectively, where the input to the short DFT producing the odd-indexed frequencies is multiplied by a so-called **twiddle factor** term $W_N^k = e^{-j \frac{2\pi k}{N}}$. This is called a **decimation in frequency** because the frequency samples are computed separately in alternating groups, and a **radix-2** algorithm because there are two groups. Figure 2.33 graphically illustrates this form of the DFT computation. This conversion of the full DFT into a series of shorter DFTs with a simple preprocessing step gives the decimation-in-frequency FFT its computational savings.

²⁷This content is available online at <<http://cnx.org/content/m12018/1.6/>>.

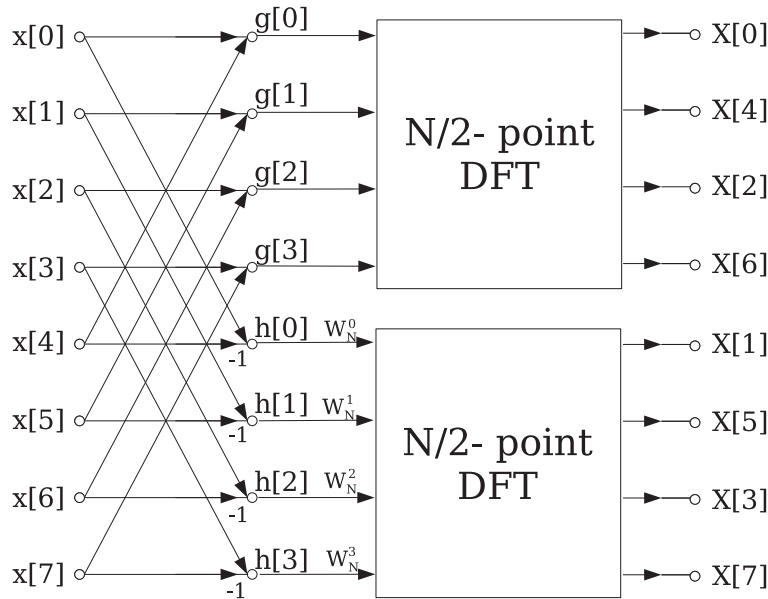


Figure 2.33: Decimation in frequency of a length- N DFT into two length- $\frac{N}{2}$ DFTs preceded by a preprocessing stage.

Whereas direct computation of all N DFT frequencies according to the DFT equation (Section 2.1.1) would require N^2 complex multiplies and $N^2 - N$ complex additions (for complex-valued data), by breaking the computation into two short-length DFTs with some preliminary combining of the data, as illustrated in Figure 2.33, the computational cost is now

New Operation Counts

- $2\left(\frac{N}{2}\right)^2 + N = \frac{N^2}{2} + \frac{N}{2}$ complex multiplies
- $2\frac{N}{2}\left(\frac{N}{2} - 1\right) + N = \frac{N^2}{2}$ complex additions

This simple manipulation has reduced the total computational cost of the DFT by almost a factor of two!

The initial combining operations for both short-length DFTs involve parallel groups of two time samples, $x(n)$ and $x\left(n + \frac{N}{2}\right)$. One of these so-called **butterfly** operations is illustrated in Figure 2.34. There are $\frac{N}{2}$ butterflies per **stage**, each requiring a complex addition and subtraction followed by one **twiddle-factor** multiplication by $W_N^n = e^{-j\frac{2\pi n}{N}}$ on the lower output branch.

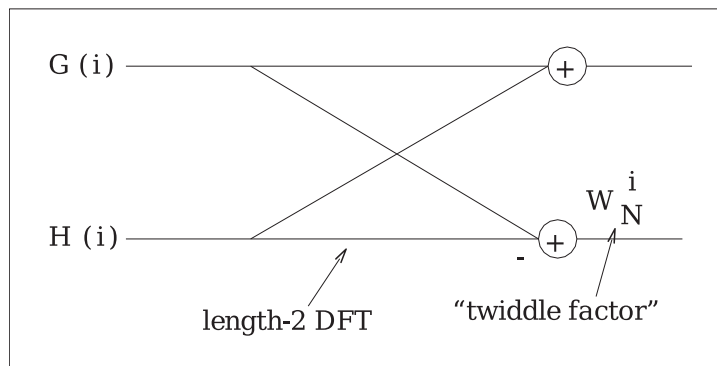


Figure 2.34: DIF butterfly: twiddle factor after length-2 DFT

It is worthwhile to note that the initial add/subtract part of the DIF butterfly is actually a length-2 DFT! The theory of multi-dimensional index maps (Section 2.3.6) shows that this must be the case, and that FFTs of any factorable length may consist of successive stages of shorter-length FFTs with twiddle-factor multiplications in between. It is also worth noting that this butterfly differs from the decimation-in-time radix-2 butterfly (Figure 2.31) in that the twiddle factor multiplication occurs *after* the combining.

2.3.4.2.2 Radix-2 decimation-in-frequency algorithm

The same radix-2 decimation in frequency can be applied recursively to the two length- $\frac{N}{2}$ DFT (Section 2.1.1)s to save additional computation. When successively applied until the shorter and shorter DFTs reach length-2, the result is the radix-2 decimation-in-frequency FFT algorithm (Figure 2.35).

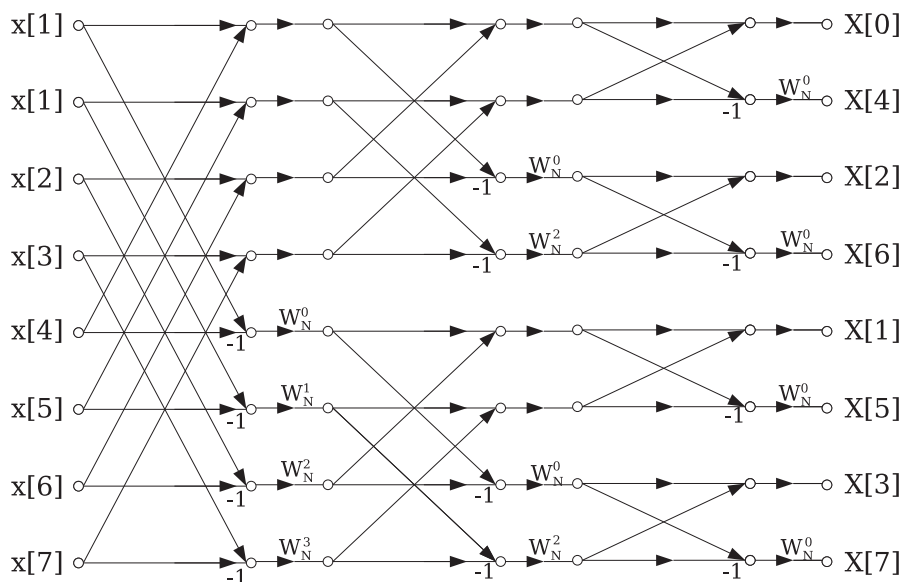


Figure 2.35: Radix-2 decimation-in-frequency FFT for a length-8 signal

The full radix-2 decimation-in-frequency decomposition illustrated in Figure 2.35 requires $M = \log_2 N$ stages, each with $\frac{N}{2}$ butterflies per stage. Each butterfly requires 1 complex multiply and 2 adds per butterfly. The total cost of the algorithm is thus

Computational cost of radix-2 DIF FFT

- $\frac{N}{2} \log_2 N$ complex multiplies
- $N \log_2 N$ complex adds

This is a remarkable savings over direct computation of the DFT. For example, a length-1024 DFT would require 1048576 complex multiplications and 1047552 complex additions with direct computation, but only 5120 complex multiplications and 10240 complex additions using the radix-2 FFT, a savings by a factor of 100 or more. The relative savings increase with longer FFT lengths, and are less for shorter lengths. Modest additional reductions in computation can be achieved by noting that certain twiddle factors, namely W_N^0 , $W_N^{\frac{N}{2}}$, $W_N^{\frac{N}{4}}$, $W_N^{\frac{N}{8}}$, $W_N^{\frac{3N}{8}}$, require no multiplications, or fewer real multiplies than other ones. By implementing special butterflies for these twiddle factors as discussed in FFT algorithm and programming tricks (Section 2.3.5), the computational cost of the radix-2 decimation-in-frequency FFT can be reduced to

- $2N \log_2 N - 7N + 12$ real multiplies
- $3N \log_2 N - 3N + 4$ real additions

The decimation-in-frequency FFT is a flow-graph reversal of the decimation-in-time (Section 2.3.4.2.1) FFT: it has the same twiddle factors (in reverse pattern) and the same operation counts.

NOTE: In a decimation-in-frequency radix-2 FFT as illustrated in Figure 2.35, the output is in **bit-reversed** order (hence "decimation-in-frequency"). That is, if the frequency-sample index n is written as a binary number, the order is that binary number reversed. The bit-reversal process is illustrated here (Example 2.3: $N=8$).

It is important to note that, if the input data are in order before beginning the FFT computations, the outputs of each butterfly throughout the computation can be placed in the same memory locations from which the inputs were fetched, resulting in an **in-place algorithm** that requires no extra memory to perform the FFT. Most FFT implementations are in-place, and overwrite the input data with the intermediate values and finally the output.

2.3.4.2.3 Alternate FFT Structures²⁸

Bit-reversing (Section 2.3.4.2.1) the input in decimation-in-time (DIT) FFTs (Section 2.3.4.2.1) or the output in decimation-in-frequency (DIF) FFTs (Section 2.3.4.2.2) can sometimes be inconvenient or inefficient. For such situations, alternate FFT structures have been developed. Such structures involve the same mathematical computations as the standard algorithms, but alter the memory locations in which intermediate values are stored or the order of computation of the FFT butterflies (Section 2.3.4.2.1).

The structure in Figure 2.36 computes a decimation-in-frequency FFT (Section 2.3.4.2.2), but remaps the memory usage so that the *input* is bit-reversed (Section 2.3.4.2.1), and the output is in-order as in the conventional decimation-in-time FFT (Section 2.3.4.2.1). This alternate structure is still considered a DIF FFT because the twiddle factors (Section 2.3.4.2.1) are applied as in the DIF FFT (Section 2.3.4.2.2). This structure is useful if for some reason the DIF butterfly is preferred but it is easier to bit-reverse the input.

²⁸This content is available online at <<http://cnx.org/content/m12012/1.6/>>.

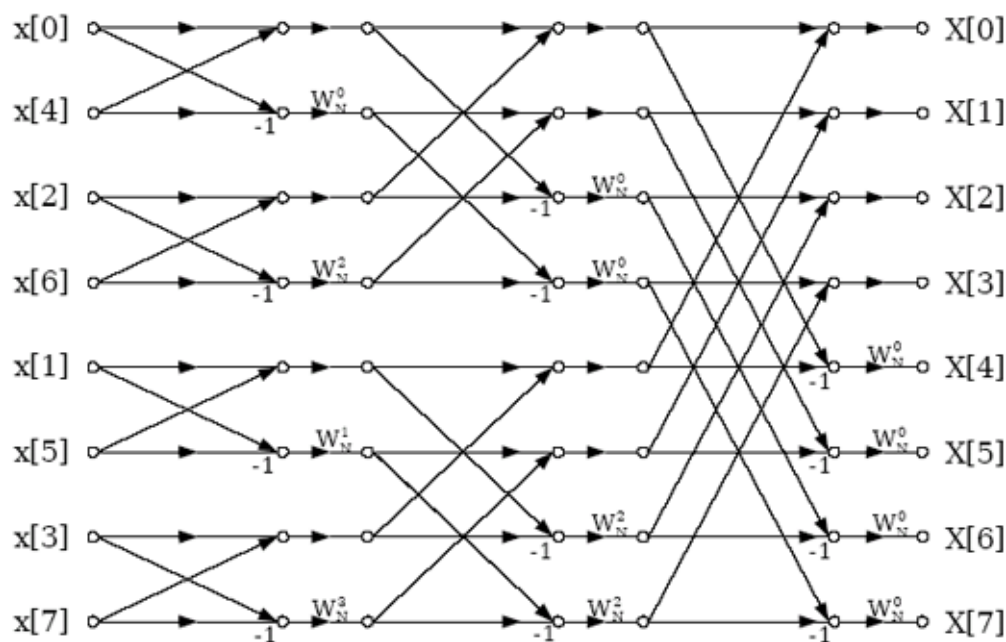


Figure 2.36: Decimation-in-frequency radix-2 FFT (Section 2.3.4.2.2) with bit-reversed *input*. This is an in-place (Section 2.3.4.2.1) algorithm in which the same memory can be reused throughout the computation.

There is a similar structure for the decimation-in-time FFT (Section 2.3.4.2.1) with in-order inputs and bit-reversed frequencies. This structure can be useful for fast convolution (Section 2.4) on machines that favor decimation-in-time algorithms because the filter can be stored in bit-reverse order, and then the inverse FFT returns an in-order result without ever bit-reversing any data. As discussed in Efficient FFT Programming Tricks (Section 2.3.5), this may save several percent of the execution time.

The structure in Figure 2.37 implements a decimation-in-frequency FFT (Section 2.3.4.2.2) that has both input and output in order. It thus avoids the need for bit-reversing altogether. Unfortunately, it destroys the in-place (Section 2.3.4.2.1) structure somewhat, making an FFT program more complicated and requiring more memory; on most machines the resulting cost exceeds the benefits. This structure can be computed in place if *two* butterflies are computed simultaneously.

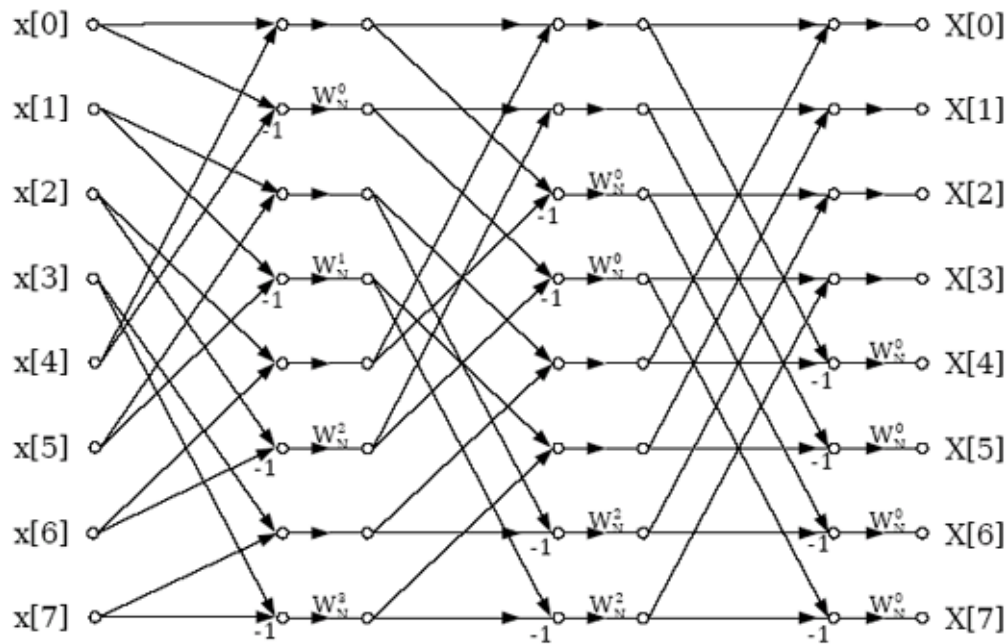


Figure 2.37: Decimation-in-frequency radix-2 FFT with in-order input and output. It can be computed in-place if two butterflies are computed simultaneously.

The structure in Figure 2.38 has a constant geometry; the connections between memory locations are identical in each FFT stage (Section 2.3.4.2.1). Since it is not in-place and requires bit-reversal, it is inconvenient for software implementation, but can be attractive for a highly parallel hardware implementation because the connections between stages can be hardwired. An analogous structure exists that has bit-reversed inputs and in-order outputs.

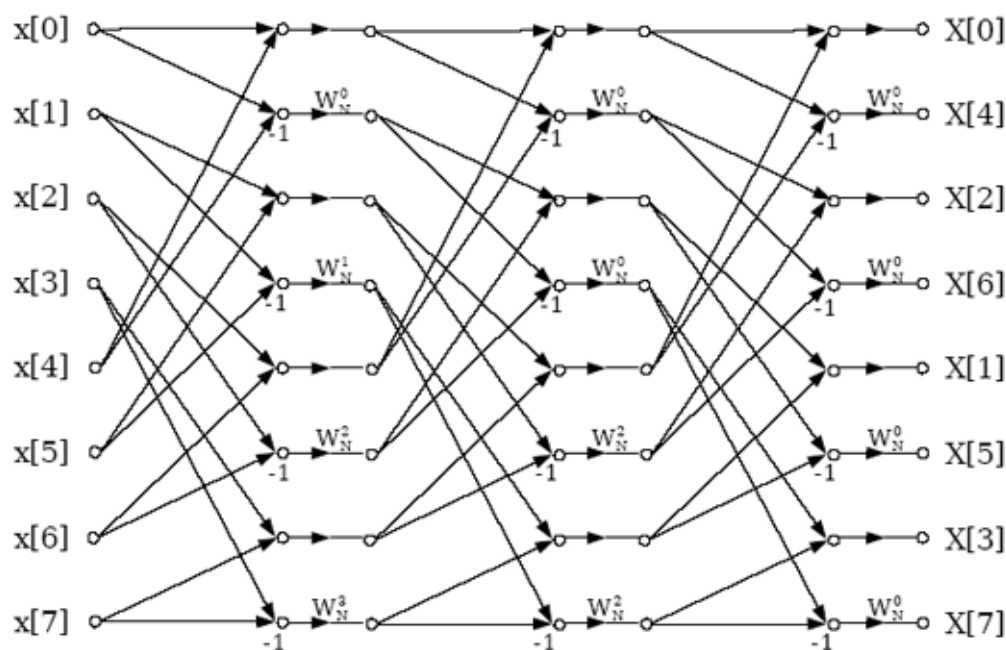


Figure 2.38: This constant-geometry structure has the same interconnect pattern from stage to stage. This structure is sometimes useful for special hardware.

2.3.4.3 Radix-4 FFT Algorithms²⁹

The radix-4 decimation-in-time (Section 2.3.4.2.1) and decimation-in-frequency (Section 2.3.4.2.2) fast Fourier transforms (FFTs) (Section 2.3.1) gain their speed by reusing the results of smaller, intermediate computations to compute multiple DFT frequency outputs. The radix-4 decimation-in-time algorithm rearranges the discrete Fourier transform (DFT) equation (Section 2.1.1) into four parts: sums over all groups of every fourth discrete-time index $n = [0, 4, 8, \dots, N - 4]$, $n = [1, 5, 9, \dots, N - 3]$, $n = [2, 6, 10, \dots, N - 2]$ and $n = [3, 7, 11, \dots, N - 1]$ as in (2.19). (This works out only when the FFT length is a multiple of four.) Just as in the radix-2 decimation-in-time FFT (Section 2.3.4.2.1), further mathematical manipulation shows that the length- N DFT can be computed as the sum of the outputs of four length- $\frac{N}{4}$ DFTs, of the even-indexed and odd-indexed discrete-time samples, respectively, where three of them are multiplied by so-called **twiddle factors** $W_N^k = e^{-j\frac{2\pi k}{N}}$, W_N^{2k} , and W_N^{3k} .

²⁹This content is available online at <<http://cnx.org/content/m12027/1.4/>>.

$$\begin{aligned}
X(k) &= \sum_{n=0}^{N-1} \left(x(n) e^{-j\frac{2\pi nk}{N}} \right) = \sum_{n=0}^{\frac{N}{4}-1} \left(x(4n) e^{-j\frac{2\pi(4n)k}{N}} \right) + \\
&\sum_{n=0}^{\frac{N}{4}-1} \left(x(4n+1) e^{-j\frac{2\pi(4n+1)k}{N}} \right) + \sum_{n=0}^{\frac{N}{4}-1} \left(x(4n+2) e^{-j\frac{2\pi(4n+2)k}{N}} \right) + \\
&\sum_{n=0}^{\frac{N}{4}-1} \left(x(4n+3) e^{-j\frac{2\pi(4n+3)k}{N}} \right) = DFT_{\frac{N}{4}}[x(4n)] + W_N^k DFT_{\frac{N}{4}}[x(4n+1)] + \\
&W_N^{2k} DFT_{\frac{N}{4}}[x(4n+2)] + W_N^{3k} DFT_{\frac{N}{4}}[x(4n+3)]
\end{aligned} \tag{2.19}$$

This is called a **decimation in time** because the time samples are rearranged in alternating groups, and a **radix-4** algorithm because there are four groups. Figure 2.39 (Radix-4 DIT structure) graphically illustrates this form of the DFT computation.

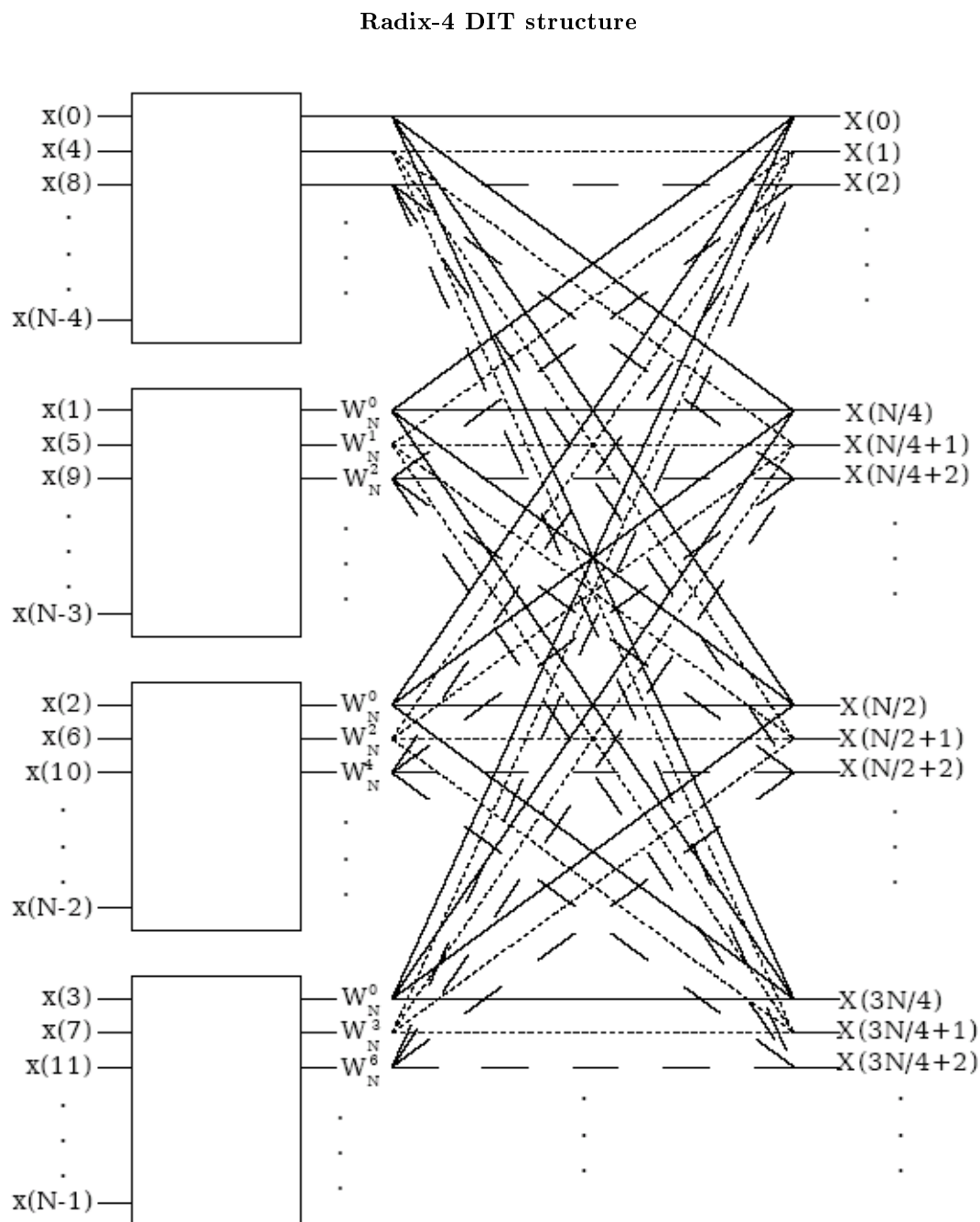


Figure 2.39: Decimation in time of a length- N DFT into four length- $\frac{N}{4}$ DFTs followed by a combining stage.

Due to the periodicity with $\frac{N}{4}$ of the short-length DFTs, their outputs for frequency-sample k are reused to compute $X(k)$, $X(k + \frac{N}{4})$, $X(k + \frac{N}{2})$, and $X(k + \frac{3N}{4})$. It is this reuse that gives the radix-4 FFT its efficiency. The computations involved with each group of four frequency samples constitute the **radix-4 butterfly**, which is shown in Figure 2.40. Through further rearrangement, it can be shown that this computation can be simplified to three twiddle-factor multiplies and a length-4 DFT! The theory of multi-dimensional index maps (Section 2.3.6) shows that this must be the case, and that FFTs of any factorable length may consist of successive stages of shorter-length FFTs with twiddle-factor multiplications in between. The length-4 DFT requires no multiplies and only eight complex additions (this efficient computation can be derived using a radix-2 FFT (Section 2.3.4.2.1)).

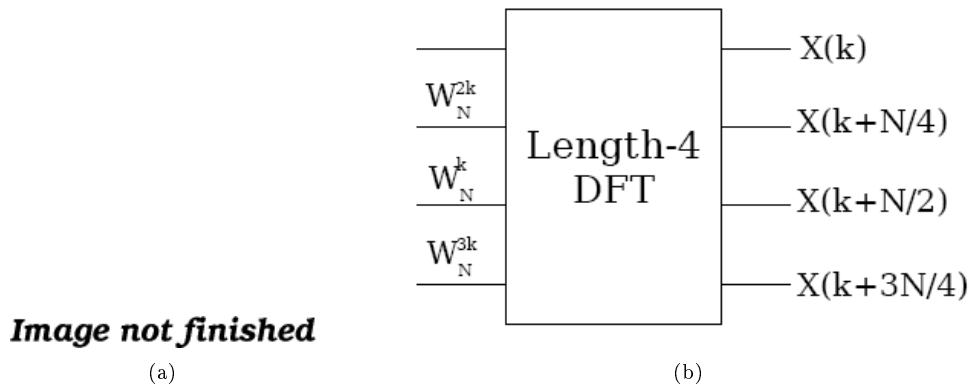


Figure 2.40: The radix-4 DIT butterfly can be simplified to a length-4 DFT preceded by three twiddle-factor multiplies.

If the FFT length $N = 4^M$, the shorter-length DFTs can be further decomposed recursively in the same manner to produce the full **radix-4 decimation-in-time FFT**. As in the radix-2 decimation-in-time FFT (Section 2.3.4.2.1), each stage of decomposition creates additional savings in computation. To determine the total computational cost of the radix-4 FFT, note that there are $M = \log_4 N = \frac{\log_2 N}{2}$ stages, each with $\frac{N}{4}$ butterflies per stage. Each radix-4 butterfly requires 3 complex multiplies and 8 complex additions. The total cost is then

Radix-4 FFT Operation Counts

- $3 \frac{N}{4} \frac{\log_2 N}{2} = \frac{3}{8} N \log_2 N$ complex multiplies (75% of a radix-2 FFT)
- $8 \frac{N}{4} \frac{\log_2 N}{2} = N \log_2 N$ complex adds (same as a radix-2 FFT)

The radix-4 FFT requires only 75% as many complex multiplies as the radix-2 (Section 2.3.4.2.1) FFTs, although it uses the same number of complex additions. These additional savings make it a widely-used FFT algorithm.

The decimation-in-time operation regroups the input samples at each successive stage of decomposition, resulting in a "digit-reversed" input order. That is, if the time-sample index n is written as a base-4 number, the order is that base-4 number reversed. The digit-reversal process is illustrated for a length- $N = 64$ example below.

Example 2.4: $N = 64 = 4^3$

Original Number	Original Digit Order	Reversed Digit Order	Digit-Reversed Number
0	000	000	0
1	001	100	16
2	002	200	32
3	003	300	48
4	010	010	4
5	011	110	20
\vdots	\vdots	\vdots	\vdots

It is important to note that, if the input signal data are placed in digit-reversed order before beginning the FFT computations, the outputs of each butterfly throughout the computation can be placed in the same memory locations from which the inputs were fetched, resulting in an **in-place algorithm** that requires no extra memory to perform the FFT. Most FFT implementations are in-place, and overwrite the input data with the intermediate values and finally the output. A slight rearrangement within the radix-4 FFT introduced by *Burrus*[5] allows the inputs to be arranged in bit-reversed (Section 2.3.4.2.1) rather than digit-reversed order.

A radix-4 decimation-in-frequency (Section 2.3.4.2.2) FFT can be derived similarly to the radix-2 DIF FFT (Section 2.3.4.2.2), by separately computing all four groups of every fourth *output* frequency sample. The DIF radix-4 FFT is a flow-graph reversal of the DIT radix-4 FFT, with the same operation counts and twiddle factors in the reversed order. The output ends up in digit-reversed order for an in-place DIF algorithm.

Exercise 2.5

(Solution on p. 189.)

How do we derive a radix-4 algorithm when $N = 4^M 2$?

2.3.4.4 Split-radix FFT Algorithms³⁰

The split-radix algorithm, first clearly described and named by *Duhamel and Hollman*[11] in 1984, required fewer total multiply and add operations than any previous power-of-two algorithm. (*Yavne*[31] first derived essentially the same algorithm in 1968, but the description was so atypical that the work was largely neglected.) For a time many FFT experts thought it to be optimal in terms of total complexity, but even more efficient variations have more recently been discovered by *Johnson and Frigo*[21].

The split-radix algorithm can be derived by careful examination of the radix-2 (Section 2.3.4.2.1) and radix-4 (Section 2.3.4.3) flowgraphs as in Figure 1 below. While in most places the radix-4 (Section 2.3.4.3) algorithm has fewer nontrivial twiddle factors, in some places the radix-2 (Section 2.3.4.2.1) actually lacks twiddle factors present in the radix-4 (Section 2.3.4.3) structure or those twiddle factors simplify to multiplication by $-j$, which actually requires only additions. By mixing radix-2 (Section 2.3.4.2.1) and radix-4 (Section 2.3.4.3) computations appropriately, an algorithm of lower complexity than either can be derived.

³⁰This content is available online at <<http://cnx.org/content/m12031/1.5/>>.

Motivation for split-radix algorithm

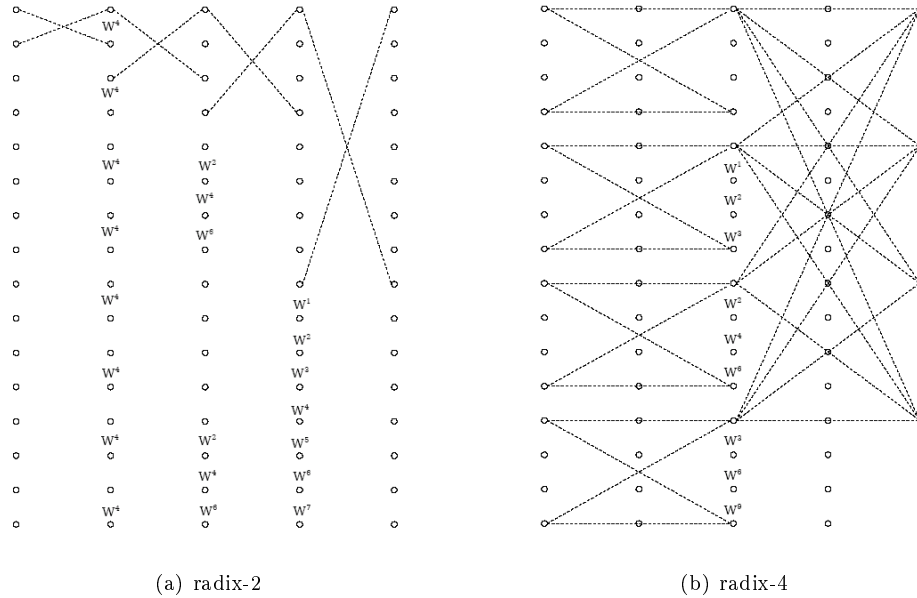


Figure 2.41: See Decimation-in-Time (DIT) Radix-2 FFT (Section 2.3.4.2.1) and Radix-4 FFT Algorithms (Section 2.3.4.3) for more information on these algorithms.

An alternative derivation notes that radix-2 butterflies of the form shown in Figure 2 can merge twiddle factors from two successive stages to eliminate one-third of them; hence, the split-radix algorithm requires only about two-thirds as many multiplications as a radix-2 FFT.

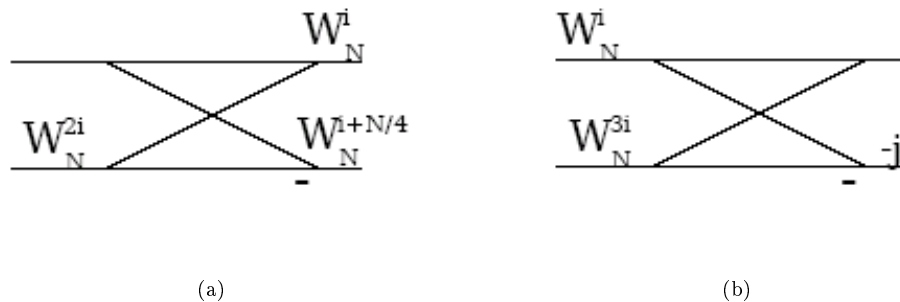


Figure 2.42: Note that these two butterflies are equivalent

The split-radix algorithm can also be derived by mixing the radix-2 (Section 2.3.4.2.1) and radix-4 (Section 2.3.4.3) decompositions.

DIT Split-radix derivation

$$\begin{aligned}
 X(k) &= \sum_{n=0}^{\frac{N}{2}-1} \left(x(2n) e^{-j\frac{2\pi(2n)k}{N}} \right) + \sum_{n=0}^{\frac{N}{4}-1} \left(x(4n+1) e^{-j\frac{2\pi(4n+1)k}{N}} \right) + \\
 &\sum_{n=0}^{\frac{N}{4}-1} \left(x(4n+3) e^{-j\frac{2\pi(4n+3)k}{N}} \right) = DFT_{\frac{N}{2}}[x(2n)] + W_N^k DFT_{\frac{N}{4}}x(4n+1) + \\
 &W_N^{3k} DFT_{\frac{N}{4}}x(4n+3)
 \end{aligned} \tag{2.20}$$

Figure 3 illustrates the resulting split-radix butterfly.

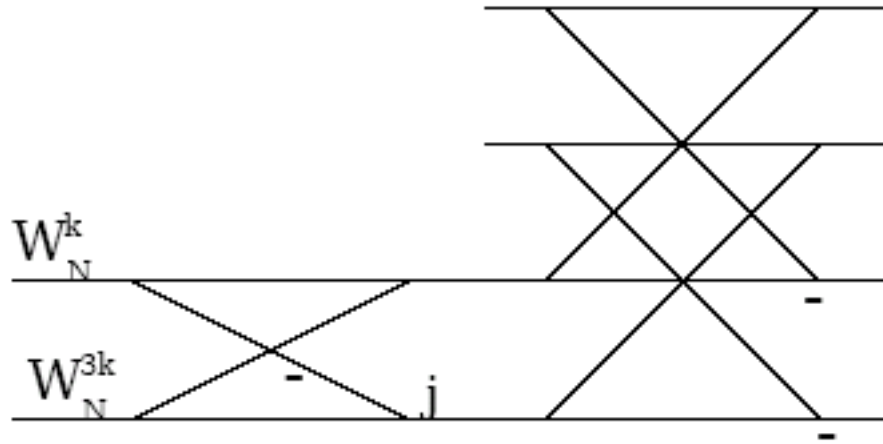


Figure 2.43: The split-radix butterfly mixes radix-2 and radix-4 decompositions and is L-shaped

Further decomposition of the half- and quarter-length DFTs yields the full split-radix algorithm. The mix of different-length FFTs in different parts of the flowgraph results in a somewhat irregular algorithm; *Sorensen et al.*[18] show how to adjust the computation such that the data retains the simpler radix-2 bit-reverse order. A decimation-in-frequency split-radix FFT can be derived analogously.

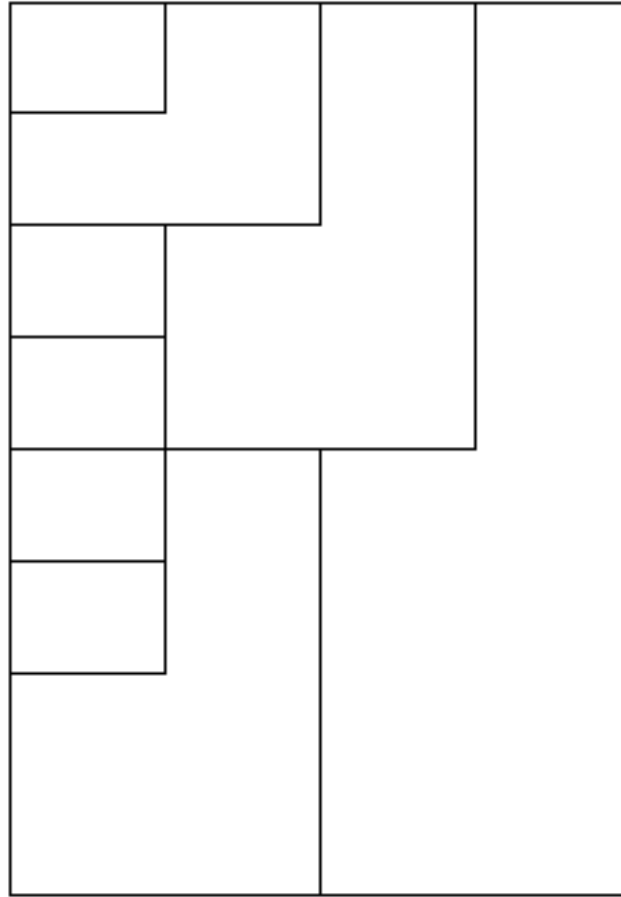


Figure 2.44: The split-radix transform has L-shaped butterflies

The multiplicative complexity of the split-radix algorithm is only about two-thirds that of the radix-2 FFT, and is better than the radix-4 FFT or any higher power-of-two radix as well. The additions within the complex twiddle-factor multiplies are similarly reduced, but since the underlying butterfly tree remains the same in all power-of-two algorithms, the butterfly additions remain the same and the overall reduction in additions is much less.

Operation Counts

	Complex M/As	Real M/As (4/2)	Real M/As (3/3)
Multiplies	$O\left[\frac{N}{3}\log_2 N\right]$	$\frac{4}{3}N\log_2 N - \frac{38}{9}N + 6 + \frac{2}{9}(-1)^M$	$N\log_2 N - 3N + 4$
Additions	$O[N\log_2 N]$	$\frac{8}{3}N\log_2 N - \frac{16}{9}N + 2 + \frac{2}{9}(-1)^M$	$3N\log_2 N - 3N + 4$

Comments

- The split-radix algorithm has a somewhat irregular structure. Successful programs have been written (Sorensen[18]) for uni-processor machines, but it may be difficult to efficiently code the split-radix algorithm for vector or multi-processor machines.
- *G. Bruun's algorithm*[2] requires only $N - 2$ more operations than the split-radix algorithm and has a regular structure, so it might be better for multi-processor or special-purpose hardware.
- The execution time of FFT programs generally depends more on compiler- or hardware-friendly software design than on the exact computational complexity. See Efficient FFT Algorithm and Programming Tricks (Section 2.3.5) for further pointers and links to good code.

2.3.5 Efficient FFT Algorithm and Programming Tricks³¹

The use of FFT algorithms (Section 2.3.1) such as the radix-2 decimation-in-time (Section 2.3.4.2.1) or decimation-in-frequency (Section 2.3.4.2.2) methods result in tremendous savings in computations when computing the discrete Fourier transform (Section 2.1.1). While most of the speed-up of FFTs comes from this, careful implementation can provide additional savings ranging from a few percent to several-fold increases in program speed.

2.3.5.1 Precompute twiddle factors

The twiddle factor (Section 2.3.4.2.1), or $W_N^k = e^{-j\frac{2\pi k}{N}}$, terms that multiply the intermediate data in the FFT algorithms (Section 2.3.1) consist of cosines and sines that each take the equivalent of several multiplies to compute. However, at most N unique twiddle factors can appear in any FFT or DFT algorithm. (For example, in the radix-2 decimation-in-time FFT (Section 2.3.4.2.1), only $\frac{N}{2}$ twiddle factors $\forall k, k = \{0, 1, 2, \dots, \frac{N}{2} - 1\} : (W_N^k)$ are used.) These twiddle factors can be precomputed once and stored in an array in computer memory, and accessed in the FFT algorithm by **table lookup**. This simple technique yields very substantial savings and is almost always used in practice.

2.3.5.2 Compiler-friendly programming

On most computers, only some of the total computation time of an FFT is spent performing the FFT butterfly computations; determining indices, loading and storing data, computing loop parameters and other operations consume the majority of cycles. Careful programming that allows the compiler to generate efficient code can make a several-fold improvement in the run-time of an FFT. The best choice of radix in terms of program speed may depend more on characteristics of the hardware (such as the number of CPU registers) or compiler than on the exact number of computations. Very often the manufacturer's library codes are carefully crafted by experts who know intimately both the hardware and compiler architecture and how to get the most performance out of them, so use of well-written FFT libraries is generally recommended. Certain freely available programs and libraries are also very good. Perhaps the best current general-purpose library is the FFTW³² package; information can be found at <http://www.fftw.org>³³. A paper by *Frigo and Johnson*[13] describes many of the key issues in developing compiler-friendly code.

2.3.5.3 Program in assembly language

While compilers continue to improve, FFT programs written directly in the assembly language of a specific machine are often several times faster than the best compiled code. This is particularly true for DSP microprocessors, which have special instructions for accelerating FFTs that compilers don't use. (I have myself seen differences of up to 26 to 1 in favor of assembly!) Very often, FFTs in the manufacturer's or high-performance third-party libraries are hand-coded in assembly. For DSP microprocessors, the codes

³¹This content is available online at <http://cnx.org/content/m12021/1.6/>.

³²<http://www.fftw.org>

³³<http://www.fftw.org>

developed by Meyer, Schuessler, and Schwarz[24] are perhaps the best ever developed; while the particular processors are now obsolete, the techniques remain equally relevant today. Most DSP processors provide special instructions and a hardware design favoring the radix-2 decimation-in-time algorithm, which is thus generally fastest on these machines.

2.3.5.4 Special hardware

Some processors have special hardware accelerators or co-processors specifically designed to accelerate FFT computations. For example, AMI Semiconductor's³⁴ Toccata³⁵ ultra-low-power DSP microprocessor family, which is widely used in digital hearing aids, have on-chip FFT accelerators; it is always faster and more power-efficient to use such accelerators and whatever radix they prefer.

In a surprising number of applications, almost all of the computations are FFTs. A number of special-purpose chips are designed to specifically compute FFTs, and are used in specialized high-performance applications such as radar systems. Other systems, such as OFDM³⁶ -based communications receivers, have special FFT hardware built into the digital receiver circuit. Such hardware can run many times faster, with much less power consumption, than FFT programs on general-purpose processors.

2.3.5.5 Effective memory management

Cache misses or excessive data movement between registers and memory can greatly slow down an FFT computation. Efficient programs such as the FFTW package³⁷ are carefully designed to minimize these inefficiencies. In-place algorithms (Section 2.3.4.2.1) reuse the data memory throughout the transform, which can reduce cache misses for longer lengths.

2.3.5.6 Real-valued FFTs

FFTs of real-valued signals require only half as many computations as with complex-valued data. There are several methods for reducing the computation, which are described in more detail in Sorensen *et al.*[19]

1. Use DFT symmetry properties (Section 2.1.1) to do two real-valued DFTs at once with one FFT program
2. Perform one stage of the radix-2 decimation-in-time (Section 2.3.4.2.1) decomposition and compute the two length- $\frac{N}{2}$ DFTs using the above approach.
3. Use a direct real-valued FFT algorithm; see H.V. Sorensen *et al.*[19]

2.3.5.7 Special cases

Occasionally only certain DFT frequencies are needed, the input signal values are mostly zero, the signal is real-valued (as discussed above), or other special conditions exist for which faster algorithms can be developed. Sorensen and Burrus[27] describe slightly faster algorithms for pruned³⁸ or zero-padded (Section 2.2.1.5: Zero-Padding) data. Goertzel's algorithm (Section 2.3.3) is useful when only a few DFT outputs are needed. The running FFT (Section 2.3.2) can be faster when DFTs of highly overlapped blocks of data are needed, as in a spectrogram (Section 2.2.3).

2.3.5.8 Higher-radix algorithms

Higher-radix algorithms, such as the radix-4 (Section 2.3.4.3), radix-8, or split-radix (Section 2.3.4.4) FFTs, require fewer computations and can produce modest but worthwhile savings. Even the split-radix FFT

³⁴<http://www.amis.com>

³⁵http://www.amis.com/products/dsp/toccata_plus.html

³⁶<http://en.wikipedia.org/wiki/OFDM>

³⁷<http://www.fftw.org>

³⁸<http://www.fftw.org/pruned.html>

(Section 2.3.4.4) reduces the multiplications by only 33% and the additions by a much lesser amount relative to the radix-2 FFTs (Section 2.3.4.2.1); significant improvements in program speed are often due to implicit loop-unrolling³⁹ or other compiler benefits than from the computational reduction itself!

2.3.5.9 Fast bit-reversal

Bit-reversing (Section 2.3.4.2.1) the input or output data can consume several percent of the total run-time of an FFT program. Several fast bit-reversal algorithms have been developed that can reduce this to two percent or less, including the method published by *D.M.W. Evans*[12].

2.3.5.10 Trade additions for multiplications

When FFTs first became widely used, hardware multipliers were relatively rare on digital computers, and multiplications generally required many more cycles than additions. Methods to reduce multiplications, even at the expense of a substantial increase in additions, were often beneficial. The prime factor algorithms (Section 2.3.7) and the Winograd Fourier transform algorithms (Section 2.6), which required fewer multiplies and considerably more additions than the power-of-two-length algorithms (Section 2.3.4.1), were developed during this period. Current processors generally have high-speed pipelined hardware multipliers, so trading multiplies for additions is often no longer beneficial. In particular, most machines now support single-cycle multiply-accumulate (MAC) operations, so balancing the number of multiplies and adds and combining them into single-cycle MACs generally results in the fastest code. Thus, the prime-factor and Winograd FFTs are rarely used today unless the application requires FFTs of a specific length.

It is possible to implement a complex multiply with 3 real multiplies and 5 real adds rather than the usual 4 real multiplies and 2 real adds:

$$(C + jS)(X + jY) = CX - SY + j(CY + SX)$$

but alternatively

$$Z = C(X - Y)$$

$$D = C + S$$

$$E = C - S$$

$$CX - SY = EY + Z$$

$$CY + SX = DX - Z$$

In an FFT, D and E come entirely from the twiddle factors, so they can be precomputed and stored in a look-up table. This reduces the cost of the complex twiddle-factor multiply to 3 real multiplies and 3 real adds, or one less and one more, respectively, than the conventional 4/2 computation.

2.3.5.11 Special butterflies

Certain twiddle factors, namely $W_N^0 = 1$, $W_N^{\frac{N}{2}}$, $W_N^{\frac{N}{4}}$, $W_N^{\frac{N}{8}}$, $W_N^{\frac{3N}{8}}$, etc., can be implemented with no additional operations, or with fewer real operations than a general complex multiply. Programs that specially implement such butterflies in the most efficient manner throughout the algorithm can reduce the computational cost by up to several N multiplies and additions in a length- N FFT.

³⁹http://en.wikipedia.org/wiki/Loop_unrolling

2.3.5.12 Practical Perspective

When optimizing FFTs for speed, it can be important to maintain perspective on the benefits that can be expected from any given optimization. The following list categorizes the various techniques by potential benefit; these will be somewhat situation- and machine-dependent, but clearly one should begin with the most significant and put the most effort where the pay-off is likely to be largest.

Methods to speed up computation of DFTs

- **Tremendous Savings** -
 1. FFT ($\frac{N}{\log_2 N}$ savings)
- **Substantial Savings** - (≥ 2)
 1. Table lookup of cosine/sine
 2. Compiler tricks/good programming
 3. Assembly-language programming
 4. Special-purpose hardware
 5. Real-data FFT for real data (factor of 2)
 6. Special cases
- **Minor Savings** -
 1. radix-4 (Section 2.3.4.3), split-radix (Section 2.3.4.4) (-10% - +30%)
 2. special butterflies
 3. 3-real-multiplication complex multiply
 4. Fast bit-reversal (up to 6%)

FACT: On general-purpose machines, computation is only part of the total run time. Address generation, indexing, data shuffling, and memory access take up much or most of the cycles.

FACT: A well-written radix-2 (Section 2.3.4.2.1) program will run much faster than a poorly written split-radix (Section 2.3.4.4) program!

2.3.6 Multidimensional Index Maps⁴⁰

2.3.6.1 Multidimensional Index Maps for DIF and DIT algorithms

2.3.6.1.1 Decimation-in-time algorithm

Radix-2 DIT (Section 2.3.4.2.1):

$$X(k) = \sum_{n=0}^{N-1} (x(n) W_N^{nk}) = \sum_{n=0}^{\frac{N}{2}-1} (x(2n) W_N^{2nk}) + \sum_{n=0}^{\frac{N}{2}-1} (x(2n+1) W_N^{(2n+1)k})$$

Formalization: Let $n = n_1 + 2n_2$: $n_1 = [0, 1]$: $n_2 = [0, 1, 2, \dots, \frac{N}{2} - 1]$

$$X(k) = \sum_{n=0}^{N-1} (x(n) W_N^{nk}) = \sum_{n_1=0}^1 \left(\sum_{n_2=0}^{\frac{N}{2}-1} (x(n_1 + 2n_2) W_N^{(n_1+2n_2)k}) \right)$$

Also, let $k = \frac{N}{2}k_1 + k_2$: $k_1 = [0, 1]$: $k_2 = [0, 1, 2, \dots, \frac{N}{2} - 1]$

NOTE: As long as there is a one-to-one correspondence between the original indices $[n, k] = [0, 1, 2, \dots, N-1]$ and the n, k generated by the index map, the computation is the *same*; only the order in which the sums are done is changed.

⁴⁰This content is available online at <<http://cnx.org/content/m12025/1.3/>>.

Rewriting the DFT (2.5) formula in terms of index map $n = n_1 + 2n_2$, $k = \frac{N}{2}k_1 + k_2$:

$$\begin{aligned}
X(k) &= X\left(\frac{N}{2}k_1 + k_2\right) \\
&= \sum_{n=0}^{N-1} \left(x(n) W_N^{n(\frac{N}{2}k_1 + k_2)} \right) \\
&= \sum_{n_1=0}^1 \left(\sum_{n_2=0}^{\frac{N}{2}-1} \left(x(n_1 + 2n_2) W_N^{(n_1+2n_2)(\frac{N}{2}k_1 + k_2)} \right) \right) \\
&= \sum_{n_1=0}^1 \left(\sum_{n_2=0}^{\frac{N}{2}-1} \left(x([n_1, n_2]) W_N^{\frac{N}{2}n_1k_1} W_N^{n_1k_2} W_N^{Nn_2k_1} W_N^{2n_2k_2} \right) \right) \\
&= \sum_{n_1=0}^1 \left(\sum_{n_2=0}^{\frac{N}{2}-1} \left(x([n_1, n_2]) W_2^{n_1k_2} W_N^{n_1k_2} 1 W_{\frac{N}{2}}^{n_2k_2} \right) \right) \\
&= \sum_{n_1=0}^1 \left(W_2^{n_1k_2} \left(W_N^{n_1k_2} \sum_{n_2=0}^{\frac{N}{2}-1} \left(x([n_1, n_2]) W_{\frac{N}{2}}^{n_2k_2} \right) \right) \right)
\end{aligned} \tag{2.21}$$

NOTE: Key to FFT is choosing index map so that one of the cross-terms disappears!

Exercise 2.6

What is an index map for a radix-4 (Section 2.3.4.3) DIT algorithm?

Exercise 2.7

What is an index map for a radix-4 (Section 2.3.4.3) DIF algorithm?

Exercise 2.8

What is an index map for a radix-3 DIT algorithm? (N a multiple of 3)

For arbitrary composite $N = N_1N_2$, we can define an index map

$$n = n_1 + N_1n_2$$

$$k = N_2k_1 + k_2$$

$$n_1 = [0, 1, 2, \dots, N_1 - 1]$$

$$k_1 = [0, 1, 2, \dots, N_1 - 1]$$

$$n_2 = [0, 1, 2, \dots, N_2 - 1]$$

$$k_2 = [0, 1, 2, \dots, N_2 - 1]$$

$$\begin{aligned}
X(k) &= X(k_1, k_2) \\
&= \sum_{n_1=0}^{N_1-1} \left(\sum_{n_2=0}^{N_2-1} \left(x(n_1, n_2) W_N^{N_2n_1k_1} W_N^{n_1k_2} W_N^{Nk_1n_2} W_N^{N_1n_2k_2} \right) \right) \\
&= \sum_{n_1=0}^{N_1-1} \left(\sum_{n_2=0}^{N_2-1} \left(x(n_1, n_2) W_{N_1}^{n_1k_1} W_N^{n_1k_2} 1 W_{N_2}^{n_2k_2} \right) \right) \\
&= DFT_{n_1, N_1} \left[W_N^{n_1k_2} DFT_{n_2, N_2} [x(n_1, n_2)] \right]
\end{aligned} \tag{2.22}$$

Computational cost in multiplies "Common Factor Algorithm (CFA)"

- N_1 length- N_2 DFTs $\Rightarrow N_1N_2^2$
- N twiddle factors $\Rightarrow N$
- N_2 length- N_1 DFTs $\Rightarrow N_2N_1^2$
- **Total** - $N_1N_2^2 + N_1N_2 + N_2N_1^2 = N(N_1 + N_2 + 1)$

"Direct": $N^2 = N(N_1N_2)$

Example 2.5

$$N_1 = 16$$

$$N_2 = 15$$

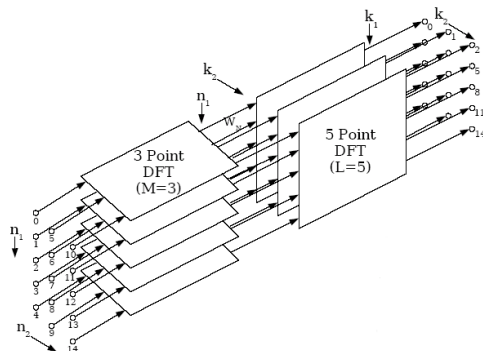
$$N = 240$$

$$\text{direct} = 240^2 = 57600$$

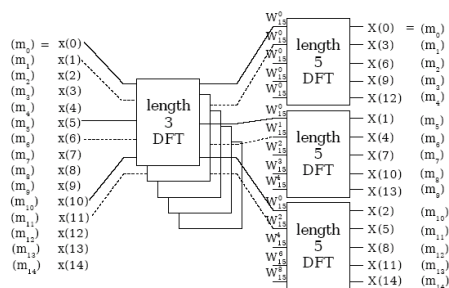
$$\text{CFA} = 7680$$

Tremendous saving for *any* composite N

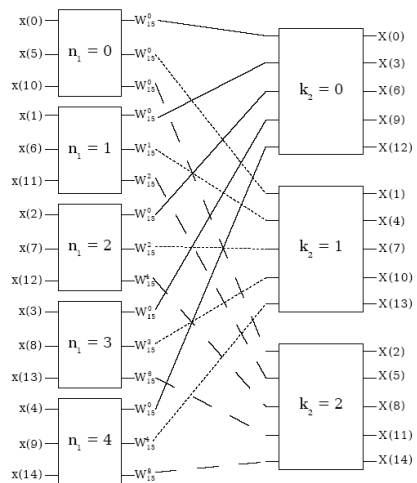
Pictorial Representations



(a) Emphasizes Multi-dimensional structure



(b) Emphasizes computer memory organization



(c) Easy to draw

Figure 2.45: $n = n_1 + 5n_2$, $k = 3k_1 + k_2$

Exercise 2.9

Can the composite CFAs be implemented in-place?

Exercise 2.10

What do we do with $N = N_1 N_2 N_3$?

2.3.7 The Prime Factor Algorithm⁴¹**2.3.7.1 General Index Maps**

$$n = (K_1 n_1 + K_2 n_2) \bmod N$$

$$n = (K_3 k_1 + K_4 k_2) \bmod N$$

$$n_1 = [0, 1, \dots, N_1 - 1]$$

$$k_1 = [0, 1, \dots, N_1 - 1]$$

$$n_2 = [0, 1, \dots, N_2 - 1]$$

$$k_2 = [0, 1, \dots, N_2 - 1]$$

The basic idea is to simply *reorder* the DFT (2.5) computation to expose the redundancies in the DFT (2.5), and exploit these to reduce computation!

Three conditions must be satisfied to make this map (p. 173) serve our purposes

1. Each map must be one-to-one from 0 to $N - 1$, because we want to do the *same* computation, just in a different order.
2. The map must be cleverly chosen so that computation is reduced
3. The map should be chosen to make the short-length transforms be DFTs (2.5). (Not essential, since fast algorithms for short-length DFT (2.5)-like computations could be developed, but it makes our work easier.)

2.3.7.1.1 Conditions for one-to-oneness of general index map**2.3.7.1.1.1 Case I**

N_1, N_2 relatively prime (greatest common denominator = 1) i.e. $\gcd(N_1, N_2) = 1$

$K_1 = aN_2$ and/or $K_2 = bN_1$ and $\gcd(K_1, N_1) = 1, \gcd(K_2, N_2) = 1$

2.3.7.1.1.2 Case II

N_1, N_2 *not* relatively prime: $\gcd(N_1, N_2) > 1$

$K_1 = aN_2$ and $K_2 \neq bN_1$ and $\gcd(a, N_1) = 1, \gcd(K_2, N_2) = 1$ or $K_1 \neq aN_2$ and $K_2 = bN_1$ and $\gcd(K_1, N_1) = 1, \gcd(b, N_2) = 1$ where $K_1, K_2, K_3, K_4, N_1, N_2, a, b$ integers

PROOF: Requires number-theory/abstract-algebra concepts. Reference: *C.S. Burrus*[3]

NOTE: Conditions of one-to-oneness must apply to both k and n

⁴¹This content is available online at <<http://cnx.org/content/m12033/1.3/>>.

2.3.7.1.2 Conditions for arithmetic savings

$$\begin{aligned}
X(k_1, k_2) &= \sum_{n_1=0}^{N_1-1} \left(\sum_{n_2=0}^{N_2-1} \left(x(n_1, n_2) W_N^{(K_1 n_1 + K_2 n_2)(K_3 k_1 + K_4 k_2)} \right) \right) \\
&= \sum_{n_1=0}^{N_1-1} \left(\sum_{n_2=0}^{N_2-1} \left(x(n_1, n_2) W_N^{K_1 K_3 n_1 k_1} W_N^{K_1 K_4 n_1 k_2} W_N^{K_2 K_3 n_2 k_1} W_N^{K_2 K_4 n_2 k_2} \right) \right)
\end{aligned} \tag{2.23}$$

- $(K_1 K_4) \bmod N = 0$ exclusive or $(K_2 K_3) \bmod N = 0 \Rightarrow$ Common Factor Algorithm (CFA). Then

$$X(k) = DFT_{N_i} [\text{twiddle factors} DFT_{N_j} [x(n_1, n_2)]]$$

- $(K_1 K_4) \bmod N$ and $(K_2 K_3) \bmod N = 0 \Rightarrow$ Prime Factor Algorithm (PFA).

$$X(k) = DFT_{N_i} [DFT_{N_j}]$$

No twiddle factors!

FACT: A PFA exists only and always for relatively prime N_1, N_2

2.3.7.1.3 Conditions for short-length transforms to be DFTs

$$(K_1 K_3) \bmod N = N_2 \text{ and } (K_2 K_4) \bmod N = N_1$$

NOTE: Convenient choice giving a PFA

$K_1 = N_2, K_2 = N_1, K_3 = N_2 ((N_2^{-1}) \bmod N_1) \bmod N_1, K_4 = N_1 ((N_1^{-1}) \bmod N_2) \bmod N_2$ where $(N_1^{-1}) \bmod N_2$ is an integer such that $(N_1 N_1^{-1}) \bmod = 1$

Example 2.6

$$N_1 = 3, N_2 = 5, N = 15$$

$$n = (5n_1 + 3n_2) \bmod 15$$

$$k = (10k_1 + 6k_2) \bmod 15$$

1. Checking Conditions for one-to-oneness -

$$5 = K_1 = aN_2 = 5a$$

$$3 = K_2 = bN_1 = 3b$$

$$\gcd(5, 3) = 1$$

$$\gcd(3, 5) = 1$$

$$10 = K_3 = aN_2 = 5a$$

$$6 = K_4 = bN_1 = 3b$$

$$\gcd(10, 3) = 1$$

$$\gcd(6, 5) = 1$$

2. Checking conditions for reduced computation -

$$(K_1 K_4) \bmod 15 = (5 \times 6) \bmod 15 = 0$$

$$(K_2 K_3) \bmod 15 = (3 \times 10) \bmod 15 = 0$$

3. Checking Conditions for making the short-length transforms be DFTS -

$$(K_1 K_3) \bmod 15 = (5 \times 10) \bmod 15 = 5 = N_2$$

$$(K_2 K_4) \bmod 15 = (3 \times 6) \bmod 15 = 3 = N_1$$

Therefore, this is a prime factor map.

2-D map

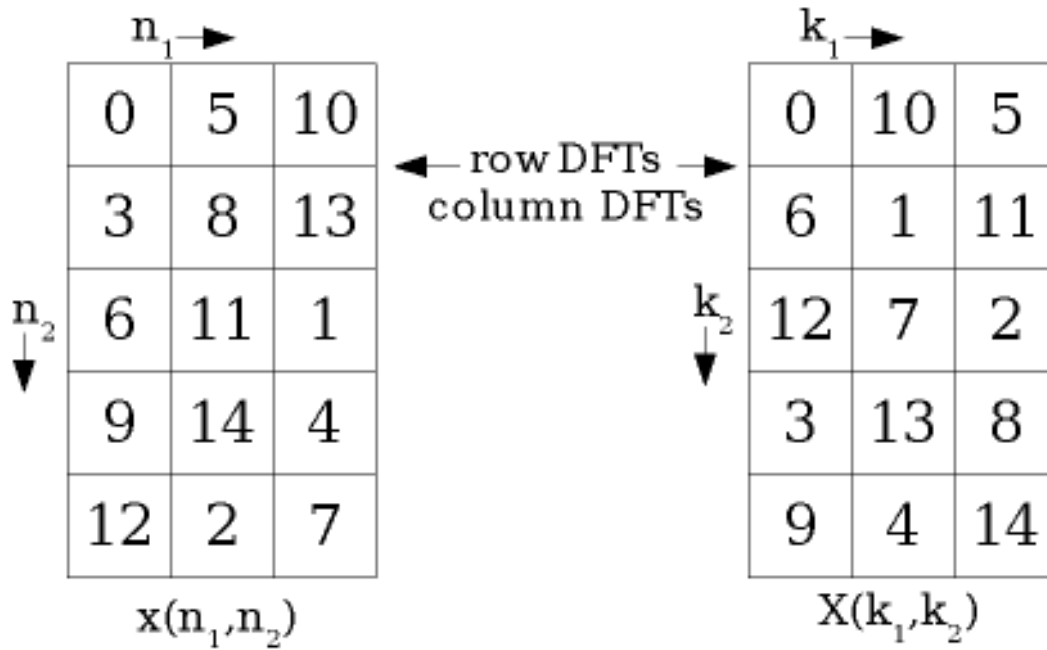


Figure 2.46: $n = (5n_1 + 3n_2) \bmod 15$ and $k = (10k_1 + 6k_2) \bmod 15$

Operation Counts

- N_2 length- N_1 DFTs + N_1 length- N_2 DFTs

$$N_2 N_1^2 + N_1 N_2^2 = N (N_1 + N_2)$$

complex multiplies

- Suppose $N = N_1 N_2 N_3 \dots N_M$

$$N(N_1 + N_2 + \dots + N_M)$$

Complex multiplies

DIFFERENT STRATEGIES: radix-2 (Section 2.3.4.2.1), radix-4 (Section 2.3.4.3) eliminate all multiplies in short-length DFTs, but have twiddle factors: PFA eliminates all twiddle factors, but ends up with multiplies in short-length DFTs (2.5). Surprisingly, total operation counts end up being very similar for similar lengths.

2.4 Fast Convolution⁴²

2.4.1 Fast Circular Convolution

Since,

$$\sum_{m=0}^{N-1} (x(m) (h(n-m)) \bmod N) = y(n) \text{ is equivalent to } Y(k) = X(k) H(k)$$

$y(n)$ can be computed as $y(n) = IDFT[DFT[x(n)] DFT[h(n)]]$

Cost

- **Direct**

- N^2 complex multiplies.
- $N(N-1)$ complex adds.

Via FFTs

- - 3 FFTs + N multiplies.
 - $N + \frac{3N}{2} \log_2 N$ complex multiplies.
 - $3(N \log_2 N)$ complex adds.

If $H(k)$ can be precomputed, cost is only 2 FFTs + N multiplies.

2.4.2 Fast Linear Convolution

DFT (2.5) produces circular convolution. For linear convolution, we must zero-pad sequences so that circular wrap-around always wraps over zeros.

⁴²This content is available online at <<http://cnx.org/content/m12022/1.5/>>.

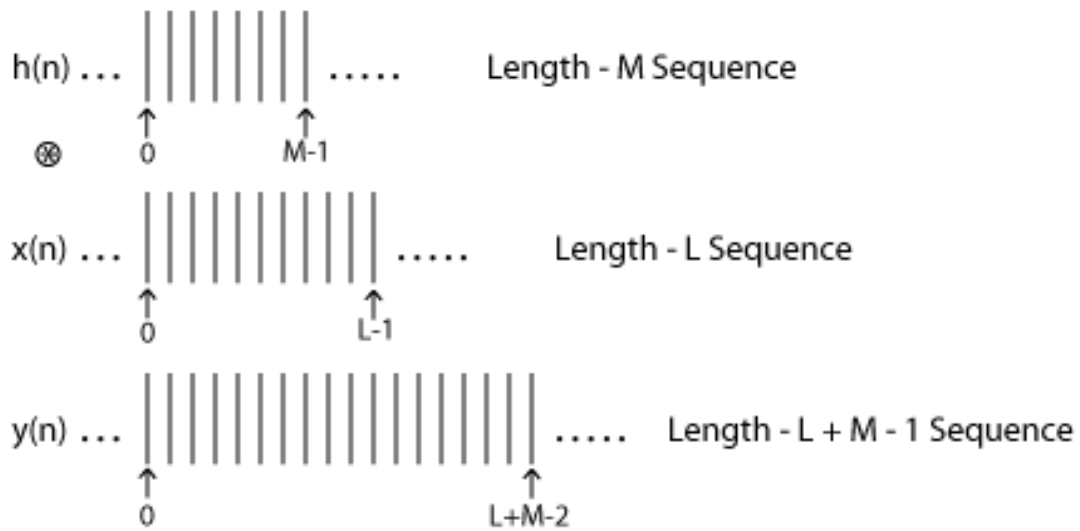


Figure 2.47

To achieve linear convolution using fast circular convolution, we must use zero-padded DFTs of length $N \geq L + M - 1$

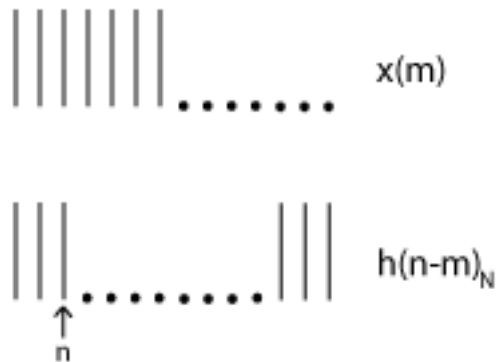


Figure 2.48

Choose shortest convenient N (usually smallest power-of-two greater than or equal to $L + M - 1$)

$$y(n) = IDFT_N [DFT_N [x(n)] DFT_N [h(n)]]$$

NOTE: There is some inefficiency when compared to circular convolution due to longer zero-padded

DFTs (2.5). Still, $O\left(\frac{N}{\log_2 N}\right)$ savings over direct computation.

2.4.3 Running Convolution

Suppose $L = \infty$, as in a real time filter application, or ($L \gg M$). There are efficient block methods for computing fast convolution.

2.4.3.1 Overlap-Save (OLS) Method

Note that if a length- M filter $h(n)$ is circularly convolved with a length- N segment of a signal $x(n)$,

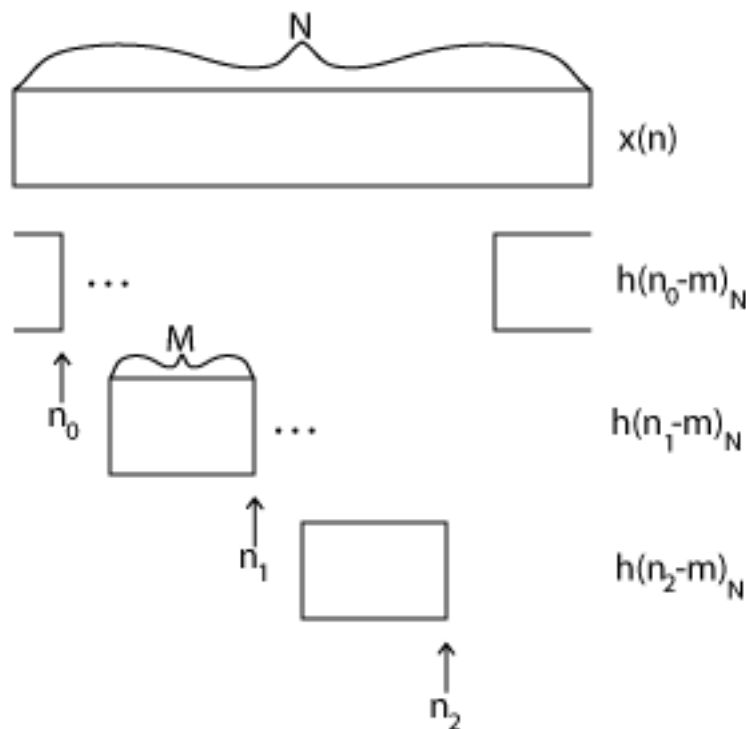


Figure 2.49

the first $M - 1$ samples are wrapped around and thus is *incorrect*. However, for $M - 1 \leq n \leq N - 1$, the convolution is linear convolution, so these samples are correct. Thus $N - M + 1$ good outputs are produced for each length- N circular convolution.

The Overlap-Save Method: Break long signal into successive blocks of N samples, each block overlapping the previous block by $M - 1$ samples. Perform circular convolution of each block with filter $h(m)$. Discard first $M - 1$ points in each output block, and concatenate the remaining points to create $y(n)$.

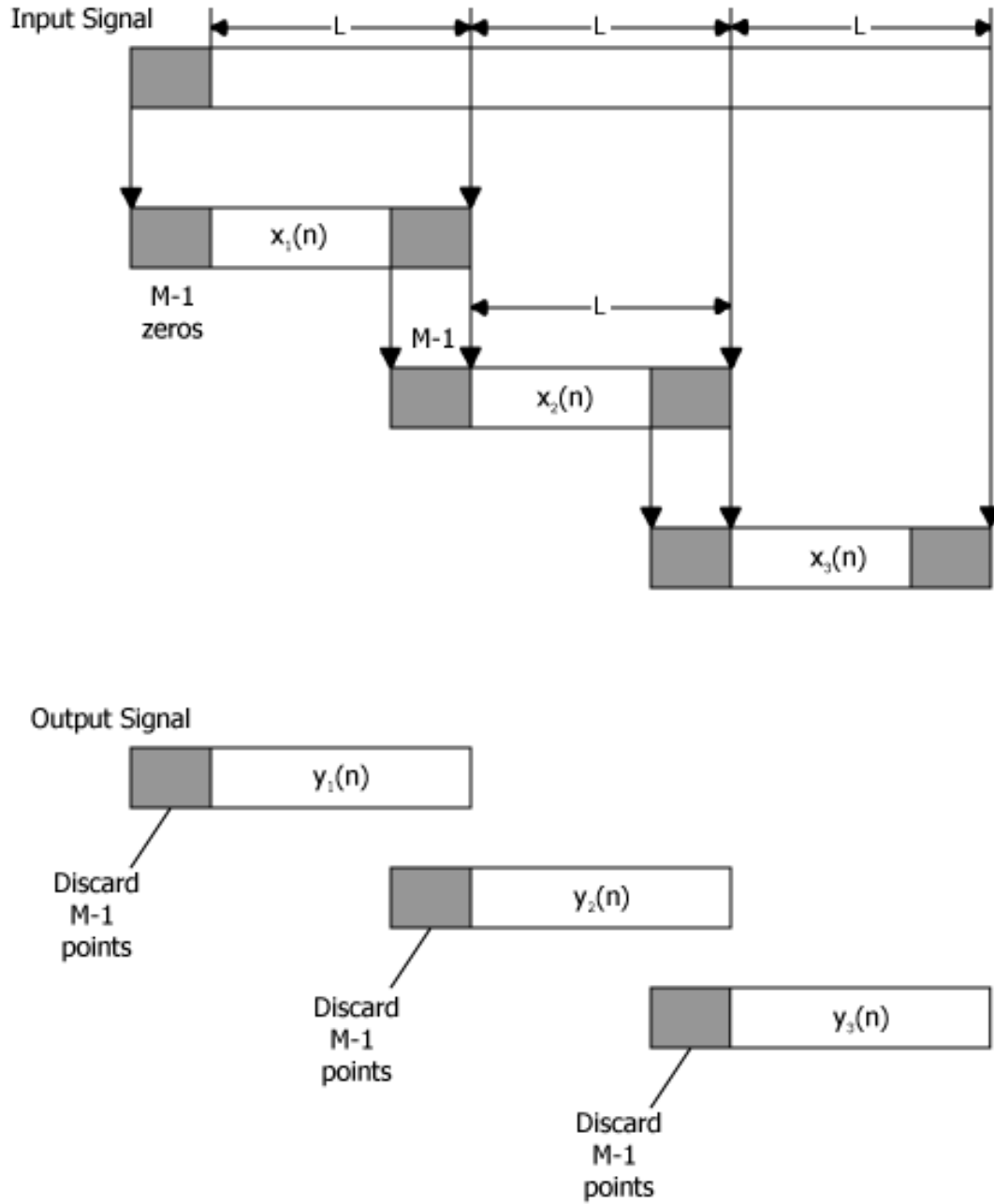


Figure 2.50

Computation cost for a length- N equals 2^n FFT per output sample is (assuming precomputed $H(k)$) 2 FFTs and N multiplies

$$\frac{2 \left(\frac{N}{2} \log_2 N \right) + N}{N - M + 1} = \frac{N (\log_2 N + 1)}{N - M + 1} \text{ complex multiplies}$$

$$\frac{2(N \log_2 N)}{N - M + 1} = \frac{2N \log_2 N}{N - M + 1} \text{complex adds}$$

Compare to M mults, $M - 1$ adds per output point for direct method. For a given M , optimal N can be determined by finding N minimizing operation counts. Usually, optimal N is $4M \leq N_{opt} \leq 8M$.

2.4.3.2 Overlap-Add (OLA) Method

Zero-pad length- L blocks by $M - 1$ samples.

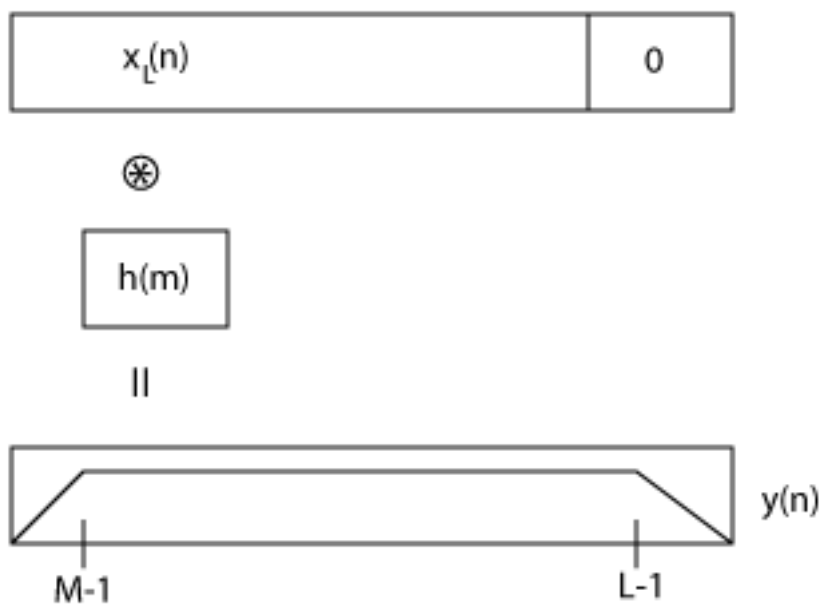


Figure 2.51

Add successive blocks, overlapped by $M - 1$ samples, so that the tails sum to produce the complete linear convolution.

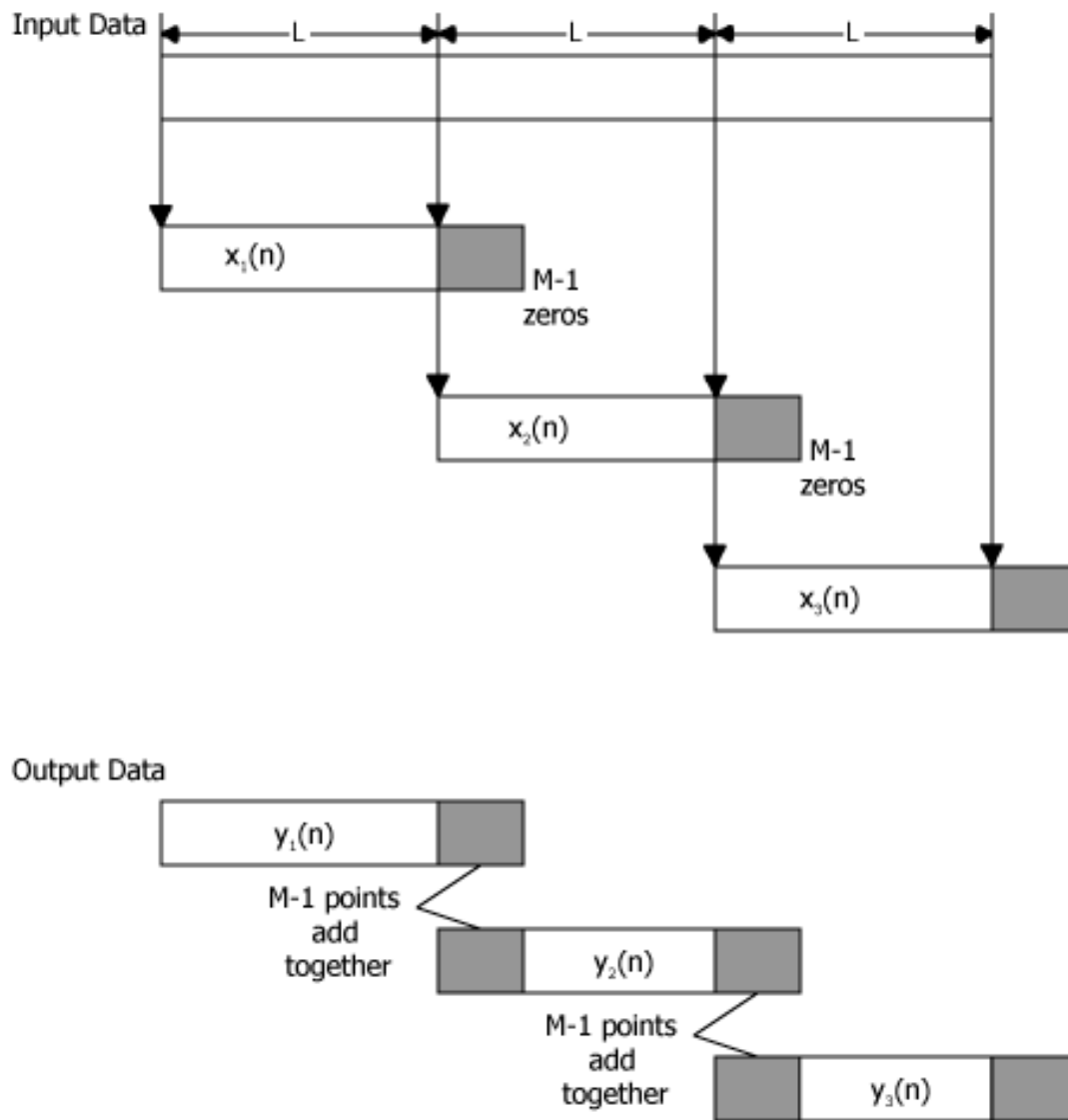


Figure 2.52

Computational Cost: Two length $N = L + M - 1$ FFTs and M mults and $M - 1$ adds per L output points; essentially the same as OLS method.

2.5 Chirp-z Transform⁴³

Let $z^k = AW^{-k}$, where $A = A_0 e^{j\theta_0}$, $W = W_0 e^{-(j\phi_0)}$.

⁴³This content is available online at <http://cnx.org/content/m12013/1.4/>.

We wish to compute M samples, $k = [0, 1, 2, \dots, M-1]$ of

$$X(z_k) = \sum_{n=0}^{N-1} (x(n) z_k^{-n}) = \sum_{n=0}^{N-1} (x(n) A^{-n} W^{nk})$$

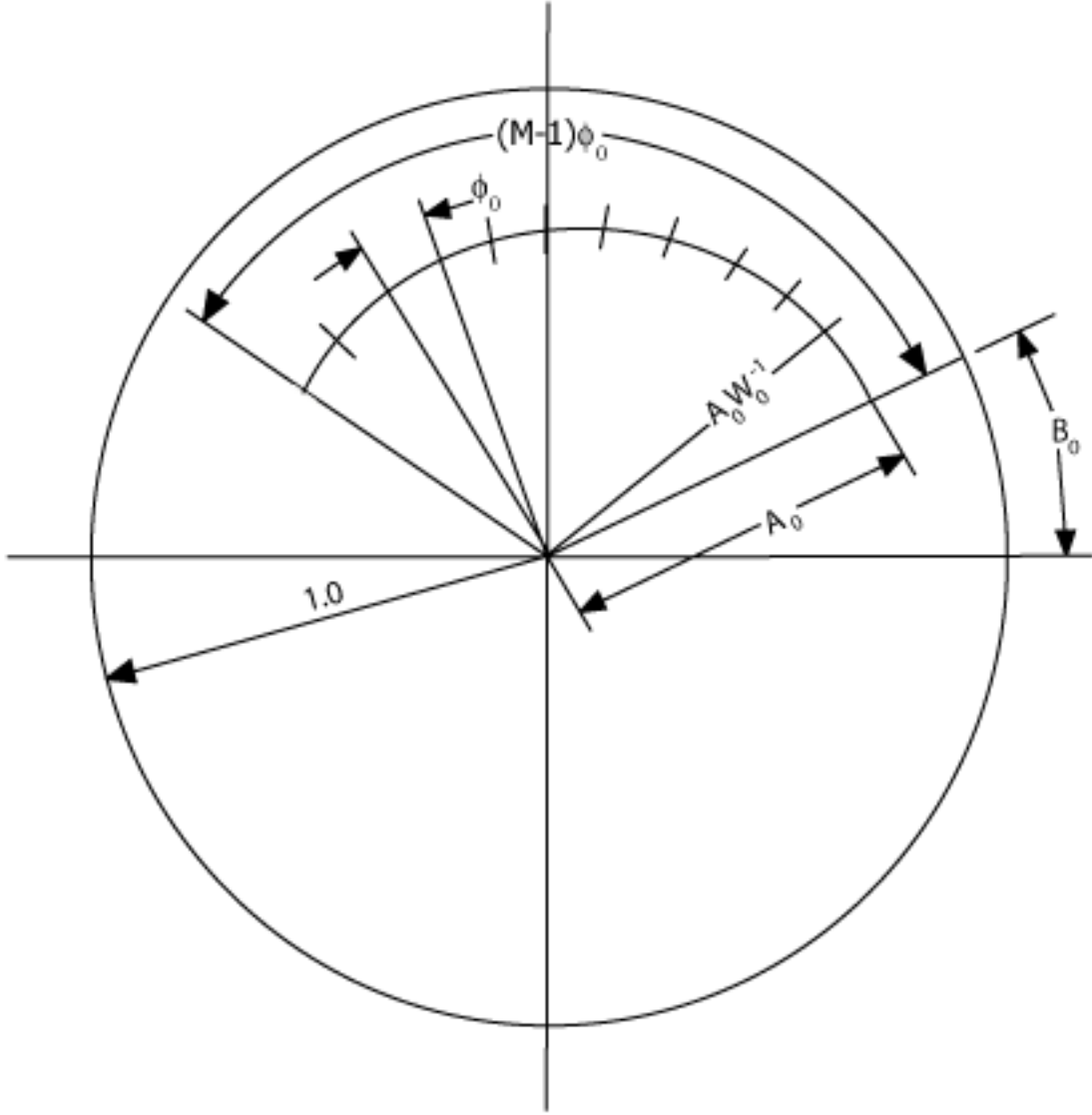


Figure 2.53

Note that $(k-n)^2 = n^2 - 2nk + k^2 \Rightarrow nk = \frac{1}{2} (n^2 + k^2 - (k-n)^2)$, So

$$X(z_k) = \sum_{n=0}^{N-1} \left(x(n) A^{-n} W^{\frac{n^2}{2}} W^{\frac{k^2}{2}} W^{-\frac{(k-n)^2}{2}} \right)$$

$$= W^{\frac{k^2}{2}} \sum_{n=0}^{N-1} \left(x(n) A^{-n} W^{\frac{n^2}{2}} W^{-\frac{(k-n)^2}{2}} \right)$$

Thus, $X(z_k)$ can be compared by

1. Premultiply $x(n)$ by $A^n W^{\frac{n^2}{2}}$, $n = [0, 1, \dots, N-1]$ to make $y(n)$
2. Linearly convolve with $W^{-\frac{(k-n)^2}{2}}$
3. Post multiply by to get $W^{\frac{k^2}{2}}$ to get $X(z_k)$.

1. (list, item 1, p. 183) and 3. (list, item 3, p. 183) require N and M operations respectively. 2. (list, item 2, p. 183) can be performed efficiently using fast convolution.

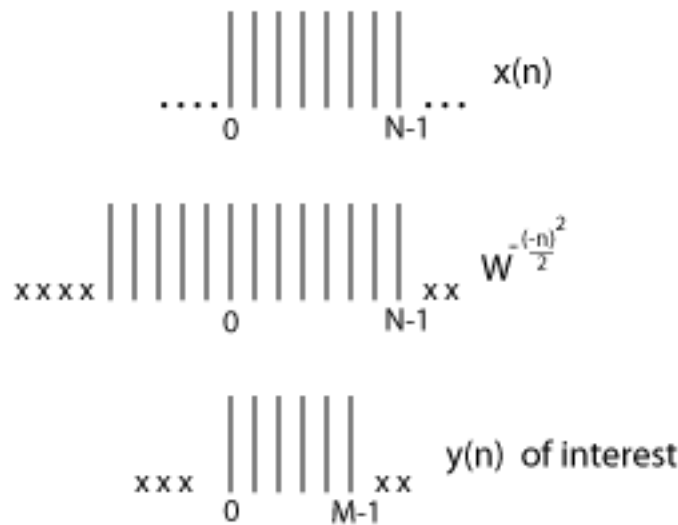


Figure 2.54

$W^{-\left(\frac{n^2}{2}\right)}$ is required only for $-(N-1) \leq n \leq M-1$, so this linear convolution can be implemented with $L \geq N + M - 1$ FFTs.

NOTE: Wrap $W^{-\left(\frac{n^2}{2}\right)}$ around L when implementing with circular convolution.

So, a weird-length DFT can be implemented relatively efficiently using power-of-two algorithms via the chirp-z transform.

Also useful for "zoom-FFTs".

2.6 FFTs of prime length and Rader's conversion⁴⁴

The power-of-two FFT algorithms (Section 2.3.4.1), such as the radix-2 (Section 2.3.4.2.1) and radix-4 (Section 2.3.4.3) FFTs, and the common-factor (Section 2.3.6) and prime-factor (Section 2.3.7) FFTs, achieve great reductions in computational complexity of the DFT (Section 2.1.1) when the length, N , is a composite

⁴⁴This content is available online at <<http://cnx.org/content/m12023/1.3/>>.

number. DFTs of prime length are sometimes needed, however, particularly for the short-length DFTs in common-factor or prime-factor algorithms. The methods described here, along with the composite-length algorithms, allow fast computation of DFTs of *any* length.

There are two main ways of performing DFTs of prime length:

1. Rader's conversion, which is most efficient, and the
2. Chirp-z transform (Section 2.5), which is simpler and more general.

Oddly enough, both work by turning prime-length DFTs into convolution! The resulting convolutions can then be computed efficiently by either

1. fast convolution (Section 2.4) via composite-length FFTs (simpler) or by
2. Winograd techniques (more efficient)

2.6.1 Rader's Conversion

Rader's conversion is a one-dimensional index-mapping (Section 2.3.6) scheme that turns a length- N DFT (2.5) (N prime) into a length- $(N-1)$ convolution and a few additions. Rader's conversion works *only* for prime-length N .

An **index map** simply rearranges the order of the sum operation in the DFT definition (Section 2.1.1). Because addition is a commutative operation, the same mathematical result is produced from any order, as long as all of the same terms are added once and only once. (This is the condition that defines an index map.) Unlike the multi-dimensional index maps (Section 2.3.6) used in deriving common factor (Section 2.3.6) and prime-factor FFTs (Section 2.3.7), Rader's conversion uses a one-dimensional index map in a finite group of N integers: $k = (r^m) \bmod N$

2.6.1.1 Fact from number theory

If N is prime, there exists an integer " r " called a **primitive root**, such that the index map $k = (r^m) \bmod N$, $m = [0, 1, 2, \dots, N-2]$, uniquely generates all elements $k = [1, 2, 3, \dots, N-1]$

Example 2.7

$$N = 5, r = 2$$

$$(2^0) \bmod 5 = 1$$

$$(2^1) \bmod 5 = 2$$

$$(2^2) \bmod 5 = 4$$

$$(2^3) \bmod 5 = 3$$

2.6.1.2 Another fact from number theory

For N prime, the inverse of r (i.e. $(r^{-1}r) \bmod N = 1$ is also a primitive root (call it r^{-1}).

Example 2.8

$$N = 5, r = 2 \quad r^{-1} = 3$$

$$(2 \times 3) \bmod 5 = 1$$

$$(3^0) \bmod 5 = 1$$

$$(3^1) \bmod 5 = 3$$

$$(3^2) \bmod 5 = 4$$

$$(3^3) \bmod 5 = 2$$

So why do we care? Because we can use these facts to turn a DFT (2.5) into a convolution!

2.6.1.3 Rader's Conversion

Let $\forall mn, m = [0, 1, \dots, N-2] \wedge n \in [1, 2, \dots, N-1] : (n = (r^{-m}) \bmod N), \forall pk, p = [0, 1, \dots, N-2] \wedge k \in [1, 2, \dots, N-1] : (k = (r^p) \bmod N)$

$$X(k) = \sum_{n=0}^{N-1} (x(n) W_N^{nk}) = \begin{cases} x(0) + \sum_{n=1}^{N-1} (x(n) W_N^{nk}) & \text{if } k \neq 0 \\ \sum_{n=0}^{N-1} (x(n)) & \text{if } k = 0 \end{cases}$$

where for convenience $W_N^{nk} = e^{-j\frac{2\pi nk}{N}}$ in the DFT equation. For $k \neq 0$

$$\begin{aligned} X((r^p) \bmod N) &= \sum_{m=0}^{N-2} \left(x((r^{-m}) \bmod N) W^{r^p r^{-m}} \right) + x(0) \\ &= \sum_{m=0}^{N-2} \left(x((r^{-m}) \bmod N) W^{r^{p-m}} \right) + x(0) \\ &= x(0) + x((r^{-l}) \bmod N) * W^{r^l} \end{aligned} \tag{2.24}$$

where $l = [0, 1, \dots, N-2]$

Example 2.9

$N = 5, r = 2, r^{-1} = 3$

$$\begin{pmatrix} X(0) \\ X(1) \\ X(2) \\ X(3) \\ X(4) \end{pmatrix} = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 2 & 3 & 4 \\ 0 & 2 & 4 & 1 & 3 \\ 0 & 3 & 1 & 4 & 2 \\ 0 & 4 & 3 & 2 & 1 \end{pmatrix} \begin{pmatrix} x(0) \\ x(1) \\ x(2) \\ x(3) \\ x(4) \end{pmatrix}$$

$$\begin{pmatrix} X(0) \\ X(1) \\ X(2) \\ X(4) \\ X(3) \end{pmatrix} = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 3 & 4 & 2 \\ 0 & 2 & 1 & 3 & 4 \\ 0 & 4 & 2 & 1 & 1 \\ 0 & 3 & 4 & 2 & 3 \end{pmatrix} \begin{pmatrix} x(0) \\ x(1) \\ x(3) \\ x(4) \\ x(2) \end{pmatrix}$$

where for visibility the matrix entries represent only the *power*, m of the corresponding DFT term W_N^m . Note that the 4-by-4 circulant matrix⁴⁵

$$\begin{pmatrix} 1 & 3 & 4 & 2 \\ 2 & 1 & 3 & 4 \\ 4 & 2 & 1 & 1 \\ 3 & 4 & 2 & 3 \end{pmatrix}$$

corresponds to a length-4 circular convolution.

Rader's conversion turns a prime-length DFT (2.5) into a few adds and a *composite-length* $(N - 1)$ circular convolution, which can be computed efficiently using either

1. fast convolution (Section 2.4) via FFT and IFFT
2. index-mapped convolution algorithms and short Winograd convolution algorithms. (Rather complicated, and trades fewer multiplies for many more adds, which may not be worthwhile on most modern processors.) See *R.C. Agarwal and J.W. Cooley*[1]

2.6.2 Winograd minimum-multiply convolution and DFT algorithms

S. Winograd has proved that a length- N circular or linear convolution or DFT (2.5) requires less than $2N$ multiplies (for real data), or $4N$ real multiplies for complex data. (This doesn't count multiplies by rational fractions, like 3 or $\frac{1}{N}$ or $\frac{5}{17}$, which can be computed with additions and one overall scaling factor.) Furthermore, Winograd showed how to construct algorithms achieving these counts. Winograd prime-length DFTs and convolutions have the following characteristics:

1. Extremely efficient for small N ($N < 20$)
2. The number of adds becomes *huge* for large N .

Thus Winograd's minimum-multiply FFT's are useful only for small N . They are very important for Prime-Factor Algorithms (Section 2.3.7), which generally use Winograd modules to implement the short-length DFTs. Tables giving the multiplies and adds necessary to compute Winograd FFTs for various lengths can be found in *C.S. Burrus (1988)*[4]. Tables and FORTRAN and TMS32010 programs for these short-length transforms can be found in *C.S. Burrus and T.W. Parks (1985)*[6]. The theory and derivation of these algorithms is quite elegant but requires substantial background in number theory and abstract algebra. Fortunately for the practitioner, all of the short algorithms one is likely to need have already been derived and can simply be looked up without mastering the details of their derivation.

2.6.3 Winograd Fourier Transform Algorithm (WFTA)

The Winograd Fourier Transform Algorithm (WFTA) is a technique that recombines the short Winograd modules in a prime-factor FFT (Section 2.3.7) into a composite- N structure with fewer multiplies but more adds. While theoretically interesting, WFTAs are complicated and different for every length, and on modern processors with hardware multipliers the trade of multiplies for many more adds is very rarely useful in practice today.

⁴⁵http://en.wikipedia.org/wiki/Circulant_matrix

2.7 Choosing the Best FFT Algorithm⁴⁶

2.7.1 Choosing an FFT length

The most commonly used FFT algorithms *by far* are the power-of-two-length FFT (Section 2.3.4.1) algorithms. The Prime Factor Algorithm (PFA) (Section 2.3.7) and Winograd Fourier Transform Algorithm (WFTA) (Section 2.6.3: Winograd Fourier Transform Algorithm (WFTA)) require somewhat fewer multiplies, but the overall difference usually isn't sufficient to warrant the extra difficulty. This is particularly true now that most processors have single-cycle pipelined hardware multipliers, so the total operation count is more relevant. As can be seen from the following table, for similar lengths the split-radix algorithm is comparable in total operations to the Prime Factor Algorithm, and is considerably better than the WFTA, although the PFA and WFTA require fewer multiplications and more additions. Many processors now support single cycle multiply-accumulate (MAC) operations; in the power-of-two algorithms all multiplies can be combined with adds in MACs, so the number of additions is the most relevant indicator of computational cost.

Representative FFT Operation Counts

	FFT length	Multiplies (real)	Adds(real)	Mults + Adds
Radix 2	1024	10248	30728	40976
Split Radix	1024	7172	27652	34824
Prime Factor Alg	1008	5804	29100	34904
Winograd FT Alg	1008	3548	34416	37964

The Winograd Fourier Transform Algorithm (Section 2.6.3: Winograd Fourier Transform Algorithm (WFTA)) is particularly difficult to program and is rarely used in practice. For applications in which the transform length is somewhat arbitrary (such as fast convolution or general spectrum analysis), the length is usually chosen to be a power of two. When a particular length is required (for example, in the USA each carrier has exactly 416 frequency channels in each band in the AMPS⁴⁷ cellular telephone standard), a Prime Factor Algorithm (Section 2.3.7) for all the relatively prime terms is preferred, with a Common Factor Algorithm (Section 2.3.6) for other non-prime lengths. Winograd's short-length modules (Section 2.6) should be used for the prime-length factors that are not powers of two. The chirp z-transform (Section 2.5) offers a universal way to compute any length DFT (Section 2.2.1) (for example, Matlab⁴⁸ reportedly uses this method for lengths other than a power of two), at a few times higher cost than that of a CFA or PFA optimized for that specific length. The chirp z-transform (Section 2.5), along with Rader's conversion (Section 2.6.1: Rader's Conversion), assure us that algorithms of $O(N \log N)$ complexity exist for *any* DFT length N .

2.7.2 Selecting a power-of-two-length algorithm

The choice of a power-of-two algorithm may not just depend on computational complexity. The latest extensions of the split-radix algorithm (Section 2.3.4.4) offer the lowest known power-of-two FFT operation counts, but the 10%-30% difference may not make up for other factors such as regularity of structure or data flow, FFT programming tricks (Section 2.3.5), or special hardware features. For example, the decimation-in-time radix-2 FFT (Section 2.3.4.2.1) is the fastest FFT on Texas Instruments'⁴⁹ TMS320C54x DSP microprocessors, because this processor family has special assembly-language instructions that accelerate this particular algorithm. On other hardware, radix-4 algorithms (Section 2.3.4.3) may be more efficient. Some

⁴⁶This content is available online at <<http://cnx.org/content/m12060/1.3/>>.

⁴⁷<http://en.wikipedia.org/wiki/AMPS>

⁴⁸<http://www.mathworks.com/products/matlab/>

⁴⁹<http://www.ti.com/>

devices, such as AMI Semiconductor's⁵⁰ Toccata⁵¹ ultra-low-power DSP microprocessor family, have on-chip FFT accelerators; it is always faster and more power-efficient to use these accelerators and whatever radix they prefer. For fast convolution (Section 2.4), the decimation-in-frequency (Section 2.3.4.2.2) algorithms may be preferred because the bit-reversing can be bypassed; however, most DSP microprocessors provide zero-overhead bit-reversed indexing hardware and prefer decimation-in-time algorithms, so this may not be true for such machines. Good, compiler- or hardware-friendly programming always matters more than modest differences in raw operation counts, so manufacturers' or good third-party FFT libraries are often the best choice. The module FFT programming tricks (Section 2.3.5) references some good, free FFT software (including the FFTW⁵² package) that is carefully coded to be compiler-friendly; such codes are likely to be considerably faster than codes written by the casual programmer.

2.7.3 Multi-dimensional FFTs

Multi-dimensional FFTs pose additional possibilities and problems. The orthogonality and separability of multi-dimensional DFTs allows them to be efficiently computed by a series of one-dimensional FFTs along each dimension. (For example, a two-dimensional DFT can quickly be computed by performing FFTs of each row of the data matrix followed by FFTs of all columns, or vice-versa.) **Vector-radix FFTs** have been developed with higher efficiency per sample than row-column algorithms. Multi-dimensional datasets, however, are often large and frequently exceed the cache size of the processor, and excessive cache misses may increase the computational time greatly, thus overwhelming any minor complexity reduction from a vector-radix algorithm. Either vector-radix FFTs must be carefully programmed to match the cache limitations of a specific processor, or a row-column approach should be used with matrix transposition in between to ensure data locality for high cache utilization throughout the computation.

2.7.4 Few time or frequency samples

FFT algorithms gain their efficiency through intermediate computations that can be reused to compute many DFT frequency samples at once. Some applications require only a handful of frequency samples to be computed; when that number is of order less than $O(\log N)$, direct computation of those values via Goertzel's algorithm (Section 2.3.3) is faster. This has the additional advantage that any frequency, not just the equally-spaced DFT frequency samples, can be selected. *Sorensen and Burrus*[28] developed algorithms for when most input samples are zero or only a block of DFT frequencies are needed, but the computational cost is of the same order.

Some applications, such as time-frequency analysis via the short-time Fourier transform (Section 2.2.3) or spectrogram (Section 2.2.3), require DFTs of overlapped blocks of discrete-time samples. When the step-size between blocks is less than $O(\log N)$, the running FFT (Section 2.3.2) will be most efficient. (Note that any window must be applied via frequency-domain convolution, which is quite efficient for sinusoidal windows such as the Hamming window.) For step-sizes of $O(\log N)$ or greater, computation of the DFT of each successive block via an FFT is faster.

⁵⁰<http://www.amis.com>

⁵¹http://www.amis.com/products/dsp/toccata_plus.html

⁵²<http://www.fftw.org/>

Solutions to Exercises in Chapter 2

Solution to Exercise 2.1 (p. 104)

In general, *NO*. The DTFT exactly corresponds to the continuous-time Fourier transform only when the signal is bandlimited and sampled at more than twice its highest frequency. The DFT frequency values exactly correspond to frequency samples of the DTFT only when the discrete-time signal is time-limited. However, a bandlimited continuous-time signal cannot be time-limited, so in general these conditions cannot both be satisfied.

It can, however, be true for a small class of analog signals which are not time-limited but happen to exactly equal zero at all *sample times* outside of the interval $n \in [0, N - 1]$. The sinc function with a bandwidth equal to the Nyquist frequency and centered at $t = 0$ is an example.

Solution to Exercise 2.2 (p. 130)

Solution to Exercise 2.3 (p. 138)

Solution to Exercise 2.4 (p. 140)

Solution to Exercise 2.5 (p. 162)

Perform a radix-2 decomposition for one stage, then radix-4 decompositions of all subsequent shorter-length DFTs.

Chapter 3

Digital Filter Design

3.1 Overview of Digital Filter Design¹

Advantages of FIR filters

1. Straight forward conceptually and simple to implement
2. Can be implemented with fast convolution
3. Always stable
4. Relatively insensitive to quantization
5. Can have linear phase (same time delay of all frequencies)

Advantages of IIR filters

1. Better for approximating analog systems
2. For a given magnitude response specification, IIR filters often require much less computation than an equivalent FIR, particularly for narrow transition bands

Both FIR and IIR filters are very important in applications.

Generic Filter Design Procedure

1. Choose a desired response, based on application requirements
2. Choose a filter class
3. Choose a quality measure
4. Solve for the filter in class 2 optimizing criterion in 3

3.1.1 Perspective on FIR filtering

Most of the time, people do L^∞ optimal design, using the Parks-McClellan algorithm (Section 3.2.4). This is probably the second most important technique in "classical" signal processing (after the Cooley-Tukey (radix-2 (Section 2.3.4.2.1)) FFT).

Most of the time, FIR filters are designed to have linear phase. The most important advantage of FIR filters over IIR filters is that they can have exactly linear phase. There are advanced design techniques for minimum-phase filters, constrained L^2 optimal designs, *etc.* (see chapter 8 of text). However, if only the *magnitude* of the response is important, IIR filters usually require much fewer operations and are typically used, so the bulk of FIR filter design work has concentrated on linear phase designs.

¹This content is available online at <<http://cnx.org/content/m12776/1.2/>>.

3.2 FIR Filter Design

3.2.1 Linear Phase Filters²

In general, for $-\pi \leq \omega \leq \pi$

$$H(\omega) = |H(\omega)|e^{-j\theta(\omega)}$$

Strictly speaking, we say $H(\omega)$ is linear phase if

$$H(\omega) = |H(\omega)|e^{-j\omega K}e^{-j\theta_0}$$

Why is this important? A linear phase response gives the *same time delay for ALL frequencies!* (Remember the shift theorem.) This is very desirable in many applications, particularly when the appearance of the time-domain waveform is of interest, such as in an oscilloscope. (see Figure 3.1)

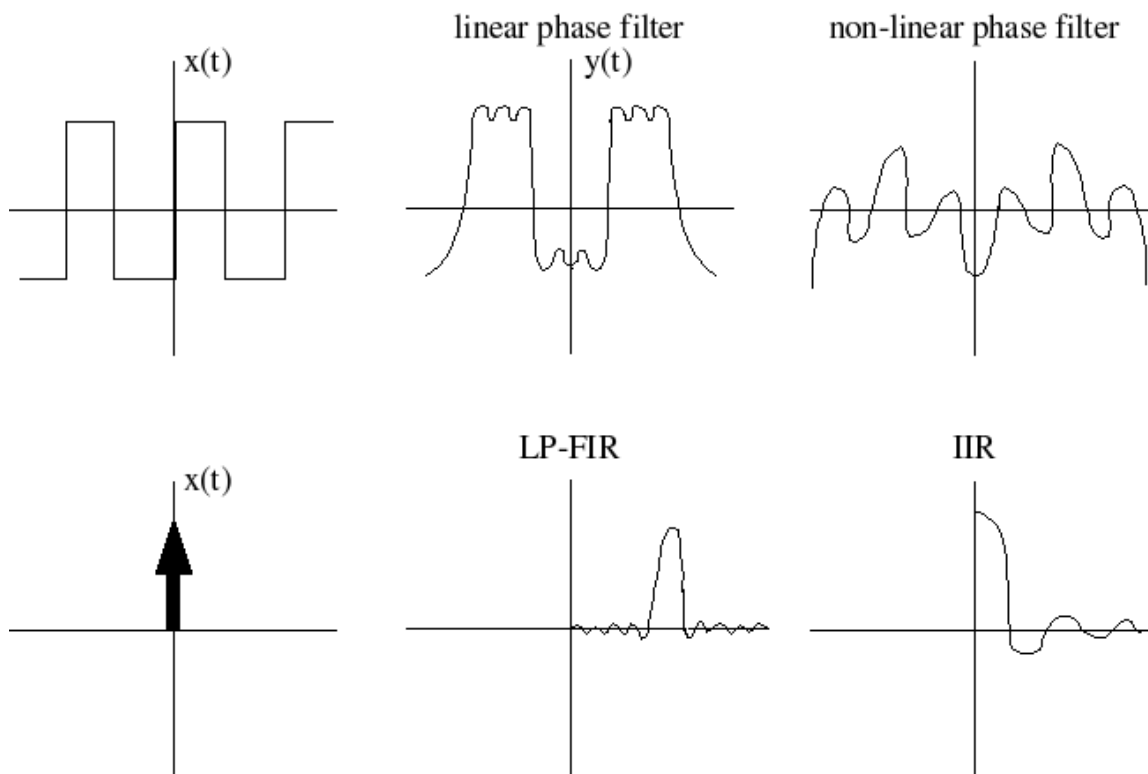


Figure 3.1

3.2.1.1 Restrictions on $h(n)$ to get linear phase

²This content is available online at <http://cnx.org/content/m12802/1.2/>.

$$\begin{aligned}
H(\omega) &= \sum_{n=0}^{M-1} (h(n) e^{-j\omega n}) = h(0) + h(1) e^{-j\omega} + h(2) e^{-j2\omega} + \dots + \\
h(M-1) e^{-j\omega(M-1)} &= e^{-j\omega \frac{M-1}{2}} \left(h(0) e^{j\omega \frac{M-1}{2}} + \dots + h(M-1) e^{-j\omega \frac{M-1}{2}} \right) = \\
e^{-j\omega \frac{M-1}{2}} &\left((h(0) + h(M-1)) \cos\left(\frac{M-1}{2}\omega\right) + (h(1) + h(M-2)) \cos\left(\frac{M-3}{2}\omega\right) + \dots + j((h(0) - h(M-1)) \sin\left(\frac{M-1}{2}\omega\right) + \dots \right)
\end{aligned} \tag{3.1}$$

For linear phase, we require the right side of (3.1) to be $e^{-j\theta_0}$ (real, positive function of ω). For $\theta_0 = 0$, we thus require

$$h(0) + h(M-1) = \text{real number}$$

$$h(0) - h(M-1) = \text{pure imaginary number}$$

$$h(1) + h(M-2) = \text{pure real number}$$

$$h(1) - h(M-2) = \text{pure imaginary number}$$

$$\vdots$$

Thus $h(k) = h^*(M-1-k)$ is a *necessary* condition for the right side of (3.1) to be real valued, for $\theta_0 = 0$.

For $\theta_0 = \frac{\pi}{2}$, or $e^{-j\theta_0} = -j$, we require

$$h(0) + h(M-1) = \text{pure imaginary}$$

$$h(0) - h(M-1) = \text{pure real number}$$

$$\vdots$$

$$\Rightarrow h(k) = -(h^*(M-1-k))$$

Usually, one is interested in filters with *real-valued* coefficients, or see Figure 3.2 and Figure 3.3.

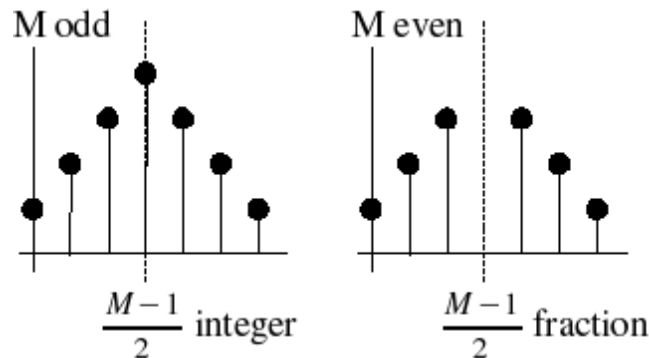


Figure 3.2: $\theta_0 = 0$ (Symmetric Filters). $h(k) = h(M-1-k)$.

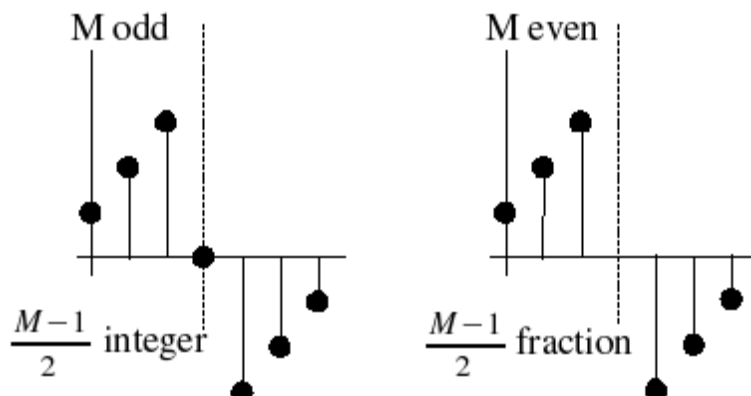


Figure 3.3: $\theta_0 = \frac{\pi}{2}$ (Anti-Symmetric Filters). $h(k) = -(h(M-1-k))$.

Filter design techniques are usually slightly different for each of these four different filter types. We will study the most common case, symmetric-odd length, in detail, and often leave the others for homework or tests or for when one encounters them in practice. Even-symmetric filters are often used; the anti-symmetric filters are rarely used in practice, except for special classes of filters, like differentiators or Hilbert transformers, in which the desired response is anti-symmetric.

So far, we have satisfied the condition that $H(\omega) = A(\omega) e^{-j\theta_0} e^{-j\omega \frac{M-1}{2}}$ where $A(\omega)$ is *real-valued*. However, we have *not* assured that $A(\omega)$ is *non-negative*. In general, this makes the design techniques much more difficult, so most FIR filter design methods actually design filters with **Generalized Linear Phase**: $H(\omega) = A(\omega) e^{-j\omega \frac{M-1}{2}}$, where $A(\omega)$ is *real-valued*, but possible negative. $A(\omega)$ is called the **amplitude of the frequency response**.

EXCUSE: $A(\omega)$ usually goes negative only in the stopband, and the stopband phase response is generally unimportant.

NOTE: $|H(\omega)| = \pm A(\omega) = A(\omega) e^{-j\pi \frac{1}{2}(1 - \text{sign} A(\omega))}$ where $\text{sign} x = \begin{cases} 1 & \text{if } x > 0 \\ -1 & \text{if } x < 0 \end{cases}$

Example 3.1 Lowpass Filter

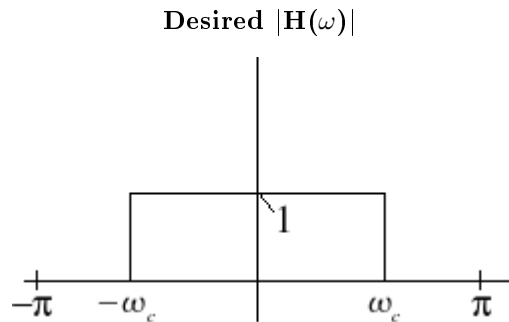
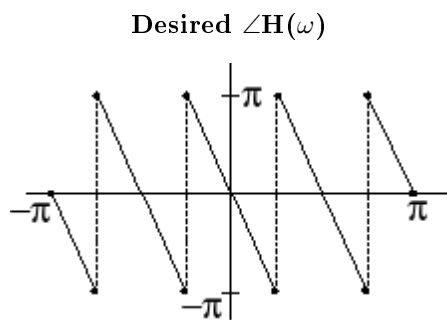
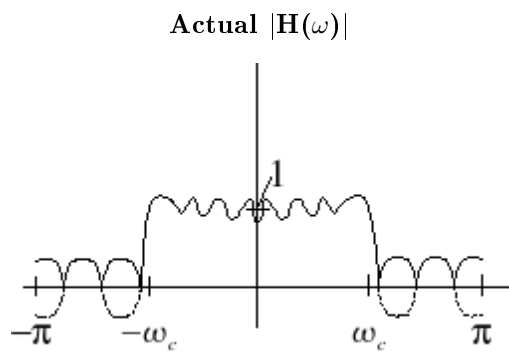


Figure 3.4

Figure 3.5: The slope of each line is $-\left(\frac{M-1}{2}\right)$.Figure 3.6: $A(\omega)$ goes negative.

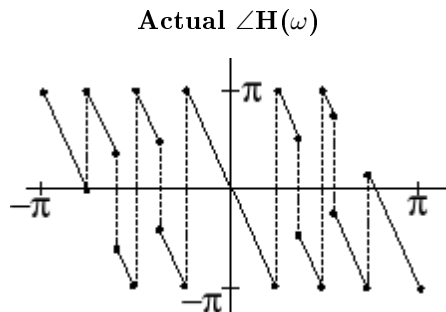


Figure 3.7: 2π phase jumps due to periodicity of phase. π phase jumps due to sign change in $A(\omega)$.

Time-delay introduces generalized linear phase.

NOTE: For odd-length FIR filters, a linear-phase design procedure is equivalent to a zero-phase design procedure followed by an $\frac{M-1}{2}$ -sample delay of the impulse response³. For even-length filters, the delay is non-integer, and the linear phase must be incorporated directly in the desired response!

3.2.2 Window Design Method⁴

The truncate-and-delay design procedure is the simplest and most obvious FIR design procedure.

Exercise 3.1

Is it any Good?

(Solution on p. 223.)

3.2.2.1 L2 optimization criterion

find $\forall n, 0 \leq n \leq M-1 : (h[n])$, maximizing the energy difference between the desired response and the actual response: i.e., find

$$\min_{h[n]} \left\{ \int_{-\pi}^{\pi} (|H_d(\omega) - H(\omega)|)^2 d\omega \right\}$$

by Parseval's relationship⁵

$$\begin{aligned} \min_{h[n]} \left\{ \int_{-\pi}^{\pi} (|H_d(\omega) - H(\omega)|)^2 d\omega \right\} &= 2\pi \sum_{n=-\infty}^{\infty} (|h_d[n] - h[n]|^2) = \\ 2\pi \left(\sum_{n=-\infty}^{-1} (|h_d[n] - h[n]|^2) + \sum_{n=0}^{M-1} (|h_d[n] - h[n]|^2) + \sum_{n=M}^{\infty} (|h_d[n] - h[n]|^2) \right) \end{aligned} \quad (3.2)$$

Since $\forall n, n < 0, n \geq M : (h[n])$ this becomes

$$\begin{aligned} \min_{h[n]} \left\{ \int_{-\pi}^{\pi} (|H_d(\omega) - H(\omega)|)^2 d\omega \right\} &= \sum_{n=-\infty}^{-1} (|h_d[n]|^2) + \\ \sum_{n=0}^{M-1} (|h[n] - h_d[n]|^2) + \sum_{n=M}^{\infty} (|h_d[n]|^2) \end{aligned}$$

³"Impulse Response of a Linear System" <<http://cnx.org/content/m12041/latest/>>

⁴This content is available online at <<http://cnx.org/content/m12790/1.2/>>.

⁵"Parseval's Theorem" <<http://cnx.org/content/m0047/latest/>>

NOTE: $h[n]$ has no influence on the first and last sums.

The best we can do is let

$$h[n] = \begin{cases} h_d[n] & \text{if } 0 \leq n \leq M-1 \\ 0 & \text{if else} \end{cases}$$

Thus $h[n] = h_d[n] w[n]$,

$$w[n] = \begin{cases} 1 & \text{if } 0 \leq n \leq M-1 \\ 0 & \text{if else} \end{cases}$$

is *optimal* in a least-total-squared-error (L_2 , or energy) sense!

Exercise 3.2

(Solution on p. 223.)

Why, then, is this design often considered undesirable?

For desired spectra with discontinuities, the least-square designs are poor in a minimax (worst-case, or L_∞) error sense.

3.2.2.2 Window Design Method

Apply a more gradual truncation to reduce "ringing" (Gibb's Phenomenon⁶)

$$\forall n, 0 \leq n \leq M-1, h[n] = h_d[n] w[n]$$

NOTE: $H(\omega) = H_d(\omega) * W(\omega)$

The window design procedure (except for the boxcar window) is ad-hoc and not optimal in any usual sense. However, it is very simple, so it is sometimes used for "quick-and-dirty" designs of if the error criterion is itself heuristic.

3.2.3 Frequency Sampling Design Method for FIR filters⁷

Given a desired frequency response, the frequency sampling design method designs a filter with a frequency response *exactly* equal to the desired response at a particular set of frequencies ω_k .

Procedure

$$\forall k, k = [0, 1, \dots, N-1] : \left(H_d(\omega_k) = \sum_{n=0}^{M-1} \left(h(n) e^{-j\omega_k n} \right) \right) \quad (3.3)$$

NOTE: Desired Response must include linear phase shift (if linear phase is desired)

Exercise 3.3

(Solution on p. 223.)

What is $H_d(\omega)$ for an ideal lowpass filter, cutoff at ω_c ?

NOTE: This set of linear equations can be written in matrix form

$$H_d(\omega_k) = \sum_{n=0}^{M-1} \left(h(n) e^{-j\omega_k n} \right) \quad (3.4)$$

⁶"Gibbs's Phenomena" <<http://cnx.org/content/m10092/latest/>>

⁷This content is available online at <<http://cnx.org/content/m12789/1.2/>>.

$$\begin{pmatrix} H_d(\omega_0) \\ H_d(\omega_1) \\ \vdots \\ H_d(\omega_{N-1}) \end{pmatrix} = \begin{pmatrix} e^{-(j\omega_0 0)} & e^{-(j\omega_0 1)} & \dots & e^{-(j\omega_0 (M-1))} \\ e^{-(j\omega_1 0)} & e^{-(j\omega_1 1)} & \dots & e^{-(j\omega_1 (M-1))} \\ \vdots & \vdots & \vdots & \vdots \\ e^{-(j\omega_{M-1} 0)} & e^{-(j\omega_{M-1} 1)} & \dots & e^{-(j\omega_{M-1} (M-1))} \end{pmatrix} \begin{pmatrix} h(0) \\ h(1) \\ \vdots \\ h(M-1) \end{pmatrix} \quad (3.5)$$

or

$$H_d = W\mathbf{h}$$

So

$$\mathbf{h} = W^{-1}H_d \quad (3.6)$$

NOTE: W is a square matrix for $N = M$, and invertible as long as $\omega_i \neq \omega_j + 2\pi l$, $i \neq j$

3.2.3.1 Important Special Case

What if the frequencies are equally spaced between 0 and 2π , i.e. $\omega_k = \frac{2\pi k}{M} + \alpha$

Then

$$H_d(\omega_k) = \sum_{n=0}^{M-1} \left(h(n) e^{-(j\frac{2\pi kn}{M})} e^{-(j\alpha n)} \right) = \sum_{n=0}^{M-1} \left(\left(h(n) e^{-(j\alpha n)} \right) e^{-(j\frac{2\pi kn}{M})} \right) = \text{DFT!}$$

so

$$h(n) e^{-(j\alpha n)} = \frac{1}{M} \sum_{k=0}^{M-1} \left(H_d(\omega_k) e^{+j\frac{2\pi kn}{M}} \right)$$

or

$$h[n] = \frac{e^{j\alpha n}}{M} \sum_{k=0}^{M-1} \left(H_d[\omega_k] e^{j\frac{2\pi kn}{M}} \right) = e^{j\alpha n} \text{IDFT}[H_d[\omega_k]]$$

3.2.3.2 Important Special Case #2

$h[n]$ symmetric, linear phase, and has real coefficients. Since $h[n] = h[-1]$, there are only $\frac{M}{2}$ degrees of freedom, and only $\frac{M}{2}$ linear equations are required.

$$\begin{aligned} H[\omega_k] &= \sum_{n=0}^{M-1} (h[n] e^{-(j\omega_k n)}) \\ &= \begin{cases} \sum_{n=0}^{\frac{M}{2}-1} (h[n] (e^{-(j\omega_k n)} + e^{-(j\omega_k (M-n-1))})) & \text{if } M \text{ even} \\ \sum_{n=0}^{M-\frac{3}{2}} \left(+h[n] (e^{-(j\omega_k n)} + e^{-(j\omega_k (M-n-1))}) \left(h\left[\frac{M-1}{2}\right] e^{-(j\omega_k \frac{M-1}{2})} \right) \right) & \text{if } M \text{ odd} \end{cases} \\ &= \begin{cases} e^{-(j\omega_k \frac{M-1}{2})} 2 \sum_{n=0}^{\frac{M}{2}-1} (h[n] \cos(\omega_k (\frac{M-1}{2} - n))) & \text{if } M \text{ even} \\ e^{-(j\omega_k \frac{M-1}{2})} 2 \sum_{n=0}^{M-\frac{3}{2}} (h[n] \cos(\omega_k (\frac{M-1}{2} - n)) + h[\frac{M-1}{2}]) & \text{if } M \text{ odd} \end{cases} \end{aligned} \quad (3.7)$$

Removing linear phase from both sides yields

$$A(\omega_k) = \begin{cases} 2 \sum_{n=0}^{\frac{M}{2}-1} (h[n] \cos(\omega_k (\frac{M-1}{2} - n))) & \text{if } M \text{ even} \\ 2 \sum_{n=0}^{M-\frac{3}{2}} (h[n] \cos(\omega_k (\frac{M-1}{2} - n)) + h[\frac{M-1}{2}]) & \text{if } M \text{ odd} \end{cases}$$

Due to symmetry of response for real coefficients, only $\frac{M}{2}$ ω_k on $\omega \in [0, \pi)$ need be specified, with the frequencies $-\omega_k$ thereby being implicitly defined also. Thus we have $\frac{M}{2}$ *real-valued* simultaneous linear equations to solve for $h[n]$.

3.2.3.2.1 Special Case 2a

$h[n]$ symmetric, odd length, linear phase, real coefficients, and ω_k equally spaced: $\forall k, 0 \leq k \leq M-1$:
 $(\omega_k = \frac{n\pi k}{M})$

$$\begin{aligned} h[n] &= IDFT[H_d(\omega_k)] \\ &= \frac{1}{M} \sum_{k=0}^{M-1} \left(A(\omega_k) e^{-j \left(\frac{2\pi k}{M} \right) \frac{M-1}{2}} e^{j \frac{2\pi n k}{M}} \right) \\ &= \frac{1}{M} \sum_{k=0}^{M-1} \left(A(k) e^{j \left(\frac{2\pi k}{M} \left(n - \frac{M-1}{2} \right) \right)} \right) \end{aligned} \quad (3.8)$$

To yield real coefficients, $A(\omega)$ must be symmetric

$$A(\omega) = A(-\omega) \Rightarrow A[k] = A[M-k]$$

$$\begin{aligned} h[n] &= \frac{1}{M} \left(A(0) + \sum_{k=1}^{\frac{M-1}{2}} \left(A[k] \left(e^{j \frac{2\pi k}{M} \left(n - \frac{M-1}{2} \right)} + e^{-j \frac{2\pi k}{M} \left(n - \frac{M-1}{2} \right)} \right) \right) \right) \\ &= \frac{1}{M} \left(A(0) + 2 \sum_{k=1}^{\frac{M-1}{2}} \left(A[k] \cos \left(\frac{2\pi k}{M} \left(n - \frac{M-1}{2} \right) \right) \right) \right) \\ &= \frac{1}{M} \left(A(0) + 2 \sum_{k=1}^{\frac{M-1}{2}} \left(A[k] (-1)^k \cos \left(\frac{2\pi k}{M} \left(n + \frac{1}{2} \right) \right) \right) \right) \end{aligned} \quad (3.9)$$

Similar equations exist for even lengths, anti-symmetric, and $\alpha = \frac{1}{2}$ filter forms.

3.2.3.3 Comments on frequency-sampled design

This method is simple conceptually and very efficient for equally spaced samples, since $h[n]$ can be computed using the IDFT.

$H(\omega)$ for a frequency sampled design goes *exactly* through the sample points, but it may be very far off from the desired response for $\omega \neq \omega_k$. This is the main problem with frequency sampled design.

Possible solution to this problem: specify more frequency samples than degrees of freedom, and minimize the total error in the frequency response at all of these samples.

3.2.3.4 Extended frequency sample design

For the samples $H(\omega_k)$ where $0 \leq k \leq M-1$ and $N > M$, find $h[n]$, where $0 \leq n \leq M-1$ minimizing $\|H_d(\omega_k) - H(\omega_k)\|$

For $\|l\|_\infty$ norm, this becomes a linear programming problem (standard packages available!)

Here we will consider the $\|l\|_2$ norm.

To minimize the $\|l\|_2$ norm; that is, $\sum_{n=0}^{N-1} (|H_d(\omega_k) - H(\omega_k)|)$, we have an overdetermined set of linear equations:

$$\begin{pmatrix} e^{-j\omega_0 0} & \dots & e^{-j\omega_0(M-1)} \\ \vdots & \vdots & \vdots \\ e^{-j\omega_{N-1} 0} & \dots & e^{-j\omega_{N-1}(M-1)} \end{pmatrix} \mathbf{h} = \begin{pmatrix} H_d(\omega_0) \\ H_d(\omega_1) \\ \vdots \\ H_d(\omega_{N-1}) \end{pmatrix}$$

or

$$W\mathbf{h} = H_d$$

The minimum error norm solution is well known to be $\mathbf{h} = (\overline{W}W)^{-1} \overline{W}H_d$; $(\overline{W}W)^{-1} \overline{W}$ is well known as the pseudo-inverse matrix.

NOTE: Extended frequency sampled design discourages radical behavior of the frequency response between samples for sufficiently closely spaced samples. However, the actual frequency response may no longer pass exactly through any of the $H_d(\omega_k)$.

3.2.4 Parks-McClellan FIR Filter Design⁸

The approximation tolerances for a filter are very often given in terms of the maximum, or worst-case, deviation within frequency bands. For example, we might wish a lowpass filter in a (16-bit) CD player to have no more than $\frac{1}{2}$ -bit deviation in the pass and stop bands.

$$H(\omega) = \begin{cases} 1 - \frac{1}{2^{17}} \leq |H(\omega)| \leq 1 + \frac{1}{2^{17}} & \text{if } |\omega| \leq \omega_p \\ \frac{1}{2^{17}} \geq |H(\omega)| & \text{if } \omega_s \leq |\omega| \leq \pi \end{cases}$$

The Parks-McClellan filter design method efficiently designs linear-phase FIR filters that are optimal in terms of worst-case (minimax) error. Typically, we would like to have the shortest-length filter achieving these specifications. Figure Figure 3.8 illustrates the amplitude frequency response of such a filter.

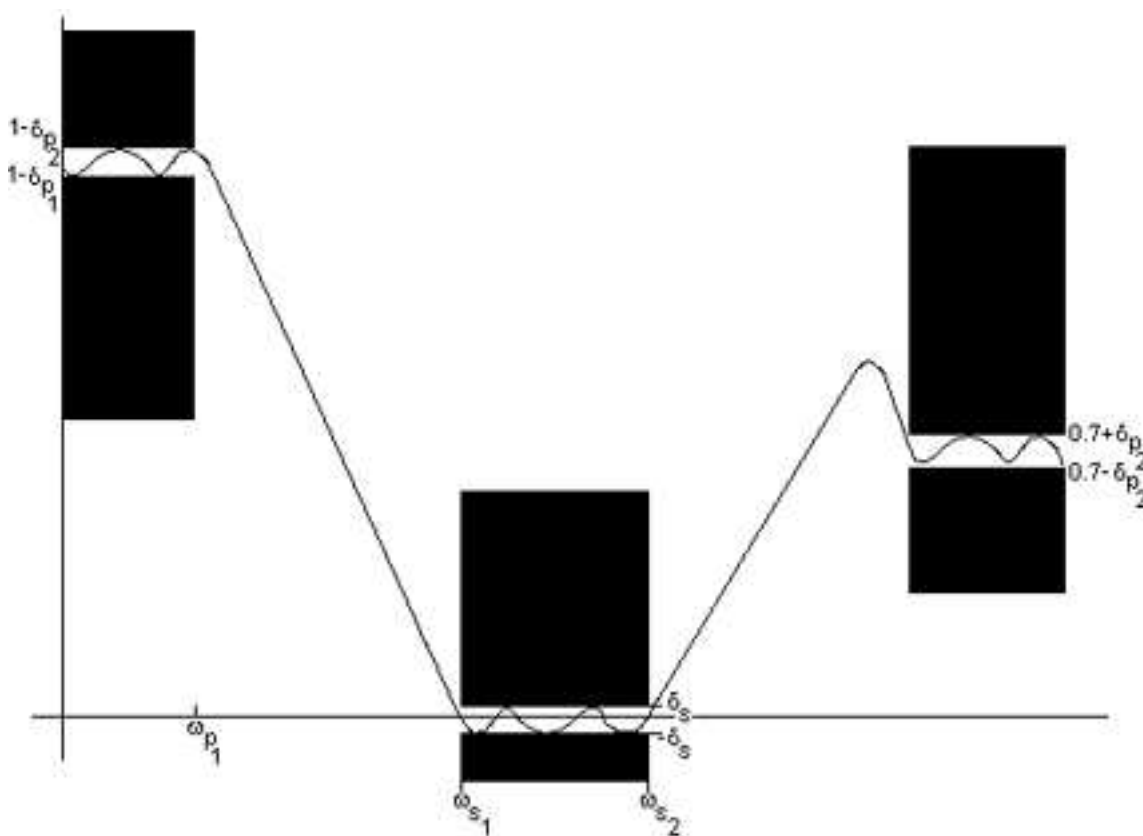


Figure 3.8: The black boxes on the left and right are the passbands, the black boxes in the middle represent the stop band, and the space between the boxes are the transition bands. Note that overshoots may be allowed in the transition bands.

Exercise 3.4

Must there be a transition band?

(Solution on p. 223.)

⁸This content is available online at <<http://cnx.org/content/m12799/1.3/>>.

3.2.4.1 Formal Statement of the L^∞ (Minimax) Design Problem

For a given filter length (M) and type (odd length, symmetric, linear phase, for example), and a relative error weighting function $W(\omega)$, find the filter coefficients minimizing the maximum error

$$\underset{\mathbf{h}}{\operatorname{argmin}} \underset{\omega \in F}{\operatorname{argmax}} |E(\omega)| = \underset{\mathbf{h}}{\operatorname{argmin}} \|E(\omega)\|_\infty$$

where

$$E(\omega) = W(\omega)(H_d(\omega) - H(\omega))$$

and F is a compact subset of $\omega \in [0, \pi]$ (i.e., all ω in the passbands and stop bands).

NOTE: Typically, we would often rather specify $\|E(\omega)\|_\infty \leq \delta$ and minimize over M and \mathbf{h} ; however, the design techniques minimize δ for a given M . One then repeats the design procedure for different M until the minimum M satisfying the requirements is found.

We will discuss in detail the design only of odd-length symmetric linear-phase FIR filters. Even-length and anti-symmetric linear phase FIR filters are essentially the same except for a slightly different implicit weighting function. For arbitrary phase, exactly optimal design procedures have only recently been developed (1990).

3.2.4.2 Outline of L^∞ Filter Design

The Parks-McClellan method adopts an indirect method for finding the minimax-optimal filter coefficients.

1. Using results from Approximation Theory, simple conditions for determining whether a given filter is L^∞ (minimax) optimal are found.
2. An iterative method for finding a filter which satisfies these conditions (and which is thus optimal) is developed.

That is, the L^∞ filter design problem is actually solved *indirectly*.

3.2.4.3 Conditions for L^∞ Optimality of a Linear-phase FIR Filter

All conditions are based on Chebyshev's "Alternation Theorem," a mathematical fact from polynomial approximation theory.

3.2.4.3.1 Alternation Theorem

Let F be a compact subset on the real axis x , and let $P(x)$ be an L th-order polynomial

$$P(x) = \sum_{k=0}^L (a_k x^k)$$

Also, let $D(x)$ be a desired function of x that is continuous on F , and $W(x)$ a positive, continuous weighting function on F . Define the error $E(x)$ on F as

$$E(x) = W(x)(D(x) - P(x))$$

and

$$\|E(x)\|_\infty = \underset{x \in F}{\operatorname{argmax}} |E(x)|$$

A necessary and sufficient condition that $P(x)$ is the unique L th-order polynomial minimizing $\|E(x)\|_\infty$ is that $E(x)$ exhibits *at least* $L+2$ "alternations;" that is, there must exist at least $L+2$ values of x , $x_k \in F$, $k = [0, 1, \dots, L+1]$, such that $x_0 < x_1 < \dots < x_{L+2}$ and such that $E(x_k) = -(E(x_{k+1})) = \pm(\|E\|_\infty)$

Exercise 3.5

(Solution on p. 223.)

What does this have to do with linear-phase filter design?

3.2.4.4 Optimality Conditions for Even-length Symmetric Linear-phase Filters

For M even,

$$A(\omega) = \sum_{n=0}^L \left(h(L-n) \cos \left(\omega \left(n + \frac{1}{2} \right) \right) \right)$$

where $L = \frac{M}{2} - 1$ Using the trigonometric identity $\cos(\alpha + \beta) = \cos(\alpha - \beta) + 2\cos(\alpha)\cos(\beta)$ to pull out the $\frac{\omega}{2}$ term and then using the other trig identities (p. 223), it can be shown that $A(\omega)$ can be written as

$$A(\omega) = \cos\left(\frac{\omega}{2}\right) \sum_{k=0}^L (\alpha_k \cos^k(\omega))$$

Again, this is a polynomial in $x = \cos(\omega)$, except for a weighting function out in front.

$$\begin{aligned} E(\omega) &= W(\omega) (A_d(\omega) - A(\omega)) \\ &= W(\omega) \left(A_d(\omega) - \cos\left(\frac{\omega}{2}\right) P(\omega) \right) \\ &= W(\omega) \cos\left(\frac{\omega}{2}\right) \left(\frac{A_d(\omega)}{\cos\left(\frac{\omega}{2}\right)} - P(\omega) \right) \end{aligned} \tag{3.10}$$

which implies

$$E(x) = W'(x) (A_d'(x) - P(x)) \tag{3.11}$$

where

$$W'(x) = W \left((\cos(x))^{-1} \right) \cos \left(\frac{1}{2} (\cos(x))^{-1} \right)$$

and

$$A_d'(x) = \frac{A_d \left((\cos(x))^{-1} \right)}{\cos \left(\frac{1}{2} (\cos(x))^{-1} \right)}$$

Again, this is a polynomial approximation problem, so the alternation theorem holds. If $E(\omega)$ has at least $L + 2 = \frac{M}{2} + 1$ alternations, the even-length symmetric filter is optimal in an L^∞ sense.

The prototypical filter design problem:

$$W = \begin{cases} 1 & \text{if } |\omega| \leq \omega_p \\ \frac{\delta_s}{\delta_p} & \text{if } |\omega_s| \leq |\omega| \end{cases}$$

See Figure 3.9.

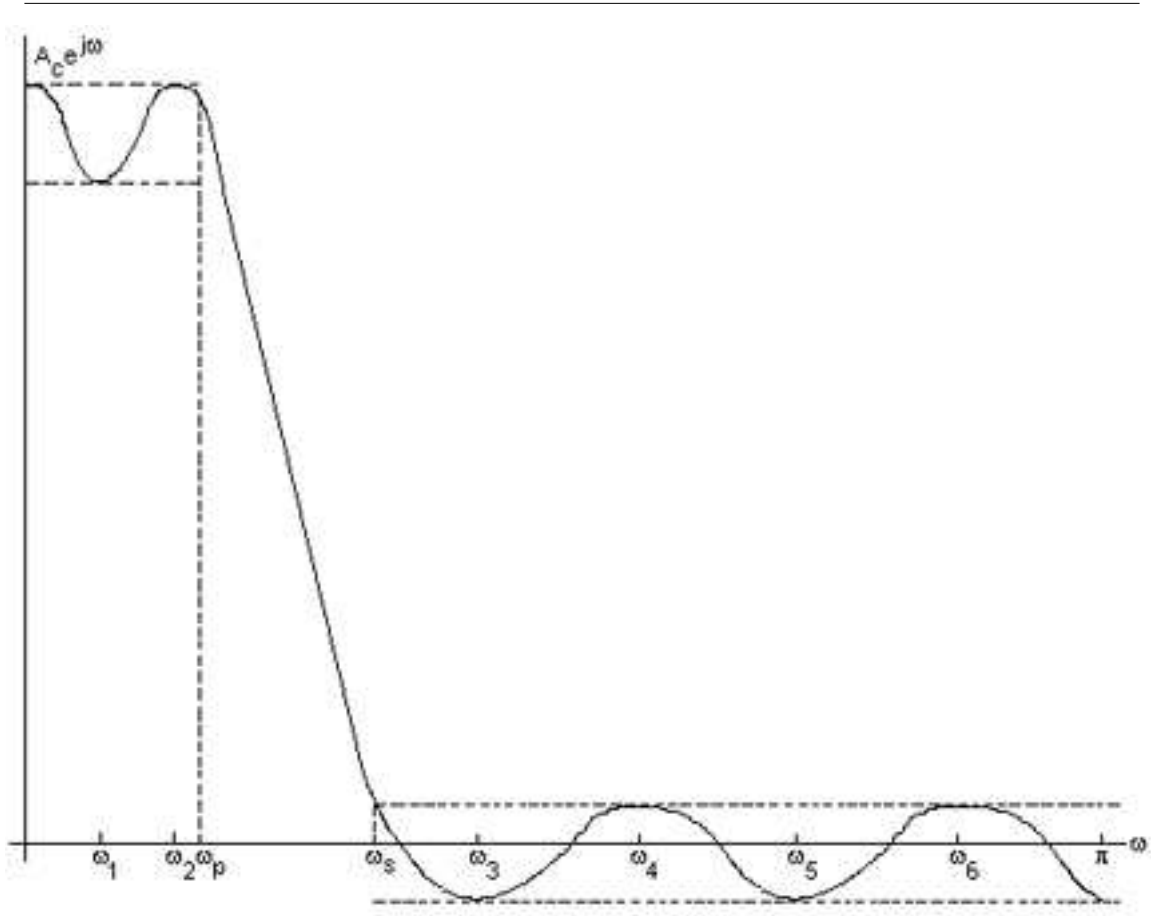


Figure 3.9

3.2.4.5 L_∞ Optimal Lowpass Filter Design Lemma

1. The maximum possible number of alternations for a lowpass filter is $L + 3$: The proof is that the extrema of a polynomial occur only where the derivative is zero: $\frac{\partial}{\partial x} P(x) = 0$. Since $P'(x)$ is an $(L - 1)$ th-order polynomial, it can have at *most* $L - 1$ zeros. *However*, the mapping $x = \cos(\omega)$ implies that $\frac{\partial}{\partial \omega} A(\omega) = 0$ at $\omega = 0$ and $\omega = \pi$, for two more possible alternation points. *Finally*, the band edges can also be alternations, for a total of $L - 1 + 2 + 2 = L + 3$ possible alternations.
2. There must be an alternation at either $\omega = 0$ or $\omega = \pi$.
3. Alternations must occur at ω_p and ω_s . See Figure 3.9.
4. The filter must be equiripple except at possibly $\omega = 0$ or $\omega = \pi$. Again see Figure 3.9.

NOTE: The alternation theorem doesn't directly suggest a method for computing the optimal filter. It simply tells us how to recognize that a filter *is* optimal, or *isn't* optimal. What we need is an intelligent way of guessing the optimal filter coefficients.

In matrix form, these $L + 2$ simultaneous equations become

$$\begin{pmatrix} 1 & \cos(\omega_0) & \cos(2\omega_0) & \dots & \cos(L\omega_0) & \frac{1}{W(\omega_0)} \\ 1 & \cos(\omega_1) & \cos(2\omega_1) & \dots & \cos(L\omega_1) & \frac{-1}{W(\omega_1)} \\ \vdots & \vdots & \ddots & \dots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \dots & \ddots & \vdots \\ 1 & \cos(\omega_{L+1}) & \cos(2\omega_{L+1}) & \dots & \cos(L\omega_{L+1}) & \frac{\pm 1}{W(\omega_{L+1})} \end{pmatrix} \begin{pmatrix} h(L) \\ h(L-1) \\ \vdots \\ h(1) \\ h(0) \\ \delta \end{pmatrix} = \begin{pmatrix} A_d(\omega_0) \\ A_d(\omega_1) \\ \vdots \\ \vdots \\ \vdots \\ A_d(\omega_{L+1}) \end{pmatrix}$$

or

$$W \begin{pmatrix} \mathbf{h} \\ \delta \end{pmatrix} = A_d$$

So, for the given set of $L + 2$ extremal frequencies, we can solve for \mathbf{h} and δ via $(\mathbf{h}, \delta)^T = W^{-1}A_d$. Using the FFT, we can compute $A(\omega)$ of $h(n)$, on a dense set of frequencies. If the old ω_k are, in fact the extremal locations of $A(\omega)$, then the alternation theorem is satisfied and $h(n)$ is **optimal**. If not, repeat the process with the new extremal locations.

3.2.4.6 Computational Cost

$O(L^3)$ for the matrix inverse and $N \log_2 N$ for the FFT ($N \geq 32L$, typically), *per iteration!*

This method is expensive computationally due to the matrix inverse.

A more efficient variation of this method was developed by Parks and McClellan (1972), and is based on the Remez exchange algorithm. To understand the Remez exchange algorithm, we first need to understand Lagrange Interpolation.

Now $A(\omega)$ is an L th-order polynomial in $x = \cos(\omega)$, so Lagrange interpolation can be used to *exactly* compute $A(\omega)$ from $L + 1$ samples of $A(\omega_k)$, $k = [0, 1, 2, \dots, L]$.

Thus, given a set of extremal frequencies and knowing δ , samples of the amplitude response $A(\omega)$ can be computed *directly* from the

$$A(\omega_k) = \frac{(-1)^{k+1}}{W(\omega_k)} \delta + A_d(\omega_k) \quad (3.12)$$

without solving for the filter coefficients!

This leads to computational savings!

Note that (3.12) is a set of $L + 2$ simultaneous equations, which can be solved for δ to obtain (Rabiner, 1975)

$$\delta = \frac{\sum_{k=0}^{L+1} (\gamma_k A_d(\omega_k))}{\sum_{k=0}^{L+1} \left(\frac{(-1)^{k+1} \gamma_k}{W(\omega_k)} \right)} \quad (3.13)$$

where

$$\gamma_k = \prod_{\substack{i=0 \\ i \neq k}}^{L+1} \left(\frac{1}{\cos(\omega_k) - \cos(\omega_i)} \right)$$

The result is the Parks-McClellan FIR filter design method, which is simply an application of the Remez exchange algorithm to the filter design problem. See Figure 3.10.

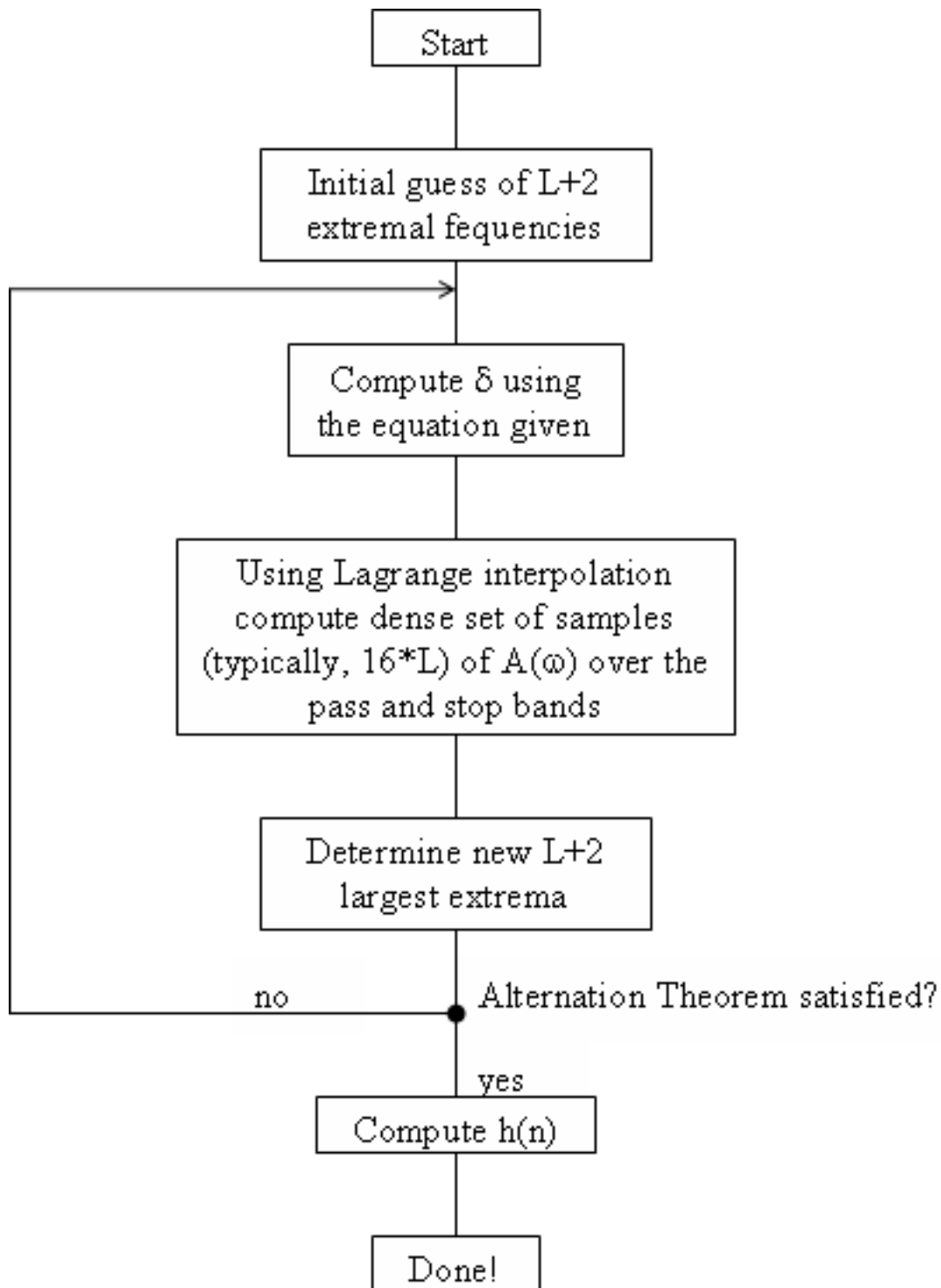


Figure 3.10: The initial guess of extremal frequencies is usually equally spaced in the band. Computing δ costs $O(L^2)$. Using Lagrange interpolation costs $O(16LL) \approx O(16L^2)$. Computing $h(n)$ costs $O(L^3)$, but it is only done once!

The cost per iteration is $O(16L^2)$, as opposed to $O(L^3)$; much more efficient for large L . Can also interpolate to DFT sample frequencies, take inverse FFT to get corresponding filter coefficients, and zeropad and take longer FFT to efficiently interpolate.

3.2.5 Lagrange Interpolation⁹

Lagrange's interpolation method is a simple and clever way of finding the unique L th-order polynomial that exactly passes through $L + 1$ distinct samples of a signal. Once the polynomial is known, its value can easily be interpolated at any point using the polynomial equation. Lagrange interpolation is useful in many applications, including Parks-McClellan FIR Filter Design (Section 3.2.4).

3.2.5.1 Lagrange interpolation formula

Given an L th-order polynomial

$$P(x) = a_0 + a_1x + \dots + a_Lx^L = \sum_{k=0}^L (a_kx^k)$$

and $L + 1$ values of $P(x_k)$ at different x_k , $k \in \{0, 1, \dots, L\}$, $x_i \neq x_j$, $i \neq j$, the polynomial can be written as

$$P(x) = \sum_{k=0}^L \left(P(x_k) \frac{(x - x_1)(x - x_2) \dots (x - x_{k-1})(x - x_{k+1}) \dots (x - x_L)}{(x_k - x_1)(x_k - x_2) \dots (x_k - x_{k-1})(x_k - x_{k+1}) \dots (x_k - x_L)} \right)$$

The value of this polynomial at other x can be computed via substitution into this formula, or by expanding this formula to determine the polynomial coefficients a_k in standard form.

3.2.5.2 Proof

Note that for each term in the Lagrange interpolation formula above,

$$\prod_{i=0, i \neq k}^L \left(\frac{x - x_i}{x_k - x_i} \right) = \begin{cases} 1 & \text{if } x = x_k \\ 0 & \text{if } x = x_j \wedge j \neq k \end{cases}$$

and that it is an L th-order polynomial in x . The Lagrange interpolation formula is thus exactly equal to $P(x_k)$ at all x_k , and as a sum of L th-order polynomials is itself an L th-order polynomial.

It can be shown that the Vandermonde matrix¹⁰

$$\begin{pmatrix} 1 & x_0 & x_0^2 & \dots & x_0^L \\ 1 & x_1 & x_1^2 & \dots & x_1^L \\ 1 & x_2 & x_2^2 & \dots & x_2^L \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_L & x_L^2 & \dots & x_L^L \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \\ a_2 \\ \vdots \\ a_L \end{pmatrix} = \begin{pmatrix} P(x_0) \\ P(x_1) \\ P(x_2) \\ \vdots \\ P(x_L) \end{pmatrix}$$

has a non-zero determinant and is thus invertible, so the L th-order polynomial passing through all $L + 1$ sample points x_j is unique. Thus the Lagrange polynomial expressions, as an L th-order polynomial passing through the $L + 1$ sample points, must be the unique $P(x)$.

⁹This content is available online at <<http://cnx.org/content/m12812/1.2/>>.

¹⁰http://en.wikipedia.org/wiki/Vandermonde_matrix

3.3 IIR Filter Design

3.3.1 Overview of IIR Filter Design¹¹

3.3.1.1 IIR Filter

$$y(n) = - \left(\sum_{k=1}^{M-1} (a_k y(n-k)) \right) + \sum_{k=0}^{M-1} (b_k x(n-k))$$

$$H(z) = \frac{b_0 + b_1 z^{-1} + b_2 z^{-2} + \dots + b_M z^{-M}}{1 + a_1 z^{-1} + a_2 z^{-2} + \dots + a_M z^{-M}}$$

3.3.1.2 IIR Filter Design Problem

Choose $\{a_i\}$, $\{b_i\}$ to best approximate some desired $|H_d(w)|$ or, (occasionally), $H_d(w)$.

As before, different design techniques will be developed for different approximation criteria.

3.3.1.3 Outline of IIR Filter Design Material

- **Bilinear Transform** - Maps $\|L\|_\infty$ optimal (and other) *analog* filter designs to $\|L\|_\infty$ optimal digital IIR filter designs.
- **Prony's Method** - Quasi- $\|L\|_2$ optimal method for time-domain fitting of a desired impulse response (*ad hoc*).
- **Lp Optimal Design** - $\|L\|_p$ optimal filter design ($1 < p < \infty$) using non-linear optimization techniques.

3.3.1.4 Comments on IIR Filter Design Methods

The bilinear transform method is used to design "typical" $\|L\|_\infty$ magnitude optimal filters. The $\|L\|_p$ optimization procedures are used to design filters for which classical analog prototype solutions don't exist. The program by Deczky (*DSP Programs Book*, IEEE Press) is widely used. Prony/Linear Prediction techniques are used often to obtain initial guesses, and are almost exclusively used in data modeling, system identification, and most applications involving the fitting of real data (for example, the impulse response of an unknown filter).

3.3.2 Prototype Analog Filter Design¹²

3.3.2.1 Analog Filter Design

Laplace transform:

$$H(s) = \int_{-\infty}^{\infty} h_a(t) e^{-(st)} dt$$

Note that the continuous-time Fourier transform (Section 1.7) is $H(j\lambda)$ (the Laplace transform evaluated on the imaginary axis).

Since the early 1900's, there has been a lot of research on designing analog filters of the form

$$H(s) = \frac{b_0 + b_1 s + b_2 s^2 + \dots + b_M s^M}{1 + a_1 s + a_2 s^2 + \dots + a_M s^M}$$

¹¹This content is available online at <<http://cnx.org/content/m12758/1.2/>>.

¹²This content is available online at <<http://cnx.org/content/m12763/1.2/>>.

A causal¹³ IIR filter *cannot* have linear phase (no possible symmetry point), and design work for analog filters has concentrated on designing filters with equiripple ($\|L\|_\infty$) *magnitude* responses. These design problems have been solved. We will not concern ourselves here with the design of the analog prototype filters, only with how these designs are mapped to discrete-time while preserving optimality.

An analog filter with *real* coefficients must have a magnitude response of the form

$$(|H(\lambda)|)^2 = B(\lambda^2)$$

$$\begin{aligned} H(j\lambda) \overline{H(j\lambda)} &= \frac{b_0 + b_1 j\lambda + b_2 (j\lambda)^2 + b_3 (j\lambda)^3 + \dots}{1 + a_1 j\lambda + a_2 (j\lambda)^2 + \dots} \overline{H(j\lambda)} \\ &= \frac{b_0 - b_2 \lambda^2 + b_4 \lambda^4 + \dots + j\lambda(b_1 - b_3 \lambda^2 + b_5 \lambda^4 + \dots)}{1 - a_2 \lambda^2 + a_4 \lambda^4 + \dots + j\lambda(a_1 - a_3 \lambda^2 + a_5 \lambda^4 + \dots)} \frac{b_0 - b_2 \lambda^2 + b_4 \lambda^4 + \dots + j\lambda(b_1 - b_3 \lambda^2 + b_5 \lambda^4 + \dots)}{1 - a_2 \lambda^2 + a_4 \lambda^4 + \dots + j\lambda(a_1 - a_3 \lambda^2 + a_5 \lambda^4 + \dots)} \\ &= \frac{(b_0 - b_2 \lambda^2 + b_4 \lambda^4 + \dots)^2 + \lambda^2 (b_1 - b_3 \lambda^2 + b_5 \lambda^4 + \dots)^2}{(1 - a_2 \lambda^2 + a_4 \lambda^4 + \dots)^2 + \lambda^2 (a_1 - a_3 \lambda^2 + a_5 \lambda^4 + \dots)^2} \\ &= B(\lambda^2) \end{aligned} \quad (3.14)$$

Let $s = j\lambda$, note that the poles and zeros of $B(-s^2)$ are symmetric around *both* the real and imaginary axes: that is, a pole at p_1 implies poles at p_1 , \bar{p}_1 , $-p_1$, and $-(\bar{p}_1)$, as seen in Figure 3.11 (s-plane).

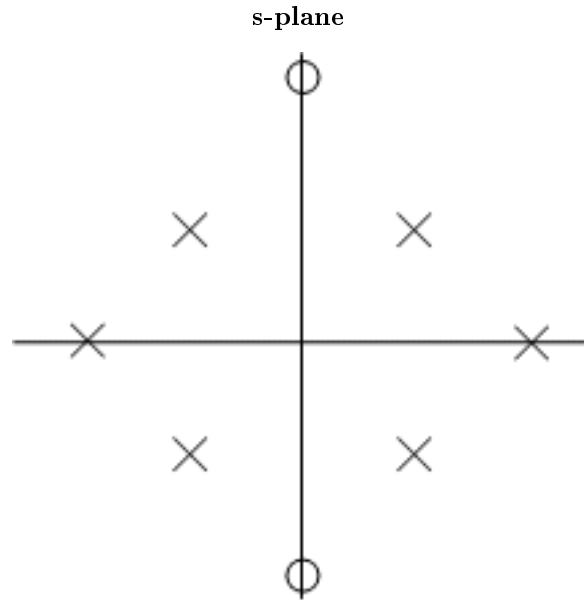


Figure 3.11

Recall that an analog filter is stable and causal if all the poles are in the left half-plane, LHP, and is **minimum phase** if all zeros and poles are in the LHP.

$s = j\lambda$: $B(\lambda^2) = B(-s^2) = H(s)H(-s) = H(j\lambda)H(-(j\lambda)) = H(j\lambda)\overline{H(j\lambda)}$ we can factor $B(-s^2)$ into $H(s)H(-s)$, where $H(s)$ has the left half plane poles and zeros, and $H(-s)$ has the RHP poles and zeros.

$(|H(s)|)^2 = H(s)H(-s)$ for $s = j\lambda$, so $H(s)$ has the magnitude response $B(\lambda^2)$. The trick to analog filter design is to design a good $B(\lambda^2)$, then factor this to obtain a filter with that *magnitude* response.

¹³"Properties of Systems": Section Causality <<http://cnx.org/content/m2102/latest/#causality>>

The traditional analog filter designs all take the form $B(\lambda^2) = (|H(\lambda)|)^2 = \frac{1}{1+F(\lambda^2)}$, where F is a rational function in λ^2 .

Example 3.2

$$B(\lambda^2) = \frac{2 + \lambda^2}{1 + \lambda^4}$$

$$B(-s^2) = \frac{2 - s^2}{1 + s^4} = \frac{(\sqrt{2} - s)(\sqrt{2} + s)}{(s + \alpha)(s - \alpha)(s + \bar{\alpha})(s - \bar{\alpha})}$$

where $\alpha = \frac{1+j}{\sqrt{2}}$.

NOTE: Roots of $1 + s^N$ are N points equally spaced around the unit circle (Figure 3.12).

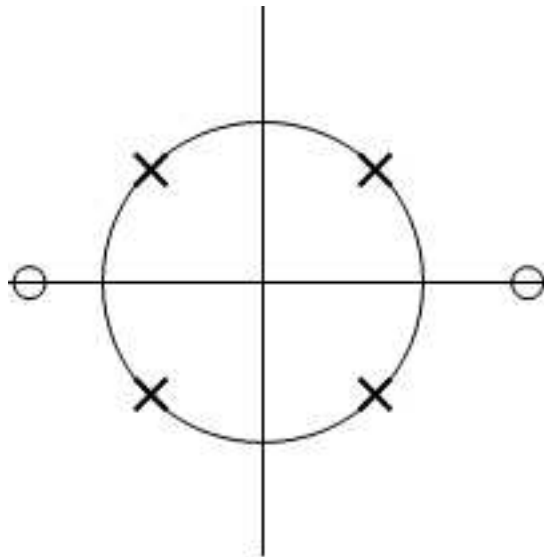


Figure 3.12

Take $H(s) = LHP$ factors:

$$H(s) = \frac{\sqrt{2} + s}{(s + \alpha)(s + \bar{\alpha})} = \frac{\sqrt{2} + s}{s^2 + \sqrt{2}s + 1}$$

3.3.2.2 Traditional Filter Designs

3.3.2.2.1 Butterworth

$$B(\lambda^2) = \frac{1}{1 + \lambda^{2M}}$$

NOTE: Remember this for homework and rest problems!

"Maximally smooth" at $\lambda = 0$ and $\lambda = \infty$ (maximum possible number of zero derivatives). Figure 3.13.

$$B(\lambda^2) = (|H(\lambda)|)^2$$

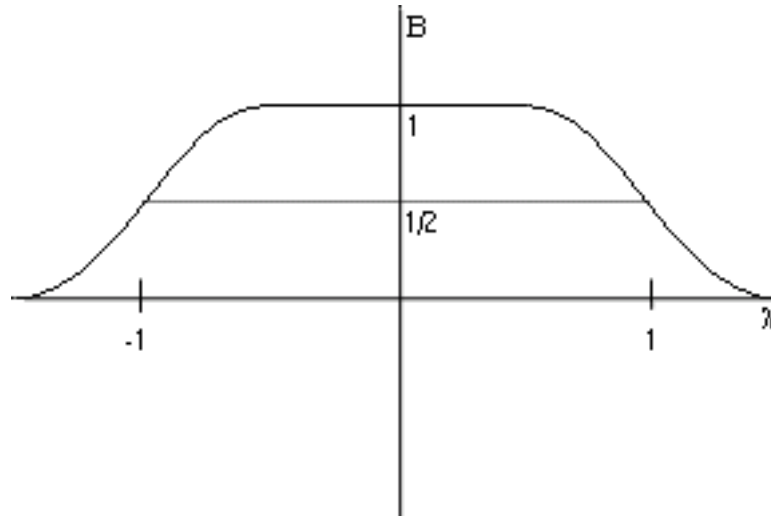


Figure 3.13

3.3.2.2.2 Chebyshev

$$B(\lambda^2) = \frac{1}{1 + \epsilon^2 C_M^2(\lambda)}$$

where $C_M^2(\lambda)$ is an M^{th} order Chebyshev polynomial. Figure 3.14.

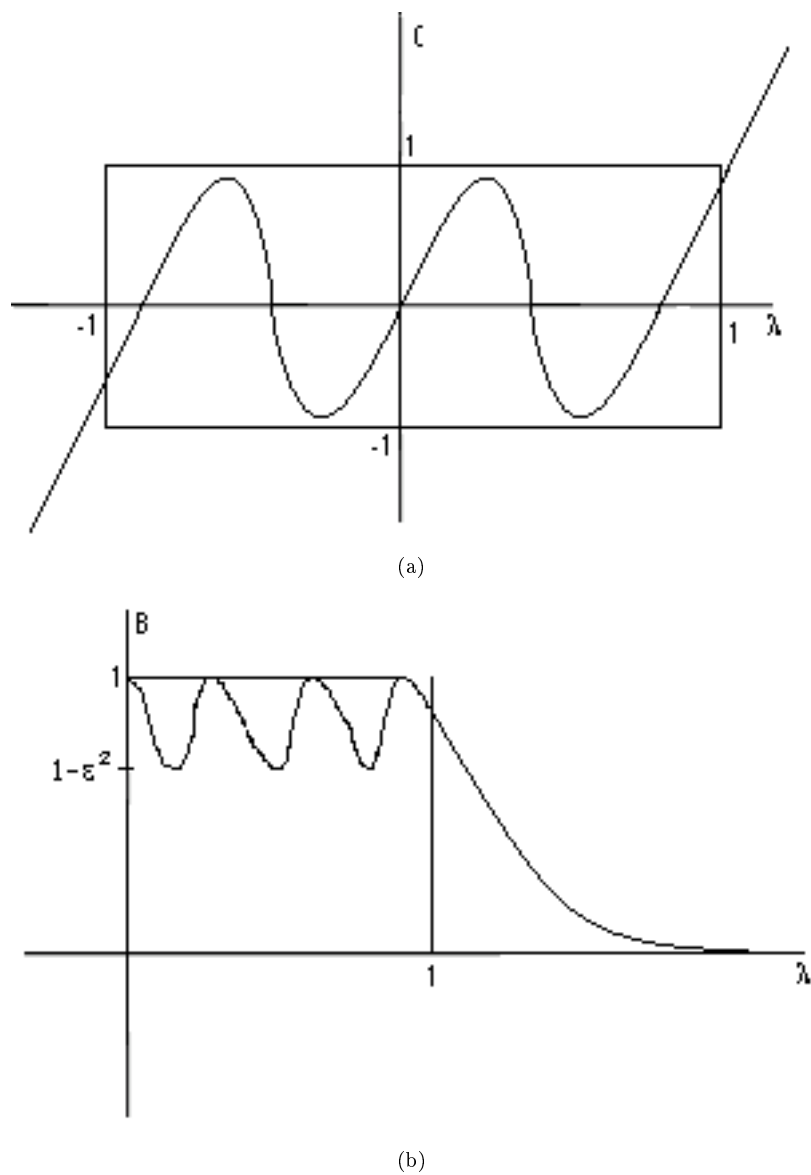


Figure 3.14

3.3.2.2.3 Inverse Chebyshev

Figure 3.15.

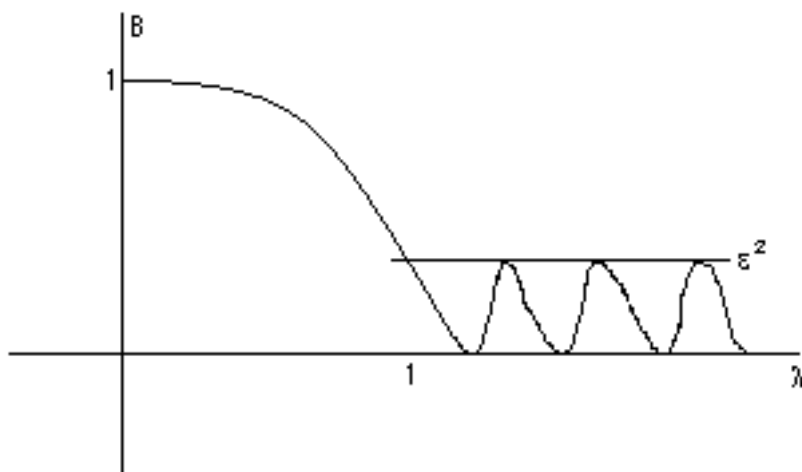


Figure 3.15

3.3.2.2.4 Elliptic Function Filter (Cauer Filter)

$$B(\lambda^2) = \frac{1}{1 + \epsilon^2 J_M^2(\lambda)}$$

where J_M is the "Jacobi Elliptic Function." Figure 3.16.

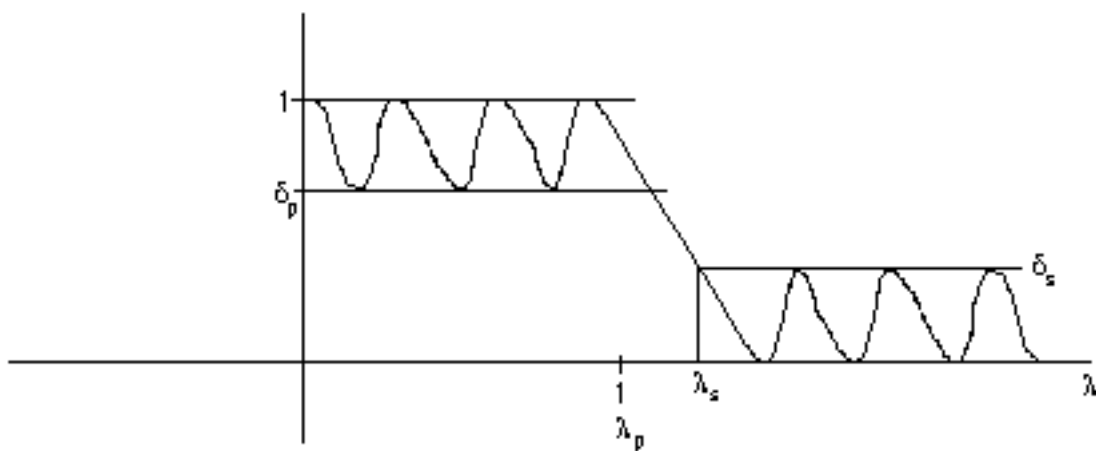


Figure 3.16

The Cauer filter is $\|\mathbf{L}\|_\infty$ optimum in the sense that for a given M , δ_p , δ_s , and λ_p , the transition bandwidth is smallest.

That is, it is $\|\mathbf{L}\|_\infty$ optimal.

3.3.3 IIR Digital Filter Design via the Bilinear Transform¹⁴

A **bilinear transform** maps an analog filter $H_a(s)$ to a discrete-time filter $H(z)$ of the same order.

If only we could somehow map these optimal analog filter designs to the digital world while preserving the magnitude response characteristics, we could make use of the already-existing body of knowledge concerning optimal analog filter design.

3.3.3.1 Bilinear Transformation

The Bilinear Transform is a nonlinear ($\mathbb{C} \rightarrow \mathbb{C}$) mapping that maps a function of the complex variable s to a function of a complex variable z . This map has the property that the LHP in s ($\Re(s) < 0$) maps to the interior of the unit circle in z , and the $j\lambda = s$ axis maps to the unit circle $e^{j\omega}$ in z .

Bilinear transform:

$$s = \alpha \frac{z - 1}{z + 1}$$

$$H(z) = H_a\left(s = \alpha \frac{z - 1}{z + 1}\right)$$

NOTE: $j\lambda = \alpha \frac{e^{j\omega} - 1}{e^{j\omega} + 1} = \alpha \frac{(e^{j\omega} - 1)(e^{-(j\omega)} + 1)}{(e^{j\omega} + 1)(e^{-(j\omega)} + 1)} = \frac{2jsin(\omega)}{2 + 2cos(\omega)} = j\alpha tan\left(\frac{\omega}{2}\right)$, so $\lambda \equiv \alpha tan\left(\frac{\omega}{2}\right)$, $\omega \equiv 2arctan\left(\frac{\lambda}{\alpha}\right)$. Figure 3.17.

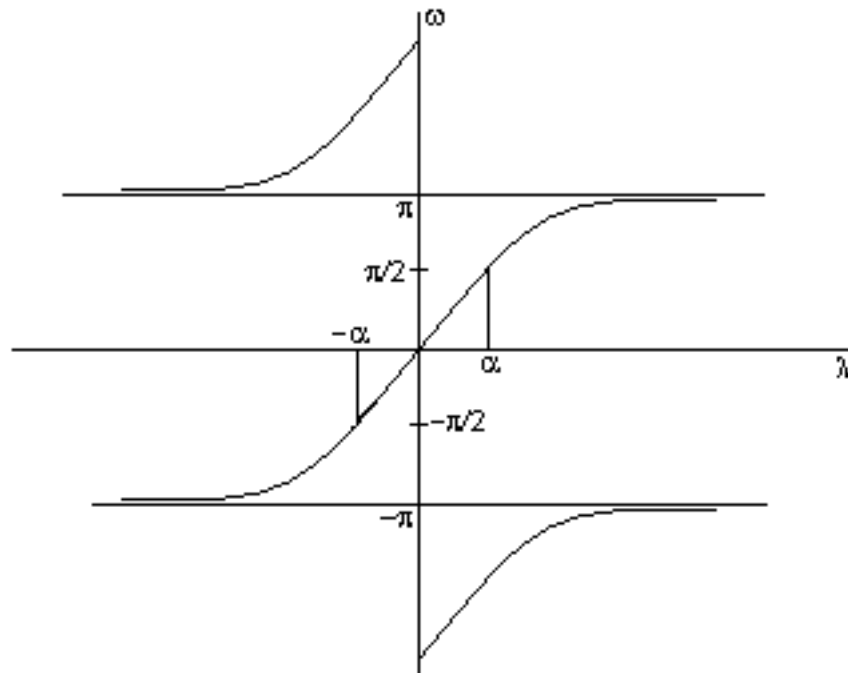


Figure 3.17

¹⁴This content is available online at <<http://cnx.org/content/m12757/1.2/>>.

The magnitude response doesn't change in the mapping from λ to ω , it is simply warped nonlinearly according to $H(\omega) = H_a\left(\alpha \tan\left(\frac{\omega}{2}\right)\right)$, Figure 3.18.

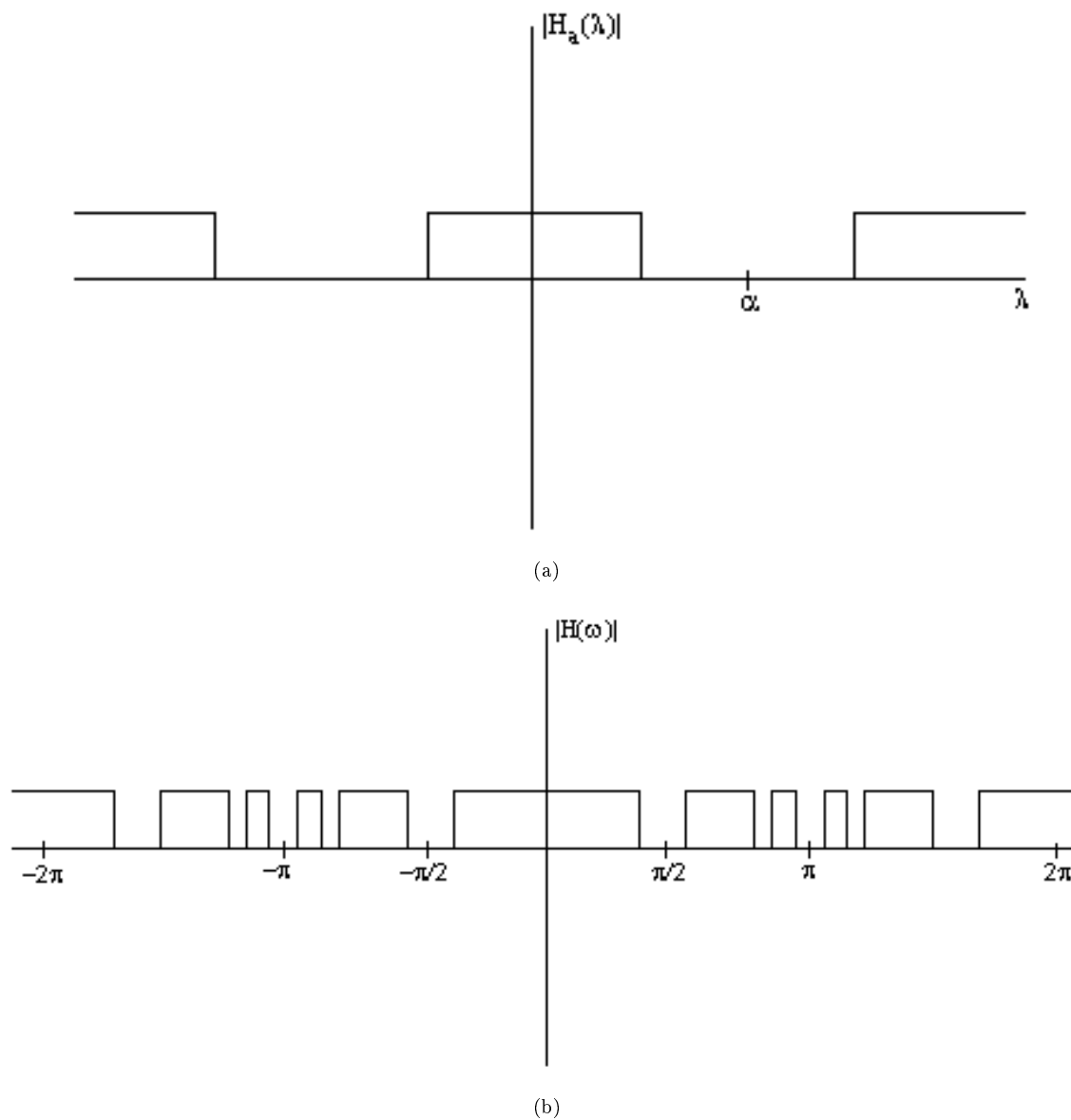


Figure 3.18: The first image implies the second one.

NOTE: This mapping preserves $\|L\|_\infty$ errors in (warped) frequency bands. Thus optimal Cauer ($\|L\|_\infty$) filters in the analog realm can be mapped to $\|L\|_\infty$ optimal discrete-time IIR filters using the bilinear transform! This is how IIR filters with $\|L\|_\infty$ optimal magnitude responses are designed.

NOTE: The parameter α provides one degree of freedom which can be used to map a single λ_0 to any desired ω_0 :

$$\lambda_0 = \alpha \tan\left(\frac{\omega_0}{2}\right)$$

or

$$\alpha = \frac{\lambda_0}{\tan\left(\frac{\omega_0}{2}\right)}$$

This can be used, for example, to map the pass-band edge of a lowpass analog prototype filter to any desired pass-band edge in ω . Often, analog prototype filters will be designed with $\lambda = 1$ as a band edge, and α will be used to locate the band edge in ω . Thus an M^{th} order optimal lowpass analog filter prototype can be used to design *any* M^{th} order discrete-time lowpass IIR filter with the same ripple specifications.

3.3.3.2 Prewarping

Given specifications on the frequency response of an IIR filter to be designed, map these to specifications in the analog frequency domain which are equivalent. Then a satisfactory analog prototype can be designed which, when transformed to discrete-time using the bilinear transformation, will meet the specifications.

Example 3.3

The goal is to design a high-pass filter, $\omega_s = \omega_s$, $\omega_p = \omega_p$, $\delta_s = \delta_s$, $\delta_p = \delta_p$; pick up some $\alpha = \alpha_0$. In Figure 3.19 the δ_i remain the same and the band edges are mapped by $\lambda_i = \alpha_0 \tan\left(\frac{\omega_i}{2}\right)$.

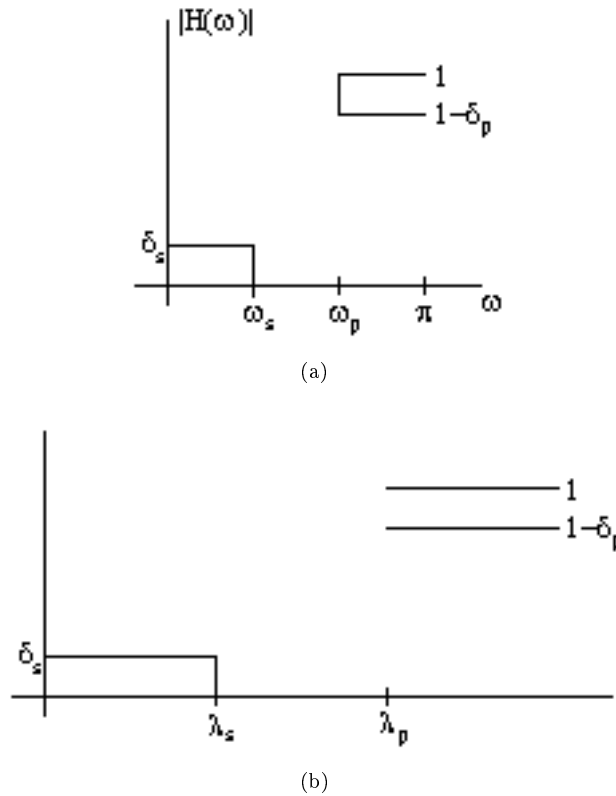


Figure 3.19: Where $\lambda_s = \alpha_0 \tan\left(\frac{\omega_s}{2}\right)$ and $\lambda_p = \alpha_0 \tan\left(\frac{\omega_p}{2}\right)$.

3.3.4 Impulse-Invariant Design¹⁵

Pre-classical, adhoc-but-easy method of converting an analog prototype filter to a digital IIR filter. Does not preserve any optimality.

Impulse invariance means that digital filter impulse response exactly equals samples of the analog prototype impulse response:

$$\forall n : (h(n) = h_a(nT))$$

How is this done?

The impulse response of a causal, stable analog filter is simply a sum of decaying exponentials:

$$H_a(s) = \frac{b_0 + b_1s + b_2s^2 + \dots + b_ps^p}{1 + a_1s + a_2s^2 + \dots + a_ps^p} = \frac{A_1}{s - s_1} + \frac{A_2}{s - s_2} + \dots + \frac{A_p}{s - s_p}$$

which implies

$$h_a(t) = (A_1e^{s_1t} + A_2e^{s_2t} + \dots + A_pe^{s_pt}) u(t)$$

For impulse invariance, we desire

$$h(n) = h_a(nT) = (A_1e^{s_1nT} + A_2e^{s_2nT} + \dots + A_pe^{s_pnT}) u(n)$$

Since

$$A_ke^{(s_kT)n}u(n) \equiv \frac{A_kz}{z - e^{s_kT}}$$

where $|z| > |e^{s_kT}|$, and

$$H(z) = \sum_{k=1}^p \left(A_k \frac{z}{z - e^{s_kT}} \right)$$

where $|z| > \max_k \{|e^{s_kT}|\}$.

This technique is used occasionally in digital simulations of analog filters.

Exercise 3.6

(Solution on p. 224.)

What is the main problem/drawback with this design technique?

3.3.5 Digital-to-Digital Frequency Transformations¹⁶

Given a prototype *digital* filter design, transformations similar to the bilinear transform can also be developed.

Requirements on such a mapping $z^{-1} = g(z^{-1})$:

1. points inside the unit circle stay inside the unit circle (condition to preserve stability)
2. unit circle is mapped to itself (preserves frequency response)

This condition (list, item 2, p. 216) implies that $e^{-(j\omega_1)} = g(e^{-(j\omega)}) = |g(\omega)|e^{j\angle(g(\omega))}$ requires that $|g(e^{-(j\omega)})| = 1$ on the unit circle!

Thus we require an **all-pass** transformation:

$$g(z^{-1}) = \prod_{k=1}^p \left(\frac{z^{-1} - \alpha_k}{1 - \alpha_k z^{-1}} \right)$$

¹⁵This content is available online at <<http://cnx.org/content/m12760/1.2/>>.

¹⁶This content is available online at <<http://cnx.org/content/m12759/1.2/>>.

where $|\alpha_K| < 1$, which is required to satisfy this condition (list, item 1, p. 216).

Example 3.4: Lowpass-to-Lowpass

$$z_1^{-1} = \frac{z^{-1} - a}{1 - az^{-1}}$$

which maps original filter with a cutoff at ω_c to a new filter with cutoff ω'_c ,

$$a = \frac{\sin\left(\frac{1}{2}(\omega_c - \omega'_c)\right)}{\sin\left(\frac{1}{2}(\omega_c + \omega'_c)\right)}$$

Example 3.5: Lowpass-to-Highpass

$$z_1^{-1} = \frac{z^{-1} + a}{1 + az^{-1}}$$

which maps original filter with a cutoff at ω_c to a frequency reversed filter with cutoff ω'_c ,

$$a = \frac{\cos\left(\frac{1}{2}(\omega_c - \omega'_c)\right)}{\cos\left(\frac{1}{2}(\omega_c + \omega'_c)\right)}$$

(Interesting and occasionally useful!)

3.3.6 Prony's Method¹⁷

Prony's Method is a quasi-least-squares time-domain IIR filter design method.

First, assume $H(z)$ is an "all-pole" system:

$$H(z) = \frac{b_0}{1 + \sum_{k=1}^M (a_k z^{-k})} \quad (3.15)$$

and

$$h(n) = - \left(\sum_{k=1}^M (a_k h(n-k)) \right) + b_0 \delta(n)$$

where $h(n) = 0$, $n < 0$ for a causal system.

NOTE: For $h = 0$, $h(0) = b_0$.

Let's attempt to fit a desired impulse response (let it be *causal*, although one can extend this technique when it isn't) $h_d(n)$.

A true least-squares solution would attempt to minimize

$$\epsilon^2 = \sum_{n=0}^{\infty} \left((|h_d(n) - h(n)|)^2 \right)$$

where $H(z)$ takes the form in (3.15). This is a difficult non-linear optimization problem which is known to be plagued by local minima in the error surface. So instead of solving this difficult non-linear problem, we solve the **deterministic linear prediction** problem, which is related to, *but not the same as*, the true least-squares optimization.

The deterministic linear prediction problem is a *linear* least-squares optimization, which is easy to solve, but it minimizes the *prediction* error, not the $(|desired - actual|)^2$ response error.

¹⁷This content is available online at <<http://cnx.org/content/m12762/1.2/>>.

Notice that for $n > 0$, with the all-pole filter

$$h(n) = - \left(\sum_{k=1}^M (a_k h(n-k)) \right) \quad (3.16)$$

the right hand side of this equation (3.16) is a **linear predictor** of $h(n)$ in terms of the M previous samples of $h(n)$.

For the desired response $h_d(n)$, one can choose the recursive filter coefficients a_k to minimize the squared prediction error

$$\epsilon_p^2 = \sum_{n=1}^{\infty} \left(\left(|h_d(n) + \sum_{k=1}^M (a_k h_d(n-k))| \right)^2 \right)$$

where, in practice, the ∞ is replaced by an N .

In matrix form, that's

$$\begin{pmatrix} h_d(0) & 0 & \dots & 0 \\ h_d(1) & h_d(0) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ h_d(N-1) & h_d(N-2) & \dots & h_d(N-M) \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_M \end{pmatrix} \approx - \begin{pmatrix} h_d(1) \\ h_d(2) \\ \vdots \\ h_d(N) \end{pmatrix}$$

or

$$H_d \mathbf{a} \approx -h_d$$

The optimal solution is

$$a_{lp} = - \left((H_d^H H_d)^{-1} H_d^H h_d \right)$$

Now suppose $H(z)$ is an M^{th} -order IIR (ARMA) system,

$$H(z) = \frac{\sum_{k=0}^M (b_k z^{-k})}{1 + \sum_{k=1}^M (a_k z^{-k})}$$

or

$$\begin{aligned} h(n) &= - \left(\sum_{k=1}^M (a_k h(n-k)) \right) + \sum_{k=0}^M (b_k \delta(n-k)) \\ &= \begin{cases} - \left(\sum_{k=1}^M (a_k h(n-k)) \right) + b_n & \text{if } 0 \leq n \leq M \\ - \left(\sum_{k=1}^M (a_k h(n-k)) \right) & \text{if } n > M \end{cases} \end{aligned} \quad (3.17)$$

For $n > M$, this is just like the all-pole case, so we can solve for the best predictor coefficients as before:

$$\begin{pmatrix} h_d(M) & h_d(M-1) & \dots & h_d(1) \\ h_d(M+1) & h_d(M) & \dots & h_d(2) \\ \vdots & \vdots & \ddots & \vdots \\ h_d(N-1) & h_d(N-2) & \dots & h_d(N-M) \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_M \end{pmatrix} \approx \begin{pmatrix} h_d(M+1) \\ h_d(M+2) \\ \vdots \\ h_d(N) \end{pmatrix}$$

or

$$\hat{H}_d \mathbf{a} \approx \hat{h}_d$$

and

$$a_{opt} = \left((\hat{H}_d)^H H_d \right)^{-1} H_d^H \hat{h}_d$$

Having determined the a 's, we can use them in (3.17) to obtain the b_n 's:

$$b_n = \sum_{k=1}^M (a_k h_d(n-k))$$

where $h_d(n-k) = 0$ for $n-k < 0$.

For $N = 2M$, \hat{H}_d is square, and we can solve *exactly* for the a_k 's with no error. The b_k 's are also chosen such that there is no error in the first $M+1$ samples of $h(n)$. Thus for $N = 2M$, the first $2M+1$ points of $h(n)$ *exactly equal* $h_d(n)$. This is called **Prony's Method**. Baron de Prony invented this in 1795.

For $N > 2M$, $h_d(n) = h(n)$ for $0 \leq n \leq M$, the prediction error is minimized for $M+1 < n \leq N$, and whatever for $n \geq N+1$. This is called the **Extended Prony Method**.

One might prefer a method which tries to minimize an overall error with the numerator coefficients, rather than just using them to exactly fit $h_d(0)$ to $h_d(M)$.

3.3.6.1 Shank's Method

1. Assume an all-pole model and fit $h_d(n)$ by minimizing the prediction error $1 \leq n \leq N$.
2. Compute $v(n)$, the impulse response of this all-pole filter.
3. Design an all-zero (MA, FIR) filter which fits $v(n) * h_z(n) \approx h_d(n)$ optimally in a least-squares sense (Figure 3.20).

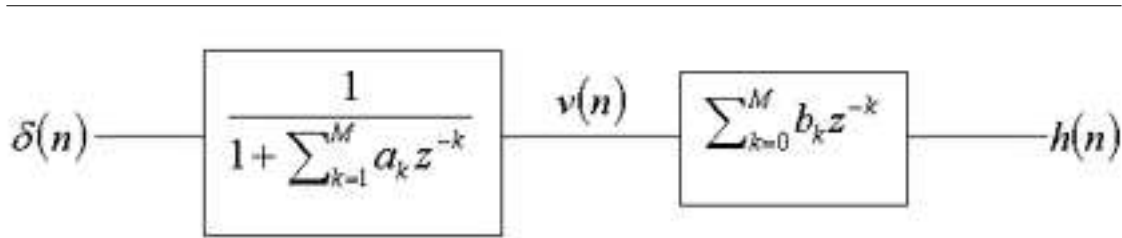


Figure 3.20: Here, $h(n) \approx h_d(n)$.

The final IIR filter is the cascade of the all-pole and all-zero filter.

This (list, item 3, p. 219) is solved by

$$\min_{b_k} \left\{ \sum_{n=0}^N \left(\left(\left| h_d(n) - \sum_{k=0}^M (b_k v(n-k)) \right| \right)^2 \right) \right\}$$

or in matrix form

$$\begin{pmatrix} v(0) & 0 & 0 & \dots & 0 \\ v(1) & v(0) & 0 & \dots & 0 \\ v(2) & v(1) & v(0) & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ v(N) & v(N-1) & v(N-2) & \dots & v(N-M) \end{pmatrix} \begin{pmatrix} b_0 \\ b_1 \\ b_2 \\ \vdots \\ b_M \end{pmatrix} \approx \begin{pmatrix} h_d(0) \\ h_d(1) \\ h_d(2) \\ \vdots \\ h_d(N) \end{pmatrix}$$

Which has solution:

$$b_{opt} = (V^H V)^{-1} V^H \mathbf{h}$$

Notice that none of these methods solve the true least-squares problem:

$$\min_{\mathbf{a}, \mathbf{b}} \left\{ \sum_{n=0}^{\infty} \left(|h_d(n) - h(n)|^2 \right) \right\}$$

which is a difficult non-linear optimization problem. The true least-squares problem can be written as:

$$\min_{\alpha, \beta} \left\{ \sum_{n=0}^{\infty} \left(\left(|h_d(n) - \sum_{i=1}^M (\alpha_i e^{\beta_i n})| \right)^2 \right) \right\}$$

since the impulse response of an IIR filter is a sum of exponentials, and non-linear optimization is then used to solve for the α_i and β_i .

3.3.7 Linear Prediction¹⁸

Recall that for the all-pole design problem, we had the overdetermined set of linear equations:

$$\begin{pmatrix} h_d(0) & 0 & \dots & 0 \\ h_d(1) & h_d(0) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ h_d(N-1) & h_d(N-2) & \dots & h_d(N-M) \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_M \end{pmatrix} \approx - \begin{pmatrix} h_d(1) \\ h_d(2) \\ \vdots \\ h_d(N) \end{pmatrix}$$

with solution $\mathbf{a} = (H_d^H H_d)^{-1} H_d^H h_d$

Let's look more closely at $H_d^H H_d = R$. r_{ij} is related to the *correlation* of h_d with itself:

$$r_{ij} = \sum_{k=0}^{N-\max\{i,j\}} (h_d(k) h_d(k + |i - j|))$$

Note also that:

$$H_d^H h_d = \begin{pmatrix} r_d(1) \\ r_d(2) \\ r_d(3) \\ \vdots \\ r_d(M) \end{pmatrix}$$

where

$$r_d(i) = \sum_{n=0}^{N-i} (h_d(n) h_d(n + i))$$

so this takes the form $a_{opt} = -(R^H r_d)$, or $\mathbf{R}\mathbf{a} = -\mathbf{r}$, where R is $M \times M$, \mathbf{a} is $M \times 1$, and \mathbf{r} is also $M \times 1$.

¹⁸This content is available online at <<http://cnx.org/content/m12761/1.2/>>.

Except for the changing endpoints of the sum, $r_{ij} \approx r(i-j) = r(j-i)$. If we tweak the problem slightly to make $r_{ij} = r(i-j)$, we get:

$$\begin{pmatrix} r(0) & r(1) & r(2) & \dots & r(M-1) \\ r(1) & r(0) & r(1) & \dots & \vdots \\ r(2) & r(1) & r(0) & \dots & \vdots \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ r(M-1) & \dots & \dots & \dots & r(0) \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \\ a_3 \\ \vdots \\ a_M \end{pmatrix} = - \begin{pmatrix} r(1) \\ r(2) \\ r(3) \\ \vdots \\ r(M) \end{pmatrix}$$

The matrix R is **Toeplitz** (diagonal elements equal), and \mathbf{a} can be solved for with $O(M^2)$ computations using Levinson's recursion.

3.3.7.1 Statistical Linear Prediction

Used very often for forecasting (e.g. stock market).

Given a time-series $y(n)$, assumed to be produced by an auto-regressive (AR) (all-pole) system:

$$y(n) = - \left(\sum_{k=1}^M (a_k y(n-k)) \right) + u(n)$$

where $u(n)$ is a white Gaussian noise sequence which is stationary and has zero mean.

To determine the model parameters $\{a_k\}$ minimizing the variance of the prediction error, we seek

$$\begin{aligned} \min_{a_k} \left\{ E \left[\left(y(n) + \sum_{k=1}^M (a_k y(n-k)) \right)^2 \right] \right\} &= \min_{a_k} \left\{ E \left[y^2(n) + 2 \sum_{k=1}^M (a_k y(n) y(n-k)) + \left(\sum_{k=1}^M (a_k y(n-k))^2 \right) \right] \right\} \\ \min_{a_k} \left\{ E[y^2(n)] + 2 \sum_{k=1}^M (a_k E[y(n) y(n-k)]) + \sum_{k=1}^M \left(\sum_{l=1}^M (a_k a_l E[y(n-k) y(n-l)]) \right) \right\} \end{aligned}$$

NOTE: The mean of $y(n)$ is zero.

$$\epsilon^2 = r(0) + 2 \begin{pmatrix} r(1) & r(2) & r(3) & \dots & r(M) \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \\ a_3 \\ \vdots \\ a_M \end{pmatrix} + \begin{pmatrix} a_1 & a_2 & a_3 & \dots & a_M \end{pmatrix} \begin{pmatrix} r(0) & r(1) & r(2) & \dots & r(M-1) \\ r(1) & r(0) & r(1) & \dots & \vdots \\ r(2) & r(1) & r(0) & \dots & \vdots \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ r(M-1) & \dots & \dots & \dots & r(0) \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \\ a_3 \\ \vdots \\ a_M \end{pmatrix} \quad (3.19)$$

$$\frac{\partial}{\partial \mathbf{a}} (\epsilon^2) = 2\mathbf{r} + 2R\mathbf{a} \quad (3.20)$$

Setting (3.20) equal to zero yields: $R\mathbf{a} = -\mathbf{r}$ These are called the **Yule-Walker** equations. In practice, given samples of a sequence $y(n)$, we estimate $r(n)$ as

$$\hat{r}(n) = \frac{1}{N} \sum_{k=0}^{N-n} (y(n) y(n+k)) \approx E[y(k) y(n+k)]$$

which is extremely similar to the deterministic least-squares technique.

Solutions to Exercises in Chapter 3

Solution to Exercise 3.1 (p. 196)

Yes; in fact it's optimal! (in a certain sense)

Solution to Exercise 3.2 (p. 197): Gibbs Phenomenon

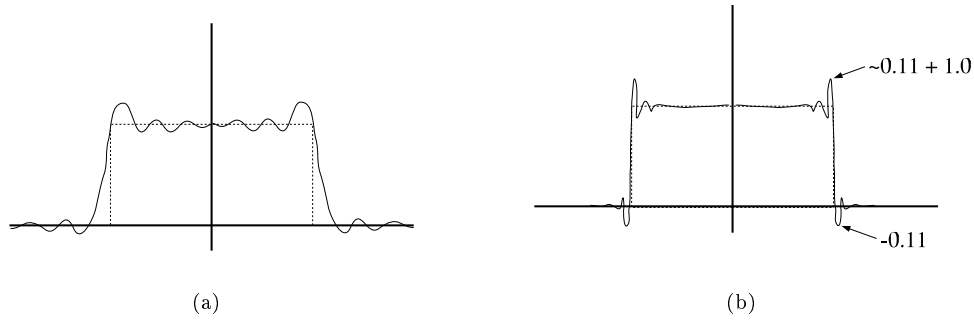


Figure 3.21: (a) $A(\omega)$, small M (b) $A(\omega)$, large M

Solution to Exercise 3.3 (p. 197)

$$\begin{cases} e^{-j\omega \frac{M-1}{2}} & \text{if } -\omega_c \leq \omega \leq \omega_c \\ 0 & \text{if } -\pi \leq \omega < -\omega_c \vee \omega_c < \omega \leq \pi \end{cases}$$

Solution to Exercise 3.4 (p. 200)

Yes, when the desired response is discontinuous. Since the frequency response of a finite-length filter must be continuous, without a transition band the worst-case error could be no less than half the discontinuity.

Solution to Exercise 3.5 (p. 201)

It's the same problem! To show that, consider an odd-length, symmetric linear phase filter.

$$\begin{aligned} H(\omega) &= \sum_{n=0}^{M-1} (h(n) e^{-j\omega n}) \\ &= e^{-j\omega \frac{M-1}{2}} \left(h\left(\frac{M-1}{2}\right) + 2 \sum_{n=1}^L \left(h\left(\frac{M-1}{2} - n\right) \cos(\omega n) \right) \right) \end{aligned} \quad (3.21)$$

$$A(\omega) = h(L) + 2 \sum_{n=1}^L (h(L-n) \cos(\omega n)) \quad (3.22)$$

Where $(L \doteq \frac{M-1}{2})$.

Using trigonometric identities (such as $\cos(n\alpha) = 2\cos((n-1)\alpha)\cos(\alpha) - \cos((n-2)\alpha)$), we can rewrite $A(\omega)$ as

$$A(\omega) = h(L) + 2 \sum_{n=1}^L (h(L-n) \cos(\omega n)) = \sum_{k=0}^L (\alpha_k \cos^k(\omega))$$

where the α_k are related to the $h(n)$ by a linear transformation. Now, let $x = \cos(\omega)$. This is a one-to-one mapping from $x \in [-1, 1]$ onto $\omega \in [0, \pi]$. Thus $A(\omega)$ is an L th-order polynomial in $x = \cos(\omega)$!

IMPLICATION: The alternation theorem holds for the L^∞ filter design problem, too!

Therefore, to determine whether or not a length- M , odd-length, symmetric linear-phase filter is optimal in an L^∞ sense, simply count the alternations in $E(\omega) = W(\omega)(A_d(\omega) - A(\omega))$ in the pass and stop bands. If there are $L + 2 = \frac{M+3}{2}$ or more alternations, $h(n)$, $0 \leq n \leq M - 1$ is the optimal filter!

Solution to Exercise 3.6 (p. 216)

Since it samples the non-bandlimited impulse response of the analog prototype filter, the frequency response **aliases**. This distorts the original analog frequency and destroys any optimal frequency properties in the resulting digital filter.

Chapter 4

Digital Filter Structures and Quantization Error Analysis

4.1 Filter Structures

4.1.1 Filter Structures¹

A realizable filter must require only a finite number of computations per output sample. For linear, causal, time-Invariant filters, this restricts one to rational transfer functions of the form

$$H(z) = \frac{b_0 + b_1 z^{-1} + \dots + b_m z^{-m}}{1 + a_1 z^{-1} + a_2 z^{-2} + \dots + a_n z^{-n}}$$

Assuming no pole-zero cancellations, $H(z)$ is FIR if $\forall i, i > 0 : (a_i = 0)$, and IIR otherwise. Filter structures usually implement rational transfer functions as difference equations.

Whether FIR or IIR, a given transfer function can be implemented with many different filter structures. With infinite-precision data, coefficients, and arithmetic, all filter structures implementing the same transfer function produce the same output. However, different filter structures may produce very different errors with quantized data and finite-precision or fixed-point arithmetic. The computational expense and memory usage may also differ greatly. Knowledge of different filter structures allows DSP engineers to trade off these factors to create the best implementation.

4.1.2 FIR Filter Structures²

Consider causal FIR filters: $y(n) = \sum_{k=0}^{M-1} (h(k) x(n-k))$; this can be realized using the following structure

¹This content is available online at <http://cnx.org/content/m11917/1.3/>.

²This content is available online at <http://cnx.org/content/m11918/1.2/>.

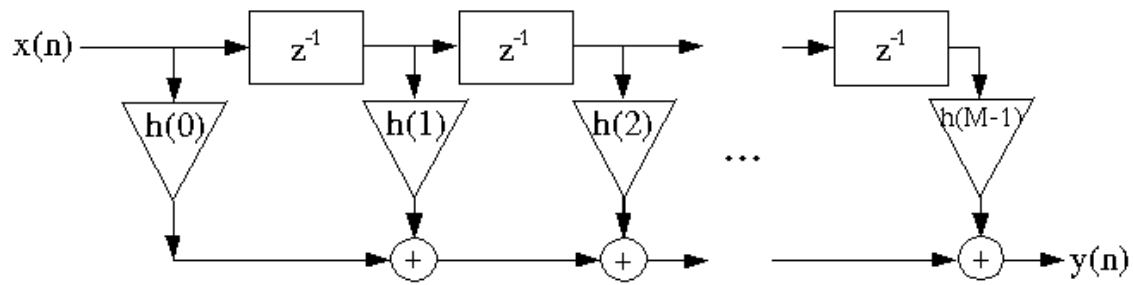


Figure 4.1

or in a different notation

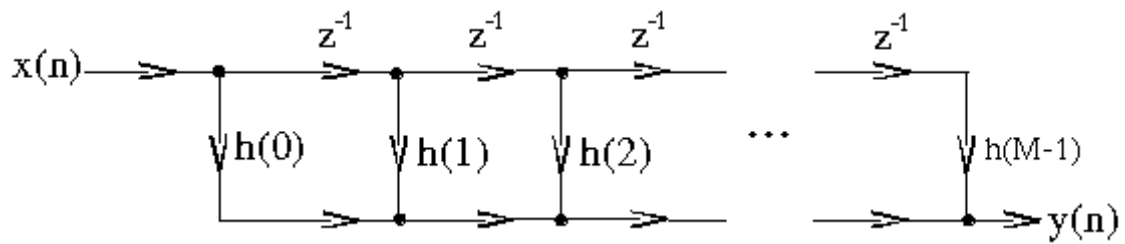
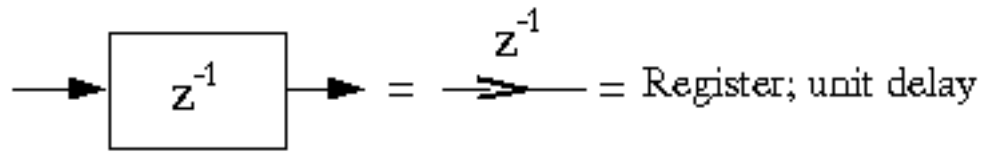
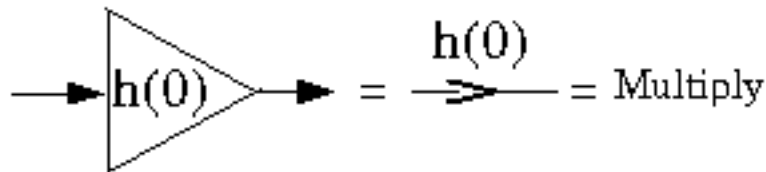


Figure 4.2



(a)



(b)



(c)

Figure 4.3

This is called the **direct-form FIR filter structure**.

There are no closed loops (no feedback) in this structure, so it is called a **non-recursive structure**. Since any FIR filter can be implemented using the direct-form, non-recursive structure, it is always possible to implement an FIR filter non-recursively. However, it is also possible to implement an FIR filter *recursively*, and for some special sets of FIR filter coefficients this is much more efficient.

Example 4.1

$$y(n) = \sum_{k=0}^{M-1} (x(n-k))$$

where

$$h(k) = \left\{ 0, 0, \underset{k=0}{\overset{[U+2303][U+FE00]}{1}}, \dots, 1, \underset{k=M-1}{\overset{[U+2303][U+FE00]}{1}}, 0, 0, 0, \dots \right\}$$

But note that

$$y(n) = y(n-1) + x(n) - x(n-M)$$

This can be implemented as

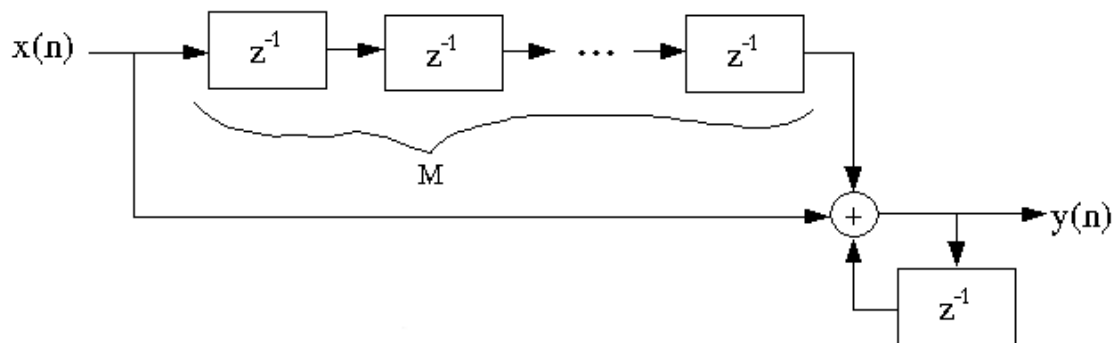


Figure 4.4

Instead of costing $M - 1$ adds/output point, this comb filter costs only two adds/output.

Exercise 4.1

Is this stable, and if not, how can it be made so?

IIR filters must be implemented with a *recursive* structure, since that's the only way a finite number of elements can generate an infinite-length impulse response in a linear, time-invariant (LTI) system. Recursive structures have the advantages of being able to implement IIR systems, and sometimes greater computational efficiency, but the disadvantages of possible instability, limit cycles, and other deleterious effects that we will study shortly.

4.1.2.1 Transpose-form FIR filter structures

The **flow-graph-reversal theorem** says that if one changes the directions of all the arrows, and inputs at the output and takes the output from the input of a reversed flow-graph, the new system has an identical input-output relationship to the original flow-graph.

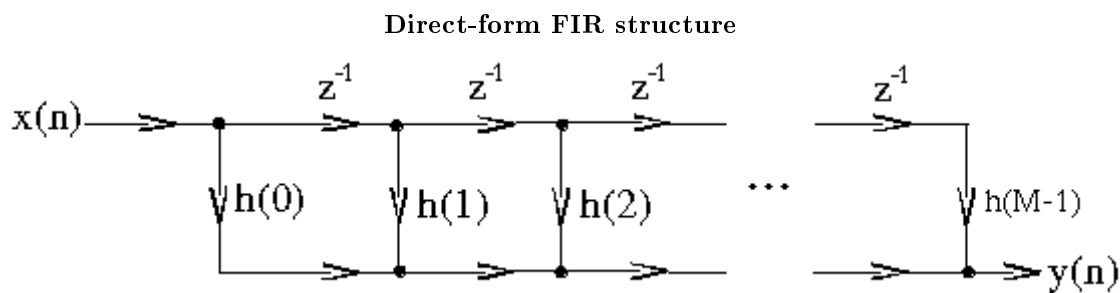


Figure 4.5

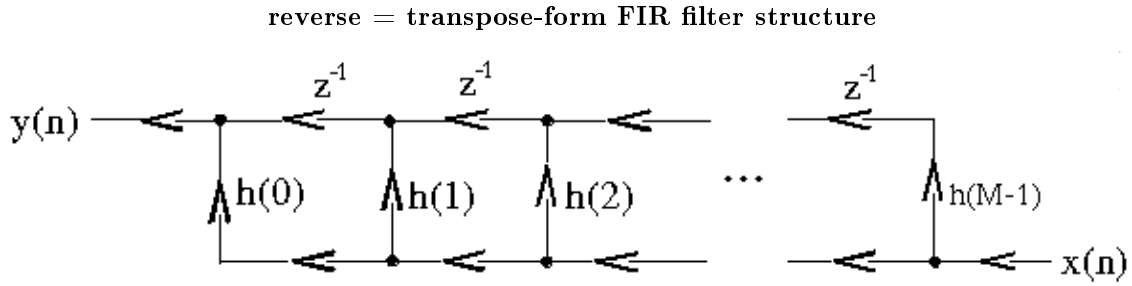


Figure 4.6

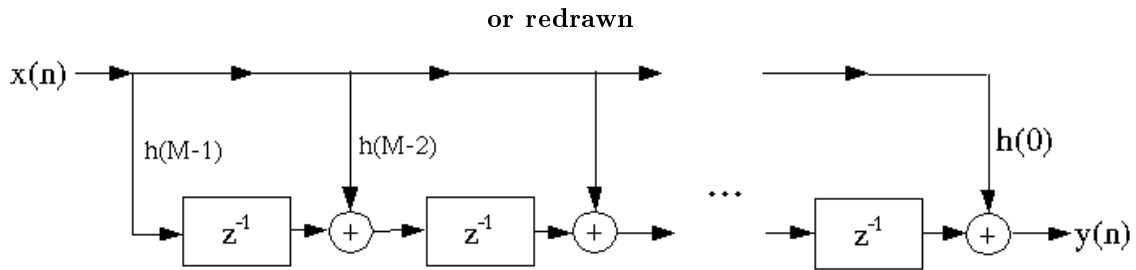


Figure 4.7

4.1.2.1.1 Cascade structures

The z -transform of an FIR filter can be factored into a cascade of short-length filters

$$b_0 + b_1 z^{-1} + b_2 z^{-2} + \dots + b_m z^{-m} = b_0 (1 - z_1 z^{-1}) (1 - z_2 z^{-1}) \dots (1 - z_m z^{-1})$$

where the z_i are the zeros of this polynomial. Since the coefficients of the polynomial are usually real, the roots are usually complex-conjugate pairs, so we generally combine $(1 - z_i z^{-1}) (1 - \bar{z}_i z^{-1})$ into one quadratic (length-2) section with *real* coefficients

$$(1 - z_i z^{-1}) (1 - \bar{z}_i z^{-1}) = 1 - 2\Re(z_i) z^{-1} + (|z_i|)^2 z^{-2} = H_i(z)$$

The overall filter can then be implemented in a **cascade** structure.

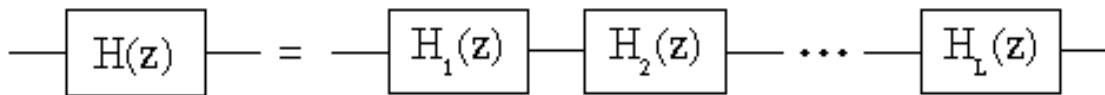


Figure 4.8

This is occasionally done in FIR filter implementation when one or more of the short-length filters can be implemented efficiently.

4.1.2.1.2 Lattice Structure

It is also possible to implement FIR filters in a lattice structure: this is sometimes used in adaptive filtering

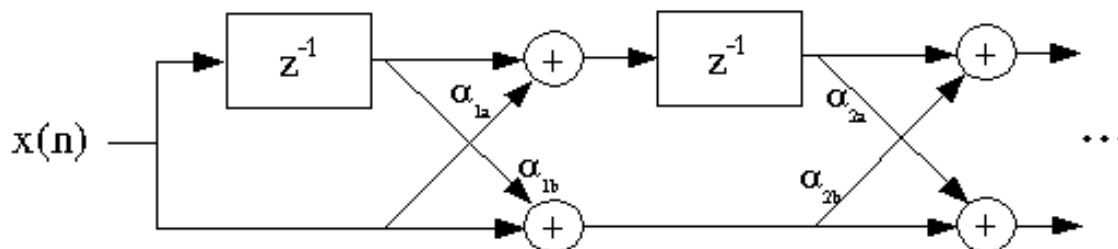


Figure 4.9

4.1.3 IIR Filter Structures³

IIR (Infinite Impulse Response) filter structures must be recursive (use feedback); an infinite number of coefficients could not otherwise be realized with a finite number of computations per sample.

$$H(z) = \frac{N(z)}{D(z)} = \frac{b_0 + b_1z^{-1} + b_2z^{-2} + \cdots + b_Mz^{-M}}{1 + a_1z^{-1} + a_2z^{-2} + \cdots + a_Nz^{-N}}$$

The corresponding time-domain difference equation is

$$y(n) = -(a_1y(n-1) + a_2y(n-2) + \cdots + a_Ny(n-N)) + b_0x(n) + b_1x(n-1) + \cdots + b_Mx(n-M)$$

4.1.3.1 Direct-form I IIR Filter Structure

The difference equation above is implemented directly as written by the Direct-Form I IIR Filter Structure.

³This content is available online at <<http://cnx.org/content/m11919/1.2/>>.

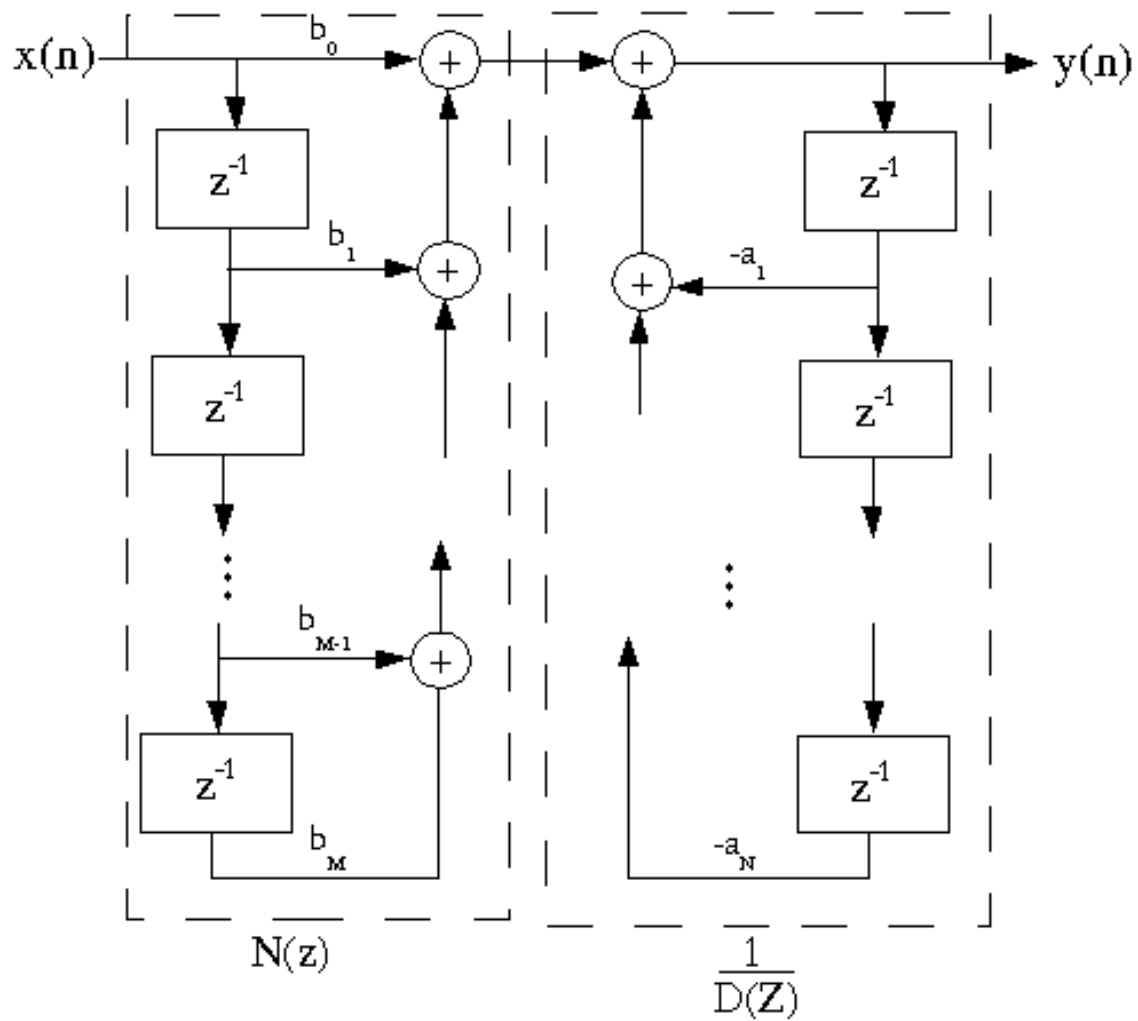


Figure 4.10

Note that this is a cascade of two systems, $N(z)$ and $\frac{1}{D(z)}$. If we reverse the order of the filters, the overall system is unchanged: The memory elements appear in the middle and store identical values, so they can be combined, to form the Direct-Form II IIR Filter Structure.

4.1.3.2 Direct-Form II IIR Filter Structure

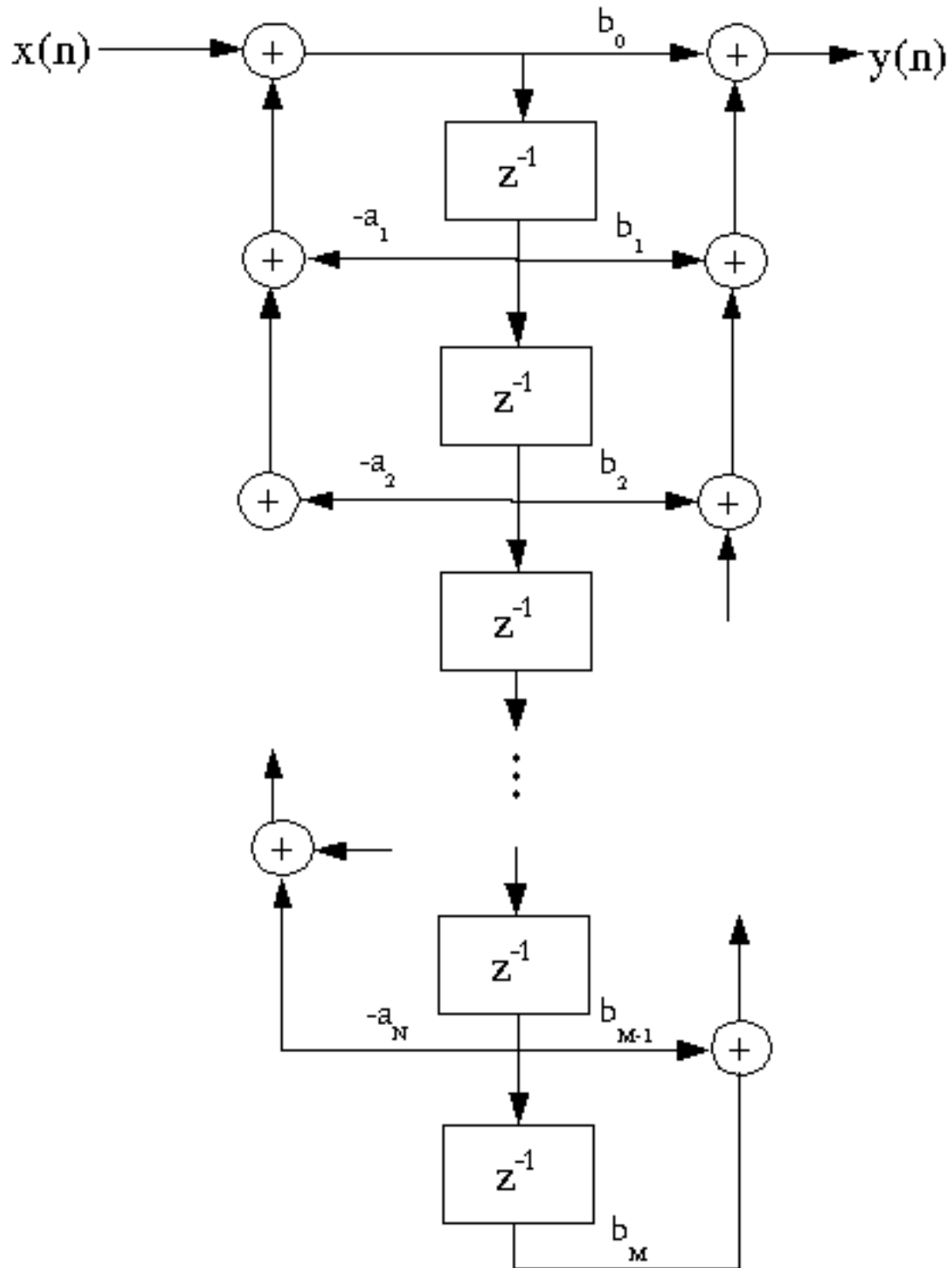


Figure 4.11

This structure is **canonic**: (i.e., it requires the minimum number of memory elements).
Flowgraph reversal gives the

4.1.3.3 Transpose-Form IIR Filter Structure

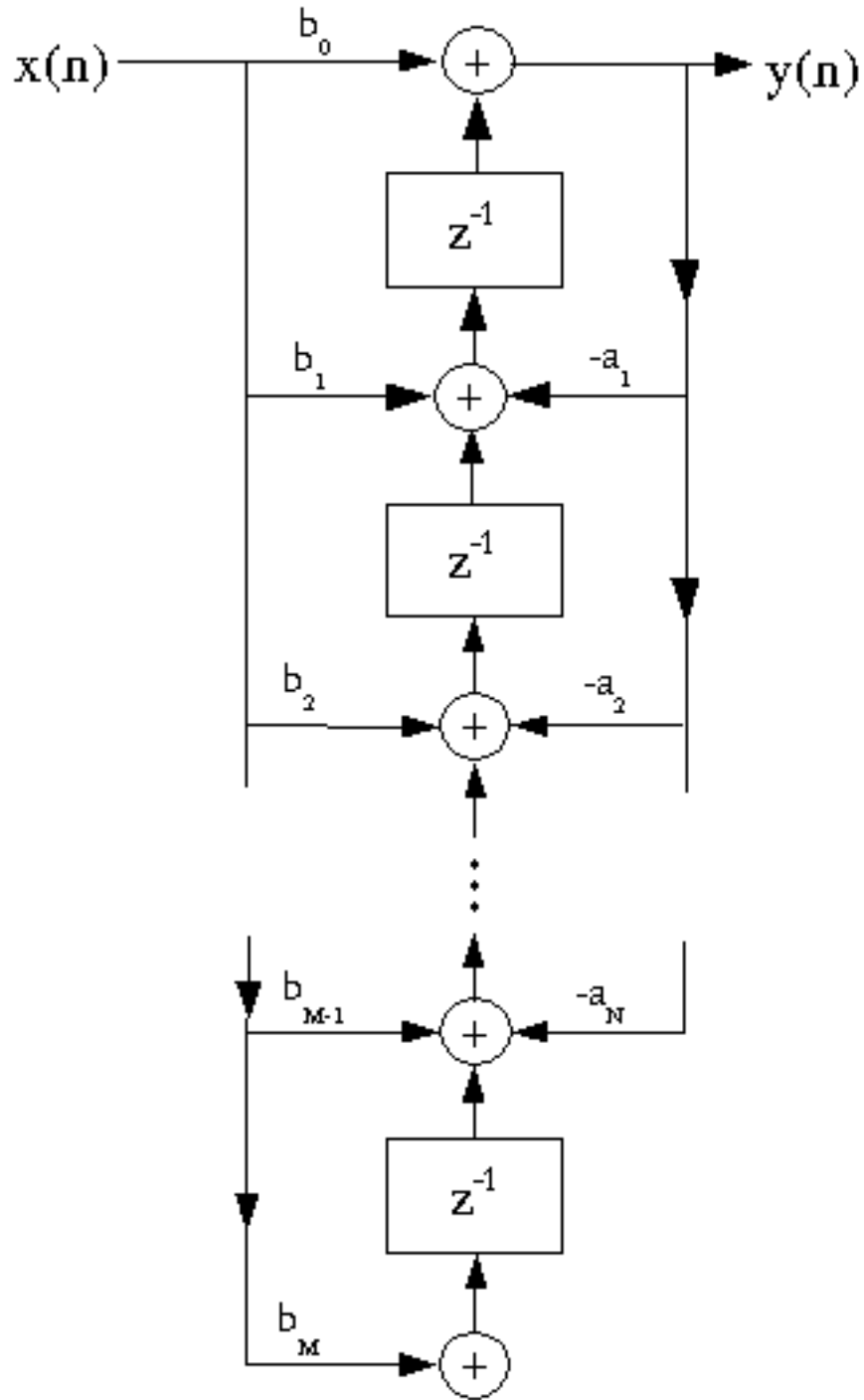


Figure 4.12

Usually we design IIR filters with $N = M$, but not always.

Obviously, since all these structures have identical frequency response, filter structures are not unique. We consider many different structures because

1. Depending on the technology or application, one might be more convenient than another
2. The response in a practical realization, in which the data and coefficients must be **quantized**, may differ substantially, and some structures behave much better than others with quantization.

The Cascade-Form IIR filter structure is one of the least sensitive to quantization, which is why it is the most commonly used IIR filter structure.

4.1.3.4 IIR Cascade Form

The numerator and denominator polynomials can be factored

$$H(z) = \frac{b_0 + b_1 z^{-1} + \dots + b_M z^{-M}}{1 + a_1 z^{-1} + \dots + a_N z^{-N}} = \frac{b_0 \prod_{k=1}^M (z - z_k)}{z^{M-N} \prod_{i=1}^N (z - p_i)}$$

and implemented as a cascade of short IIR filters.

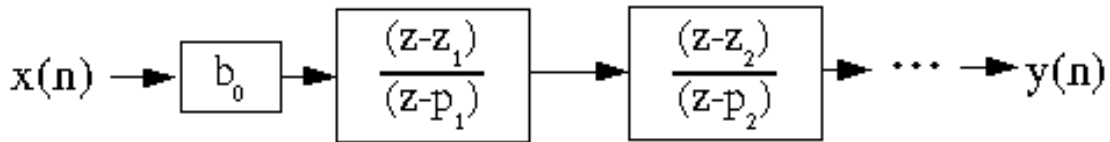


Figure 4.13

Since the filter coefficients are usually real yet the roots are mostly complex, we actually implement these as second-order sections, where complex-conjugate pole and zero pairs are combined into second-order sections with real coefficients. The second-order sections are usually implemented with either the Direct-Form II or Transpose-Form structure.

4.1.3.5 Parallel form

A rational transfer function can also be written as

$$\frac{b_0 + b_1 z^{-1} + \dots + b_M z^{-M}}{1 + a_1 z^{-1} + \dots + a_N z^{-N}} = c_0' + c_1' z^{-1} + \dots + \frac{A_1}{z - p_1} + \frac{A_2}{z - p_2} + \dots + \frac{A_N}{z - p_N}$$

which by linearity can be implemented as

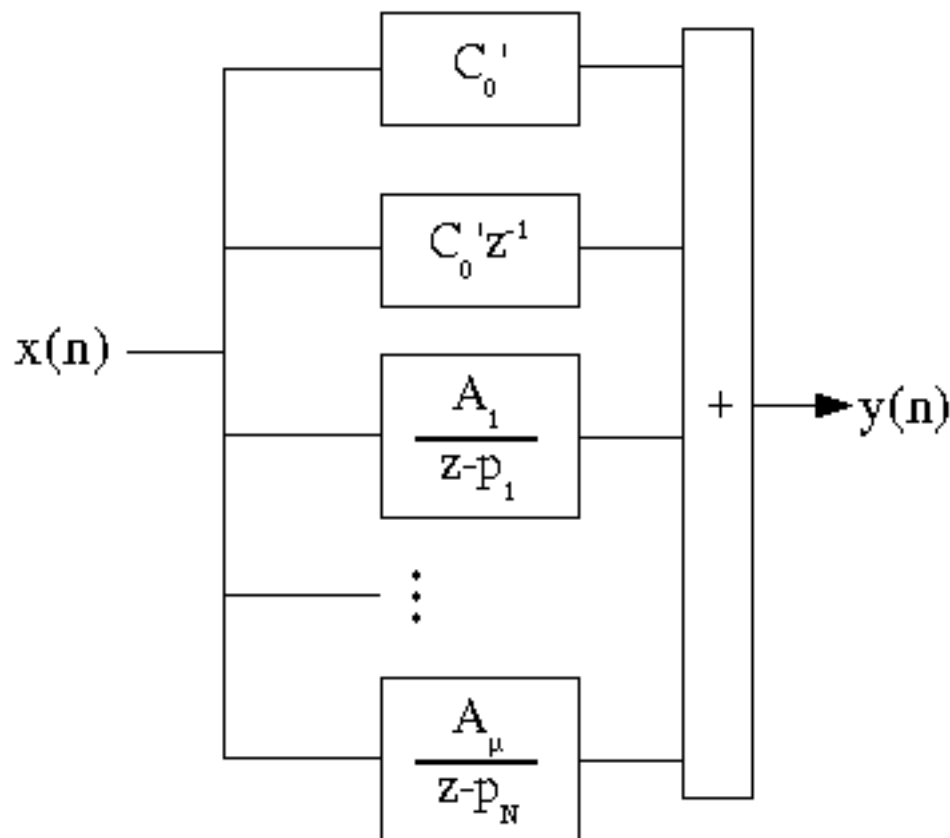


Figure 4.14

As before, we combine complex-conjugate pole pairs into second-order sections with real coefficients.

The cascade and parallel forms are of interest because they are much less sensitive to coefficient quantization than higher-order structures, as analyzed in later modules in this course.

4.1.3.6 Other forms

There are many other structures for IIR filters, such as wave digital filter structures, lattice-ladder, all-pass-based forms, and so forth. These are the result of extensive research to find structures which are computationally efficient *and* insensitive to quantization error. They all represent various tradeoffs; the best choice in a given context is not yet fully understood, and may never be.

4.1.4 State-Variable Representation of Discrete-Time Systems⁴

4.1.4.1 State and the State-Variable Representation

Definition 10: State

the minimum additional information at time n , which, along with all current and future input values, is necessary to compute all future outputs.

⁴This content is available online at <<http://cnx.org/content/m11920/1.2/>>.

Essentially, the state of a system is the information held in the delay registers in a filter structure or signal flow graph.

FACT: Any LTI (linear, time-invariant) system of finite order M can be represented by a state-variable description

$$\mathbf{x}(n+1) = \mathbf{A}\mathbf{x}(n) + \mathbf{B}u(n)$$

$$\mathbf{y}(n) = \mathbf{C}\mathbf{x}(n) + Du(n)$$

where \mathbf{x} is an $(M \times 1)$ "state vector," $u(n)$ is the input at time n , $y(n)$ is the output at time n ; \mathbf{A} is an $(M \times M)$ matrix, \mathbf{B} is an $(M \times 1)$ vector, \mathbf{C} is a $(1 \times M)$ vector, and D is a (1×1) scalar.

One can always obtain a state-variable description of a signal flow graph.

Example 4.2: 3rd-Order IIR

$$y(n) = -(a_1y(n-1)) - a_2y(n-2) - a_3y(n-3) + b_0x(n) + b_1x(n-1) + b_2x(n-2) + b_3x(n-3)$$

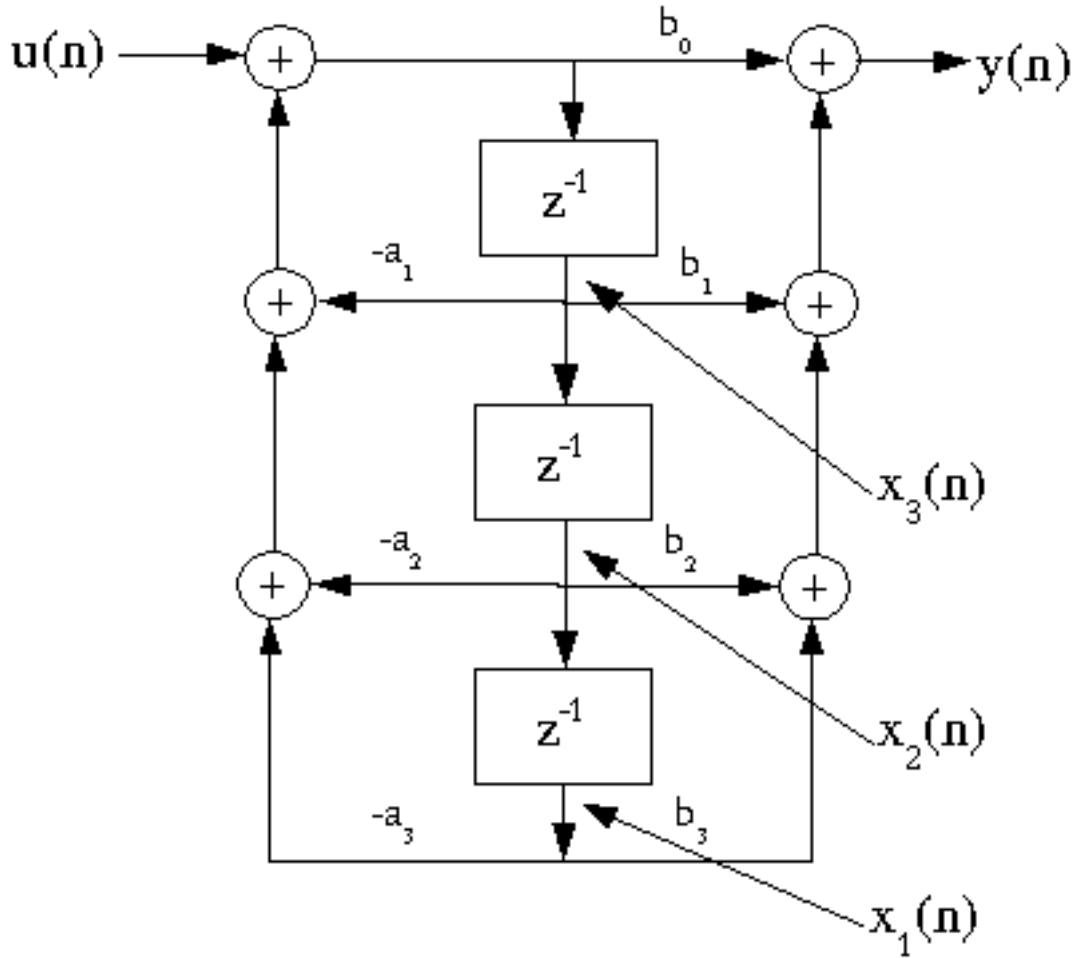


Figure 4.15

$$\begin{pmatrix} x_1(n+1) \\ x_2(n+1) \\ x_3(n+1) \end{pmatrix} = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ -a_3 & -a_2 & -a_1 \end{pmatrix} \begin{pmatrix} x_1(n) \\ x_2(n) \\ x_3(n) \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} u(n)$$

$$y(n) = \begin{pmatrix} -(a_3b_0) & -(a_2b_0) & -(a_1b_0) \end{pmatrix} \begin{pmatrix} x_1(n) \\ x_2(n) \\ x_3(n) \end{pmatrix} + \begin{pmatrix} b_0 \end{pmatrix} u(n)$$

Exercise 4.2

Is the state-variable description of a filter $H(z)$ unique?

Exercise 4.3

Does the state-variable description fully describe the signal flow graph?

4.1.4.2 State-Variable Transformation

Suppose we wish to define a new set of state variables, related to the old set by a linear transformation: $\mathbf{q}(n) = T\mathbf{x}(n)$, where T is a nonsingular ($M \times M$) matrix, and $\mathbf{q}(n)$ is the new state vector. We wish the overall system to remain the same. Note that $\mathbf{x}(n) = T^{-1}\mathbf{q}(n)$, and thus

$$(\mathbf{x}(n+1) = A\mathbf{x}(n) + Bu(n) \Rightarrow T^{-1}\mathbf{q}(n) = AT^{-1}\mathbf{q}(n) + Bu(n) \Rightarrow \mathbf{q}(n) = TAT^{-1}\mathbf{q}(n) + TBu(n))$$

$$(y(n) = C\mathbf{x}(n) + Du(n) \Rightarrow y(n) = CT^{-1}\mathbf{q}(n) + Du(n))$$

This defines a new state system with an input-output behavior identical to the old system, but with different internal memory contents (states) and state matrices.

$$\mathbf{q}(n) = \hat{A}\mathbf{q}(n) + \hat{B}u(n)$$

$$y(n) = \hat{C}\mathbf{q}(n) + \hat{D}u(n)$$

$$\hat{A} = TAT^{-1}, \hat{B} = TB, \hat{C} = CT^{-1}, \hat{D} = D$$

These transformations can be used to generate a wide variety of alternative structures or implementations of a filter.

4.1.4.3 Transfer Function and the State-Variable Description

Taking the z transform of the state equations

$$Z[\mathbf{x}(n+1)] = Z[A\mathbf{x}(n) + Bu(n)]$$

$$Z[y(n)] = Z[C\mathbf{x}(n) + Du(n)]$$

$$\Downarrow$$

$$z\mathbf{X}(z) = A\mathbf{X}(z) + BU(z)$$

NOTE: $\mathbf{X}(z)$ is a vector of scalar z -transforms $(\mathbf{X}(z))^T = \begin{pmatrix} X_1(z) & X_2(z) & \dots \end{pmatrix}$

$$\mathbf{Y}(z) = C\mathbf{X}(z) + DU(z)$$

$$\left((zI - A)\mathbf{X}(z) = BU(z) \Rightarrow \mathbf{X}(z) = (zI - A)^{-1}BU(z)\right)$$

so

$$\begin{aligned} Y(z) &= C(zI - A)^{-1}BU(z) + DU(z) \\ &= \left(C(-(zI))^{-1}B + D\right)U(z) \end{aligned} \quad (4.1)$$

and thus

$$H(z) = C(zI - A)^{-1}B + D$$

Note that since $(zI - A)^{-1} = \frac{(\pm \det(zI - A)_{red})^T}{\det(zI - A)}$, this transfer function is an M th-order rational fraction in z . The denominator polynomial is $D(z) = \det(zI - A)$. A discrete-time state system is thus stable if the M roots of $\det(zI - A)$ (i.e., the poles of the digital filter) are all inside the unit circle.

Consider the transformed state system with $\hat{A} = TAT^{-1}$, $\hat{B} = TB$, $\hat{C} = CT^{-1}$, $\hat{D} = D$:

$$\begin{aligned} H(z) &= \hat{C}(zI - \hat{A})^{-1}\hat{B} + \hat{D} \\ &= CT^{-1}(zI - TAT^{-1})^{-1}TB + D \\ &= CT^{-1}(T(zI - A)T^{-1})^{-1}TB + D \\ &= CT^{-1}(T^{-1})^{-1}(zI - A)^{-1}T^{-1}TB + D \\ &= C(zI - A)^{-1}B + D \end{aligned} \quad (4.2)$$

This proves that state-variable transformation doesn't change the transfer function of the underlying system. However, it can provide alternate forms that are less sensitive to coefficient quantization or easier to analyze, understand, or implement.

State-variable descriptions of systems are useful because they provide a fairly general tool for analyzing all systems; they provide a more detailed description of a signal flow graph than does the transfer function (although not a full description); and they suggest a large class of alternative implementations. They are even more useful in control theory, which is largely based on state descriptions of systems.

4.2 Fixed-Point Numbers

4.2.1 Fixed-Point Number Representation⁵

Fixed-point arithmetic is generally used when hardware cost, speed, or complexity is important. Finite-precision quantization issues usually arise in fixed-point systems, so we concentrate on fixed-point quantization and error analysis in the remainder of this course. For basic signal processing computations such as digital filters and FFTs, the magnitude of the data, the internal states, and the output can usually be scaled to obtain good performance with a fixed-point implementation.

4.2.1.1 Two's-Complement Integer Representation

As far as the hardware is concerned, fixed-point number systems represent data as B -bit integers. The two's-complement number system is usually used:

$$k = \begin{cases} \text{binary integer representation if } 0 \leq k \leq 2^{B-1} - 1 \\ \text{bit-by-bit inverse } (-k) + 1 \text{ if } -(2^{B-1}) \leq k \leq 0 \end{cases}$$

⁵This content is available online at <<http://cnx.org/content/m11930/1.2/>>.

	Quantized representation	Quantization error
0.10	-0.12	-0.02
+ 0.11	+ 0.12	+ 0.01
1.01	-0.04	-0.99

Figure 4.19

Note the occurrence of wraparound **overflow**; this only happens with *addition*. Obviously, it can be a bad problem.

There are thus two types of fixed-point error: roundoff error, associated with data quantization and multiplication, and overflow error, associated with data quantization and additions. In fixed-point systems, one must strike a balance between these two error sources; by scaling down the data, the occurrence of overflow errors is reduced, but the relative size of the roundoff error is increased.

NOTE: Since multiplies require a number of additions, they are especially expensive in terms of hardware (with a complexity proportional to $B_x B_h$, where B_x is the number of bits in the data, and B_h is the number of bits in the filter coefficients). Designers try to minimize both B_x and B_h , and often choose $B_x \neq B_h$!

4.2.2 Fixed-Point Quantization⁶

The fractional B -bit two's complement number representation evenly distributes 2^B quantization levels between -1 and $1 - 2^{-(B-1)}$. The spacing between quantization levels is then

$$\left(\frac{2}{2^B} = 2^{-(B-1)} \doteq \Delta_B \right)$$

Any signal value falling between two levels is assigned to one of the two levels.

$X_Q = Q[x]$ is our notation for quantization. $e = Q[x] - x$ is then the quantization error.

One method of quantization is **rounding**, which assigns the signal value to the *nearest* level. The maximum error is thus $\frac{\Delta_B}{2} = 2^{-B}$.

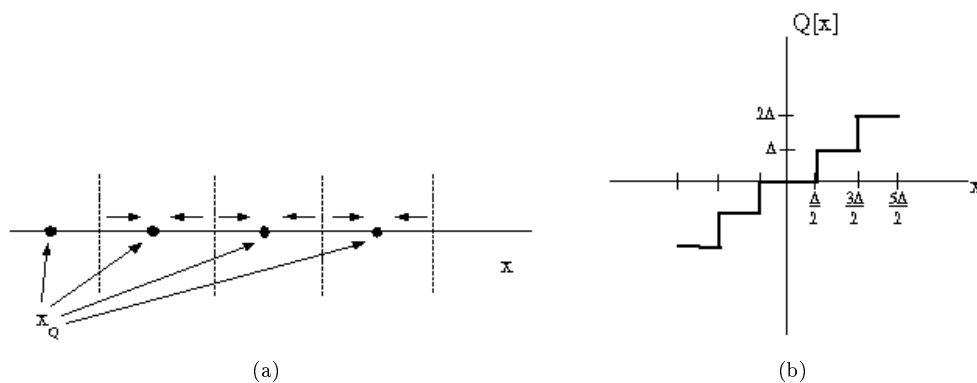


Figure 4.20

⁶This content is available online at <<http://cnx.org/content/m11921/1.2/>>.

Another common scheme, which is often easier to implement in hardware, is **truncation**. $Q[x]$ assigns x to the next lowest level.

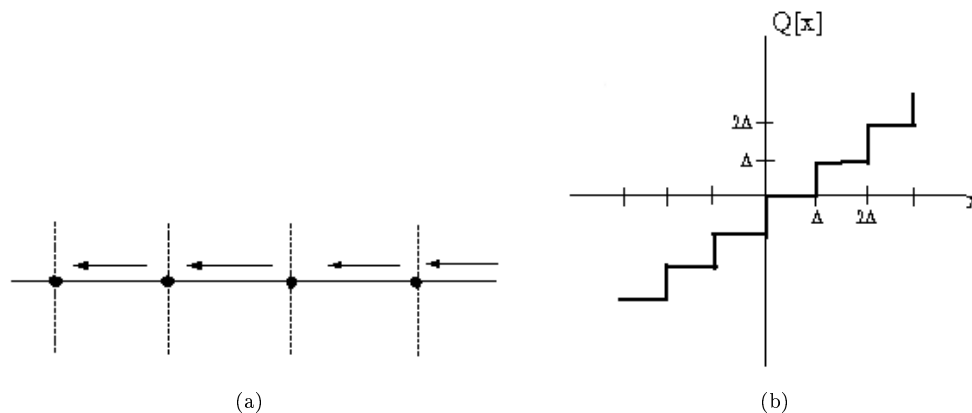


Figure 4.21

The worst-case error with truncation is $\Delta = 2^{-(B-1)}$, which is twice as large as with rounding. Also, the error is always negative, so on average it may have a non-zero mean (i.e., a bias component).

Overflow is the other problem. There are two common types: two's complement (or **wraparound**) overflow, or **saturation** overflow.

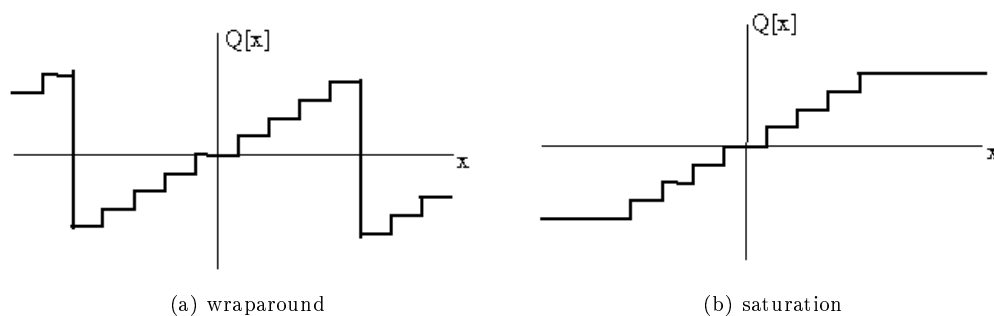


Figure 4.22

Obviously, overflow errors are bad because they are typically *large*; two's complement (or wraparound) overflow introduces more error than saturation, but is easier to implement in hardware. It also has the advantage that if the *sum* of several numbers is between $[-1, 1)$, the final answer will be correct even if intermediate sums overflow! However, wraparound overflow leaves IIR systems susceptible to zero-input large-scale limit cycles, as discussed in another module. As usual, there are many tradeoffs to evaluate, and no one right answer for all applications.

4.3 Quantization Error Analysis

4.3.1 Finite-Precision Error Analysis⁷

4.3.1.1 Fundamental Assumptions in finite-precision error analysis

Quantization is a highly nonlinear process and is very difficult to analyze precisely. Approximations and assumptions are made to make analysis tractable.

4.3.1.1.1 Assumption #1

The roundoff or truncation errors at any point in a system at each time are *random*, *stationary*, and *statistically independent* (white and independent of all other quantizers in a system).

That is, the error autocorrelation function is $r_e[k] = E[e_n e_{n+k}] = \sigma_q^2 \delta[k]$. Intuitively, and confirmed experimentally in some (but not all!) cases, one expects the quantization error to have a uniform distribution over the interval $[-(\frac{\Delta}{2}), (\frac{\Delta}{2})$ for rounding, or $(-\Delta, 0]$ for truncation.

In this case, rounding has zero mean and variance

$$E[Q[x_n] - x_n] = 0$$

$$\sigma_Q^2 = E[e_n^2] = \frac{\Delta_B^2}{12}$$

and truncation has the statistics

$$E[Q[x_n] - x_n] = -\left(\frac{\Delta}{2}\right)$$

$$\sigma_Q^2 = \frac{\Delta_B^2}{12}$$

Please note that the independence assumption may be very bad (for example, when quantizing a sinusoid with an integer period N). There is another quantizing scheme called **dithering**, in which the values are randomly assigned to nearby quantization levels. This can be (and often is) implemented by adding a small (one- or two-bit) random input to the signal before a truncation or rounding quantizer.

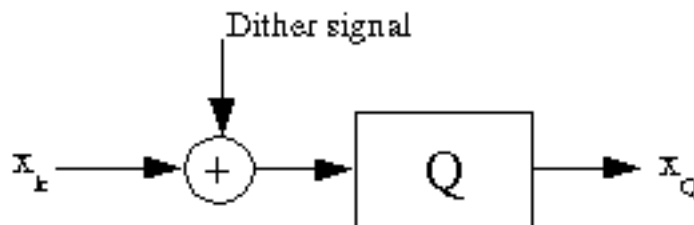


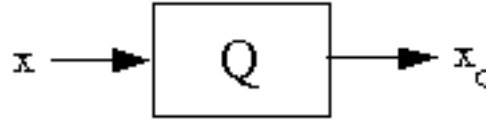
Figure 4.23

This is used extensively in practice. Although the overall error is somewhat higher, it is spread evenly over all frequencies, rather than being concentrated in spectral lines. This is very important when quantizing sinusoidal or other periodic signals, for example.

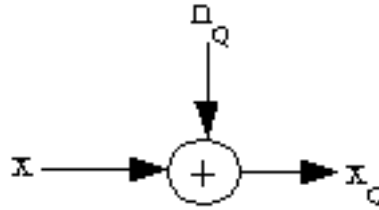
⁷This content is available online at <<http://cnx.org/content/m11922/1.2/>>.

4.3.1.1.2 Assumption #2

Pretend that the quantization error is really additive *Gaussian* noise with the same mean and variance as the uniform quantizer. That is, model



(a)



(b) as

Figure 4.24

This model is a *linear* system, which our standard theory can handle easily. We model the noise as Gaussian because it remains Gaussian after passing through filters, so analysis in a system context is tractable.

4.3.1.2 Summary of Useful Statistical Facts

- **correlation function** - ($r_x[k] \doteq E[x_n x_{n+k}]$)
- **power spectral density** - ($S_x(w) \doteq DTFT[r_x[n]]$)
- Note $r_x[0] = \sigma_x^2 = \frac{1}{2\pi} \int_{-\pi}^{\pi} S_x(w) dw$
- ($r_{xy}[k] \doteq E[x^*[n] y[n+k]]$)
- **cross-spectral density** - $S_{xy}(w) = DTFT[r_{xy}[n]]$
- For $y = h * x$:

$$S_{yx}(w) = H(w) S_x(w)$$

$$S_{yy}(w) = (|H(w)|)^2 S_x(w)$$

- Note that the *output* noise level after filtering a noise sequence is

$$\sigma_y^2 = r_{yy}[0] = \frac{1}{\pi} \int_{-\pi}^{\pi} (|H(w)|)^2 S_x(w) dw$$

so postfiltering quantization noise alters the noise power spectrum and may change its variance!

- For x_1, x_2 statistically independent

$$r_{x_1+x_2}[k] = r_{x_1}[k] + r_{x_2}[k]$$

$$S_{x_1+x_2}(w) = S_{x_1}(w) + S_{x_2}(w)$$

- For independent random variables

$$\sigma_{x_1+x_2}^2 = \sigma_{x_1}^2 + \sigma_{x_2}^2$$

4.3.2 Input Quantization Noise Analysis⁸

All practical analog-to-digital converters (A/D) must quantize the input data. This can be modeled as an ideal sampler followed by a B -bit quantizer.

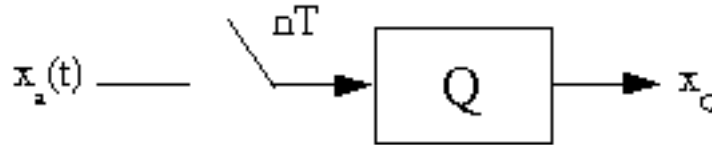


Figure 4.25

The signal-to-noise ratio (SNR) of an A/D is

$$\begin{aligned} SNR &= 10\log_{10}\left(\frac{P_x}{P_n}\right) \\ &= 10\log_{10}P_x - 10\log_{10}\left(\frac{\Delta_B^2}{12}\right) \\ &= 10\log_{10}P_x + 4.77 + 6.02B \end{aligned} \quad (4.3)$$

where P_x is the power in the signal and P_n is the power of the quantization noise, which equals its variance if it has a zero mean. The SNR increases by 6dB with each additional bit.

4.3.3 Quantization Error in FIR Filters⁹

In digital filters, both the data at various places in the filter, which are continually varying, and the coefficients, which are fixed, must be quantized. The effects of quantization on data and coefficients are quite different, so they are analyzed separately.

4.3.3.1 Data Quantization

Typically, the input and output in a digital filter are quantized by the analog-to-digital and digital-to-analog converters, respectively. Quantization also occurs at various points in a filter structure, usually after a multiply, since multiplies increase the number of bits.

4.3.3.1.1 Direct-form Structures

There are two common possibilities for quantization in a direct-form FIR filter structure: after each multiply, or only once at the end.

⁸This content is available online at <<http://cnx.org/content/m11923/1.2/>>.

⁹This content is available online at <<http://cnx.org/content/m11924/1.2/>>.

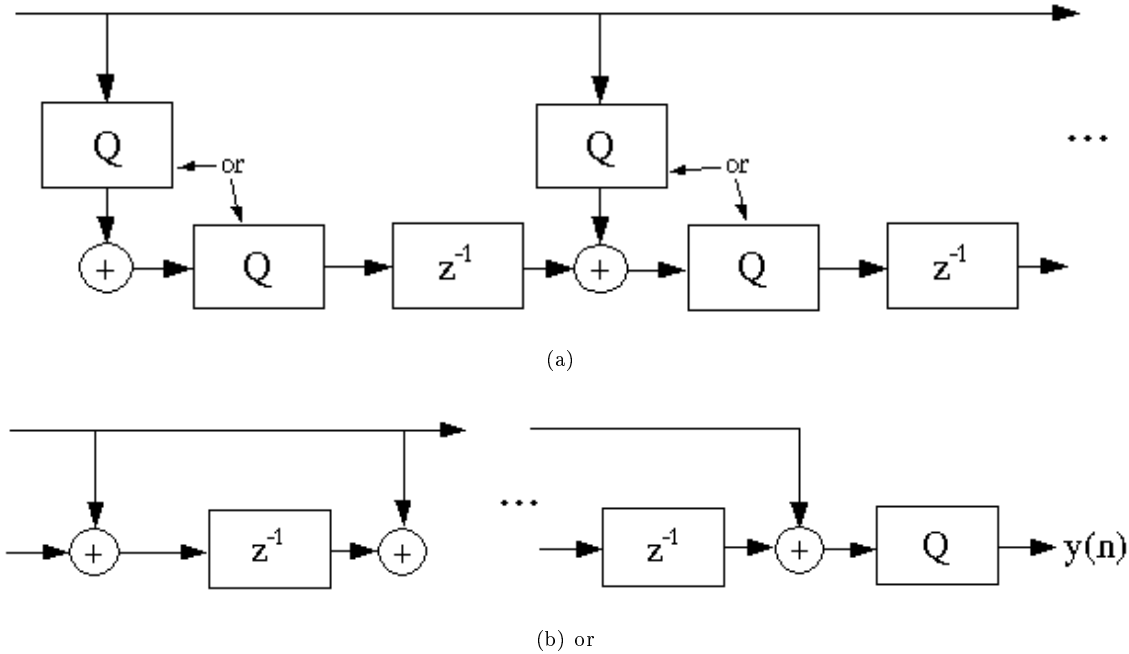


Figure 4.27: (a) Quantize at each stage before storing intermediate sum. Output variance = $M \frac{\Delta^2}{12}$
 (b) Store double-precision partial sums. Costs more memory, but variance = $\frac{\Delta^2}{12}$

The transpose form is not as convenient in terms of supporting double-precision accumulation, which is a significant disadvantage of this structure.

4.3.3.2 Coefficient Quantization

Since a quantized coefficient is fixed for all time, we treat it differently than data quantization. The fundamental question is: how much does the quantization affect the frequency response of the filter?

The quantized filter frequency response is

$$DTFT[h_Q] = DTFT[h_{inf.prec.} + e] = H_{inf.prec.}(w) + H_e(w)$$

Assuming the quantization model is correct, $H_e(w)$ should be fairly random and white, with the error spread fairly equally over all frequencies $w \in [-\pi, \pi]$; however, the randomness of this error destroys any equiripple property or any infinite-precision optimality of a filter.

Exercise 4.4

What quantization scheme minimizes the L_2 quantization error in frequency (minimizes $\int_{-\pi}^{\pi} (|H(w) - H_Q(w)|)^2 dw$)? On average, how big is this error?

Ideally, if one knows the coefficients are to be quantized to B bits, one should incorporate this directly into the filter design problem, and find the M B -bit binary fractional coefficients minimizing the maximum deviation (L_∞ error). This can be done, but it is an integer program, which is known to be np-hard (i.e., requires almost a brute-force search). This is so expensive computationally that it's rarely done. There are some sub-optimal methods that are much more efficient and usually produce pretty good results.

4.3.4 Data Quantization in IIR Filters¹⁰

Finite-precision effects are much more of a concern with IIR filters than with FIR filters, since the effects are more difficult to analyze and minimize, coefficient quantization errors can cause the filters to become unstable, and disastrous things like large-scale limit cycles can occur.

4.3.4.1 Roundoff noise analysis in IIR filters

Suppose there are several quantization points in an IIR filter structure. By our simplifying assumptions about quantization error and Parseval's theorem, the quantization noise variance $\sigma_{y,i}^2$ at the output of the filter from the i th quantizer is

$$\begin{aligned}\sigma_{y,i}^2 &= \frac{1}{2\pi} \int_{-\pi}^{\pi} (|H_i(w)|)^2 S_{n_i}(w) dw \\ &= \frac{\sigma_{n_i}^2}{2\pi} \int_{-\pi}^{\pi} (|H_i(w)|)^2 dw \\ &= \sigma_{n_i}^2 \sum_{n=-\infty}^{\infty} (h_i^2(n))\end{aligned}\tag{4.4}$$

where $\sigma_{n_i}^2$ is the variance of the quantization error at the i th quantizer, $S_{n_i}(w)$ is the power spectral density of that quantization error, and $H_i(w)$ is the transfer function from the i th quantizer to the output point. Thus for P independent quantizers in the structure, the total quantization noise variance is

$$\sigma_y^2 = \frac{1}{2\pi} \sum_{i=1}^P \left(\sigma_{n_i}^2 \int_{-\pi}^{\pi} (|H_i(w)|)^2 dw \right)$$

Note that in general, each $H_i(w)$, and thus the variance at the output due to each quantizer, is different; for example, the system as seen by a quantizer at the input to the first delay state in the Direct-Form II IIR filter structure to the output, call it n_4 , is

¹⁰This content is available online at <<http://cnx.org/content/m11925/1.2/>>.

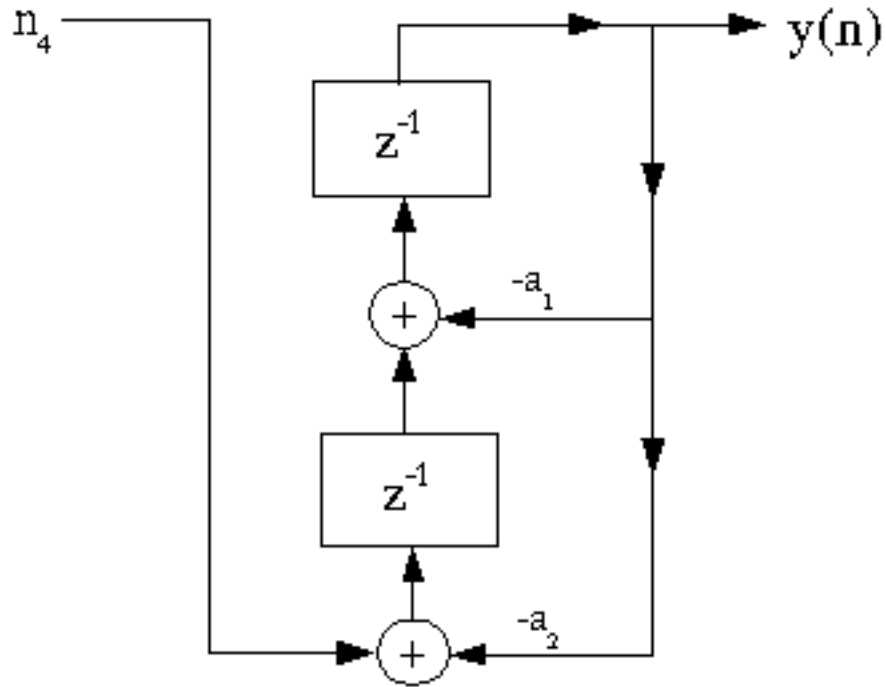


Figure 4.28

with a transfer function

$$H_4(z) = \frac{z^{-2}}{1 + a_1 z^{-1} + a_2 z^{-2}}$$

which can be evaluated at $z = e^{j\omega}$ to obtain the frequency response.

A general approach to find $H_i(w)$ is to write state equations for the equivalent structure as seen by n_i , and to determine the transfer function according to $H(z) = C(zI - A)^{-1}B + d$.

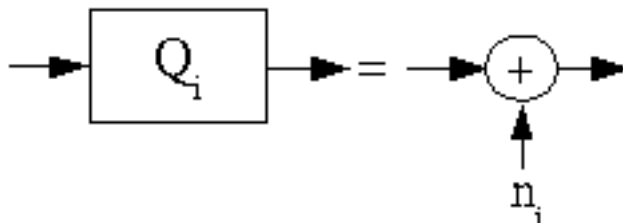


Figure 4.30

the variance at the output is the sum of the variances at the output due to each noise source:

$$\sigma_y^2 = \sum_{i=1}^4 (\sigma_{y,i}^2)$$

The variance due to each noise source at the output can be determined from $\frac{1}{2\pi} \int_{-\pi}^{\pi} (|H_i(w)|)^2 S_{n_i}(w) dw$; note that $S_{n_i}(w) = \sigma_{n_i}^2$ by our assumptions, and $H_i(w)$ is the transfer function *from the noise source to the output*.

4.3.5 IIR Coefficient Quantization Analysis¹¹

Coefficient quantization is an important concern with IIR filters, since straightforward quantization often yields poor results, and because quantization can produce unstable filters.

4.3.5.1 Sensitivity analysis

The performance and stability of an IIR filter depends on the pole locations, so it is important to know how quantization of the filter coefficients a_k affects the pole locations p_j . The denominator polynomial is

$$D(z) = 1 + \sum_{k=1}^N (a_k z^{-k}) = \prod_{i=1}^N (1 - p_i z^{-1})$$

We wish to know $\frac{\partial}{\partial a_k}(p_i)$, which, for small deviations, will tell us that a δ change in a_k yields an $\epsilon = \delta \frac{\partial}{\partial a_k}(p_i)$ change in the pole location. $\frac{\partial}{\partial a_k}(p_i)$ is the **sensitivity** of the pole location to quantization of a_k . We can find $\frac{\partial}{\partial a_k}(p_i)$ using the chain rule.

$$\frac{\partial}{\partial a_k} A(z) \big|_{z=p_i} = \frac{\partial}{\partial z} A(z) \frac{\partial}{\partial a_k} (z) \big|_{z=p_i}$$

\Downarrow

$$\frac{\partial}{\partial a_k} (p_i) = \frac{\frac{\partial}{\partial a_k} A(z_i) \big|_{z=p_i}}{\frac{\partial}{\partial z} A(z_i) \big|_{z=p_i}}$$

¹¹This content is available online at <<http://cnx.org/content/m11926/1.2/>>.

which is

$$\begin{aligned}\frac{\partial}{\partial a_k}(p_i) &= \frac{z^{-k}}{-\left(z^{-1} \prod_{\substack{j=1 \\ j \neq i}}^N (1 - p_j z^{-1})\right)} \Big|_{z=p_i} \\ &= \frac{-(p_i^{N-k})}{\prod_{\substack{j=1 \\ j \neq i}}^N (p_j - p_i)}\end{aligned}\quad (4.5)$$

Note that as the poles get closer together, the sensitivity increases greatly. So as the filter order increases and more poles get stuffed closer together inside the unit circle, the error introduced by coefficient quantization in the pole locations grows rapidly.

How can we reduce this high sensitivity to IIR filter coefficient quantization?

4.3.5.1.1 Solution

Cascade (Section 4.1.3.4: IIR Cascade Form) or parallel form (Section 4.1.3.5: Parallel form) implementations! The numerator and denominator polynomials can be factored off-line at very high precision and grouped into second-order sections, which are then quantized section by section. The sensitivity of the quantization is thus that of second-order, rather than N -th order, polynomials. This yields major improvements in the frequency response of the overall filter, and is almost always done in practice.

Note that the numerator polynomial faces the same sensitivity issues; the *cascade* form also improves the sensitivity of the zeros, because they are also factored into second-order terms. However, in the *parallel* form, the zeros are globally distributed across the sections, so they suffer from quantization of all the blocks. Thus the *cascade* form preserves zero locations much better than the parallel form, which typically means that the stopband behavior is better in the cascade form, so it is most often used in practice.

NOTE ON FIR FILTERS: On the basis of the preceding analysis, it would seem important to use cascade structures in FIR filter implementations. However, most FIR filters are linear-phase and thus symmetric or anti-symmetric. As long as the quantization is implemented such that the filter coefficients retain symmetry, the filter retains linear phase. Furthermore, since all zeros off the unit circle must appear in groups of four for symmetric linear-phase filters, zero pairs can leave the unit circle only by joining with another pair. This requires relatively severe quantizations (enough to completely remove or change the sign of a ripple in the amplitude response). This "reluctance" of pole pairs to leave the unit circle tends to keep quantization from damaging the frequency response as much as might be expected, enough so that cascade structures are rarely used for FIR filters.

Exercise 4.6

(Solution on p. 259.)

What is the worst-case pole pair in an IIR digital filter?

4.3.5.2 Quantized Pole Locations

In a direct-form (Section 4.1.3.1: Direct-form I IIR Filter Structure) or transpose-form (Section 4.1.3.3: Transpose-Form IIR Filter Structure) implementation of a second-order section, the filter coefficients are quantized versions of the polynomial coefficients.

$$D(z) = z^2 + a_1 z + a_2 = (z - p)(z - \bar{p})$$

$$p = \frac{(-a_1 \pm \sqrt{a_1^2 - 4a_2})}{2}$$

$$p = r e^{j\theta}$$

$$D(z) = z^2 - 2r \cos(\theta) z + r^2$$

So

$$a_1 = -(2r \cos(\theta))$$

$$a_2 = r^2$$

Thus the quantization of a_1 and a_2 to B bits restricts the radius r to $r = \sqrt{k\Delta_B}$, and $a_1 = -(2\Re(p)) = k\Delta_B$. The following figure shows all stable pole locations after four-bit two's-complement quantization.

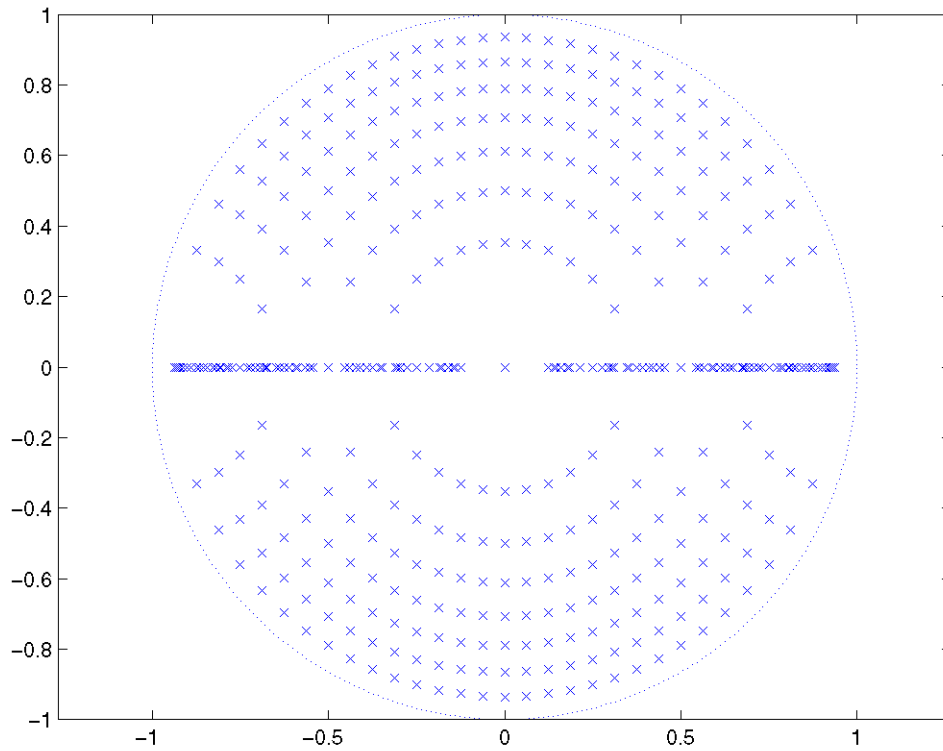
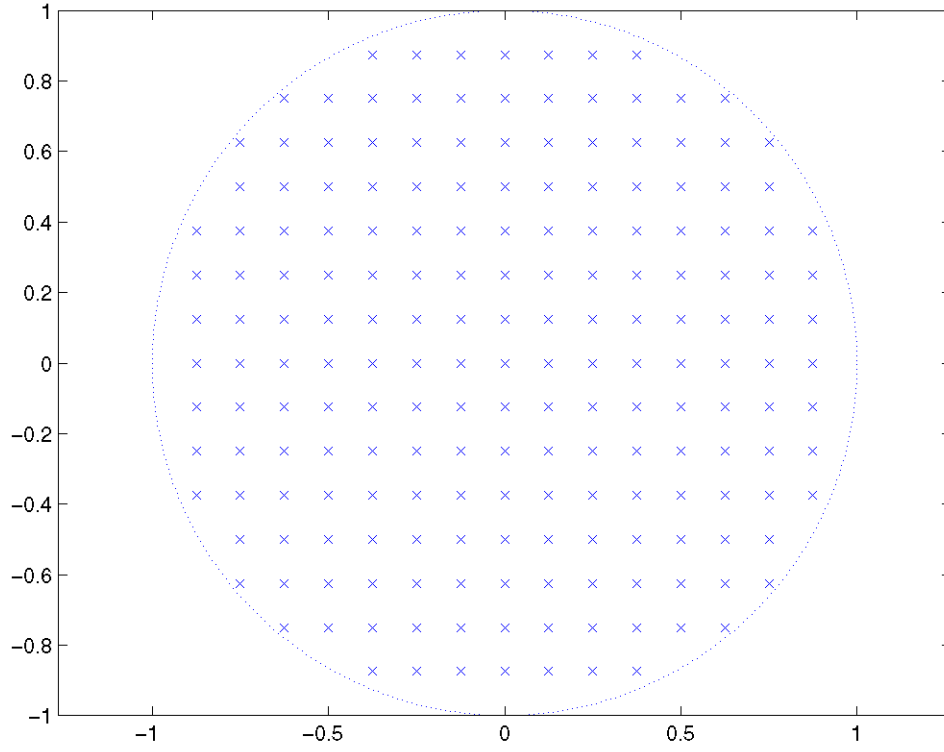


Figure 4.31

Note the nonuniform distribution of possible pole locations. This might be *good* for poles near $r = 1$, $\theta = \frac{\pi}{2}$, but not so good for poles near the origin or the Nyquist frequency.

In the "normal-form" structures, a state-variable (Definition: "State", p. 236) based realization, the poles are uniformly spaced.

**Figure 4.32**

This can only be accomplished if the coefficients to be quantized equal the real and imaginary parts of the pole location; that is,

$$\alpha_1 = r \cos(\theta) = \Re(r)$$

$$\alpha_2 = r \sin(\theta) = \Im(p)$$

This is the case for a 2nd-order system with the state matrix (Section 4.1.4.1: State and the State-Variable Representation) $A = \begin{pmatrix} \alpha_1 & \alpha_2 \\ -\alpha_1 & \alpha_1 \end{pmatrix}$: The denominator polynomial is

$$\begin{aligned} \det(zI - A) &= (z - \alpha_1)^2 + \alpha_2^2 \\ &= z^2 - 2\alpha_1 z + \alpha_1^2 + \alpha_2^2 \\ &= z^2 - 2r \cos(\theta) z + r^2 (\cos^2(\theta) + \sin^2(\theta)) \\ &= z^2 - 2r \cos(\theta) z + r^2 \end{aligned} \tag{4.6}$$

Given any second-order filter coefficient set, we can write it as a state-space system (Section 4.1.4.1: State and the State-Variable Representation), find a transformation matrix (Section 4.1.4.2: State-Variable Trans-

formation) T such that $\hat{A} = T^{-1}AT$ is in normal form, and then implement the second-order section using a structure corresponding to the state equations.

The normal form has a number of other advantages; both eigenvalues are equal, so it minimizes the norm of Ax , which makes overflow less likely, and it minimizes the output variance due to quantization of the state values. It is sometimes used when minimization of finite-precision effects is critical.

Exercise 4.7

(Solution on p. 259.)

What is the disadvantage of the normal form?

4.4 Overflow Problems and Solutions

4.4.1 Limit Cycles¹²

4.4.1.1 Large-scale limit cycles

When overflow occurs, even otherwise stable filters may get stuck in a **large-scale limit cycle**, which is a short-period, almost full-scale persistent filter output caused by overflow.

Example 4.3

Consider the second-order system

$$H(z) = \frac{1}{1 - z^{-1} + \frac{1}{2}z^{-2}}$$

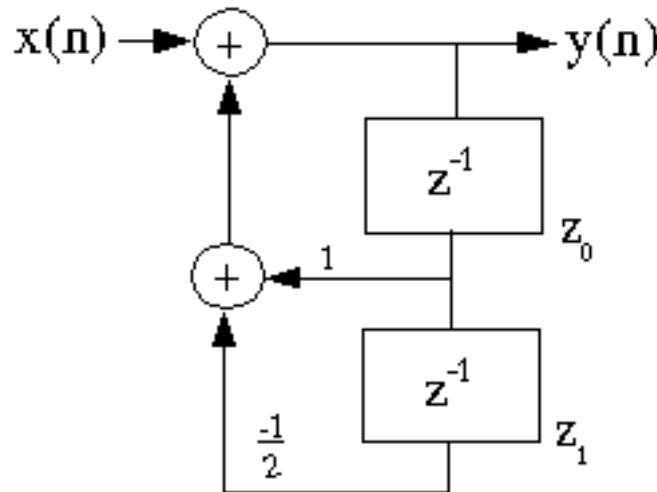


Figure 4.33

with zero input and initial state values $z_0[0] = 0.8$, $z_1[0] = -0.8$. Note $y[n] = z_0[n+1]$.

The filter is obviously stable, since the magnitude of the poles is $\frac{1}{\sqrt{2}} = 0.707$, which is well inside the unit circle. However, with wraparound overflow, note that $y[0] = z_0[1] = \frac{4}{5} - \frac{1}{2}(-\frac{4}{5}) = \frac{6}{5} =$

¹²This content is available online at <<http://cnx.org/content/m11928/1.2/>>.

$-\left(\frac{4}{5}\right)$, and that $z_0[2] = y[1] = -\left(\frac{4}{5}\right) - \frac{1}{2}\frac{4}{5} = -\left(\frac{6}{5}\right) = \frac{4}{5}$, so $y[n] = -\left(\frac{4}{5}\right), \frac{4}{5}, -\left(\frac{4}{5}\right), \frac{4}{5}, \dots$ even with zero input.

Clearly, such behavior is intolerable and must be prevented. Saturation arithmetic has been proved to prevent **zero-input limit cycles**, which is one reason why all DSP microprocessors support this feature. In many applications, this is considered sufficient protection. Scaling to prevent overflow is another solution, if as well the initial state values are never initialized to limit-cycle-producing values. The normal-form structure also reduces the chance of overflow.

4.4.1.2 Small-scale limit cycles

Small-scale limit cycles are caused by quantization. Consider the system

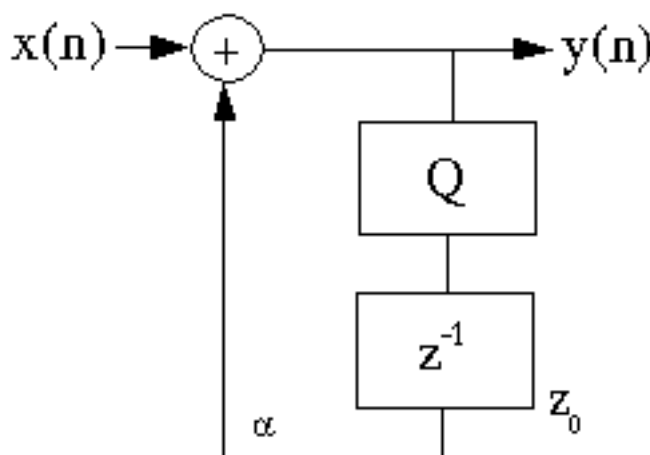


Figure 4.34

Note that when $\alpha z_0 > z_0 - \frac{\Delta_B}{2}$, rounding will quantize the output to the current level (with zero input), so the output will remain at this level forever. Note that the maximum amplitude of this "small-scale limit cycle" is achieved when

$$\alpha z_0 = z_0 - \frac{\Delta_B}{2} \Rightarrow z_{max} = \frac{\Delta_B}{2(1-\alpha)}$$

In a higher-order system, the small-scale limit cycles are oscillatory in nature. Any quantization scheme that never increases the magnitude of any quantized value prevents small-scale limit cycles.

NOTE: Two's-complement truncation does *not* do this; it increases the magnitude of negative numbers.

However, this introduces greater error and bias. Since the level of the limit cycles is proportional to Δ_B , they can be reduced by increasing the number of bits. Poles close to the unit circle increase the magnitude and likelihood of small-scale limit cycles.

4.4.2 Scaling¹³

Overflow is clearly a serious problem, since the errors it introduces are very large. As we shall see, it is also responsible for large-scale limit cycles, which cannot be tolerated. One way to prevent overflow, or to render

¹³This content is available online at <<http://cnx.org/content/m11927/1.2/>>.

it acceptably unlikely, is to **scale** the input to a filter such that overflow cannot (or is sufficiently unlikely to) occur.

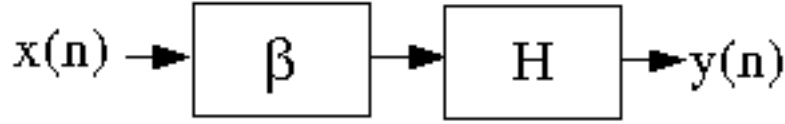


Figure 4.35

In a fixed-point system, the range of the input signal is limited by the fractional fixed-point number representation to $|x[n]| \leq 1$. If we scale the input by multiplying it by a value β , $0 < \beta < 1$, then $|\beta x[n]| \leq \beta$.

Another option is to incorporate the scaling directly into the filter coefficients.

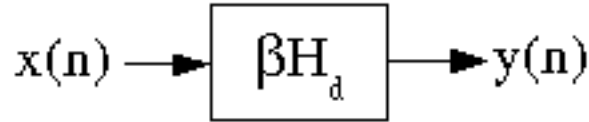


Figure 4.36

4.4.2.1 FIR Filter Scaling

What value of β is required so that the output of an FIR filter cannot overflow ($\forall n : (|y(n)| \leq 1)$, $\forall n : (|x(n)| \leq 1)$)?

$$|y(n)| = \left| \sum_{k=0}^{M-1} (h(k) \beta x(n-k)) \right| \leq \sum_{k=0}^{M-1} (|h(k)| |\beta| |x(n-k)|) \leq \beta \sum_{k=0}^{M-1} (|h(k)| \cdot 1)$$

\Downarrow

$$\beta < \sum_{k=0}^{M-1} (|h(k)|)$$

Alternatively, we can incorporate the scaling directly into the filter, and require that

$$\sum_{k=0}^{M-1} (|h(k)|) < 1$$

to prevent overflow.

4.4.2.2 IIR Filter Scaling

To prevent the output from overflowing in an IIR filter, the condition above still holds: ($M = \infty$)

$$|y(n)| < \sum_{k=0}^{\infty} (|h(k)|)$$

so an initial scaling factor $\beta < \frac{1}{\sum_{k=0}^{\infty} (|h(k)|)}$ can be used, or the filter itself can be scaled.

However, it is also necessary to prevent the *states* from overflowing, and to prevent overflow at any point in the signal flow graph where the arithmetic hardware would thereby produce errors. To prevent the states from overflowing, we determine the transfer function from the input to all states i , and scale the filter such that $\forall i : (\sum_{k=0}^{\infty} (|h_i(k)|) \leq 1)$

Although this method of scaling guarantees no overflows, it is often too conservative. Note that a worst-case signal is $x(n) = \text{sign}(h(-n))$; this input may be extremely unlikely. In the relatively common situation in which the input is expected to be mainly a single-frequency sinusoid of unknown frequency and amplitude less than 1, a scaling condition of

$$\forall w : (|H(w)| \leq 1)$$

is sufficient to guarantee no overflow. This scaling condition is often used. If there are several potential overflow locations i in the digital filter structure, the scaling conditions are

$$\forall i, w : (|H_i(w)| \leq 1)$$

where $H_i(w)$ is the frequency response from the input to location i in the filter.

Even this condition may be excessively conservative, for example if the input is more-or-less random, or if occasional overflow can be tolerated. In practice, experimentation and simulation are often the best ways to optimize the scaling factors in a given application.

For filters implemented in the cascade form, rather than scaling for the entire filter at the beginning, (which introduces lots of quantization of the input) the filter is usually scaled so that each stage is just prevented from overflowing. This is best in terms of reducing the quantization noise. The scaling factors are incorporated either into the previous or the next stage, whichever is most convenient.

Some heuristic rules for grouping poles and zeros in a cascade implementation are:

1. Order the poles in terms of decreasing radius. Take the pole pair closest to the unit circle and group it with the zero pair closest to that pole pair (to minimize the gain in that section). Keep doing this with all remaining poles and zeros.
2. Order the section with those with highest gain ($\text{argmax} |H_i(w)|$) in the middle, and those with lower gain on the ends.

Leland B. Jackson[20] has an excellent intuitive discussion of finite-precision problems in digital filters. The book by *Roberts and Mullis*[26] is one of the most thorough in terms of detail.

Solutions to Exercises in Chapter 4

Solution to Exercise 4.6 (p. 252)

The pole pair closest to the real axis in the z -plane, since the complex-conjugate poles will be closest together and thus have the highest sensitivity to quantization.

Solution to Exercise 4.7 (p. 255)

It requires more computation. The general state-variable equation (Definition: "State", p. 236) requires nine multiplies, rather than the five used by the Direct-Form II (Section 4.1.3.2: Direct-Form II IIR Filter Structure) or Transpose-Form (Section 4.1.3.3: Transpose-Form IIR Filter Structure) structures.

Chapter 5

Adaptive Filters and Applications

5.1 Introduction to Adaptive Filters¹

In many applications requiring filtering, the necessary frequency response may not be known beforehand, or it may vary with time. (Example; suppression of engine harmonics in a car stereo.) In such applications, an adaptive filter which can automatically design itself and which can track system variations in time is extremely useful. Adaptive filters are used extensively in a wide variety of applications, particularly in telecommunications.

Outline of adaptive filter material

1. **Wiener Filters** - L_2 optimal (FIR) filter design in a statistical context
2. **LMS algorithm** - simplest and by-far-the-most-commonly-used adaptive filter algorithm
3. **Stability and performance of the LMS algorithm** - When and how well it works
4. **Applications of adaptive filters** - Overview of important applications
5. **Introduction to advanced adaptive filter algorithms** - Techniques for special situations or faster convergence

5.2 Wiener Filter Algorithm

5.2.1 Discrete-Time, Causal Wiener Filter²

Stochastic L_2 optimal (least squares) FIR filter design problem: Given a wide-sense stationary (WSS) input signal x_k and desired signal d_k ($\text{WSS} \Leftrightarrow E[y_k] = E[y_{k+d}], r_{yz}(l) = E[y_k z_{k+l}], \forall k, l : (r_{yy}(0) < \infty)$)

¹This content is available online at <<http://cnx.org/content/m11535/1.3/>>.

²This content is available online at <<http://cnx.org/content/m11825/1.1/>>.

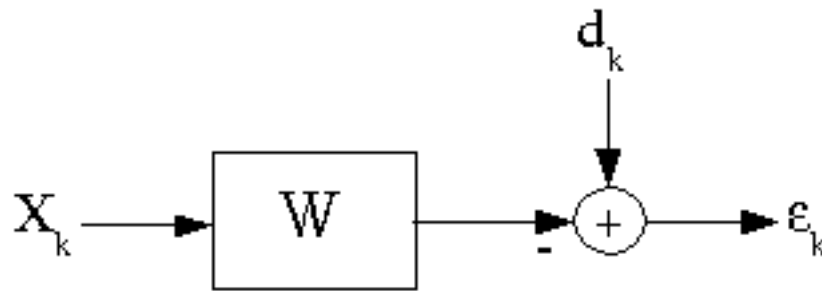


Figure 5.1

The Wiener filter is the linear, time-invariant filter minimizing $E[\epsilon^2]$, the variance of the error.

As posed, this problem seems slightly silly, since d_k is already available! However, this idea is useful in a wide variety of applications.

Example 5.1

active suspension system design

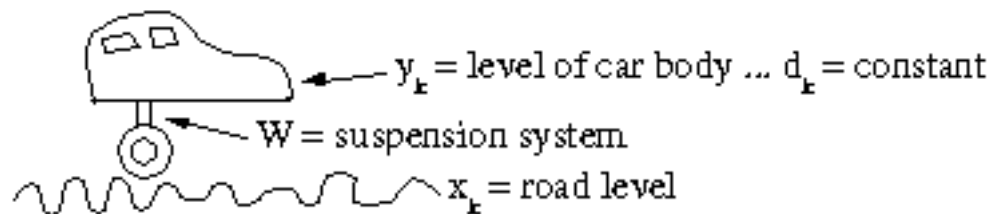


Figure 5.2

NOTE: optimal system may change with different road conditions or mass in car, so an **adaptive** system might be desirable.

Example 5.2

System identification (radar, non-destructive testing, adaptive control systems)

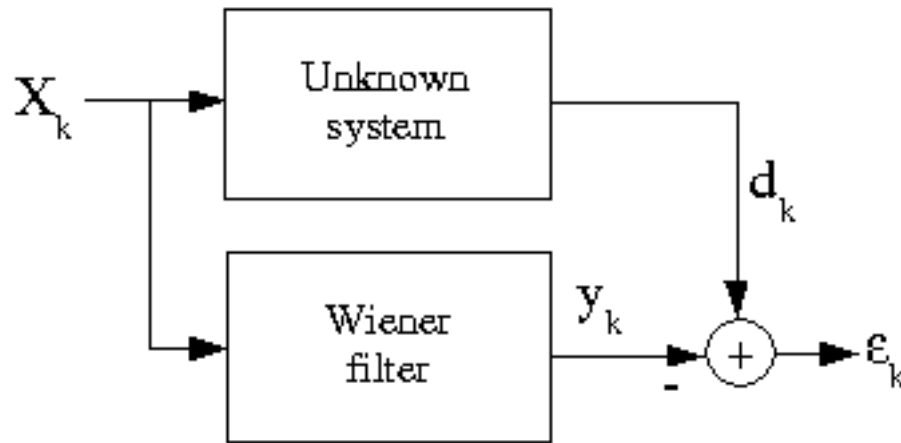


Figure 5.3

Exercise 5.1

Usually one desires that the input signal x_k be "persistently exciting," which, among other things, implies non-zero energy in all frequency bands. Why is this desirable?

5.2.1.1 Determining the optimal length-N causal FIR Wiener filter

NOTE: for convenience, we will analyze only the causal, real-data case; extensions are straightforward.

$$y_k = \sum_{l=0}^{M-1} (w_l x_{k-l})$$

$$\underset{w_l}{\operatorname{argmin}} E[\epsilon^2] = E[(d_k - y_k)^2] = E\left[\left(d_k - \sum_{l=0}^{M-1} (w_l x_{k-l})\right)^2\right] = E[d_k^2] - 2 \sum_{l=0}^{M-1} (w_l E[d_k x_{k-l}]) + \sum_{l=0}^{M-1} \left(\sum_{m=0}^{M-1} ((w_l w_m E[x_{k-l} x_{k-m}])) \right)$$

$$E[\epsilon^2] = r_{dd}(0) - 2 \sum_{l=0}^{M-1} (w_l r_{dx}(l)) + \sum_{l=0}^{M-1} \left(\sum_{m=0}^{M-1} (w_l w_m r_{xx}(l-m)) \right)$$

where

$$r_{dd}(0) = E[d_k^2]$$

$$r_{dx}(l) = E[d_k x_{k-l}]$$

$$r_{xx}(l-m) = E[x_k x_{k+l-m}]$$

This can be written in matrix form as

$$E[\epsilon^2] = r_{dd}(0) - 2\mathbf{P}\mathbf{W}^T + \mathbf{W}^T R \mathbf{W}$$

where

$$\mathbf{P} = \begin{pmatrix} r_{dx}(0) \\ r_{dx}(1) \\ \vdots \\ r_{dx}(M-1) \end{pmatrix}$$

$$R = \begin{pmatrix} r_{xx}(0) & r_{xx}(1) & \dots & \dots & r_{xx}(M-1) \\ r_{xx}(1) & r_{xx}(0) & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & r_{xx}(0) & r_{xx}(1) \\ r_{xx}(M-1) & \dots & \dots & r_{xx}(1) & r_{xx}(0) \end{pmatrix}$$

To solve for the optimum filter, compute the gradient with respect to the top weights vector W

$$\nabla \doteq \begin{pmatrix} \frac{\partial}{\partial w_0}(\epsilon^2) \\ \frac{\partial}{\partial w_1}(\epsilon^2) \\ \vdots \\ \frac{\partial}{\partial w_{M-1}}(\epsilon^2) \end{pmatrix}$$

$$\nabla = -(2\mathbf{P}) + 2R\mathbf{W}$$

(recall $\frac{d}{d\mathbf{W}}(\mathbf{A}^T \mathbf{W}) = \mathbf{A}^T$, $\frac{d}{d\mathbf{W}}(\mathbf{W}^T \mathbf{M} \mathbf{W}) = 2\mathbf{M}\mathbf{W}$ for symmetric M) setting the gradient equal to zero \Rightarrow

$$(W_{opt} R = \mathbf{P} \Rightarrow W_{opt} = R^{-1} \mathbf{P})$$

Since R is a correlation matrix, it must be non-negative definite, so this is a minimizer. For R positive definite, the minimizer is unique.

5.2.2 Practical Issues in Wiener Filter Implementation³

The weiner-filter, $W_{opt} = R^{-1} \mathbf{P}$, is ideal for many applications. But several issues must be addressed to use it in practice.

Exercise 5.2

(Solution on p. 288.)

In practice one usually won't know exactly the statistics of x_k and d_k (i.e. R and \mathbf{P}) needed to compute the Wiener filter.

How do we surmount this problem?

Exercise 5.3

(Solution on p. 288.)

In many applications, the statistics of x_k , d_k vary slowly with time.

How does one develop an **adaptive** system which tracks these changes over time to keep the system near optimal at all times?

Exercise 5.4

(Solution on p. 288.)

How can $r_{xx}^k(l)$ be computed efficiently?

³This content is available online at <<http://cnx.org/content/m11824/1.1/>>.

Exercise 5.5

how does one choose N ?

5.2.2.1 Tradeoffs

Larger $N \rightarrow$ more accurate estimates of the correlation values \rightarrow better \hat{W}_{opt} . However, larger N leads to slower adaptation.

NOTE: The success of adaptive systems depends on x, d being roughly stationary over at least N samples, $N > M$. That is, all adaptive filtering algorithms require that the underlying system varies slowly with respect to the sampling rate and the filter length (although they can tolerate occasional step discontinuities in the underlying system).

5.2.2.2 Computational Considerations

As presented here, an adaptive filter requires computing a matrix inverse at each sample. Actually, since the matrix R is Toeplitz, the linear system of equations can be solved with $O(M^2)$ computations using Levinson's algorithm, where M is the filter length. However, in many applications this may be too expensive, especially since computing the filter output itself requires $O(M)$ computations. There are two main approaches to resolving the computation problem

1. Take advantage of the fact that R^{k+1} is only slightly changed from R^k to reduce the computation to $O(M)$; these algorithms are called Fast Recursive Least Squares algorithms; all methods proposed so far have stability problems and are dangerous to use.
2. Find a different approach to solving the optimization problem that doesn't require explicit inversion of the correlation matrix.

NOTE: Adaptive algorithms involving the correlation matrix are called **Recursive least Squares** (RLS) algorithms. Historically, they were developed after the LMS algorithm, which is the simplest and most widely used approach $O(M)$. $O(M^2)$ RLS algorithms are used in applications requiring very fast adaptation.

5.2.3 Quadratic Minimization and Gradient Descent⁴**5.2.3.1 Quadratic minimization problems**

The least squares optimal filter design problem is quadratic in the filter coefficients:

$$E[\epsilon^2] = r_{dd}(0) - 2\mathbf{P}^T \mathbf{W} + \mathbf{W}^T R \mathbf{W}$$

If R is positive definite, the error surface $E[\epsilon^2](w_0, w_1, \dots, w_{M-1})$ is a unimodal "bowl" in \mathbb{R}^N .

⁴This content is available online at <<http://cnx.org/content/m11826/1.2/>>.

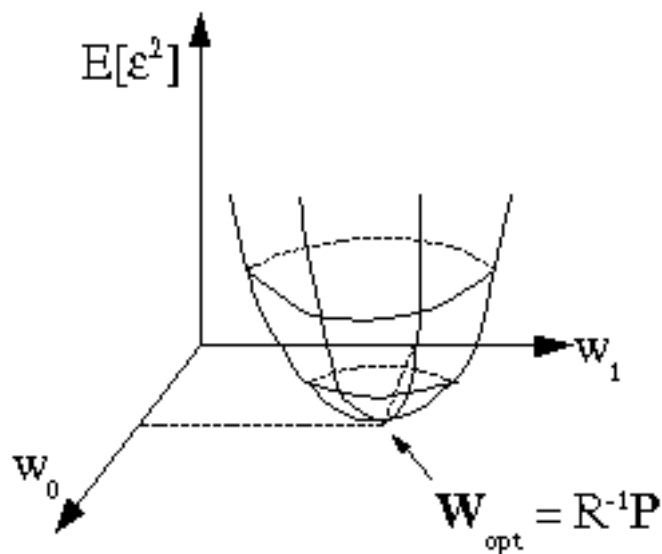


Figure 5.4

The problem is to find the bottom of the bowl. In an adaptive filter context, the shape and bottom of the bowl may drift slowly with time; hopefully slow enough that the adaptive algorithm can track it.

For a quadratic error surface, the bottom of the bowl can be found in one step by computing $\mathbf{R}^{-1}\mathbf{P}$. Most modern nonlinear optimization methods (which are used, for example, to solve the L^P optimal IIR filter design problem!) locally approximate a nonlinear function with a second-order (quadratic) Taylor series approximation and step to the bottom of this quadratic approximation on each iteration. However, an older and simpler approach to nonlinear optimization exists, based on **gradient descent**.

Contour plot of ϵ -squared

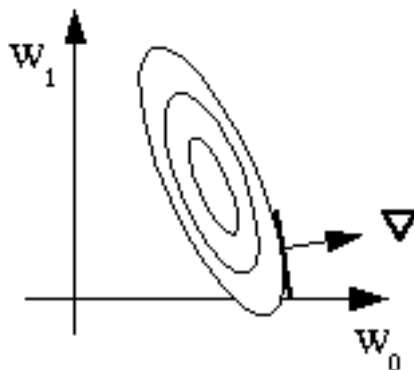


Figure 5.5

The idea is to iteratively find the minimizer by computing the gradient of the error function: $E\nabla =$

$\frac{\partial}{\partial w_i} (E [\epsilon^2])$. The gradient is a vector in \mathbb{R}^M pointing in the steepest uphill direction on the error surface at a given point \mathbf{W}^i , with ∇ having a magnitude proportional to the slope of the error surface in this steepest direction.

By updating the coefficient vector by taking a step *opposite* the gradient direction : $\mathbf{W}^{i+1} = \mathbf{W}^i - \mu \nabla^i$, we go (locally) "downhill" in the steepest direction, which seems to be a sensible way to iteratively solve a nonlinear optimization problem. The performance obviously depends on μ ; if μ is too large, the iterations could bounce back and forth up out of the bowl. However, if μ is too small, it could take many iterations to approach the bottom. We will determine criteria for choosing μ later.

In summary, the gradient descent algorithm for solving the Wiener filter problem is:

Guess W^0

do $i = 1, \infty$

$$\nabla^i = -(2P) + 2RW^i$$

$$W^{i+1} = W^i - \mu \nabla^i$$

repeat

$$W_{opt} = W^\infty$$

The gradient descent idea is used in the LMS adaptive filter algorithm. As presented, this algorithm costs $O(M^2)$ computations per iteration and doesn't appear very attractive, but LMS only requires $O(M)$ computations and is stable, so it is very attractive when computation is an issue, even though it converges more slowly than the RLS algorithms we have discussed so far.

5.3 The LMS Adaptive Filter Algorithm

5.3.1 The LMS Adaptive Filter Algorithm⁵

Recall the Wiener filter problem

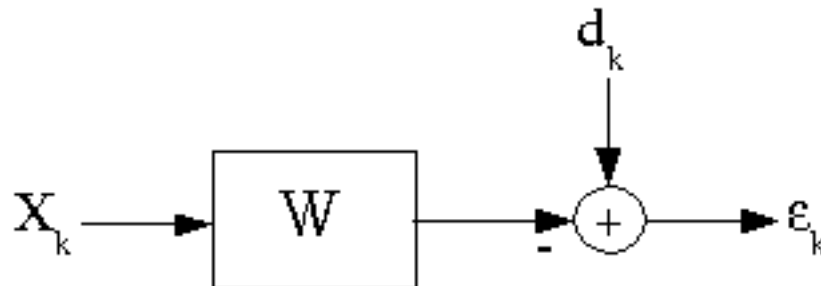


Figure 5.6

⁵This content is available online at <<http://cnx.org/content/m11829/1.1/>>.

$\{x_k\}, \{d_k\}$ jointly wide sense stationary
Find W minimizing $E[\epsilon_k^2]$

$$\epsilon_k = d_k - y_k = d_k - \sum_{i=0}^{M-1} (w_i x_{k-i}) = d_k - X^k T W^k$$

$$X^k = \begin{pmatrix} x_k \\ x_{k-1} \\ \vdots \\ x_{k-M+1} \end{pmatrix}$$

$$W^k = \begin{pmatrix} w_0^k \\ w_1^k \\ \vdots \\ w_{M-1}^k \end{pmatrix}$$

The superscript denotes absolute time, and the subscript denotes time or a vector index.
the solution can be found by setting the gradient = 0

$$\begin{aligned} \nabla^k &= \frac{\partial}{\partial W} (E[\epsilon_k^2]) \\ &= E[2\epsilon_k (-X^k)] \\ &= E[-2(d_k - X^k T W_k) X^k] \\ &= -(2E[d_k X^k]) + (E[X^k T]) W \\ &= -2P + 2RW \end{aligned} \tag{5.1}$$

$$\Rightarrow (W_{opt} = R^{-1}P)$$

Alternatively, W_{opt} can be found iteratively using a gradient descent technique

$$W^{k+1} = W^k - \mu \nabla^k$$

In practice, we don't know R and P exactly, and in an adaptive context they may be slowly varying with time.

To find the (approximate) Wiener filter, some approximations are necessary. As always, the key is to make the *right* approximations!

GOOD IDEA: Approximate R and P : \Rightarrow RLS methods, as discussed last time.

BETTER IDEA: Approximate the gradient!

$$\hat{\nabla}^k = \frac{\partial}{\partial W} (E[\epsilon_k^2])$$

Note that ϵ_k^2 itself is a very noisy approximation to $E[\epsilon_k^2]$. We can get a noisy approximation to the gradient by finding the gradient of ϵ_k^2 ! Widrow and Hoff first published the LMS algorithm, based on this clever idea, in 1960.

$$\hat{\nabla}^k = \frac{\partial}{\partial W} (\epsilon_k^2) = 2\epsilon_k \frac{\partial}{\partial W} (d_k - W^k T X^k) = 2\epsilon_k (-X^k) = -(2\epsilon_k X^k)$$

This yields the LMS adaptive filter algorithm

Example 5.3: The LMS Adaptive Filter Algorithm

1. $y_k = W^k X^k = \sum_{i=0}^{M-1} (w_i^k x_{k-i})$
2. $\epsilon_k = d_k - y_k$
3. $W^{k+1} = W^k - \mu \hat{\nabla}^k = W^k - \mu (-2\epsilon_k X^k) = W^k + 2\mu\epsilon_k X^k$ ($w_i^{k+1} = w_i^k + 2\mu\epsilon_k x_{k-i}$)

The LMS algorithm is often called a **stochastic gradient** algorithm, since $\hat{\nabla}^k$ is a noisy gradient. This is by far the most commonly used adaptive filtering algorithm, because

1. it was the first
2. it is very simple
3. in practice it works well (except that sometimes it converges slowly)
4. it requires relatively little computation
5. it updates the tap weights every sample, so it continually adapts the filter
6. it tracks slow changes in the signal statistics well

5.3.1.1 Computational Cost of LMS

To Compute \Rightarrow	y_k	ϵ_k	W^{k+1}	= Total
multiplies	M	0	$M + 1$	$2M + 1$
adds	$M - 1$	1	M	$2M$

So the LMS algorithm is $O(M)$ per sample. In fact, it is nicely balanced in that the filter computation and the adaptation require the same amount of computation.

Note that the parameter μ plays a very important role in the LMS algorithm. It can also be varied with time, but usually a constant μ ("convergence weight factor") is used, chosen after experimentation for a given application.

5.3.1.1.1 Tradeoffs

large μ : fast convergence, fast adaptivity

small μ : accurate $W \rightarrow$ less misadjustment error, stability

5.3.2 First Order Convergence Analysis of the LMS Algorithm⁶

5.3.2.1 Analysis of the LMS algorithm

It is important to analyze the LMS algorithm to determine under what conditions it is stable, whether or not it converges to the Wiener solution, to determine how quickly it converges, how much degradation is suffered due to the noisy gradient, etc. In particular, we need to know how to choose the parameter μ .

5.3.2.1.1 Mean of W

does W^k , $k \rightarrow \infty$ approach the Wiener solution? (since W^k is always somewhat random in the approximate gradient-based LMS algorithm, we ask whether the expected value of the filter coefficients converge to the Wiener solution)

$$\begin{aligned}
 E[W^{k+1}] &= \overline{W^{k+1}} \\
 &= E[W^k + 2\mu\epsilon_k X^k] \\
 &= \overline{W^k} + 2\mu E[d_k X^k] + 2\mu E\left[-\left((W^k X^k) X^k\right)\right] \\
 &= \overline{W^k} + 2\mu P + \left(-\left(2\mu E\left[(W^k X^k) X^k\right]\right)\right)
 \end{aligned} \tag{5.2}$$

⁶This content is available online at <<http://cnx.org/content/m11830/1.1/>>.

5.3.2.1.1 Patently False Assumption

X^k and X^{k-i} , X^k and d^{k-i} , and d_k and d_{k-i} are statistically independent, $i \neq 0$. This assumption is obviously false, since X^{k-1} is the same as X^k except for shifting down the vector elements one place and adding one new sample. We make this assumption because otherwise it becomes extremely difficult to analyze the LMS algorithm. (First good analysis not making this assumption: *Macchi and Eweda*[23]) Many simulations and much practical experience has shown that the results one obtains with analyses based on the patently false assumption above are quite accurate in most situations

With the independence assumption, W^k (which depends only on previous X^{k-i} , d^{k-i}) is statistically independent of X^k , and we can simplify $E[(W^{kT} X^k) X^k]$

Now $(W^{kT} X^k) X^k$ is a vector, and

$$\begin{aligned}
 E[(W^{kT} X^k) X^k] &= E \left[\begin{pmatrix} \vdots \\ (\sum_{i=0}^{M-1} (w_i^k x_{k-i})) x_{k-j} \\ \vdots \end{pmatrix} \right] \\
 &= \begin{pmatrix} \vdots \\ \sum_{i=0}^{M-1} (E[w_i^k x_{k-i} x_{k-j}]) \\ \vdots \end{pmatrix} \\
 &= \begin{pmatrix} \vdots \\ \sum_{i=0}^{M-1} ((w_i^k) E[x_{k-i} x_{k-j}]) \\ \vdots \end{pmatrix} \\
 &= \begin{pmatrix} \vdots \\ \sum_{i=0}^{M-1} (\overline{w_i^k} r_{xx}(i-j)) \\ \vdots \end{pmatrix} \\
 &= R \overline{W^k}
 \end{aligned} \tag{5.3}$$

where $R = E[X^k X^{kT}]$ is the data correlation matrix.

Putting this back into our equation

$$\begin{aligned}
 \overline{W^{k+1}} &= \overline{W^k} + 2\mu P + \left(- \left(2\mu R \overline{W^k} \right) \right) \\
 &= (I - 2\mu R) \overline{W^k} + 2\mu P
 \end{aligned} \tag{5.4}$$

Now if $\overline{W^{k \rightarrow \infty}}$ converges to a vector of finite magnitude ("convergence in the mean"), what does it converge to?

If $\overline{W^k}$ converges, then as $k \rightarrow \infty$, $\overline{W^{k+1}} \approx \overline{W^k}$, and

$$\overline{W^\infty} = (I - 2\mu R) \overline{W^\infty} + 2\mu P$$

$$2\mu R \overline{W^\infty} = 2\mu P$$

$$R \overline{W^\infty} = P$$

or

$$\overline{W_{opt}} = R^{-1}P$$

the Wiener solution!

So the LMS algorithm, *if* it converges, gives filter coefficients which on average are the Wiener coefficients! This is, of course, a desirable result.

5.3.2.1.2 First-order stability

But does $\overline{W^k}$ converge, or under what conditions?

Let's rewrite the analysis in term of $\overline{V^k}$, the "mean coefficient error vector" $\overline{V^k} = \overline{W^k} - W_{opt}$, where W_{opt} is the Wiener filter

$$\overline{W^{k+1}} = \overline{W^k} - 2\mu R \overline{W^k} + 2\mu P$$

$$\overline{W^{k+1}} - W_{opt} = \overline{W^k} - W_{opt} + \left(- \left(2\mu R \overline{W^k} \right) \right) + 2\mu R W_{opt} - 2\mu R W_{opt} + 2\mu P$$

$$\overline{V^{k+1}} = \overline{V^k} - 2\mu R \overline{V^k} + (- (2\mu R W_{opt})) + 2\mu P$$

Now $W_{opt} = R^{-1}P$, so

$$\overline{V^{k+1}} = \overline{V^k} - 2\mu R \overline{V^k} + (- (2\mu R R^{-1}P)) + 2\mu P = (I - 2\mu R) \overline{V^k}$$

We wish to know under what conditions $\overline{V^{k \rightarrow \infty}} \rightarrow \vec{0}$?

5.3.2.1.2.1 Linear Algebra Fact

Since R is positive definite, real, and symmetric, all the eigenvalues are real and positive. Also, we can write R as $(Q^{-1}\Lambda Q)$, where Λ is a diagonal matrix with diagonal entries λ_i equal to the eigenvalues of R , and Q is a unitary matrix with rows equal to the eigenvectors corresponding to the eigenvalues of R .

Using this fact,

$$V^{k+1} = (I - 2\mu (Q^{-1}\Lambda Q)) V^k$$

multiplying both sides through on the left by Q : we get

$$Q \overline{V^{k+1}} = (Q - 2\mu \Lambda Q) \overline{V^k} = (1 - 2\mu \Lambda) Q \overline{V^k}$$

Let $V' = QV$:

$$V'^{k+1} = (1 - 2\mu \Lambda) V'^k$$

Note that V' is simply V in a rotated coordinate set in \mathbb{R}^m , so convergence of V' implies convergence of V .

Since $1 - 2\mu \Lambda$ is diagonal, all elements of V' evolve independently of each other. Convergence (stability) boils down to whether all M of these scalar, first-order difference equations are stable, and thus $\rightarrow 0$.

$$\forall i, i = [1, 2, \dots, M] : (V_i'^{k+1} = (1 - 2\mu \lambda_i) V_i'^k)$$

These equations converge to zero if $|1 - 2\mu \lambda_i| < 1$, or $\forall i : (|\mu \lambda_i| < 1)$ μ and λ_i are positive, so we require $\forall i : \left(\mu < \frac{1}{\lambda_i} \right)$ so for convergence in the mean of the LMS adaptive filter, we require

$$\mu < \frac{1}{\lambda_{max}} \quad (5.5)$$

This is an elegant theoretical result, but in practice, we may not know λ_{max} , it may be time-varying, and we certainly won't want to compute it. However, another useful mathematical fact comes to the rescue...

$$tr(R) = \sum_{i=1}^M r_{ii} = \sum_{i=1}^M \lambda_i \geq \lambda_{max}$$

Since the eigenvalues are all positive and real.

For a correlation matrix, $\forall i, i \in \{1, M\} : (r_{ii} = r(0))$. So $tr(R) = Mr(0) = ME[x_k x_k]$. We can easily estimate $r(0)$ with $O(1)$ computations/sample, so in practice we might require

$$\mu < \frac{1}{Mr(\hat{0})}$$

as a conservative bound, and perhaps adapt μ accordingly with time.

5.3.2.1.3 Rate of convergence

Each of the modes decays as

$$(1 - 2\mu\lambda_i)^k$$

GOOD NEWS: The *initial* rate of convergence is dominated by the fastest mode $1 - 2\mu\lambda_{max}$. This is not surprising, since a gradient descent method goes "downhill" in the steepest direction

BAD NEWS: The *final* rate of convergence is dominated by the slowest mode $1 - 2\mu\lambda_{min}$. For small λ_{min} , it can take a long time for LMS to converge.

Note that the convergence behavior depends on the data (via R). LMS converges relatively quickly for roughly equal eigenvalues. Unequal eigenvalues slow LMS down a lot.

5.3.3 Second-order Convergence Analysis of the LMS Algorithm and Misadjustment Error⁷

Convergence of the mean (first-order analysis) is insufficient to guarantee desirable behavior of the LMS algorithm; the variance could still be infinite. It is important to show that the variance of the filter coefficients is finite, and to determine how close the average squared error is to the minimum possible error using an exact Wiener filter.

$$\begin{aligned} E[\epsilon_k^2] &= E\left[\left(d_k - W^k T X^k\right)^2\right] \\ &= E\left[d_k^2 - 2d_k X^k T W^k - W^k T X^k X^k T W^k\right] \\ &= r_{dd}(0) - 2W^k T P + W^k T R W^k \end{aligned} \quad (5.6)$$

The minimum error is obtained using the Wiener filter

$$W_{opt} = R^{-1}P$$

$$\begin{aligned} \epsilon_{min}^2 &= E[\epsilon^2] \\ &= (r_{dd}(0) - 2P^T R^{-1}P + P^T R^{-1} R R^{-1}P) \\ &= r_{dd}(0) - P^T R^{-1}P \end{aligned} \quad (5.7)$$

To analyze the average error in LMS, write (5.6) in terms of $V = Q[W - W_{opt}]$, where $Q^{-1}\Lambda Q = R$

$$\begin{aligned} E[\epsilon_k^2] &= r_{dd}(0) - 2W^k T P + W^k T R W^k + -\left(W^k T R W_{opt}\right) - W_{opt}^T R W^k + \\ &W_{opt}^T R W_{opt} + W^k T R W_{opt} + W_{opt}^T R W^k - W_{opt}^T R W_{opt} = r_{dd}(0) + V^k T R V^k - \\ &P^T R^{-1}P = \epsilon_{min}^2 + V^k T R V^k = \epsilon_{min}^2 + V^k T Q^{-1} Q R Q^{-1} Q V^k = \epsilon_{min}^2 + V^k T \Lambda V^k \end{aligned} \quad (5.8)$$

⁷This content is available online at <<http://cnx.org/content/m11831/1.2/>>.

$$E[\epsilon_k^2] = \epsilon_{min}^2 + \sum_{j=0}^{N-1} \left(\lambda_j E[v_j^{k2}] \right)$$

So we need to know $E[v_j^{k2}]$, which are the diagonal elements of the covariance matrix of V'^k , or $E[V'^k V'^k T]$.

From the LMS update equation

$$W^{k+1} = W^k + 2\mu\epsilon_k X^k$$

we get

$$V'^{k+1} = W'^k + 2\mu\epsilon_k QX^k$$

$$\begin{aligned} \mathcal{V}^{k+1} &= E[V'^{k+1} V'^{k+1 T}] \\ &= E[4\mu^2 \epsilon_k^2 QX^k X^{kT} Q^T] \\ &= \mathcal{V}^k + 2\mu \left(\epsilon_k QX^k V'^k T \right) + 2\mu \left(\epsilon_k V'^k X^{kT} Q^T \right) + 4\mu^2 E[\epsilon_k^2 QX^k X^{kT} Q^T] \end{aligned} \quad (5.9)$$

Note that

$$\epsilon_k = d_k - W^{kT} X^k = d_k - W_{opt}^T - V'^k T QX^k$$

so

$$\begin{aligned} E[\epsilon_k^2 QX^k V'^k T] &= E[d_k QX^k V'^k T - W_{opt}^T X^k QX^k V'^k T - V'^k T QX^k V'^k T] \\ &= 0 + 0 - \left(QX^k X^{kT} Q^T V'^k V'^k T \right) \\ &= - \left(QE[X^k X^{kT}] Q^T E[V'^k V'^k T] \right) \\ &= -(\Lambda \mathcal{V}^k) \end{aligned} \quad (5.10)$$

Note that the Patently False independence Assumption was invoked here.

To analyze $E[\epsilon_k^2 QX^k X^{kT} Q^T]$, we make yet another obviously false assumption that ϵ_k^2 and X^k are statistically independent. This is obviously false, since $\epsilon_k = d_k - W^{kT} X^k$. Otherwise, we get 4th-order terms in X in the product. These can be dealt with, at the expense of a more complicated analysis, if a particular type of distribution (such as Gaussian) is assumed. See, for example *Gardner*[14]. A questionable justification for this assumption is that as $W^k \approx W_{opt}$, W^k becomes uncorrelated with X^k (if we invoke the original independence assumption), which tends to randomize the error signal relative to X^k . With this assumption,

$$E[\epsilon_k^2 QX^k X^{kT} Q^T] = E[\epsilon_k^2] E[QX^k X^{kT} Q^T] = E[\epsilon_k^2] \Lambda$$

Now

$$\epsilon_k^2 = \epsilon_{min}^2 + V'^k T \Lambda V'^k$$

so

$$\begin{aligned} E[\epsilon_k^2] &= \epsilon_{min}^2 + E\left[\sum (\lambda_j V_j^{k2})\right] \\ &= \epsilon_{min}^2 + \sum (\lambda_j \mathcal{V}_{jj}^k) \end{aligned} \quad (5.11)$$

Thus, (5.9) becomes

$$\mathcal{V}^{k+1} = (I - 4\mu\Lambda) \mathcal{V}^k + 4\mu^2 \sum (\lambda_j \mathcal{V}_{jj}^k \Lambda) + 4\mu^2 \epsilon_{min}^2 \Lambda \quad (5.12)$$

Now if this system is stable and converges, it converges to $\mathcal{V}^\infty = \mathcal{V}^{\infty+1}$

$$\begin{aligned} \Rightarrow (4\mu\Lambda \mathcal{V}^\infty &= 4\mu^2 \left(\sum (\lambda_j \mathcal{V}_{jj}^\infty) + \epsilon_{min}^2 \right) \Lambda) \\ \Rightarrow (\mathcal{V}^\infty &= \mu \left(\sum (\lambda_j \mathcal{V}_{jj}^\infty) + \epsilon_{min}^2 \right) I) \end{aligned}$$

So it is a diagonal matrix with all elements on the diagonal equal:

Then

$$\mathcal{V}_{ii}^{\infty} = \mu \left(\mathcal{V}_{ii}^{\infty} \sum \lambda_j + \epsilon_{min}^2 \right)$$

$$\mathcal{V}_{ii}^{\infty} \left(1 - \mu \sum \lambda_j \right) = \mu \epsilon_{min}^2$$

$$\mathcal{V}_{ii}^{\infty} = \frac{\mu \epsilon_{min}^2}{1 - \mu \sum \lambda_j}$$

Thus the error in the LMS adaptive filter after convergence is

$$\begin{aligned} E[\epsilon_{\infty}^2] &= \epsilon_{min}^2 + E[V^{\infty} \lambda V^{\infty}] \\ &= \epsilon_{min}^2 + \frac{\mu \epsilon_{min}^2 \sum (\lambda_j)}{1 - \mu \sum \lambda_j} \\ &= \epsilon_{min}^2 \frac{1}{1 - \mu \sum \lambda_j} \\ &= \epsilon_{min}^2 \frac{1}{1 - \mu \text{tr}(R)} \\ &= \epsilon_{min}^2 \frac{1}{1 - \mu r_{xx}(0)N} \end{aligned} \tag{5.13}$$

$$E[\epsilon_{\infty}^2] = \epsilon_{min}^2 \frac{1}{1 - \mu N \sigma_x^2} \tag{5.14}$$

$1 - \mu N \sigma_x^2$ is called the **misadjustment factor**. Often, one chooses μ to select a desired misadjustment factor, such as an error 10% higher than the Wiener filter error.

5.3.3.1 2nd-Order Convergence (Stability)

To determine the range for μ for which (5.12) converges, we must determine the μ for which the matrix difference equation converges.

$$\mathcal{V}^{k+1} = (I - 4\mu\Lambda) \mathcal{V}^k + 4\mu^2 \sum (\lambda_j \mathcal{V}_{jj}^k \Lambda) + 4\mu^2 \epsilon_{min}^2 \Lambda$$

The off-diagonal elements each evolve independently according to $\mathcal{V}_{ij}^{k+1} = 1 - 4\mu\lambda_i \mathcal{V}_{ij}^k$. These terms will decay to zero if $\forall i : (4\mu\lambda_i < 2)$, or $\mu < \frac{1}{2\lambda_{max}}$

The diagonal terms evolve according to

$$\mathcal{V}_{ii}^{k+1} = (1 - 4\mu\lambda_i) \mathcal{V}_{ii}^k + 4\mu^2 \lambda_i \left(\sum \lambda_j \right) \mathcal{V}_{jj}^k + 4\mu^2 \epsilon_{min}^2 \lambda_i$$

For the homogeneous equation

$$\mathcal{V}_{ii}^{k+1} = (1 - 4\mu\lambda_i) \mathcal{V}_{ii}^k + 4\mu^2 \lambda_i \left(\sum \lambda_j \right) \mathcal{V}_{jj}^k$$

for $1 - 4\mu\lambda_i$ positive,

$$\mathcal{V}_{ii}^{k+1} \leq (1 - 4\mu\lambda_i) \mathcal{V}_{iimax}^k + 4\mu^2 \lambda_i \left(\sum \lambda_j \right) \mathcal{V}_{jjmax}^k = \left(1 - 4\mu\lambda_i + 4\mu^2 \lambda_i \sum \lambda_j \right) \mathcal{V}_{jjmax}^k \tag{5.15}$$

\mathcal{V}_{ii}^{k+1} will be strictly less than \mathcal{V}_{jjmax}^k for

$$1 - 4\mu\lambda_i + 4\mu^2 \lambda_i \sum \lambda_j < 1$$

or

$$4\mu^2 \lambda_i \sum \lambda_j < 4\mu\lambda_i$$

or

$$\begin{aligned}\mu < \frac{1}{\sum \lambda_j} &= \frac{1}{\text{tr}(R)} \\ &= \frac{1}{Nr_{xx}(0)} \\ &= \frac{1}{N\sigma_x^2}\end{aligned}\tag{5.16}$$

This is a more rigorous bound than the first-order bounds. Often engineers choose μ a few times smaller than this, since more rigorous analyses yield a slightly smaller bound. $\mu = \frac{\mu}{3N\sigma_x^2}$ is derived in some analyses assuming Gaussian x_k , d_k .

5.4 Applications of Adaptive Filters

5.4.1 Applications of Adaptive Filters⁸

Adaptive filters can be configured in a surprisingly large variety of ways for application in a large number of applications. The same LMS algorithm can be used to adapt the filter coefficients in most of these configurations; only the "plumbing" changes.

5.4.2 Adaptive System Identification⁹

GOAL: To approximate an unknown system (or the behavior of that system) as closely as possible

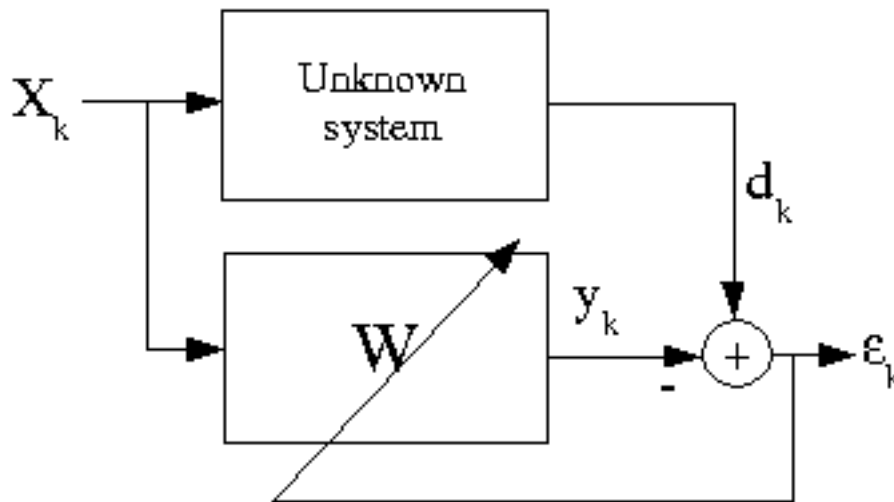


Figure 5.7

The optimal solution is $R^{-1}P = W$

⁸This content is available online at <http://cnx.org/content/m11536/1.1/>.

⁹This content is available online at <http://cnx.org/content/m11906/1.1/>.

Suppose the unknown system is a causal, linear time-invariant filter:

$$d_k = x_k * h_k = \sum_{i=0}^{\infty} (x_{k-i} h_i)$$

Now

$$\begin{aligned} P &= \begin{pmatrix} E[d_k x_{k-j}] \end{pmatrix} = \begin{pmatrix} E[(\sum_{i=0}^{\infty} (x_{k-i} h_i)) x_{k-j}] \end{pmatrix} = \quad (5.17) \\ &\begin{pmatrix} \sum_{i=0}^{\infty} (h_i E[x_{k-i} x_{k-j}]) \end{pmatrix} = \begin{pmatrix} \sum_{i=0}^{\infty} (r_{xx}(j-i)) \end{pmatrix} = \\ &\begin{pmatrix} r_{xx}(0) & r(1) & \dots & \dots & r(M-1) & | & r(M) & r(M+1) & \dots \\ r(1) & r(0) & \ddots & \ddots & \vdots & | & \vdots & \vdots & \dots \\ r(2) & r(1) & \ddots & \ddots & \vdots & | & \vdots & \vdots & \dots \\ \vdots & \vdots & \dots & r(0) & r(1) & | & r(2) & r(3) & \dots \\ r(M-1) & r(M-2) & \dots & r(1) & r(0) & | & r(1) & r(2) & \dots \end{pmatrix} \begin{pmatrix} h(0) \\ h(1) \\ h(2) \\ \vdots \end{pmatrix} \end{aligned}$$

If the adaptive filter H is a length- M FIR filter ($h(m) = h(m+1) = \dots = 0$), this reduces to

$$P = Rh^{-1}$$

and

$$W_{opt} = R^{-1}P = R^{-1}(Rh) = \mathbf{h}$$

FIR adaptive system identification thus converges in the mean to the corresponding M samples of the impulse response of the unknown system.

5.4.3 Adaptive Equalization¹⁰

GOAL: Design an approximate inverse filter to cancel out as much distortion as possible.

¹⁰This content is available online at <<http://cnx.org/content/m11907/1.1/>>.

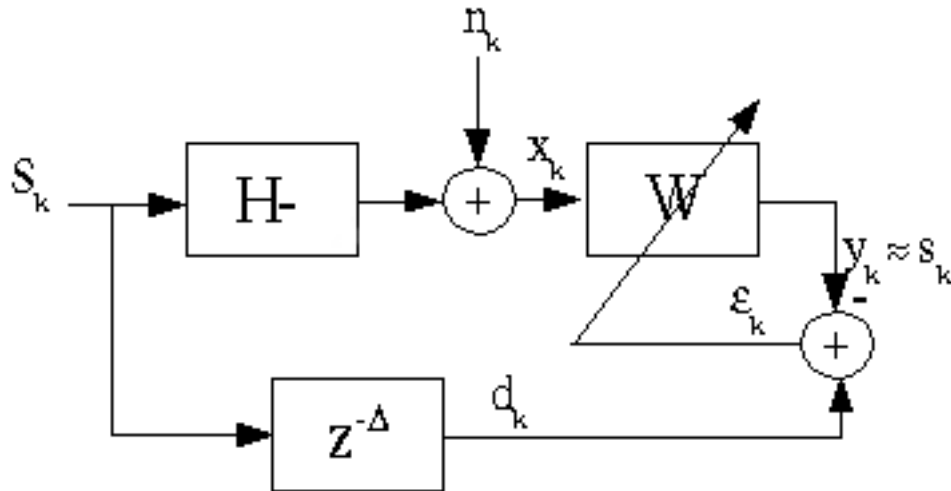


Figure 5.8

In principle, $WH \approx z^{-\Delta}$, or $W \approx \frac{z^{-\Delta}}{H}$, so that the overall response of the top path is approximately $\delta(n - \Delta)$. However, limitations on the form of W (FIR) and the presence of noise cause the equalization to be imperfect.

5.4.3.1 Important Application

Channel equalization in a digital communication system.



Figure 5.9

If the channel distorts the pulse shape, the matched filter will no longer be matched, intersymbol interference may increase, and the system performance will degrade.

An adaptive filter is often inserted in front of the matched filter to compensate for the channel.

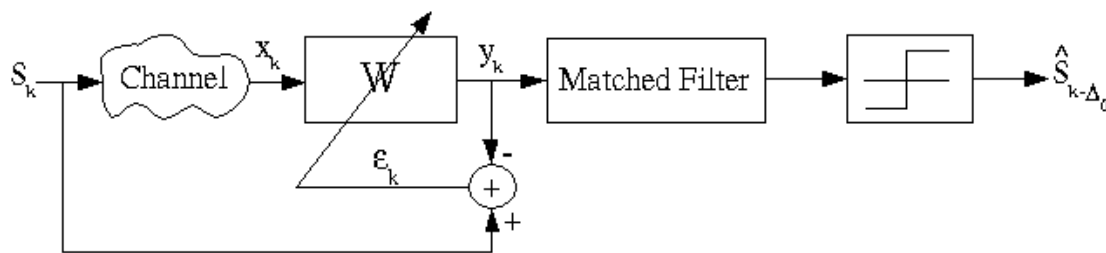


Figure 5.10

This is, of course, unrealizable, since we do not have access to the original transmitted signal, s_k . There are two common solutions to this problem:

1. Periodically broadcast a known **training signal**. The adaptation is switched on only when the training signal is being broadcast and thus s_k is known.
2. Decision-directed feedback: If the overall system is working well, then the output $\hat{s}_{k-\Delta_0}$ should almost always equal $s_{k-\Delta_0}$. We can thus use our received digital communication signal as the desired signal, since it has been cleaned of noise (we hope) by the non-linear threshold device!

Decision-directed equalizer

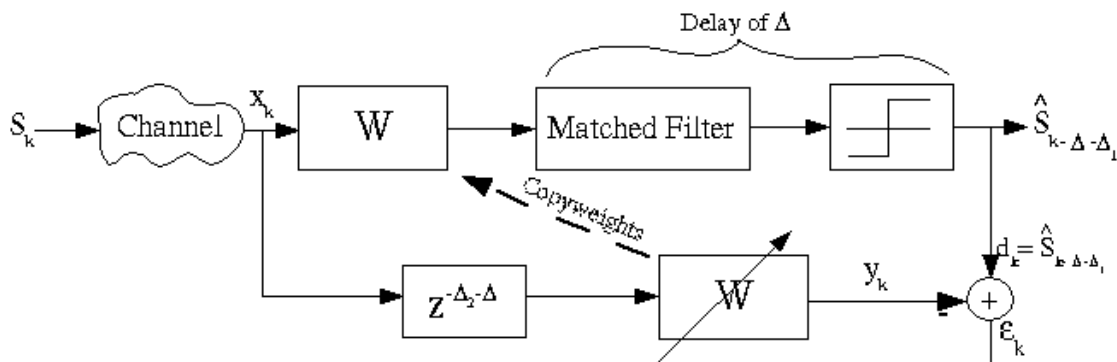


Figure 5.11

As long as the error rate in \hat{s}_k is not too high (say $< 75\%$), this method works. Otherwise, d_k is so inaccurate that the adaptive filter can never find the Wiener solution. This method is widely used in the telephone system and other digital communication networks.

5.4.4 Adaptive Interference (Noise) Cancellation¹¹

GOAL: Automatically eliminate unwanted interference in a signal.

¹¹This content is available online at <http://cnx.org/content/m11835/1.1/>.

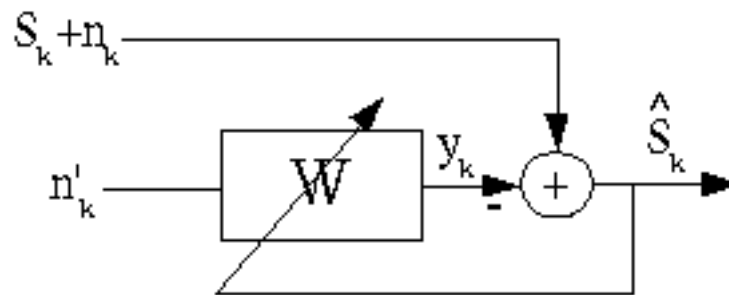


Figure 5.12

The object is to subtract out as much of the noise as possible.

Example 5.4: Engine noise cancellation in automobiles

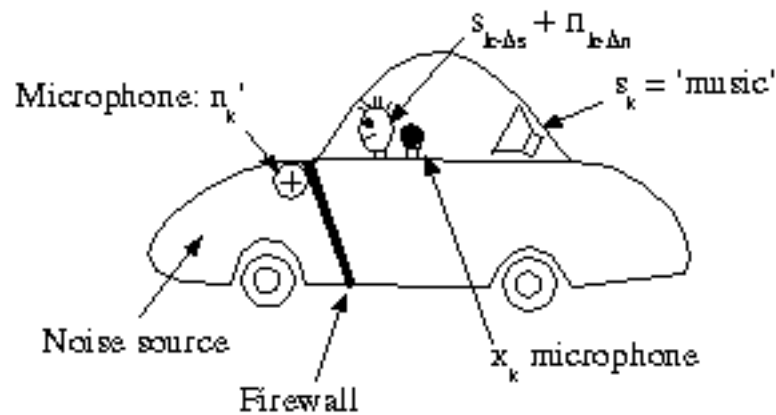


Figure 5.13

The firewall attenuates and filters the noise reaching the listener's ear, so it is not the same as n'_k . There is also a delay due to acoustic propagation in the air. For maximal cancellation, an adaptive filter is thus needed to make n'_k as similar as possible to the delayed n_k .

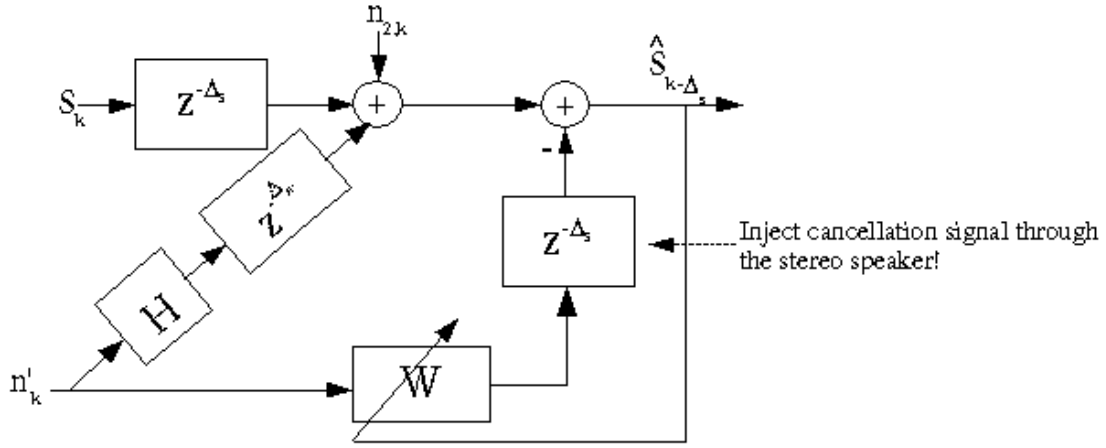


Figure 5.14

Exercise 5.6

What conditions must we impose upon the microphone locations for this to work? (Think causality and physics!)

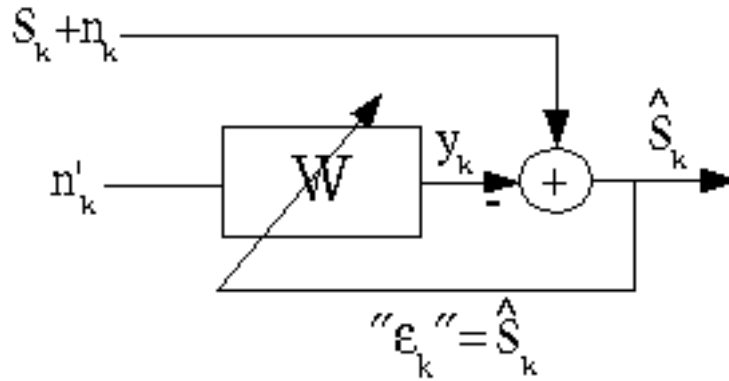
5.4.4.1 Analysis of the interference cancellor

Figure 5.15

$$E[\epsilon_k^2] = E[(s_k + n_k - y_k)^2] = E[s_k^2] + 2E[s_k(n_k - y_k)] + E[(n_k - y_k)^2]$$

We assume s_k , n_k , and n'_k are zero-mean signals, and that s_k is independent of n_k and n'_k . Then

$$E[s_k(n_k - y_k)] = E[s_k]E[n_k - y_k] = 0$$

$$E[\epsilon_k^2] = E[s_k^2] + E[(n_k - y_k)^2]$$

Since the input signal has no information about s_k in it, minimizing $E[\epsilon_k^2]$ can only affect the second term, which is the standard Wiener filtering problem, with solution

$$W = R_{n'n'}^{-1} P_{nn'}$$

5.4.5 Adaptive Echo Cancellation¹²

An adaptive echo canceller is a specific type of adaptive interference canceller that removes echos. Adaptive echo cancellers are found in all modern telephone systems.

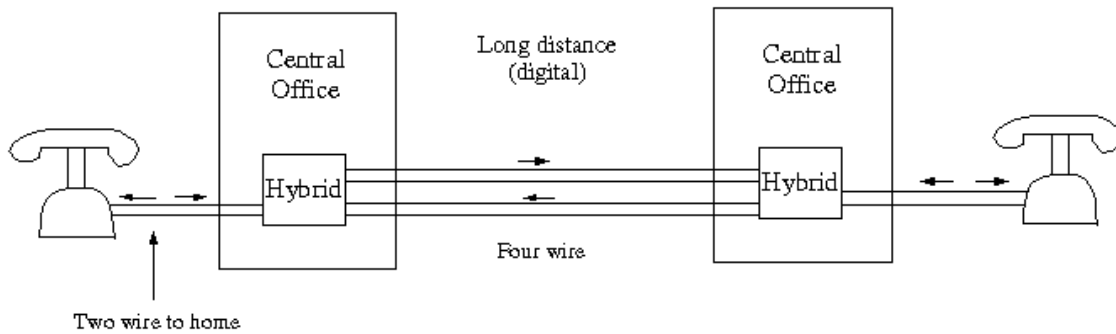


Figure 5.16

The hybrid is supposed to split the opposite-going waves, but typically achieves only about 15dB of suppression. This signal will eventually reach the other end and be coupled back, with a long delay, to the original source, which gives a very annoying echo.

¹²This content is available online at <<http://cnx.org/content/m11909/1.1/>>.

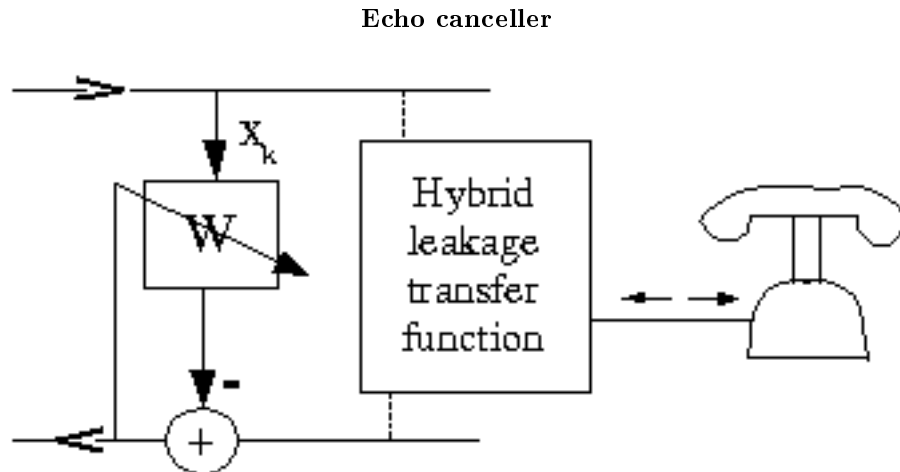


Figure 5.17

Because the input to the adaptive echo canceller contains only the signal from the far end that will echo off of the hybrid, it cancels the echo while passing the near-end signal as desired.

5.4.5.1 Narrowband interference canceller

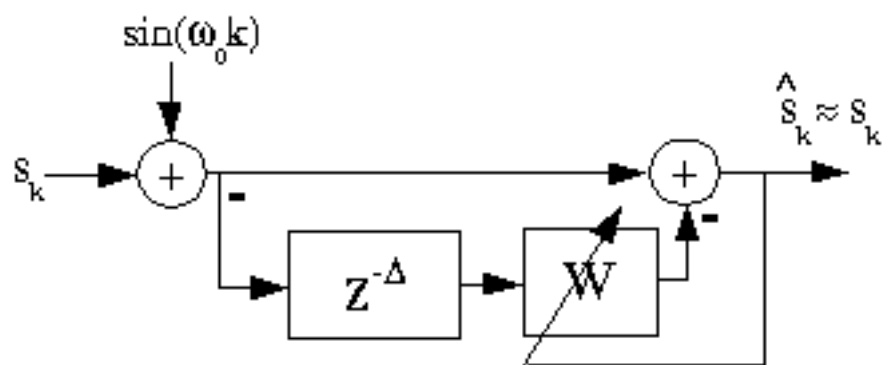


Figure 5.18

A sinusoid is predictable Δ samples ahead, whereas s_k may not be, so the sinusoid can be cancelled using the adaptive system in the Figure. This is another special case of the adaptive interference canceller in which the noise reference input is a delayed version of the primary (signal plus noise) input. Note that Δ must be large enough so that s_k and $s_{k-\Delta}$ are uncorrelated, or some of the signal will be cancelled as well!

Exercise 5.7

How would you construct an "adaptive line enhancer" that preserves the sinusoids but cancels the uncorrelated noise?

Other Applications

- Adaptive array processing
- Adaptive control
- etc...

5.5 Other Adaptive Filter Algorithms

5.5.1 Beyond LMS: an overview of other adaptive filter algorithms¹³

5.5.1.1 RLS algorithms

FIR adaptive filter algorithms with faster convergence. Since the Wiener solution can be obtained on one step by computing $W_{opt} = R^{-1}P$, most RLS algorithms attempt to estimate R^{-1} and P and compute W_{opt} from these.

There are a number of $O(N^2)$ algorithms which are stable and converge quickly. A number of $O(N)$ algorithms have been proposed, but these are all unstable except for the lattice filter method. This is described to some extent in the text. The adaptive lattice filter converges quickly and is stable, but reportedly has a very high noise floor.

Many of these approaches can be thought of as attempting to "orthogonalize" R , or to rotate the data or filter coefficients to a domain where R is diagonal, then doing LMS in each dimension *separately*, so that a fast-converging step size can be chosen in all directions.

5.5.1.2 Frequency-domain methods

Frequency-domain methods implicitly attempt to do this:

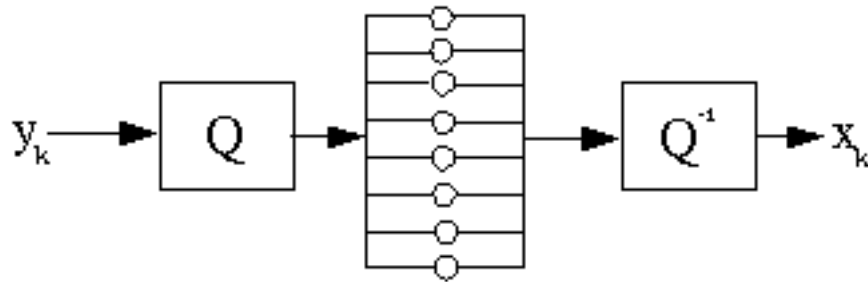


Figure 5.19

If QRQ^{-1} is a diagonal matrix, this yields a fast algorithm. If Q is chosen as an FFT matrix, each channel becomes a different frequency bin. Since R is Toeplitz and not a circulant, the FFT matrix will not exactly diagonalize R , but in many cases it comes very close and frequency domain methods converge very quickly. However, for some R they perform no better than LMS. By using an FFT, the transformation Q becomes inexpensive $O(N \log N)$. If one only updates on a block-by-block basis (once per N samples), the frequency domain methods only cost $O(\log N)$ computations per sample, which can be important for some applications with large N . (Say 16,000,000)

¹³This content is available online at <<http://cnx.org/content/m11911/1.1/>>.

5.5.2 Adaptive IIR filters¹⁴

Adaptive IIR filters are attractive for the same reasons that IIR filters are attractive: many fewer coefficients may be needed to achieve the desired performance in some applications. However, it is more difficult to develop stable IIR algorithms, they can converge very slowly, and they are susceptible to local minima. Nonetheless, adaptive IIR algorithms are used in some applications (such as low frequency noise cancellation) in which the need for IIR-type responses is great. In some cases, the exact algorithm used by a company is a tightly guarded trade secret.

Most adaptive IIR algorithms minimize the **prediction** error, to linearize the estimation problem, as in deterministic or block linear prediction.

$$y_k = \sum_{n=1}^L (v_n^k y_{k-n}) + \sum_{n=0}^L (w_n^k x_{k-n})$$

Thus the coefficient vector is

$$W_k = \begin{pmatrix} v_1^k \\ v_2^k \\ \vdots \\ v_L^k \\ w_0^k \\ w_1^k \\ \vdots \\ w_L^k \end{pmatrix}$$

and the "signal" vector is

$$U_k = \begin{pmatrix} y_{k-1} \\ y_{k-2} \\ \vdots \\ y_{k-L} \\ x_k \\ x_{k-1} \\ \vdots \\ x_{k-L} \end{pmatrix}$$

The error is

$$\epsilon_k = d_k - y_k = d_k - W_k^T U_k$$

¹⁴This content is available online at <<http://cnx.org/content/m11912/1.1/>>.

An LMS algorithm can be derived using the approximation $E[\epsilon_k^2] = \epsilon_k^2$ or

$$\hat{\nabla}_k = \frac{\partial}{\partial W_k} (\epsilon_k^2) = 2\epsilon_k \frac{\partial}{\partial W_k} (\epsilon_k) = 2\epsilon_k \begin{pmatrix} \frac{\partial}{\partial v_1^k} (\epsilon_k) \\ \vdots \\ \frac{\partial}{\partial \epsilon_k} (w_1^k) \\ \vdots \end{pmatrix} = -2\epsilon_k \begin{pmatrix} \frac{\partial}{\partial v_1^k} (y_k) \\ \vdots \\ \frac{\partial}{\partial v_L^k} (y_k) \\ \frac{\partial}{\partial w_0^k} (y_k) \\ \vdots \\ \frac{\partial}{\partial w_L^k} (y_k) \end{pmatrix}$$

Now

$$\begin{aligned} \frac{\partial}{\partial v_i^k} (y_k) &= \frac{\partial}{\partial v_i^k} \left(\sum_{n=1}^L (v_n^k y_{k-n}) + \sum_{n=0}^L (w_n^k x_{k-n}) \right) = y_{k-n} + \sum_{n=1}^L \left(v_n^k \frac{\partial}{\partial v_i^k} (y_{k-n}) \right) + 0 \\ \frac{\partial}{\partial w_i^k} (y_k) &= \frac{\partial}{\partial w_i^k} \left(\sum_{n=1}^L (v_n^k y_{k-n}) + \sum_{n=0}^L (w_n^k x_{k-n}) \right) = \sum_{n=1}^L \left(v_n^k \frac{\partial}{\partial w_i^k} (y_{k-n}) \right) + x_{k-n} \end{aligned}$$

Note that these are difference equations in $\frac{\partial}{\partial v_i^k} (y_k)$, $\frac{\partial}{\partial w_i^k} (y_k)$: call them $\alpha_i^k = \frac{\partial}{\partial v_i^k} (y_k)$, $\beta_i^k = \frac{\partial}{\partial w_i^k} (y_k)$, then

$\hat{\nabla}_k = \begin{pmatrix} \beta_1^k & \beta_2^k & \dots & \beta_L^k & \alpha_0^k & \dots & \alpha_L^k \end{pmatrix}^T$, and the IIR LMS algorithm becomes

$$y_k = W_k^T U_k$$

$$\alpha_i^k = x_{k-i} + \sum_{j=1}^L (v_j^k \alpha_i^{k-j})$$

$$\beta_i^k = y_{k-i} + \sum_{j=1}^L (v_j^k \beta_i^{k-j})$$

$$\hat{\nabla}_k = -2\epsilon_k \begin{pmatrix} \beta_1^k & \beta_2^k & \dots & \alpha_0^k & \alpha_1^k & \dots & \alpha_L^k \end{pmatrix}^T$$

and finally

$$W_{k+1} = W_k - U \hat{\nabla}_k$$

where the μ may be *different* for the different IIR coefficients. Stability and convergence rate depends on these choices, of course. There are a number of variations on this algorithm.

Due to the slow convergence and the difficulties in tweaking the algorithm parameters to ensure stability, IIR algorithms are used only if there is an overriding need for an IIR-type filter.

5.5.3 The Constant-Modulus Algorithm and the Property-Restoral Principle¹⁵

The adaptive filter configurations that we have examined so far require a "desired signal" d_k . There are many clever ways to obtain such a signal, but in some potential applications a desired signal is simply not available. However, a "property-restoral algorithm" can sometimes circumvent this problem.

If the uncorrupted signal has special properties that are characteristic of the signal and not of the distortion or interference, an algorithm can be constructed which attempts to cause the output of the adaptive filter to exhibit that property. Hopefully, the adapting filter will restore that property by removing the distortion or interference!

¹⁵This content is available online at <<http://cnx.org/content/m11913/1.1/>>.

Example 5.5: the Constant-Modulus Algorithm (CMA)

Certain communication modulation schemes, such as PSK and FSK, transmit a sinusoid of a constant analytic magnitude. Only the frequency or phase change with time. The constant modulus algorithm tries to drive the output signal to one having a constant amplitude:

$$\epsilon_k = (|y_k|)^2 - A^2$$

One can derive an LMS (or other) algorithm that seeks a Wiener filter minimizing this error. In practice, this works very well for equalization of PSK and FSK communication channels.

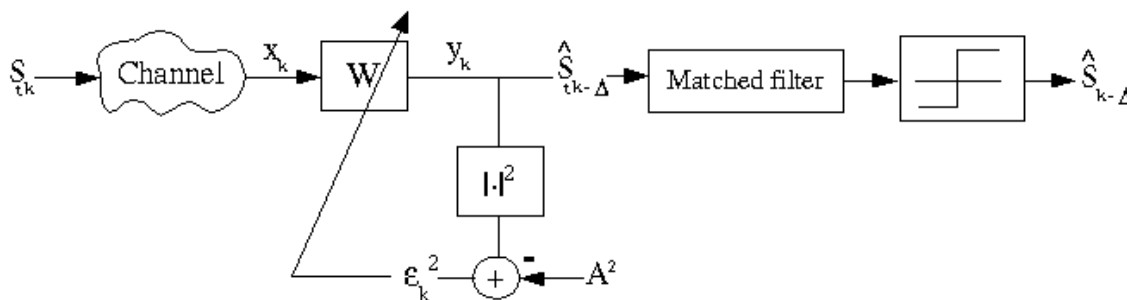


Figure 5.20

CMA is simpler than decision-directed feedback, and can work for high initial error rates!

This property-restoral idea can be used in any context in which a property-related error can be defined.

5.5.4 Complex LMS¹⁶

LMS for complex data and coefficients (such as quadrature communication systems) takes the form

$$y_k = W_k^H X_k$$

$$e_k = d_k - y_k$$

$$W_{k+1} = W_k + 2\mu e_k^* X_k$$

It is derived in exactly the same way as LMS, using the following complex vector differentiation formulas

$$\frac{d}{d\mathbf{W}} (P^H \mathbf{W}) = 0$$

$$\frac{d}{d\mathbf{W}} (W^H \mathbf{P}) = 2\mathbf{P}$$

$$\frac{d}{d\mathbf{W}} (W^H R \mathbf{W}) = 2R\mathbf{W}$$

or by differentiating with respect to the real and imaginary parts separately and recombining the results.

¹⁶This content is available online at <<http://cnx.org/content/m11914/1.3/>>.

5.5.5 Normalized LMS¹⁷

In "normalized" LMS, the gradient step factor μ is normalized by the energy of the data vector:

$$\mu_{NLMS} = \frac{\alpha}{X_k^H X_k + \sigma}$$

where α is usually $\frac{1}{2}$ and σ is a very small number introduced to prevent division by zero if $X_k^H X_k$ is very small.

$$W_{k+1} = W_k + \frac{1}{X_k^H X_k} e_k X_k$$

The normalization has several interpretations

1. corresponds to the 2nd-order convergence bound
2. makes the algorithm independent of signal scalings
3. adjusts W_{k+1} to give zero error with current input: $W_{k+1} X_k = d_k$
4. minimizes mean effort at time $k + 1$

NLMS usually converges *much* more quickly than LMS at very little extra cost; NLMS is very commonly used. In some applications, normalization is so universal that "we use the LMS algorithm" implies normalization as well.

5.6 Summary of Adaptive Filtering Methods¹⁸

1. **LMS** - remains the simplest and best algorithm when slow convergence is not a serious issue (typically used) $O(N)$
2. **NLMS** - simple extension of the LMS with much faster convergence in many cases (very commonly used) $O(N)$
3. **Frequency-domain methods** - offer computational savings ($O(\log N)$) for long filters and usually offer faster convergence, too (sometimes used; very commonly used when there are already FFTs in the system)
4. **Lattice methods** - are stable and converge quickly, but cost substantially more than LMS and have higher residual EMSE than many methods (very occasionally used) $O(N)$
5. **RLS** - algorithms that converge quickly and are stable exist. However, they are considerably more expensive than LMS. (almost never used) $O(N)$
6. **Block RLS** - (least squares) methods exist and can be pretty efficient in some cases. (occasionally used) $O(\log N)$, $O(N)$, $O(N^2)$
7. **IIR** - methods are difficult to implement successfully and pose certain difficulties, but are sometimes used in some applications, for example noise cancellation of low frequency noise (very occasionally used)
8. **CMA** - very useful when applicable (blind equalization); CMA is *the* method for blind equalizer initialization (commonly used in a few specific equalization applications) $O(N)$

NOTE: In general, getting adaptive filters to work well in an application is much more challenging than, say, FFTs or IIR filters; they generally require lots of tweaking!

¹⁷This content is available online at <<http://cnx.org/content/m11915/1.2/>>.

¹⁸This content is available online at <<http://cnx.org/content/m11916/1.1/>>.

Solutions to Exercises in Chapter 5

Solution to Exercise 5.2 (p. 264)

Estimate the statistics

$$r_{xx}^{\wedge}(l) \approx \frac{1}{N} \sum_{k=0}^{N-1} (x_k x_{k+l})$$

$$r_{xd}^{\wedge}(l) \approx \frac{1}{N} \sum_{k=0}^{N-1} (d_k x_{k-l})$$

then solve $\hat{W}_{opt} = \hat{R}^{-1} = \hat{P}$

Solution to Exercise 5.3 (p. 264)

Use short-time windowed estimates of the correlation functions.

EQUATION IN QUESTION:

$$\left(r_{xx}^{\wedge}(l)\right)^k = \frac{1}{N} \sum_{m=0}^{N-1} (x_{k-m} x_{k-m-l})$$

$$\left(r_{dx}^{\wedge}(l)\right)^k = \frac{1}{N} \sum_{m=0}^{N-1} (x_{k-m-l} d_{k-m})$$

and $W_{opt}^k \approx \left(\hat{R}_k\right)^{-1} \hat{P}_k$

Solution to Exercise 5.4 (p. 264)

Recursively!

$$r_{xx}^k(l) = r_{xx}^{k-1}(l) + x_k x_{k-l} - x_{k-N} x_{k-N-l}$$

This is critically stable, so people usually do

$$(1 - \alpha) r_{xx}^k(l) = \alpha r_{xx}^{k-1}(l) + x_k x_{k-l}$$

Chapter 6

Multirate Signal Processing

6.1 Overview of Multirate Signal Processing¹

Digital transformation of the sampling rate of signals, or signal processing with different sampling rates in the system.

6.1.1 Applications

1. **Sampling-rate conversion** - CD to DAT format change, for example.
2. **Improved D/A, A/D conversion** - oversampling converters; which reduce performance requirements on anti-aliasing or reconstruction filters
3. **FDM channel modulation and processing** - bandwidth of individual channels is much less than the overall bandwidth
4. **Subband coding of speech and images** - Eyes and ears are not as sensitive to errors in higher frequency bands, so many coding schemes split signals into different frequency bands and quantize higher-frequency bands with much less precision.

6.1.2 Outline of Multirate DSP material

1. General Rate-changing System (Section 6.1.3: General Rate-Changing Procedure)
2. Integer-factor Interpolation and Decimation and Rational-factor Rate Changing (Section 6.2)
3. Efficient Multirate Filter Structures (Section 6.3)
4. Optimal Filter Design for Multirate Systems (Section 6.4)
5. Multi-stage Multirate Systems (Section 6.5)
6. Oversampling D/As (Section 6.6)
7. Perfect-Reconstruction Filter Banks and Quadrature Mirror Filters (Section 6.7)

6.1.3 General Rate-Changing Procedure

This procedure is motivated by an analog-based method: one conceptually simple method to change the sampling rate is to simply convert a digital signal to an analog signal and resample it! (Figure 6.1)

¹This content is available online at <<http://cnx.org/content/m12777/1.2/>>.

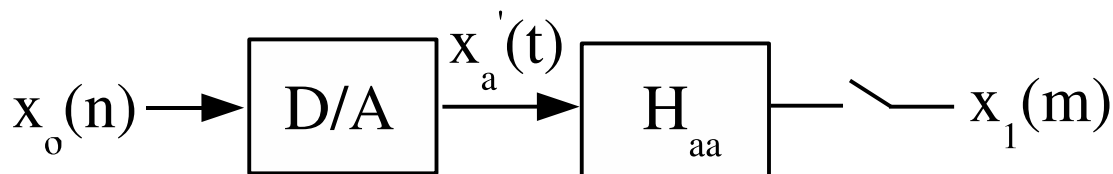


Figure 6.1

$$H_{aa}(\Omega) = \begin{cases} 1 & \text{if } |\Omega| < \frac{\pi}{T_1} \\ 0 & \text{otherwise} \end{cases}$$

$$h_{aa}(t) = \frac{\sin\left(\frac{\pi}{T_1}t\right)}{\frac{\pi}{T_1}t}$$

Recall the ideal D/A:

$$x'_a(t) = \sum_{n=-\infty}^{\infty} \left(x_0(n) \frac{\sin\left(\frac{\pi(t-nT_0)}{T_0}\right)}{\frac{\pi(t-nT_0)}{T_0}} \right) \quad (6.1)$$

The problems with this scheme are:

1. A/D, D/A, filters cost money
2. imperfections in these devices introduce errors

Digital implementation of rate-changing according to this formula has three problems:

1. Infinite sum: The solution is to truncate. Consider $\text{sinc}t \approx 0$ for $t < t_1$, $t > t_2$: Then $mT_1 - nT_0 < t_1$ and $mT_1 - nT_0 > t_2$ which implies

$$N_1 = \lceil \frac{mT_1 - t_2}{T_0} \rceil$$

$$N_2 = \lfloor \frac{mT_1 - t_1}{T_0} \rfloor$$

$$x_1(m) = \sum_{n=N_1}^{N_2} (x_0(n) \text{sinc}_{T'}(mT_1 - nT_0))$$

NOTE: This is essentially lowpass filter design using a boxcar window: other finite-length filter design methods could be used for this.

2. Lack of causality²: The solution is to delay by $\max\{|N|\}$ samples. The mathematics of the analog portions of this system can be implemented digitally.

$$\begin{aligned} x_1(m) &= h_{aa}(t) * x'_a(t) \big|_{t=mT_1} \\ &= \int_{-\infty}^{\infty} \left(\sum_{n=-\infty}^{\infty} \left(x_0(n) \frac{\sin\left(\frac{\pi(mT_1 - \tau - nT_0)}{T_0}\right)}{\frac{\pi(mT_1 - \tau - nT_0)}{T_0}} \right) \right) \frac{\sin\left(\frac{\pi\tau}{T_1}\right)}{\frac{\pi\tau}{T_1}} d\tau \end{aligned} \quad (6.2)$$

²"Properties of Systems": Section Causality <<http://cnx.org/content/m2102/latest/#causality>>

$$\begin{aligned}
 x_1(m) &= \sum_{n=-\infty}^{\infty} \left(x_0(n) \frac{\sin\left(\frac{\pi}{T'}(mT_1 - nT_0)\right)}{\frac{\pi}{T'}(mT_1 - nT_0)} \right) \Big|_{T'=\max\{T_0, T_1\}} \\
 &= \sum_{n=-\infty}^{\infty} (x_0(n) \operatorname{sinc}_{T'}(mT_1 - nT_0))
 \end{aligned} \tag{6.3}$$

So we have an all-digital formula for *exact* digital-to-digital rate changing!

3. Cost of computing $\operatorname{sinc}_{T'}(mT_1 - nT_0)$: The solution is to precompute the table of $\operatorname{sinc}(t)$ values. However, if $\frac{T_1}{T_0}$ is not a rational fraction, an infinite number of samples will be needed, so some approximation will have to be tolerated.

NOTE: Rate transformation of any rate to any other rate can be accomplished digitally with arbitrary precision (if some delay is acceptable). This method is used in practice in many cases. We will examine a number of special cases and computational improvements, but in some sense everything that follows are details; the above idea is the central idea in multirate signal processing.

Useful references for the traditional material (everything except PRFBs) are *Crochiere and Rabiner, 1981*[9] and *Crochiere and Rabiner, 1983*[10]. A more recent tutorial is *Vaidyanathan*[29]; see also *Rioul and Vetterli*[25]. References to most of the original papers can be found in these tutorials.

6.2 Interpolation, Decimation, and Rate Changing by Integer Fractions³

6.2.1 Interpolation: by an integer factor L

Interpolation means increasing the sampling rate, or filling in in-between samples. Equivalent to sampling a bandlimited analog signal L times faster. For the ideal interpolator,

$$X_1(\omega) = \begin{cases} X_0(L\omega) & \text{if } |\omega| < \frac{\pi}{L} \\ 0 & \text{if } \frac{\pi}{L} \leq |\omega| \leq \pi \end{cases} \tag{6.4}$$

We wish to accomplish this digitally. Consider (6.5) and Figure 6.2.

$$y(m) = \begin{cases} X_0\left(\frac{m}{L}\right) & \text{if } m = \{0, \pm L, \pm (2L), \dots\} \\ 0 & \text{otherwise} \end{cases} \tag{6.5}$$

³This content is available online at <<http://cnx.org/content/m12801/1.2/>>.

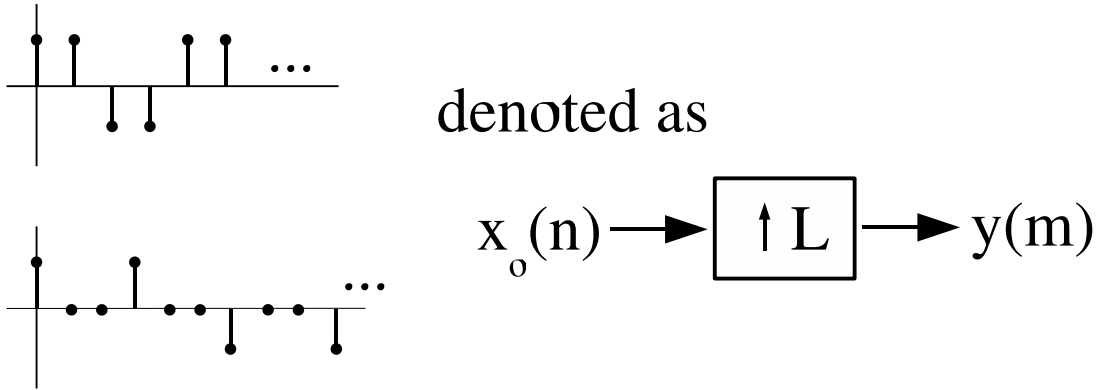


Figure 6.2

The DTFT of $y(m)$ is

$$\begin{aligned}
 Y(\omega) &= \sum_{m=-\infty}^{\infty} (y(m) e^{-j\omega m}) \\
 &= \sum_{n=-\infty}^{\infty} (x_0(n) e^{-j\omega Ln}) \\
 &= \sum_{n=-\infty}^{\infty} (x(n) e^{-j\omega Ln}) \\
 &= X_0(\omega L)
 \end{aligned} \tag{6.6}$$

Since $X_0(\omega')$ is periodic with a period of 2π , $X_0(L\omega) = Y(\omega)$ is periodic with a period of $\frac{2\pi}{L}$ (see Figure 6.3).

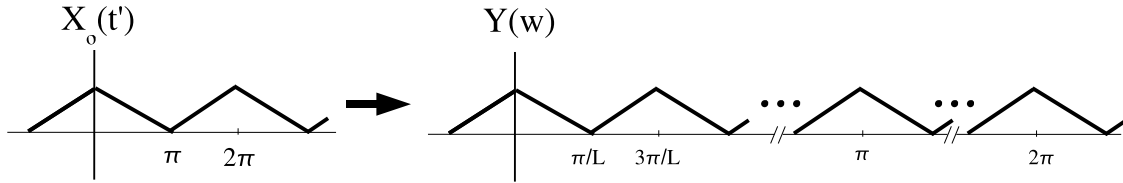


Figure 6.3

By inserting zero samples between the samples of $x_0(n)$, we obtain a signal with a scaled frequency response that simply replicates $X_0(\omega')$ L times over a 2π interval!

Obviously, the desired $x_1(m)$ can be obtained simply by lowpass filtering $y(m)$ to remove the replicas.

$$x_1(m) = y(m) * h_L(m) \tag{6.7}$$

Given

$$H_L(m) = \begin{cases} 1 & \text{if } |\omega| < \frac{\pi}{L} \\ 0 & \text{if } \frac{\pi}{L} \leq |\omega| \leq \pi \end{cases}$$

In practice, a finite-length lowpass filter is designed using any of the methods studied so far (Figure 6.4 (Interpolator Block Diagram)).

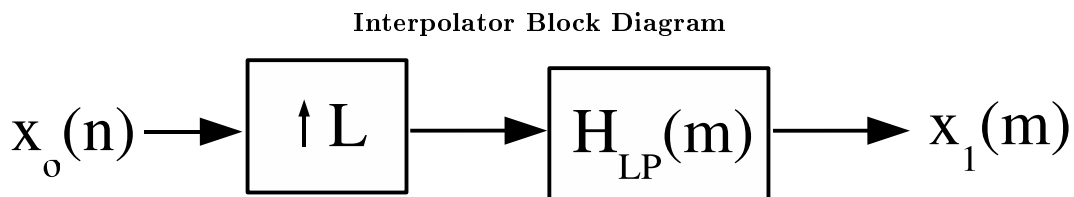


Figure 6.4

6.2.2 Decimation: sampling rate reduction (by an integer factor M)

Let $y(m) = x_0(Lm)$ (Figure 6.5)

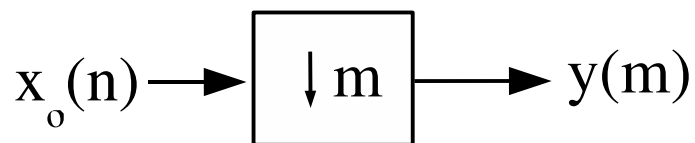


Figure 6.5

That is, keep only every L th sample (Figure 6.6)

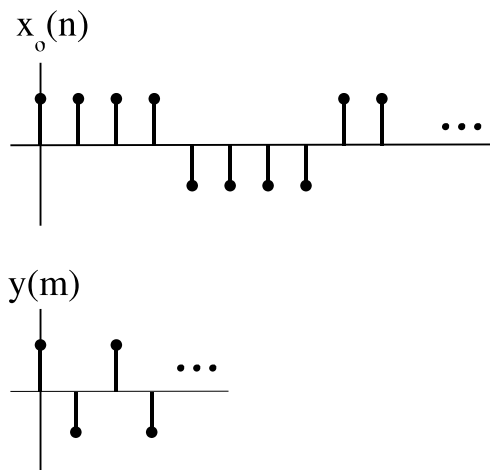


Figure 6.6

In frequency (DTFT):

$$\begin{aligned}
 Y(\omega) &= \sum_{m=-\infty}^{\infty} (y(m) e^{-j\omega m}) \\
 &= \sum_{m=-\infty}^{\infty} (x_0(Mm) e^{-j\omega m}) \\
 &= \sum_{n=-\infty}^{\infty} \left(x_0(n) \left(\sum_{k=-\infty}^{\infty} (\delta(n - Mk)) \right) e^{-j\omega \frac{n}{M}} \right) \Big|_{n=Mm} \\
 &= \sum_{n=-\infty}^{\infty} \left(x_0(n) \left(\sum_{k=-\infty}^{\infty} (\delta(n - Mk)) \right) e^{-j\omega' n} \right) \Big|_{\omega' = \frac{\omega}{M}} \\
 &= DTFT[x_0(n)] * DTFT[\sum (\delta(n - Mk))]
 \end{aligned} \tag{6.8}$$

Now $DTFT[\sum (\delta(n - Mk))] = 2\pi \sum_{k=0}^{M-1} (X(k) \delta(\omega, -\frac{2\pi k}{M}))$ for $|\omega| < \pi$ as shown in homework #1, where $X(k)$ is the DFT of one period of the periodic sequence. In this case, $X(k) = 1$ for $k \in \{0, 1, \dots, M-1\}$ and $DTFT[\sum (\delta(n - Mk))] = 2\pi \sum_{k=0}^{M-1} (\delta(\omega, -\frac{2\pi k}{M}))$.

$$\begin{aligned}
 DTFT[x_0(n)] * DTFT[\sum (\delta(n - Mk))] &= X_0(\omega') * 2\pi \sum_{k=0}^{M-1} (\delta(\omega, -\frac{2\pi k}{M})) \\
 &= \frac{1}{2\pi} \int_{-\pi}^{\pi} X_0(\mu') \left(2\pi \sum_{k=0}^{M-1} (\delta(\omega, -\mu' - \frac{2\pi k}{M})) \right) d\mu' \\
 &= \sum_{k=0}^{M-1} (X_0(\omega - \frac{2\pi k}{M}))
 \end{aligned} \tag{6.9}$$

so $Y(\omega) = \sum_{k=0}^{M-1} (X_0(\frac{\omega}{M} - \frac{2\pi k}{M}))$ i.e., we get **digital aliasing**. (Figure 6.7)

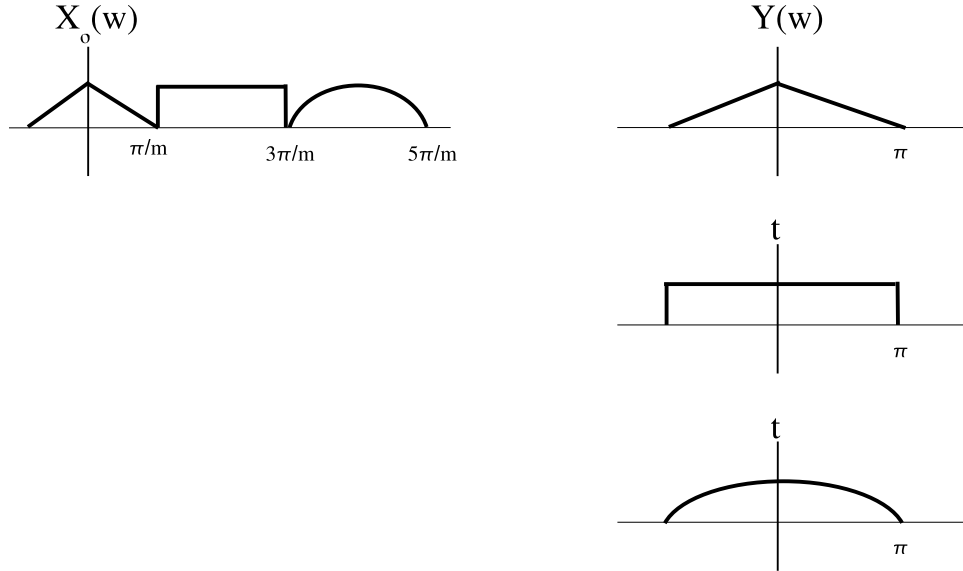


Figure 6.7

Usually, we prefer not to have aliasing, so the downsampler is preceded by a lowpass filter to remove all frequency components above $|\omega| < \frac{\pi}{M}$ (Figure 6.8).

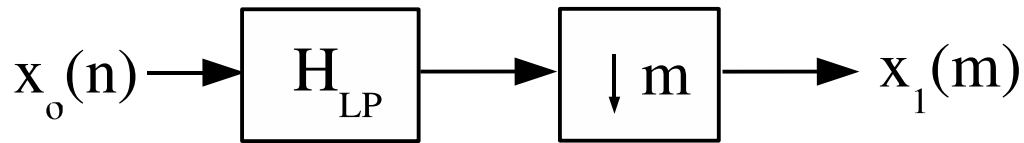


Figure 6.8

6.2.3 Rate-Changing by a Rational Fraction L/M

This is easily accomplished by interpolating by a factor of L , then decimating by a factor of M (Figure 6.9).

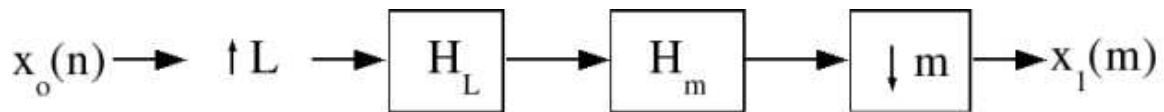


Figure 6.9

The two lowpass filters can be combined into one LP filter with the lower cutoff,

$$H(\omega) = \begin{cases} 1 & \text{if } |\omega| < \frac{\pi}{\max\{L, M\}} \\ 0 & \text{if } \frac{\pi}{\max\{L, M\}} \leq |\omega| \leq \pi \end{cases}$$

Obviously, the computational complexity and simplicity of implementation will depend on $\frac{L}{M}$: $2/3$ will be easier to implement than $1061/1060$!

6.3 Efficient Multirate Filter Structures⁴

Rate-changing appears expensive computationally, since for both decimation and interpolation the lowpass filter is implemented at the higher rate. However, this is not necessary.

6.3.1 Interpolation

For the interpolator, most of the samples in the upsampled signal are zero, and thus require no computation. (Figure 6.10)

⁴This content is available online at <http://cnx.org/content/m12800/1.2/>.

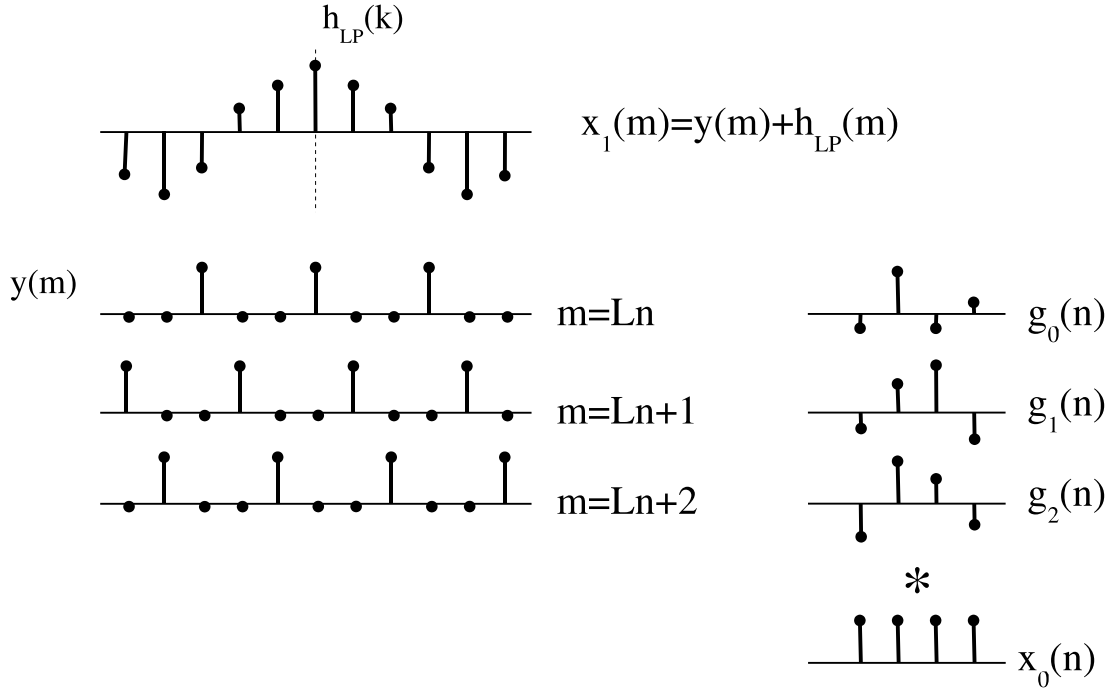


Figure 6.10

For $m = L \lfloor \frac{m}{L} \rfloor + m \bmod L$ and $p = m \bmod L$,

$$\begin{aligned}
 x_1(m) &= \sum_{m=N_1}^{N_2} (h_{LP}(m) y(m)) \\
 &= \sum_{k=\frac{N_1}{L}}^{\frac{N_2}{L}} (g_p(k) x_0(\lfloor \frac{m}{L} \rfloor - k))
 \end{aligned} \tag{6.10}$$

$$g_p(n) = h(Ln + p)$$

Pictorially, this can be represented as in Figure 6.11.

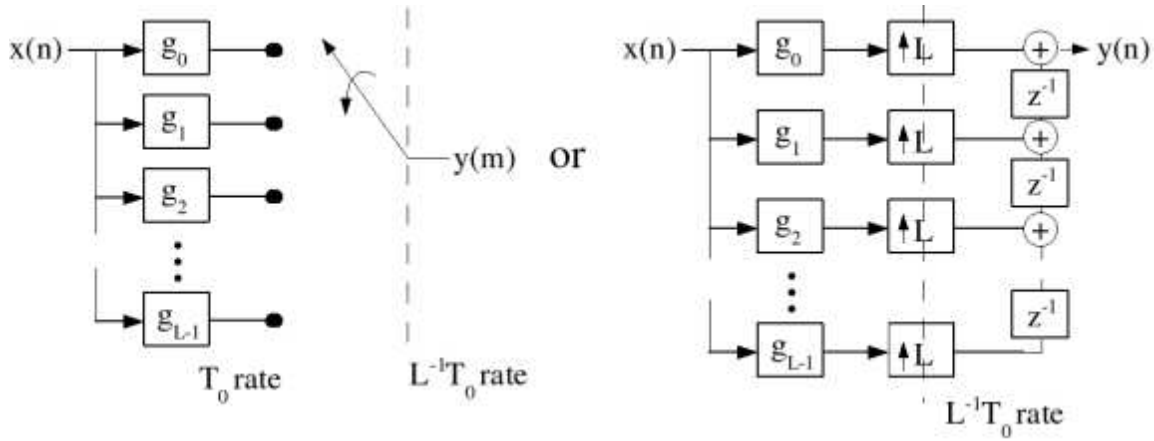


Figure 6.11

These are called **polyphase structures**, and the $g_p(n)$ are called **polyphase filters**.

Computational cost

If $h(m)$ is a length- N filter:

- No simplification: $\frac{N}{T_1} = \frac{LN}{T_0} \frac{\text{computations}}{\text{sec}}$
- Polyphase structure: $\left(L \frac{L}{N} \frac{1}{T_0}\right) \frac{\text{computations}}{\text{sec}} = \frac{N}{T_0}$ where L is the number of filters, $\frac{N}{L}$ is the taps/filter, and $\frac{1}{T_0}$ is the rate.

Thus we save a factor of L by not being dumb.

NOTE: For a given precision, N is proportional to L , (why?), so the computational cost does increase with the interpolation rate.

QUESTION: Can similar computational savings be obtained with IIR structures?

6.3.2 Efficient Decimation Structures

We only want every M th output, so we compute only the outputs of interest. (Figure 6.12 (Polyphase Decimation Structure))

$$x_1(m) = \sum_{k=N_1}^{N_2} (x_0(Lm - k) h(k))$$

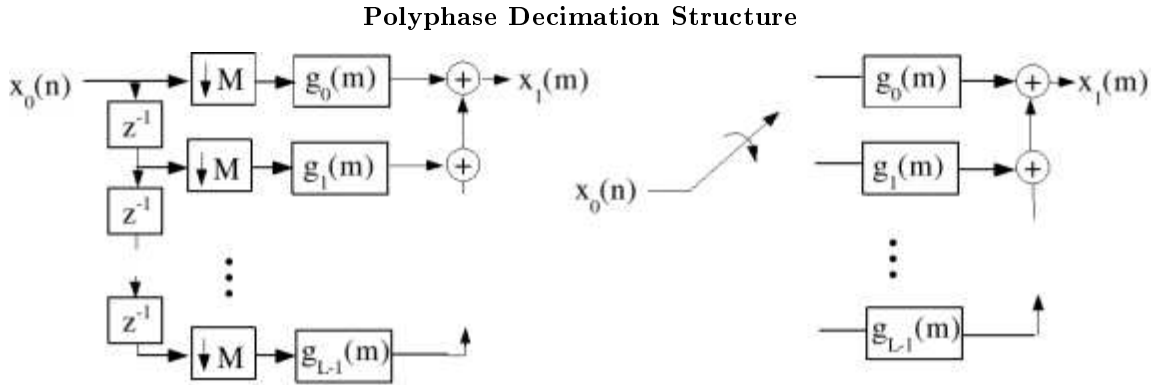


Figure 6.12

The decimation structures are flow-graph reversals of the interpolation structure. Although direct implementation of the full filter for every M th sample is obvious and straightforward, these polyphase structures give some idea as to how one might evenly partition the computation over M cycles.

6.3.3 Efficient L/M rate changers

Interpolate by L and decimate by M (Figure 6.13).

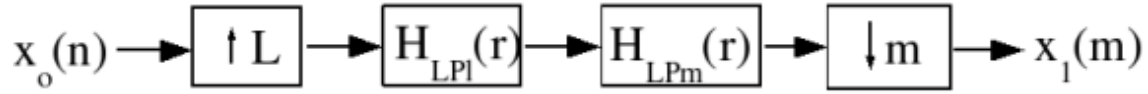


Figure 6.13

Combine the lowpass filters (Figure 6.14).

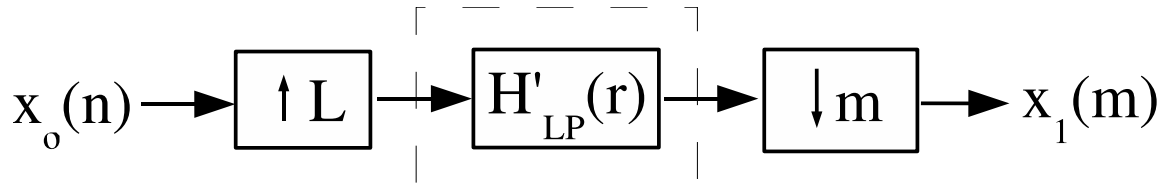


Figure 6.14

We can couple the lowpass filter either to the interpolator or the decimator to implement it efficiently (Figure 6.15).

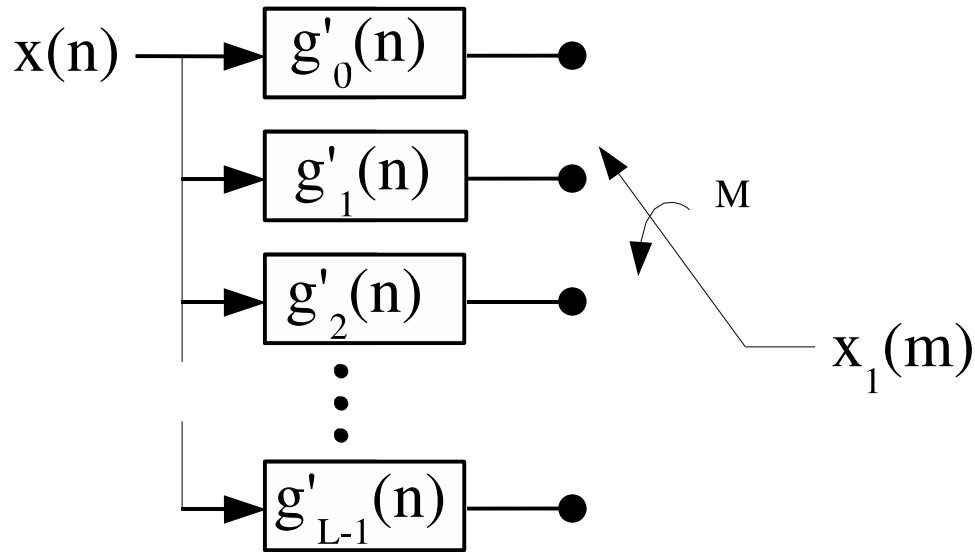


Figure 6.15

Of course we only compute the polyphase filter output selected by the decimator.

Computational Cost

Every $T_1 = \frac{M}{L}T_0$ seconds, compute one polyphase filter of length $\frac{N}{L}$, or

$$\frac{\frac{N}{L}}{T_1} = \frac{\frac{N}{L}}{\frac{M}{L}T_0} = \frac{N}{MT_0} \frac{\text{multiplies}}{\text{second}}$$

However, note that N is proportional to $\max\{L, M\}$.

6.4 Filter Design for Multirate Systems⁵

The filter design techniques⁶ learned earlier can be applied to the design of filters in multirate systems, with a few twists.

Example 6.1

Design a factor-of- L interpolator for use in a CD player, we might wish that the out-of-band error be below the least significant bit, or 96dB down, and $< 0.05\%$ error in the passband, so these specifications could be used for optimal L_∞ filter design.

In a CD player, the sampling rate is 44.1kHz, corresponding to a Nyquist frequency of 22.05kHz, but the sampled signal is bandlimited to 20kHz. This leaves a small transition band, from 20kHz to 24.1kHz. However, note that in any case where the signal spectrum is zero over some band, this introduces *other* zero bands in the scaled, replicated spectrum (Figure 6.16).

⁵This content is available online at <http://cnx.org/content/m12773/1.2/>.

⁶Digital Filter Design <http://cnx.org/content/col10285/latest/>

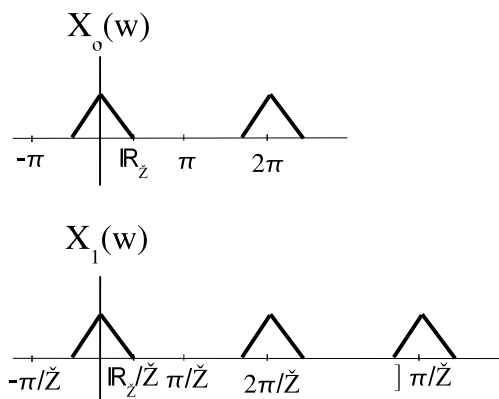


Figure 6.16

So we need only control the filter response in the stopbands over the frequency regions with nonzero energy. (Figure 6.17)

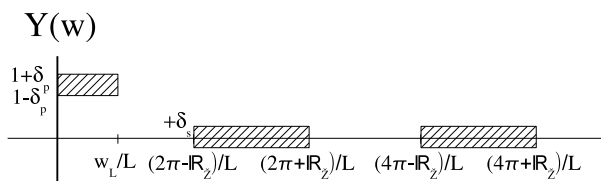


Figure 6.17

The extra "don't care" bands allow a given set of specifications to be satisfied with a shorter-length filter.

6.4.1 Direct polyphase filter design

Note that in an integer-factor interpolator, each set of output samples $x_1(Ln + p)$, $p = \{0, 1, \dots, L-1\}$, is created by a different polyphase filter $g_p(n)$, which has no interaction with the other polyphase filters except in that they each interpolate the same signal. We can thus treat the design of each polyphase filter *independently*, as an $\frac{N}{L}$ -length filter design problem. (Figure 6.18)

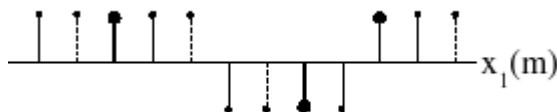


Figure 6.18

Each $g_p(n)$ produces samples $x_1(Ln + p) = x_0(n + \frac{p}{L})$, where $n + \frac{p}{L}$ is not an integer. That is, $g_p(n)$ is to produce an output signal (at a T_0 rate) that is $x_0(n)$ time-advanced by a non-integer advance $\frac{p}{L}$.

The desired response of this polyphase filter is thus

$$H_{Dp}(\omega) = e^{\frac{j\omega p}{L}}$$

for $|\omega| \leq \pi$, an all-pass filter with a linear, non-integer, phase. Each polyphase filter can be designed independently to approximate this response according to any of the design criteria developed so far.

Exercise 6.1

(Solution on p. 312.)

What should the polyphase filter for $p = 0$ be?

Example 6.2: Least-squares Polyphase Filter Design

Deterministic $x(n)$ - Minimize

$$\sum_{n=-\infty}^{\infty} \left((|x(n) - x_d(n)|)^2 \right)$$

Given $x(n) = x(n) * h(n)$ and $x_d(n) = x(n) * h_d(n)$. Using Parseval's theorem, this becomes

$$\begin{aligned} \min \left\{ \sum_{n=-\infty}^{\infty} \left((|x(n) - x_d(n)|)^2 \right) \right\} &= \min \left\{ \frac{1}{2\pi} \int_{-\pi}^{\pi} (|X(\omega)H(\omega) - X(\omega)H_d(\omega)|)^2 d\omega \right\} \\ &= \min \left\{ \frac{1}{2\pi} \int_{-\pi}^{\pi} |H(\omega) - H_d(\omega)|^2 |X(\omega)|^2 d\omega \right\} \end{aligned} \quad (6.11)$$

This is simply weighted least squares design, with $(|X(\omega)|)^2$ as the weighting function.

stochastic $X(\omega)$ -

$$\begin{aligned} \min \left\{ E \left[(|x(n) - x_d(n)|)^2 \right] \right\} &= E \left[(|x(n) * (h(n) - h_d(n))|)^2 \right] \\ &= \min \left\{ \frac{1}{2\pi} \int_{-\pi}^{\pi} |H_d(\omega) - H(\omega)|^2 S_{xx}(\omega) d\omega \right\} \end{aligned} \quad (6.12)$$

$S_{xx}(\omega)$ is the power spectral density of x .

$$S_{xx}(\omega) = DTFT[r_{xx}(k)]$$

$$r_{xx}(k) = E \left[x(k+l) \overline{x(l)} \right]$$

Again, a weighted least squares filter design problem.

Problem

Is it feasible to use IIR polyphase filters?

Solution

The recursive feedback of previous outputs means that portions of each IIR polyphase filter must be computed for every input sample; this usually makes IIR filters more expensive than FIR implementations.

6.5 Multistage Multirate Systems⁷

Multistage multirate systems are often more efficient. Suppose one wishes to decimate a signal by an integer factor M , where M is a composite integer $M = M_1 M_2 \dots M_p = \prod_{i=1}^p M_i$. A decimator can be implemented in a multistage fashion by first decimating by a factor M_1 , then decimating this signal by a factor M_2 , etc. (Figure 6.19 (Multistage decimator))

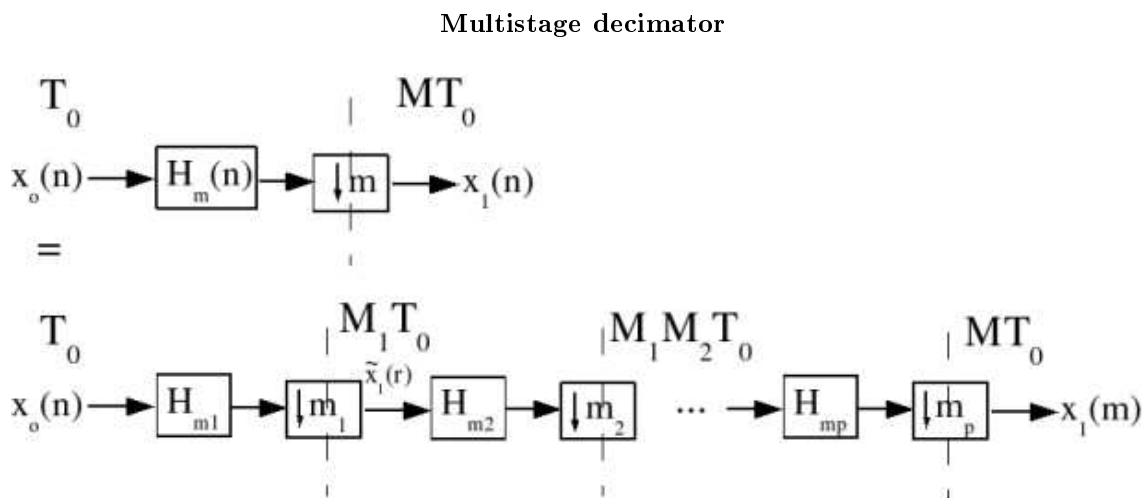


Figure 6.19

Multistage implementations are of significant practical interest only if they offer significant computational savings. In fact, they often do!

The computational cost of a *single-stage* interpolator is:

$$\frac{N}{MT_0} \frac{\text{taps}}{\text{sec}}$$

The computational cost of a *multistage* interpolator is:

$$\frac{N_1}{M_1 T_0} + \frac{N_2}{M_1 M_2 T_0} + \dots + \frac{N_p}{M T_0}$$

The first term is the most significant, since the rate is highest. Since $(N_i \propto M_i)$ for a lowpass filter, it is not immediately clear that a multistage system should require less computation. However, the multistage structure relaxes the requirements on the filters, which reduces their length and makes the overall computation less.

6.5.1 Filter design for Multi-stage Structures

Ostensibly, the first-stage filter must be a lowpass filter with a cutoff at $\frac{\pi}{M_1}$, to prevent aliasing after the downsampler. However, note that aliasing outside the *final overall* passband $|\omega| < \frac{\pi}{M}$ is of no concern, since it will be removed by later stages. We only need prevent aliasing into the band $|\omega| < \frac{\pi}{M}$; thus we need only specify the passband over the interval $|\omega| < \frac{\pi}{M}$, and the stopband over the intervals $\omega \in \left[\frac{2\pi k}{M_1} - \frac{\pi}{M}, \frac{2\pi k}{M_1} + \frac{\pi}{M} \right]$, for $k \in \{1, \dots, M-1\}$. (Figure 6.20)

⁷This content is available online at <http://cnx.org/content/m12803/1.2/>.

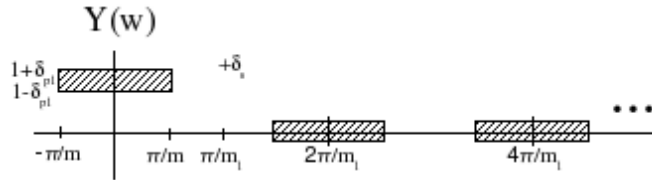


Figure 6.20

Of course, we don't want *gain* in the transition bands, since this would need to be suppressed later, but otherwise we don't care about the response in those regions. Since the transition bands are so large, the required filter turns out to be quite short. The final stage has no "don't care" regions; however, it is operating at a low rate, so it is relatively unimportant if the final filter turns out to be rather long!

6.5.2 L-infinity Tolerances on the Pass and Stopbands

The overall response is a cascade of multiple filters, so the worst-case overall passband deviation, assuming all the peaks just happen to line up, is

$$1 + \delta_{p_{ov}} = \prod_{i=1}^p (1 + \delta_{p_i})$$

$$1 - \delta_{p_{ov}} = \prod_{i=1}^p (1 - \delta_{p_i})$$

So one could choose all filters to have equal specifications and require for each-stage filter. For $(\delta_{p_{ov}} \ll 1)$,

$$1 + \delta_{p_i}^+ \leq \sqrt[p]{1 + \delta_{p_{ov}}} \approx 1 + p^{-1} \delta_{p_{ov}}$$

$$1 - \delta_{p_i}^- \leq \sqrt[p]{1 - \delta_{p_{ov}}} \approx 1 - p^{-1} \delta_{p_{ov}}$$

Alternatively, one can design later stages (at lower computation rates) to compensate for the passband ripple in earlier stages to achieve exceptionally accurate passband response.

δ_s remains essentially unchanged, since the worst-case scenario is for the error to alias into the passband and undergo no further suppression in subsequent stages.

6.5.3 Interpolation

Interpolation is the flow-graph reversal of the multi-stage decimator. The first stage has a cutoff at $\frac{\pi}{L}$ (Figure 6.21):

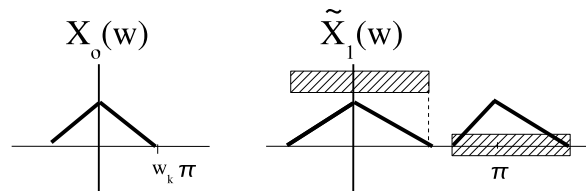


Figure 6.21

However, all subsequent stages have large bands without signal energy, due to the earlier stages (Figure 6.22):

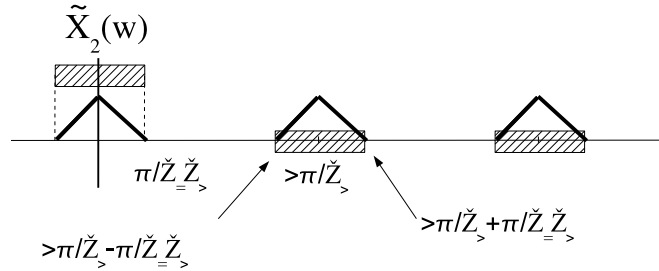


Figure 6.22

The order of the filters is reversed, but otherwise the filters are identical to the decimator filters.

6.5.4 Efficient Narrowband Lowpass Filtering

A very narrow lowpass filter requires a very long FIR filter to achieve reasonable resolution in the frequency response. However, were the input sampled at a lower rate, the cutoff frequency would be correspondingly higher, and the filter could be shorter!

The transition band is also broader, which helps as well. Thus, Figure 6.23 can be implemented as Figure 6.24.



Figure 6.23



Figure 6.24

and in practice the inner lowpass filter can be coupled to the decimator or interpolator filters. If the decimator and interpolator are implemented as multistage structures, the overall algorithm can be dramatically more efficient than direct implementation!

6.6 DFT-Based Filterbanks⁸

One common application of multirate processing arises in multirate, multi-channel filter banks (Figure 6.25).

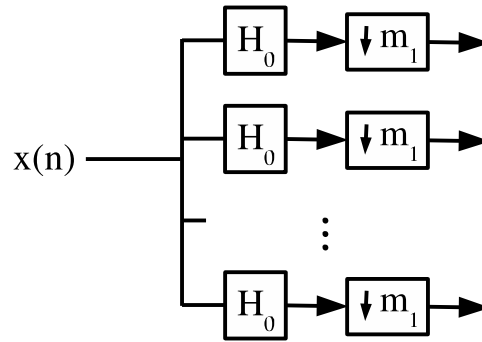


Figure 6.25

One application is separating frequency-division-multiplexed channels. If the filters are narrowband, the output channels can be decimated without significant aliasing.

Such structures are especially attractive when they can be implemented efficiently. For example, if the filters are simply frequency modulated (by $e^{-j\frac{2\pi k}{L}n}$) versions of each other, they can be efficiently implemented using FFTs!

Furthermore, there are classes of filters called **perfect reconstruction filters** which are of finite length but from which the signal can be reconstructed exactly (using all M channels), even though the output of each channel experiences aliasing in the decimation step. These types of filterbanks have received a lot of research attention, culminating in wavelet theory and techniques.

6.6.1 Uniform DFT Filter Banks

Suppose we wish to split a digital input signal into N frequency bands, uniformly spaced at center frequencies

$$\omega_k = \frac{2\pi k}{N}, \text{ for } 0 \leq k \leq N-1. \text{ Consider also a lowpass filter } h(n), H(\omega) \approx \begin{cases} 1 & \text{if } |\omega| < \frac{\pi}{N} \\ 0 & \text{otherwise} \end{cases}. \text{ Bandpass}$$

filters can be constructed which have the frequency response

$$H_k(\omega) = H\left(\omega + \frac{2\pi k}{N}\right)$$

from

$$h_k(n) = h(n) e^{-j\frac{2\pi kn}{N}}$$

The output of the k th bandpass filter is simply (assume $h(n)$ are FIR)

$$\begin{aligned} x(n) * h_k(n) &= \sum_{m=0}^{M-1} \left(x(n-m) h(m) e^{-j\frac{2\pi km}{N}} \right) \\ &= y_k(n) \end{aligned} \quad (6.13)$$

This looks suspiciously like a DFT, except that $M \neq N$, in general. However, if we fix $M = N$, then we can compute *all* $y_k(n)$ outputs simultaneously using an FFT of $x(n-m)h(m)$: The

⁸This content is available online at <<http://cnx.org/content/m12771/1.2/>>.

kth FFT frequency output = $y_k(n)$! So the cost of computing all of these filter banks outputs is $O[N \log N]$, rather than N^2 , per a given n . This is very useful for efficient implementation of **transmultiplexors** (FDM to TDM).

Exercise 6.2*(Solution on p. 312.)*

How would we implement this efficiently if we wanted to decimate the individual channels $y_k(n)$ by a factor of N , to their approximate Nyquist bandwidth?

Exercise 6.3*(Solution on p. 312.)*

Do you expect significant aliasing? If so, how do you propose to combat it? Efficiently?

Exercise 6.4*(Solution on p. 312.)*

How might one convert from N input channels into an FDM signal efficiently? (Figure 6.26)

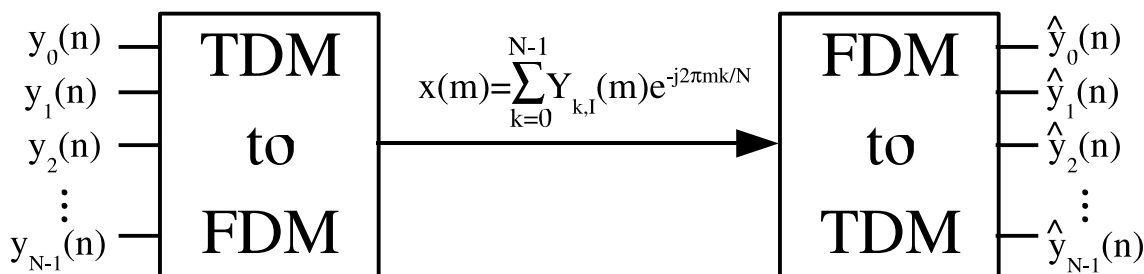


Figure 6.26

NOTE: Such systems are used throughout the telephone system, satellite communication links, etc.

6.7 Quadrature Mirror Filterbanks (QMF)⁹

Although the DFT filterbanks are widely used, there is a problem with aliasing in the decimated channels. At first glance, one might think that this is an insurmountable problem and must simply be accepted. Clearly, with FIR filters and maximal decimation, aliasing will occur. However, a simple example will show that it is possible to *exactly cancel out aliasing* under certain conditions!!!

Consider the following trivial filterbank system, with two channels. (Figure 6.27)

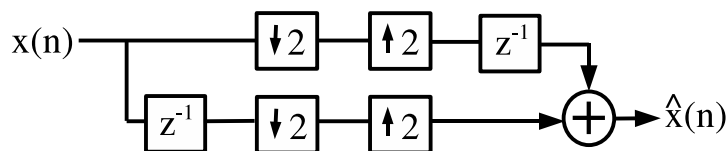


Figure 6.27

⁹This content is available online at <http://cnx.org/content/m12770/1.2/>.

Note $\hat{x}(n) = x(n)$ with no error whatsoever, although clearly aliasing occurs in both channels! Note that the overall data rate is still the Nyquist rate, so there are clearly enough degrees of freedom available to reconstruct the data, if the filterbank is designed carefully. However, this isn't splitting the data into separate frequency bands, so one questions whether something other than this trivial example could work.

Let's consider a general two-channel filterbank, and try to determine conditions under which aliasing can be cancelled, and the signal can be reconstructed perfectly (Figure 6.28).

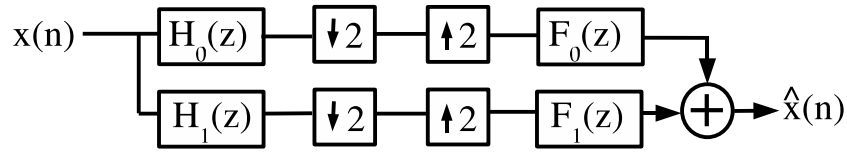


Figure 6.28

Let's derive $\hat{x}(n)$, using z-transforms, in terms of the components of this system. Recall (Figure 6.29) is equivalent to

$$Y(z) = H(z)X(z)$$

$$Y(\omega) = H(\omega)X(\omega)$$

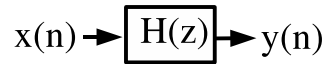


Figure 6.29

and note that (Figure 6.30) is equivalent to

$$Y(z) = \sum_{m=-\infty}^{\infty} \left(x(m) z^{-(Lm)} \right) = x(z^L)$$

$$Y(\omega) = X(L\omega)$$

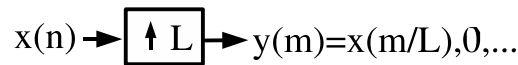


Figure 6.30

and (Figure 6.31) is equivalent to

$$Y(z) = \frac{1}{M} \sum_{k=0}^{M-1} \left(X \left(z^{\frac{1}{M}} W_M^k \right) \right)$$

$$Y(\omega) = \frac{1}{M} \sum_{k=0}^{M-1} \left(X \left(\frac{\omega}{M} + \frac{2\pi k}{M} \right) \right)$$

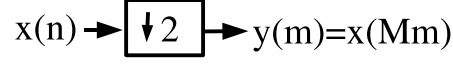


Figure 6.31

$Y(z)$ is derived in the downsampler as follows:

$$Y(z) = \sum_{m=-\infty}^{\infty} (x(Mm) z^{-m})$$

Let $n = Mm$ and $m = \frac{n}{M}$, then

$$Y(z) = \sum_{n=-\infty}^{\infty} \left(x(n) \left(\sum_{p=-\infty}^{\infty} (\delta(n - Mp)) \right) z^{-\left(\frac{n}{M}\right)} \right)$$

Now

$$\begin{aligned} x(n) \sum_{p=-\infty}^{\infty} (\delta(n - Mp)) &= IDFT \left[x(\omega) * \frac{2\pi}{M} \sum_{k=0}^{M-1} (\delta(\omega - \frac{2\pi k}{M})) \right] \\ &= IDFT \left[\frac{2\pi}{M} \sum_{k=0}^{M-1} (X(\omega - \frac{2\pi k}{M})) \right] \\ &= \frac{1}{M} \sum_{k=0}^{M-1} (X(n) W_M^{-nk}) \Big|_{W_M = e^{-j\frac{2\pi}{M}}} \end{aligned} \quad (6.14)$$

so

$$\begin{aligned} Y(z) &= \sum_{n=-\infty}^{\infty} \left(\left(\frac{1}{M} \sum_{k=0}^{M-1} (x(n) W_M^{-nk}) \right) z^{-\left(\frac{n}{M}\right)} \right) \\ &= \frac{1}{M} \sum_{k=0}^{M-1} \left(x(n) \left(W_M^{+k} z^{+\frac{1}{M}} \right)^{-n} \right) \\ &= \frac{1}{M} \sum_{k=0}^{M-1} \left(X \left(z^{\frac{1}{M}} W_M^k \right) \right) \end{aligned} \quad (6.15)$$

Armed with these results, let's determine $(\hat{X}(z) \Leftrightarrow \hat{x}(n))$. (Figure 6.32)

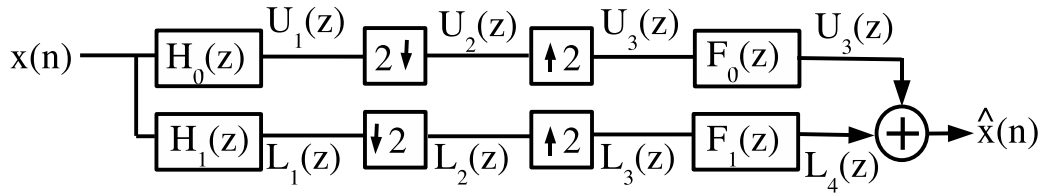


Figure 6.32

Note

$$U_1(z) = X(z) H_0(z)$$

$$U_2(z) = \frac{1}{2} \sum_{k=0}^1 \left(X \left(z^{\frac{1}{2}} e^{-\left(\frac{j2\pi k}{2}\right)} \right) H_0 \left(z^{\frac{1}{2}} e^{-\left(\frac{j\pi k}{2}\right)} \right) \right) = \frac{1}{2} X \left(z^{\frac{1}{2}} \right) H_0 \left(z^{\frac{1}{2}} \right) + \frac{1}{2} X \left(- \left(z^{\frac{1}{2}} \right) \right) H_0 \left(- \left(z^{\frac{1}{2}} \right) \right)$$

$$U_3(z) = \frac{1}{2} X(z) H_0(z) + \frac{1}{2} X(-z) H_0(-z)$$

$$U_4(z) = \frac{1}{2} F_0(z) H_0(z) X(z) + \frac{1}{2} F_0(z) H_0(-z) X(-z)$$

and

$$L_4(z) = \frac{1}{2} F_1(z) H_1(z) X(z) + \frac{1}{2} F_1(z) H_1(-z) X(-z) = \frac{1}{2} F_1(z) H_1(z) X(z) + \frac{1}{2} F_1(z) H_1(-z) X(-z)$$

Finally then,

$$\hat{X}(z) = U_4(z) + L_4(z) = \frac{1}{2} (H_0(z) F_0(z) X(z) + H_0(-z) F_0(z) X(-z) + H_1(z) F_1(z) X(z) + H_1(-z) F_1(z) X(-z)) + \frac{1}{2} (H_0(z) F_0(z) + H_1(z) F_1(z)) X(z) + \frac{1}{2} (H_0(-z) F_0(z) + H_1(-z) F_1(z)) X(-z)$$

Note that the $(X(-z) \rightarrow X(\omega + \pi))$ corresponds to the aliasing terms!

There are four things we would like to have:

1. No aliasing distortion
2. No phase distortion (overall linear phase \rightarrow simple time delay)
3. No amplitude distortion
4. FIR filters

No aliasing distortion

By insisting that $H_0(-z) F_0(z) + H_1(-z) F_1(z) = 0$, the $X(-z)$ component of $\hat{X}(z)$ can be removed, and all aliasing will be eliminated! There may be many choices for H_0, H_1, F_0, F_1 that eliminate aliasing, but most research has focused on the choice

$$F_0(z) = H_1(-z) : F_1(z) = - (H_0(-z))$$

We will consider only this choice in the following discussion.

Phase distortion

The transfer function of the filter bank, with aliasing cancelled, becomes $T(z) = \frac{1}{2} (H_0(z) F_0(z) + H_1(z) F_1(z))$, which with the above choice becomes $T(z) = \frac{1}{2} (H_0(z) H_1(-z) - H_1(z) H_0(-z))$. We would like $T(z)$ to correspond to a linear-phase filter to eliminate phase distortion: Call

$$P(z) = H_0(z) H_1(-z)$$

Note that

$$T(z) = \frac{1}{2} (P(z) - P(-z))$$

Note that $(P(-z) \Leftrightarrow (-1)^n p(n))$, and that if $p(n)$ is a linear-phase filter, $(-1)^n p(n)$ is also (perhaps of the opposite symmetry). Also note that the sum of two linear-phase filters of the same symmetry (i.e., length of $p(n)$ must be *odd*) is also linear phase, so if $p(n)$ is an odd-length linear-phase filter, there will be no phase distortion. Also note that

$$Z^{-1} (p(z) - p(-z)) = p(n) - (-1)^n p(n) = \begin{cases} 2p(n) & \text{if } n \text{ is odd} \\ 0 & \text{if } n \text{ is even} \end{cases}$$

means $p(n) = 0$, when n is even. If we choose $h_0(n)$ and $h_1(n)$ to be linear phase, $p(n)$ will also be linear phase. Thus by choosing $h_0(n)$ and $h_1(n)$ to be FIR linear phase, we eliminate phase distortion and get FIR filters as well (condition 4).

Amplitude distortion

Assuming aliasing cancellation and elimination of phase distortion, we might also desire no amplitude distortion ($|T(\omega)| = 1$). All of these conditions require

$$T(z) = \frac{1}{2} (H_0(z) H_1(-z) - H_1(z) H_0(-z)) = cz^{-D}$$

where c is some constant and D is a linear phase delay. $c = 1$ for $|T(\omega)| = 1$. It can be shown by considering that the following can be satisfied!

$$\left(T(z) = P(z) - P(-z) = 2cz^{-D} \Leftrightarrow \begin{cases} 2p(z) = 2c\delta(n-D) & \text{if } n \text{ is odd} \\ p(n) = \text{anything} & \text{if } n \text{ is even} \end{cases} \right)$$

Thus we require

$$P(z) = \sum_{n=0}^{N'} \left(p(2n) z^{-(2n)} \right) + z^{-D}$$

Any factorization of a $P(z)$ of this form, $P(z) = A(z) B(z)$ can lead to a Perfect Reconstruction filter bank of the form

$$H_0(z) = A(z)$$

$$H_1(-z) = B(z)$$

[This result is attributed to Vetterli.] A well-known special case (Smith and Barnwell)

$$H_1(z) = - \left(z^{-(2D)+1} H_0(-z^{-1}) \right)$$

Design techniques exist for optimally choosing the coefficients of these filters, under all of these constraints.

Quadrature Mirror Filters

$$(H_1(z) = H_0(-z) \Leftrightarrow H_1(\omega) = H_0(\pi + \omega) = H_0^*(\pi - \omega)) \quad (6.17)$$

for real-valued filters. The frequency response is "mirrored" around $\omega = \frac{\pi}{2}$. This choice leads to $T(z) = H_0^2(z) - H_0^2(-z)$: it can be shown that this can be a perfect reconstruction system only if

$$H_0(z) = c_0 z^{-(2n_0)} + c_1 z^{-(2n_1)}$$

which isn't a very flexible choice of filters, and not a very good lowpass! The Smith and Barnwell approach is more commonly used today.

6.8 M-Channel Filter Banks¹⁰

The theory of M-band QMFBs and PRFBs has been investigated recently. Some results are available.

¹⁰This content is available online at <<http://cnx.org/content/m12775/1.2/>>.

6.8.1 Tree-structured filter banks

Once we have a two-band PRFB, we can continue to split the subbands with similar systems! (Figure 6.33)

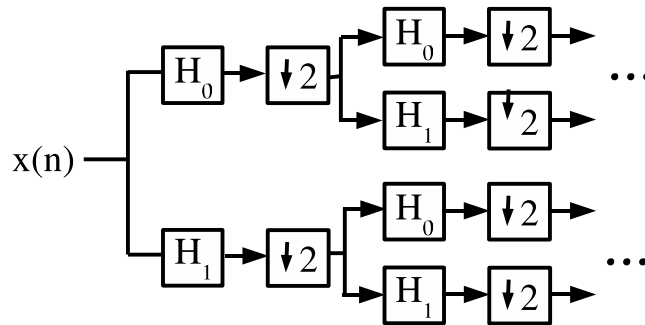


Figure 6.33

Thus we can recursively decompose a signal into 2^p bands, each sampled at 2^p th the rate of the original signal, and reconstruct exactly! Due to the tree structure, this can be quite efficient, and in fact close to the efficiency of an FFT filter bank, which does *not* have perfect reconstruction.

6.8.2 Wavelet decomposition

We need not split both the upper-frequency and lower-frequency bands identically. (Figure 6.34)

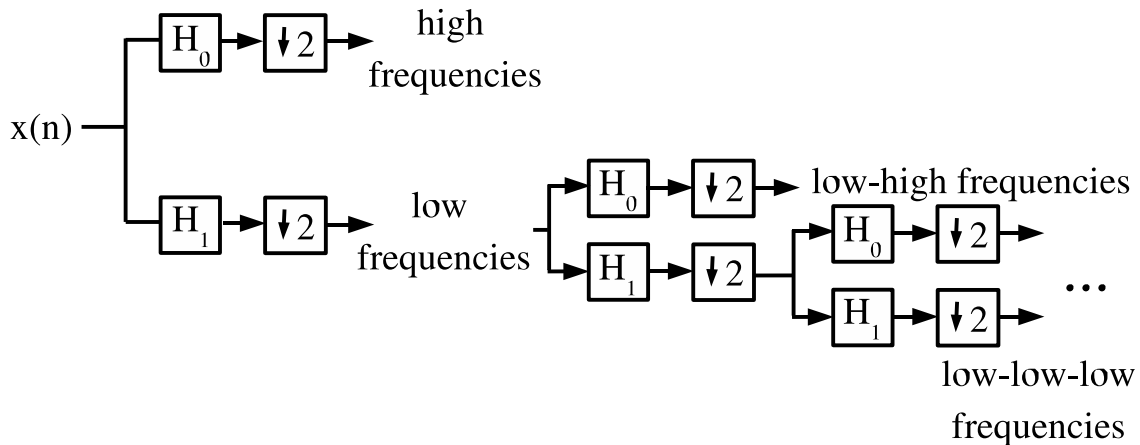


Figure 6.34

This is good for image coding, because the energy tends to be distributed such that after a wavelet decomposition, each band has roughly equal energy.

Solutions to Exercises in Chapter 6

Solution to Exercise 6.1 (p. 301)

A delta function: $h_0(n) = \delta(n')$

Solution to Exercise 6.2 (p. 306)

Simply step by N time samples between FFTs.

Solution to Exercise 6.3 (p. 306)

Aliasing should be expected. There are two ways to reduce it:

1. Decimate by less ("oversample" the individual channels) such as decimating by a factor of $\frac{N}{2}$. This is efficiently done by time-stepping by the appropriate factor.
2. Design better (and thus longer) filters, say of length LN . These can be efficiently computed by producing only N (every L th) FFT outputs using simplified FFTs.

Solution to Exercise 6.4 (p. 306)

Use an FFT and an inverse FFT for the modulation (TDM to FDM) and demodulation (FDM to TDM), respectively.

Glossary

A Autocorrelation

the expected value of the product of a random variable or signal realization with a time-shifted version of itself

C Correlation

A measure of how much one random variable depends upon the other.

Covariance

A measure of how much the deviations of two or more variables or processes match.

Crosscorrelation

if two processes are wide sense stationary, the expected value of the product of a random variable from one random process with a time-shifted, random variable from a different random process

D difference equation

An equation that shows the relationship between consecutive values of a sequence and the differences among them. They are often rearranged as a recursive formula so that a systems output can be computed from the input signal and past outputs.

Example:

$$y[n] + 7y[n - 1] + 2y[n - 2] = x[n] - 4x[n - 1] \quad ()$$

P poles

1. The value(s) for z where $Q(z) = 0$.
2. The complex frequencies that make the overall gain of the filter transfer function infinite.

R random process

A family or ensemble of signals that correspond to every possible outcome of a certain signal measurement. Each signal in this collection is referred to as a **realization** or **sample function** of the process.

Example: As an example of a random process, let us look at the Random Sinusoidal Process below. We use $f[n] = A\sin(\omega n + \phi)$ to represent the sinusoid with a given amplitude and phase. Note that the phase and amplitude of each sinusoid is based on a random number, thus making this a random process.

S State

the minimum additional information at time n , which, along with all current and future input values, is necessary to compute all future outputs.

stationary process

a random process where all of its statistical properties do not vary with time

Z zeros

1. The value(s) for z where $P(z) = 0$.
2. The complex frequencies that make the overall gain of the filter transfer function zero.

Bibliography

- [1] R.C. Agarwal and J.W. Cooley. New algorithms for digital convolution. IEEE Trans. on Acoustics, Speech, and Signal Processing, 25:392–410, Oct 1977.
- [2] G. Bruun. Z-transform dft filters and ffts. IEEE Transactions on Signal Processing, 26:56–63, February 1978.
- [3] C.S. Burrus. Index mappings for multidimensional formulation of the dft and convolution. ASSP, 25:239–242, June 1977.
- [4] C.S. Burrus. Efficient fourier transform and convolution algorithms. In J.S. Lin and A.V. Oppenheim, editors, Advanced Topics in Signal Processing, chapter Chapter 4. Prentice-Hall, 1988.
- [5] C.S. Burrus. Unscrambling for fast dft algorithms. IEEE Transactions on Acoustics, Speech, and Signal Processing, ASSP-36(7):1086–1089, July 1988.
- [6] C.S. Burrus and T.W. Parks. DFT/FFT and Convolution Algorithms. Wiley-Interscience, 1985.
- [7] Jr. C.G. Boncelet. A rearranged dft algorithm requiring $n^2/6$ multiplications. IEEE Trans. on Acoustics, Speech, and Signal Processing, ASSP-34(6):1658–1659, Dec 1986.
- [8] C.F.N. Cowan and P.M. Grant. Adaptive Filters. Prentice-Hall, 1985. Good overview of lots of topics.
- [9] R.E. Crochiere and L.R. Rabiner. Interpolation and decimation of digital signals: A tutorial review. Proc. IEEE, 69(3):300–331, March 1981.
- [10] R.E. Crochiere and L.R. Rabiner. Multirate Digital Signal Processing. Prentice-Hall, Englewood Cliffs, NJ, 1983.
- [11] P. Duhamel and H. Hollman. Split-radix fft algorithms. Electronics Letters, 20:14–16, Jan 5 1984.
- [12] D.M.W. Evans. An improved digit-reversal permutation algorithm for the fast fourier and hartley transforms. IEEE Transactions on Signal Processing, 35(8):1120–1125, August 1987.
- [13] M. Frigo and S.G. Johnson. The design and implementation of fftw3. Proceedings of the IEEE, 93(2):216–231, February 2005.
- [14] W.A. Gardner. Learning characteristics of stochastic-gradient-descent algorithms: A general study, analysis, and critique. Signal Processing, 6:113–133, 1984.
- [15] G Goertzel. An algorithm for the evaluation of finite trigonometric series. The American Mathematical Monthly, 1958.
- [16] S. Haykin. Adaptive Filters Theory. Prentice-Hall, 1986. Nice coverage of adaptive filter theory; Good reference.
- [17] M.L. Honig and D.G. Messerschmidt. Adaptive Filters: Structures, Algorithms, and Applications. Kluwer, 1984. Good coverage of lattice algorithms.

- [18] M.T. Heideman H.V. Sorensen and C.S. Burrus. On computing the split-radix fft. IEEE Transactions on Signal Processing, 34(1):152–156, 1986.
- [19] M.T. Heideman H.V. Sorensen, D.L Jones and C.S. Burrus. Real-valued fast fourier transform algorithms. IEEE Transactions on Signal Processing, 35(6):849–863, June 1987.
- [20] Leland B. Jackson. Digital Filters and Signal Processing. Kluwer Academic Publishers, 2nd edition edition, 1989.
- [21] S.G Johnson and M. Frigo. A modified split-radix fft with fewer arithmetic operations. IEEE Transactions on Signal Processing, 54, 2006.
- [22] C.R. Johnson J.R. Treichler and M.G. Larimore. Theory and Design of Adaptive Filters. Wiley-Interscience, 1987. Good introduction to adaptive filtering, CMA; nice coverage of hardware.
- [23] O. Macchi and E. Eweda. Second-order convergence analysis of stochastic adaptive linear filtering. IEEE Trans. on Automatic Controls, AC-28 #1:76–85, Jan 1983.
- [24] H.W. Schuessler R. Meyer and K. Schwarz. Fft implmentation on dsp chips - theory and practice. IEEE International Conference on Acoustics, Speech, and Signal Processing, 1990.
- [25] O. Rioul and M. Vetterli. Wavelets and signal processing. IEEE Signal Processing Magazine, 8(4):14–38, October 1991.
- [26] Richard A. Roberts and Clifford T. Mullis. Digital Signal Processing. Prentice Hall, 1987.
- [27] H.V. Sorensen and C.S. Burrus. Efficient computation of the dft with only a subset of input or output points. IEEE Transactions on Signal Processing, 41(3):1184–1200, March 1993.
- [28] H.V. Sorensen and C.S. Burrus. Efficient computation of the dft with only a subset of input or output points. IEEE Transactions on Signal Processing, 41(3):1184–1200, 1993.
- [29] P.P Vaidyanathan. Multirate digital filters, filter banks, polyphase networks, and applications: A tutorial. Proc. IEEE, 78(1):56–93, January 1990.
- [30] B. Widrow and S.D. Stearns. Adaptive Signal Processing. Prentice-Hall, 1985. Good on applications, LMS.
- [31] R. Yavne. An economical method for calculating the discrete fourier transform. Proc. AFIPS Fall Joint Computer Conf., 33:115–125, 1968.

Index of Keywords and Terms

Keywords are listed by the section with that keyword (page numbers are in parentheses). Keywords do not necessarily appear in the text of the page. They are merely associated with that section. *Ex.* apples, § 1.1 (1) **Terms** are referenced by the page they appear on. *Ex.* apples, 1

- A** A/D, § 1.10.6(49), § 1.10.7(58)
 adaptive, 262, 264
 adaptive interference cancellation, § 5.4.4(278)
 adaptive interference canceller, § 5.4.4(278)
 adaptive noise cancellation, § 5.4.4(278)
 adaptive noise canceller, § 5.4.4(278)
 alias, § 3.3.4(216)
 aliases, 224
 Aliasing, § 1.10.3(40)
 all-pass, 216
 alphabet, 5
 amplitude of the frequency response, 194
 analog, § 1.8(29), § 1.10.6(49), § 1.10.7(58)
 Applet, § 1.10.3(40)
 assembly language, § 2.3.5(166)
 auto-correlation, § 2.2.2(123), 127
 autocorrelation, § 1.12.2(81), § 1.12.5(90), 91, 91
 average, § 1.12.3(83)
 average power, 84
- B** bandlimited, 63
 basis, § 1.4(12), 16, § 1.5(22), 23
 basis matrix, § 1.5(22), 24
 bilateral z-transform, 68
 bilinear transform, 213
 bit reverse, § 2.3.5(166)
 bit-reversed, 149, 155
 boxcar, 113
 Bruun, § 2.3.4.4(162)
 butterfly, 147, 153
- C** canonic, 233
 cascade, 229
 cascade form, § 4.1.3(230)
 cascade-form structure, § 4.3.5(251)
 causal, 77
 characteristic polynomial, 67
 chirp z-transform, § 2.6(183), § 2.7(187)
 circular convolution, § 2.1.1(99)
 coefficient quantization, § 4.3.3(245), § 4.3.5(251)
 coefficient vector, § 1.5(22), 24
 commutative, 6
 compiler, § 2.3.5(166)
 complement, § 1.4(12), 17
 complex, § 1.11.4(74)
 complex exponential sequence, 4
 continuous frequency, § 1.7(27)
 continuous random process, 81
 continuous time, § 1.6(26), § 1.7(27)
 Continuous-Time Fourier Transform, 27
 convolution, § 1.3(6), 7, § 2.4(176)
 convolution property, § 2.1.1(99)
 convolution sum, 6
 Cooley-Tukey, § 2.3.4.2.1(146), § 2.3.4.2.2(151)
 correlation, 88, 88, § 1.12.5(90)
 correlation coefficient, 88
 correlation function, § 4.3.1(243)
 correlation functions, 89
 covariance, 87, 87
 Crosscorrelation, 94
 crosscorrelation function, 94
 CT, § 1.10.5(46)
 CTFT, § 1.7(27), § 1.10.7(58)
- D** D/A, § 1.10.6(49), § 1.10.7(58)
 data quantization, § 4.3.3(245)
 decimation in frequency, § 2.3.4.2.2(151), 152, § 2.3.4.2.3(155)
 decimation in time, § 2.3.4.2.1(146), 146, § 2.3.4.2.3(155), 159
 decimation-in-frequency, § 2.3.4.3(158)
 decimation-in-time, § 2.3.4.3(158)
 decompose, § 1.5(22), 23
 delayed, 5
 density function, § 1.12.2(81)
 design, § 3.2.4(199)
 deterministic, § 1.12.1(78)
 deterministic linear prediction, 217
 deterministic signals, 79
 DFT, § 1.9(34), § 1.10.6(49), § 2.1.1(99), § 2.2.1(102), § 2.3.1(141), § 2.3.2(142), § 2.3.3(144), § 2.3.4.3(158), § 2.3.5(166),

- § 2.5(181)
- DFT even symmetric, 101
- DFT odd symmetric, 101
- DFT properties, § 2.1.1(99)
- DFT symmetry, § 2.1.1(99)
- DFT-symmetric, 119
- difference equation, § 1.2(5), 6, 64, 65
- digital, § 1.8(29), § 1.10.6(49), § 1.10.7(58), § 3.2.4(199)
- digital aliasing, 294
- digital signal processing, § (1), § 1.8(29)
- direct form, § 4.1.3(230), § 4.3.3(245)
- direct method, 67
- direct sum, § 1.4(12), 17
- direct-form FIR filter structure, 227
- direct-form structure, § 4.3.5(251)
- discrete fourier transform, § 1.10.6(49), § 2.1.1(99)
- Discrete Fourier Transform (DFT), 50
- discrete random process, 80
- discrete time, § 1.3(6), § 1.6(26), § 1.11.3(73)
- discrete time fourier transform, § 1.10.6(49)
- discrete-time, § 1.8(29)
- discrete-time convolution, 6
- discrete-time Fourier transform, § 1.8(29)
- discrete-time sinc function, 33
- discrete-time systems, § 4.1.4(236)
- distribution function, § 1.12.2(81)
- dithering, § 4.3.1(243), 243
- DSP, § 1.8(29), § 1.11.3(73), § 1.12.2(81), § 1.12.3(83), § 1.12.5(90), § 3.2.2(196), § 3.2.3(197), § 3.2.4(199)
- DT, § 1.3(6)
- DTFT, § 1.10.6(49), § 1.10.7(58), § 2.2.1(102)
- DTMF, § 2.3.2(142)
- E** ergodic, 85
- Examples, § 1.10.3(40)
- Extended Prony Method, 219
- F** fast, § 2.4(176)
- fast Fourier transform, § 2.3.4.2.1(146), § 2.3.4.2.2(151)
- FFT, § 1.9(34), § 2.3.1(141), § 2.3.2(142), § 2.3.3(144), § 2.3.4.1(146), § 2.3.4.2.1(146), § 2.3.4.2.2(151), § 2.3.4.2.3(155), § 2.3.4.3(158), § 2.3.4.4(162), § 2.3.5(166), § 2.4(176), § 2.7(187)
- FFT structures, § 2.3.4.2.3(155)
- filter, § 3.2.3(197), § 3.2.4(199)
- filter structures, § 4.1.1(225)
- filters, § 3.2.2(196)
- finite, 16
- finite dimensional, § 1.4(12), 17
- FIR, § 3.2.2(196), § 3.2.3(197), § 3.2.4(199), § 4.1.2(225)
- FIR filter, § 3.1(191)
- FIR filter design, § 3.2.1(192)
- FIR filters, § 4.3.3(245), § 4.4.2(256)
- first order stationary, § 1.12.2(81)
- first-order stationary, 82
- fixed point, § 4.2.1(239)
- flights, 150
- flow-graph-reversal theorem, 228
- Fourier, § 2.2.1(102)
- fourier series, § 1.6(26)
- fourier transform, § 1.6(26), § 1.7(27), § 1.8(29), § 1.10.7(58), § 1.11.2(68), 68
- fourier transforms, § 1.10.5(46)
- fractional number representation, § 4.2.1(239)
- frequency, § 1.10.5(46)
- FT, § 1.10.5(46)
- G** Generalized Linear Phase, 194
- geometric series, 30
- Goertzel, § 2.3.3(144)
- gradient descent, 266
- H** Hamming window, 117
- Hann window, 115
- hanning, 115
- hilbert, § 1.5(22)
- Hilbert Space, § 1.4(12), 20
- hilbert spaces, § 1.5(22)
- Hold, § 1.10.4(44)
- homogeneous solution, 67
- I** identity matrix, 25
- IDFT, § 2.2.1(102)
- IIR Filter, § 3.1(191)
- IIR filter quantization, § 4.3.5(251)
- IIR Filters, § 4.1.3(230), § 4.3.4(247), § 4.4.2(256)
- Illustrations, § 1.10.3(40)
- implementation, § 4.1.2(225)
- impulse response, § 1.3(6)
- in-order FFT, § 2.3.4.2.3(155)
- in-place algorithm, 149, 155, 162
- independent, 85
- index map, § 2.6(183), 184
- indirect method, 67
- initial conditions, 65
- interpolation, § 3.2.5(206), 291
- invertible, § 1.4(12), 22

- J** Java, § 1.10.3(40)
 joint density function, § 1.12.2(81), 81
 joint distribution function, 81
- L** lag, 128
 lagrange, § 3.2.5(206)
 laplace transform, § 1.6(26)
 large-scale limit cycle, 255
 large-scale limit cycles, § 4.4.1(255)
 limit cycles, § 4.4.1(255)
 linear algebra, § 1.9(34)
 linear and time-invariant, 6
 Linear discrete-time systems, 5
 linear predictor, 218
 linear transformation, § 1.4(12), 20
 linearly dependent, § 1.4(12), 14
 linearly independent, § 1.4(12), 14, 83
 live, 40
 LTI, 6
 LTI Systems, § 1.10.7(58)
- M** matrix, § 1.9(34)
 matrix representation, § 1.4(12), 21
 mean, § 1.12.3(83), 83
 mean-square value, 84
 minimum phase, 208
 misadjustment factor, 274
 moment, 84
 multirate signal processing, § 6.1(289),
 § 6.2(291), § 6.3(295), § 6.4(299), § 6.5(302),
 § 6.6(305), § 6.7(306), § 6.8(310)
- N** narrow-band spectrogram, 130
 non-recursive structure, 227
 nonstationary, § 1.12.2(81), 81
 normal form, § 4.3.5(251)
 nyquist, § 1.10.5(46)
 Nyquist frequency, § 1.8(29)
- O** optimal, 204
 order, 65
 orthogonal, § 1.4(12), 19, 22
 orthogonal compliment, § 1.4(12), 20
 orthonormal, § 1.4(12), 19, § 1.5(22)
 orthonormal basis, § 1.5(22), 23
 overflow, § 4.2.1(239), 241, § 4.2.2(241),
 § 4.4.1(255), § 4.4.2(256)
 Overview, § 1.10.1(36)
- P** parallel form, § 4.1.3(230)
 Parseval's Theorem, § 1.8(29), § 2.1.1(99)
 particular solution, 67
 pdf, § 1.12.2(81)
 Pearson's Correlation Coefficient, 88
 perfect reconstruction filters, 305
 periodogram, § 2.2.2(123)
 picket-fence effect, 107
 pole, § 1.11.4(74)
 poles, 74
 polynomial, § 3.2.5(206)
 polyphase filters, 297
 polyphase structures, 297
 power series, 69
 power spectral density, § 2.2.2(123),
 § 4.3.1(243)
 power spectral density (PSD), 124
 power spectrum, § 2.2.2(123)
 pre-classical, § 3.3.4(216), 216
 prediction, 284
 prime factor algorithm, § 2.6(183), § 2.7(187)
 prime length FFTs, § 2.6(183)
 prime-factor algorithm, § 2.3.1(141)
 primitive root, 184
 probability, § 1.12.2(81)
 probability density function (pdf), 81
 probability distribution function, 81
 probability function, § 1.12.2(81)
 Prony's Method, 219
 Proof, § 1.10.2(38)
- Q** quantization, § 4.2.1(239), § 4.3.1(243),
 § 4.3.4(247)
 quantization error, 240, § 4.2.2(241),
 § 4.3.3(245)
 quantization noise, § 4.3.2(245)
 quantized, 235
 quantized pole locations, § 4.3.5(251)
- R** Rader's conversion, § 2.6(183), 184
 radix-2, § 2.3.4.2.1(146), 147, § 2.3.4.2.2(151),
 152
 radix-2 algorithm, § 2.3.1(141)
 radix-2 FFT, § 2.3.4.1(146)
 radix-4, § 2.3.4.3(158), 159
 radix-4 butterfly, 161
 radix-4 decimation-in-time FFT, 161
 radix-4 FFT, § 2.3.4.1(146)
 raised-cosine windows, 117
 random, § 1.12.1(78), § 1.12.3(83), § 1.12.5(90)
 random process, § 1.12.1(78), 80, 80, 81,
 § 1.12.2(81), § 1.12.3(83), § 1.12.5(90)
 random sequence, 80
 random signal, § 1.12.1(78), § 1.12.3(83)
 random signals, § 1.12.1(78), 79, § 1.12.3(83),
 § 1.12.5(90)

- real FFT, § 2.3.5(166)
 - realization, 313
 - Reconstruction, § 1.10.2(38), § 1.10.4(44)
 - rectangle, 113
 - recursive FFT, § 2.3.2(142)
 - Recursive least Squares, 265
 - ROC, § 1.11.2(68), 69
 - rounding, § 4.2.2(241), 241
 - roundoff error, § 4.2.1(239), 240, § 4.2.2(241), § 4.3.1(243)
 - roundoff noise analysis, § 4.3.4(247)
 - running FFT, § 2.3.2(142), 142
- S**
- sample function, 313
 - Sampling, § 1.10.1(36), § 1.10.2(38), § 1.10.3(40), § 1.10.4(44), § 1.10.5(46), § 1.10.6(49)
 - saturation, § 4.2.2(241), 242
 - scale, 257
 - scaling, § 4.4.2(256)
 - scaloping loss, 107
 - second order stationary, § 1.12.2(81)
 - second-order stationary, 82
 - sensitivity, 251
 - sensitivity analysis, § 4.3.5(251)
 - Shannon, § 1.10.2(38)
 - shift property, § 2.1.1(99)
 - shift-invariant, § 1.2(5), 5
 - short time fourier transform, § 2.2.3(128)
 - side lobe, 113
 - sign bit, 240
 - Signal-to-noise ratio, § 4.3.2(245)
 - signals, § 1.3(6), § 1.6(26)
 - signals and systems, § 1.3(6)
 - small-scale limit cycles, § 4.4.1(255), 256
 - SNR, § 4.3.2(245)
 - span, § 1.4(12), 15
 - spectrogram, 130
 - split-radix, § 2.3.4.4(162)
 - split-radix FFT, § 2.3.4.1(146)
 - SSS, § 1.12.2(81)
 - stable, 77
 - stage, 147, 153
 - standard basis, § 1.4(12), 21, § 1.5(22)
 - state, § 4.1.4(236), 236
 - state space, § 4.1.4(236)
 - state variable representation, § 4.1.4(236)
 - state variable transformation, § 4.1.4(236)
 - state variables, § 4.1.4(236)
 - stationarity, § 1.12.2(81)
 - stationary, § 1.12.2(81), § 1.12.5(90)
 - stationary process, 81
 - stationary processes, 81
 - statistical spectrum estimation, § 2.2.2(123)
 - stft, § 2.2.3(128)
 - stochastic, § 1.12.1(78)
 - stochastic gradient, 269
 - stochastic signals, 79
 - strict sense stationary, § 1.12.2(81)
 - strict sense stationary (SSS), 82
 - structures, § 4.1.2(225), § 4.1.3(230)
 - subspace, § 1.4(12), 14
 - superposition, § 1.2(5)
 - System, § 1.10.4(44)
 - systems, § 1.6(26)
- T**
- table lookup, § 2.3.5(166), 166
 - The stagecoach effect, 43
 - time, § 1.10.5(46)
 - time reversal, § 2.1.1(99)
 - time-varying behavior, 61
 - Toeplitz, 221
 - training signal, 278
 - transfer function, 65, § 4.1.1(225)
 - transform pairs, § 1.11.3(73)
 - transforms, 26
 - transmultiplexors, 306
 - transpose form, § 4.1.3(230), § 4.3.3(245)
 - transpose-form structure, § 4.3.5(251)
 - truncation, 113, 113, § 4.2.1(239), § 4.2.2(241), 242
 - truncation error, 240, § 4.2.2(241), § 4.3.1(243)
 - twiddle factor, 142, § 2.3.4.2.1(146), 146, § 2.3.4.2.2(151), 152
 - twiddle factors, 158
 - twiddle-factor, 146, 147, 153
 - two's complement, § 4.2.1(239)
- U**
- uncorrelated, 83
 - unilateral, § 1.11.3(73)
 - unilateral z-transform, 68
 - unique, 23
 - unit sample, 4, 5
 - unitary, § 1.4(12), 22
- V**
- variance, § 1.12.3(83), 84
 - vector, § 1.9(34)
 - Vector-radix FFTs, 188
 - vectors, § 1.4(12), 13
- W**
- well-defined, § 1.4(12), 17
 - WFTA, § 2.6(183)
 - wide sense stationary, § 1.12.2(81)
 - wide-band spectrogram, 130
 - wide-sense stationary (WSS), 83

- window, 113
 - Winograd, § 2.6(183)
 - Winograd Fourier Transform Algorithm,
§ 2.6(183), § 2.7(187)
 - wraparound, 242
 - wraparound error, § 4.2.2(241)
 - WSS, § 1.12.2(81)
- Y**
- Yavne, § 2.3.4.4(162)
 - Yule-Walker, 222
- Z**
- z transform, § 1.6(26), § 1.11.3(73)
 - z-plane, 68, § 1.11.4(74)
 - z-tranform, § 2.5(181)
 - z-transform, § 1.11.2(68), 68, § 1.11.3(73)
 - z-transforms, 73
 - zero, § 1.11.4(74)
 - zero-input limit cycles, § 4.4.1(255), 256
 - zero-padding, 104
 - zeros, 74

Attributions

Collection: *Digital Signal Processing: A User's Guide*

Edited by: Douglas L. Jones

URL: <http://cnx.org/content/col10372/1.2/>

License: <http://creativecommons.org/licenses/by/2.0/>

Module: "Preface for Digital Signal Processing: A User's Guide"

By: Douglas L. Jones

URL: <http://cnx.org/content/m13782/1.1/>

Pages: 1-2

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/2.0/>

Module: "Discrete-Time Signals and Systems"

By: Don Johnson

URL: <http://cnx.org/content/m10342/2.13/>

Pages: 3-5

Copyright: Don Johnson

License: http://creativecommons.org/licenses/by/1.0

Module: "Systems in the Time-Domain"

By: Don Johnson

URL: <http://cnx.org/content/m0508/2.7/>

Pages: 5-6

Copyright: Don Johnson

License: http://creativecommons.org/licenses/by/1.0

Module: "Discrete-Time Convolution"

By: Ricardo Radaelli-Sanchez, Richard Baraniuk

URL: <http://cnx.org/content/m10087/2.18/>

Pages: 6-12

Copyright: Ricardo Radaelli-Sanchez, Richard Baraniuk

License: http://creativecommons.org/licenses/by/1.0

Module: "Review of Linear Algebra"

By: Clayton Scott

URL: <http://cnx.org/content/m11948/1.2/>

Pages: 12-22

Copyright: Clayton Scott

License: http://creativecommons.org/licenses/by/1.0

Module: "Orthonormal Basis Expansions"

By: Michael Haag, Justin Romberg

URL: <http://cnx.org/content/m10760/2.4/>

Pages: 22-26

Copyright: Michael Haag, Justin Romberg

License: http://creativecommons.org/licenses/by/1.0

Module: "Fourier Analysis"

By: Richard Baraniuk

URL: <http://cnx.org/content/m10096/2.10/>

Pages: 26-27

Copyright: Richard Baraniuk

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Continuous-Time Fourier Transform (CTFT)"

By: Richard Baraniuk, Melissa Selik

URL: <http://cnx.org/content/m10098/2.9/>

Pages: 27-29

Copyright: Richard Baraniuk, Melissa Selik

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Discrete-Time Fourier Transform (DTFT)"

By: Don Johnson

URL: <http://cnx.org/content/m10247/2.28/>

Pages: 29-34

Copyright: Don Johnson

License: <http://creativecommons.org/licenses/by/1.0>

Module: "DFT as a Matrix Operation"

By: Robert Nowak

URL: <http://cnx.org/content/m10962/2.5/>

Pages: 34-36

Copyright: Robert Nowak

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Introduction"

By: Anders Gjendemsj 

URL: <http://cnx.org/content/m11419/1.29/>

Pages: 36-38

Copyright: Anders Gjendemsj 

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Proof"

By: Anders Gjendemsj 

URL: <http://cnx.org/content/m11423/1.27/>

Pages: 38-40

Copyright: Anders Gjendemsj 

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Illustrations"

By: Anders Gjendemsj 

URL: <http://cnx.org/content/m11443/1.33/>

Pages: 40-44

Copyright: Anders Gjendemsj 

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Systems view of sampling and reconstruction"

By: Anders Gjendemsjø

URL: <http://cnx.org/content/m11465/1.20/>

Pages: 44-46

Copyright: Anders Gjendemsjø

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Sampling CT Signals: A Frequency Domain Perspective"

By: Robert Nowak

URL: <http://cnx.org/content/m10994/2.2/>

Pages: 46-49

Copyright: Robert Nowak

License: <http://creativecommons.org/licenses/by/1.0>

Module: "The DFT: Frequency Domain with a Computer Analysis"

By: Robert Nowak

URL: <http://cnx.org/content/m10992/2.3/>

Pages: 49-58

Copyright: Robert Nowak

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Discrete-Time Processing of CT Signals"

By: Robert Nowak

URL: <http://cnx.org/content/m10993/2.2/>

Pages: 58-64

Copyright: Robert Nowak

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Difference Equation"

By: Michael Haag

URL: <http://cnx.org/content/m10595/2.5/>

Pages: 64-68

Copyright: Michael Haag

License: <http://creativecommons.org/licenses/by/1.0>

Module: "The Z Transform: Definition"

By: Benjamin Fite

URL: <http://cnx.org/content/m10549/2.9/>

Pages: 68-73

Copyright: Benjamin Fite

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Table of Common z-Transforms"

By: Melissa Selik, Richard Baraniuk

URL: <http://cnx.org/content/m10119/2.13/>

Pages: 73-74

Copyright: Melissa Selik, Richard Baraniuk

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Understanding Pole/Zero Plots on the Z-Plane"

By: Michael Haag

URL: <http://cnx.org/content/m10556/2.8/>

Pages: 74-78

Copyright: Michael Haag

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Introduction to Random Signals and Processes"

By: Michael Haag

URL: <http://cnx.org/content/m10649/2.2/>

Pages: 78-81

Copyright: Michael Haag

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Stationary and Nonstationary Random Processes"

By: Michael Haag

URL: <http://cnx.org/content/m10684/2.2/>

Pages: 81-83

Copyright: Michael Haag

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Random Processes: Mean and Variance"

By: Michael Haag

URL: <http://cnx.org/content/m10656/2.3/>

Pages: 83-87

Copyright: Michael Haag

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Correlation and Covariance of a Random Signal"

By: Michael Haag

URL: <http://cnx.org/content/m10673/2.3/>

Pages: 87-90

Copyright: Michael Haag

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Autocorrelation of Random Processes"

By: Michael Haag

URL: <http://cnx.org/content/m10676/2.4/>

Pages: 90-93

Copyright: Michael Haag

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Crosscorrelation of Random Processes"

By: Michael Haag

URL: <http://cnx.org/content/m10686/2.2/>

Pages: 93-95

Copyright: Michael Haag

License: <http://creativecommons.org/licenses/by/1.0>

Module: "DFT Definition and Properties"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12019/1.5/>

Pages: 99-102

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Spectrum Analysis Using the Discrete Fourier Transform"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12032/1.6/>

Pages: 102-123

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Classical Statistical Spectral Estimation"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12014/1.3/>

Pages: 123-128

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Short Time Fourier Transform"

By: Ivan Selesnick

URL: <http://cnx.org/content/m10570/2.4/>

Pages: 128-141

Copyright: Ivan Selesnick

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Overview of Fast Fourier Transform (FFT) Algorithms"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12026/1.3/>

Pages: 141-142

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Running FFT"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12029/1.5/>

Pages: 142-144

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Goertzel's Algorithm"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12024/1.5/>

Pages: 144-145

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Power-of-two FFTs"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12059/1.2/>

Pages: 146-146

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Decimation-in-time (DIT) Radix-2 FFT"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12016/1.7/>

Pages: 146-151

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Decimation-in-Frequency (DIF) Radix-2 FFT"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12018/1.6/>

Pages: 151-155

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Alternate FFT Structures"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12012/1.6/>

Pages: 155-158

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Radix-4 FFT Algorithms"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12027/1.4/>

Pages: 158-162

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Split-radix FFT Algorithms"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12031/1.5/>

Pages: 162-166

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Efficient FFT Algorithm and Programming Tricks"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12021/1.6/>

Pages: 166-169

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Multidimensional Index Maps"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12025/1.3/>

Pages: 169-173

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "The Prime Factor Algorithm"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12033/1.3/>

Pages: 173-176

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Fast Convolution"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12022/1.5/>

Pages: 176-181

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Chirp-z Transform"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12013/1.4/>

Pages: 181-183

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "FFTs of prime length and Rader's conversion"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12023/1.3/>

Pages: 183-186

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Choosing the Best FFT Algorithm"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12060/1.3/>

Pages: 187-188

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Overview of Digital Filter Design"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12776/1.2/>

Pages: 191-191

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/2.0/>

Module: "Linear Phase Filters"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12802/1.2/>

Pages: 192-196

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/2.0/>

Module: "Window Design Method"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12790/1.2/>

Pages: 196-197

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/2.0/>

Module: "Frequency Sampling Design Method for FIR filters"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12789/1.2/>

Pages: 197-199

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/2.0/>

Module: "Parks-McClellan FIR Filter Design"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12799/1.3/>

Pages: 199-206

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/2.0/>

Module: "Lagrange Interpolation"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12812/1.2/>

Pages: 206-206

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/2.0/>

Module: "Overview of IIR Filter Design"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12758/1.2/>

Pages: 207-207

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/2.0/>

Module: "Prototype Analog Filter Design"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12763/1.2/>

Pages: 207-213

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/2.0/>

Module: "IIR Digital Filter Design via the Bilinear Transform"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12757/1.2/>

Pages: 213-216

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/2.0/>

Module: "Impulse-Invariant Design"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12760/1.2/>

Pages: 216-216

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/2.0/>

Module: "Digital-to-Digital Frequency Transformations"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12759/1.2/>

Pages: 216-217

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/2.0/>

Module: "Prony's Method"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12762/1.2/>

Pages: 217-220

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/2.0/>

Module: "Linear Prediction"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12761/1.2/>

Pages: 220-222

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/2.0/>

Module: "Filter Structures"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11917/1.3/>

Pages: 225-225

Copyright: Douglas L. Jones

License: http://creativecommons.org/licenses/by/1.0

Module: "FIR Filter Structures"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11918/1.2/>

Pages: 225-230

Copyright: Douglas L. Jones

License: http://creativecommons.org/licenses/by/1.0

Module: "IIR Filter Structures"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11919/1.2/>

Pages: 230-236

Copyright: Douglas L. Jones

License: http://creativecommons.org/licenses/by/1.0

Module: "State-Variable Representation of Discrete-Time Systems"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11920/1.2/>

Pages: 236-239

Copyright: Douglas L. Jones

License: http://creativecommons.org/licenses/by/1.0

Module: "Fixed-Point Number Representation"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11930/1.2/>

Pages: 239-241

Copyright: Douglas L. Jones

License: http://creativecommons.org/licenses/by/1.0

Module: "Fixed-Point Quantization"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11921/1.2/>

Pages: 241-242

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Finite-Precision Error Analysis"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11922/1.2/>

Pages: 243-245

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Input Quantization Noise Analysis"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11923/1.2/>

Pages: 245-245

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Quantization Error in FIR Filters"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11924/1.2/>

Pages: 245-247

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Data Quantization in IIR Filters"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11925/1.2/>

Pages: 247-251

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "IIR Coefficient Quantization Analysis"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11926/1.2/>

Pages: 251-255

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Limit Cycles"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11928/1.2/>

Pages: 255-256

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Scaling"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11927/1.2/>

Pages: 256-258

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Introduction to Adaptive Filters"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11535/1.3/>

Pages: 261-261

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Discrete-Time, Causal Wiener Filter"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11825/1.1/>

Pages: 261-264

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Practical Issues in Wiener Filter Implementation"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11824/1.1/>

Pages: 264-265

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Quadratic Minimization and Gradient Descent"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11826/1.2/>

Pages: 265-267

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "The LMS Adaptive Filter Algorithm"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11829/1.1/>

Pages: 267-269

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "First Order Convergence Analysis of the LMS Algorithm"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11830/1.1/>

Pages: 269-272

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Second-order Convergence Analysis of the LMS Algorithm and Misadjustment Error"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11831/1.2/>

Pages: 272-275

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Applications of Adaptive Filters"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11536/1.1/>

Pages: 275-275

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Adaptive System Identification"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11906/1.1/>

Pages: 275-276

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Adaptive Equalization"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11907/1.1/>

Pages: 276-278

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Adaptive Interference (Noise) Cancellation"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11835/1.1/>

Pages: 278-281

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Adaptive Echo Cancellation"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11909/1.1/>

Pages: 281-283

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Beyond LMS: an overview of other adaptive filter algorithms"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11911/1.1/>

Pages: 283-283

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Adaptive IIR filters"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11912/1.1/>

Pages: 283-285

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "The Constant-Modulus Algorithm and the Property-Restoral Principle"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11913/1.1/>

Pages: 285-286

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Complex LMS"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11914/1.3/>

Pages: 286-286

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Normalized LMS"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11915/1.2/>

Pages: 286-287

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Summary of Adaptive Filtering Methods"

By: Douglas L. Jones

URL: <http://cnx.org/content/m11916/1.1/>

Pages: 287-287

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/1.0>

Module: "Overview of Multirate Signal Processing"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12777/1.2/>

Pages: 289-291

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/2.0/>

Module: "Interpolation, Decimation, and Rate Changing by Integer Fractions"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12801/1.2/>

Pages: 291-295

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/2.0/>

Module: "Efficient Multirate Filter Structures"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12800/1.2/>

Pages: 295-299

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/2.0/>

Module: "Filter Design for Multirate Systems"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12773/1.2/>

Pages: 299-301

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/2.0/>

Module: "Multistage Multirate Systems"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12803/1.2/>

Pages: 302-304

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/2.0/>

Module: "DFT-Based Filterbanks"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12771/1.2/>

Pages: 305-306

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/2.0/>

Module: "Quadrature Mirror Filterbanks (QMF)"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12770/1.2/>

Pages: 306-310

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/2.0/>

Module: "M-Channel Filter Banks"

By: Douglas L. Jones

URL: <http://cnx.org/content/m12775/1.2/>

Pages: 310-311

Copyright: Douglas L. Jones

License: <http://creativecommons.org/licenses/by/2.0/>

Digital Signal Processing: A User's Guide

Digital Signal Processing: A User's Guide is intended both for the practicing engineer with a basic knowledge of DSP and for a second course in signal processing at the senior or first-year postgraduate level. FFTs, digital filter design, adaptive filters, and multirate signal processing are covered with an emphasis on the techniques that have found wide use in practice.

About Connexions

Since 1999, Connexions has been pioneering a global system where anyone can create course materials and make them fully accessible and easily reusable free of charge. We are a Web-based authoring, teaching and learning environment open to anyone interested in education, including students, teachers, professors and lifelong learners. We connect ideas and facilitate educational communities.

Connexions's modular, interactive courses are in use worldwide by universities, community colleges, K-12 schools, distance learners, and lifelong learners. Connexions materials are in many languages, including English, Spanish, Chinese, Japanese, Italian, Vietnamese, French, Portuguese, and Thai. Connexions is part of an exciting new information distribution system that allows for **Print on Demand Books**. Connexions has partnered with innovative on-demand publisher QOOP to accelerate the delivery of printed course materials and textbooks into classrooms worldwide at lower prices than traditional academic publishers.