Deep learning primer

### The archetype



### What is deep learning

What can it do

Flexibility capable of fitting complex functions

How to train it

Computable gradient function *largely* smooth

• Universal representation theorem:

Any function in finite dimensions can be approximated arbitrarily well with a two-layer neural network with finite number of hidden unit



• Universal representation theorem (improved):

There exists a two-layer neural network with ReLU activations and 2n+d weights that can represent any function on a sample of size n in d dimensions.

Zhang, et al. 2016, Understanding deep learning requires rethinking generalization



• Depth efficiency hypothesis (widely held belief + proof for certain models):

Some functions expressed in multi-layer models requires superpolynomial sized units to express in shallow models



• Flexible model does not generalize?

Rademacher complexity-based generalization bound

$$\mathrm{Rad}_S(F) = rac{1}{m} \mathbb{E}_\sigma \left[ \sup_{f \in F} \left| \sum_{i=1}^m \sigma_i f(z_i) 
ight| 
ight] \qquad \qquad L_P(h) - L_S(h) \leq 2 \, \mathrm{Rad}(F \circ S) + 4 \sqrt{rac{2 \ln(4/\delta)}{m}}$$

with probability at least 1- sigma

## Fun fact: neural network usually has the capacity to memorize random labels perfectly

Zhang et al. Understanding deep learning requires rethinking generalization.

• Flexible model does not generalize?

In practice, models are never trained to obtain the minimal training loss



Notion of generalization based on the 'length' of training path?

## Gradient

 Implicit assumption is that deep learning models can be learned by simply gradient descent

It will be interesting to know theoretically when this assumption fails (e.g. I guess, prime factorization)

Computation of Gradient: Automatic differentiation

Allow trivial solution to complex models / changing model structure dynamically (data-dependent)

• The basics:

 $\frac{dy}{dx} = \frac{dy}{dw}\frac{dw}{dx}$ 

• Computational graph:



### Computation of Gradient: Automatic differentiation

Allow trivial solution to complex models / changing model structure dynamically (data-dependent)

- The basics:  $\frac{dy}{dx} = \frac{dy}{dw}\frac{dw}{dx}$
- Two modes: forward mode and backward mode

(optimal traversal path for arbitrary computational graph is NP-complete)

- The future?
  - Compiler for mathematical expressions that achieves acceleration and numeric stability
  - Mixing programing language (conditionals, loops, etc with mathematical functions)
  - Higher-order derivative (e.g. Hessian)

Use gradient efficiently: Stochastic gradient descent

 $1/\sqrt{t}$  error rate (stochastic) vs 1/t error rate (batch)



'High optimization error' is tolerable:

No need to optimize beyond the statistical limits

Is SGD adaptive to the data uncertainty?

### Connection between Stochastic Gradient Descent and Bayesian inference

SGD as MCMC

Stochastic gradient Langevin dynamics, Welling and Teh, 2011

SGD

$$\Delta \theta_t = \frac{\epsilon_t}{2} \left( \nabla \log p(\theta_t) + \frac{N}{n} \sum_{i=1}^n \nabla \log p(x_{ti} | \theta_t) \right)$$

MCMC by Stochastic gradient Langevin dynamics

$$\Delta \theta_t = \frac{\epsilon_t}{2} \left( \nabla \log p(\theta_t) + \frac{N}{n} \sum_{i=1}^n \nabla \log p(x_{ti}|\theta_t) \right) + \eta_t$$
$$\eta_t \sim N(0, \epsilon_t) \tag{4}$$

MCMC by Langevin dynamics  

$$\Delta \theta_t = \frac{\epsilon}{2} \left( \nabla \log p(\theta_t) + \sum_{i=1}^N \nabla \log p(x_i | \theta_t) \right) + \eta_t$$

$$\eta_t \sim N(0, \epsilon)$$
(3)

SGD as VI

Stochastic Gradient Descent as Approximate Bayesian Inference, Mandt, 2017

### Optimization: scale invariance



The exact way: compute Hessian matrix (second order derivatives) / Newton's method

The cheap way : approximation using the history of gradients



### Optimization: variance reduction and scale invariance



RMSprop 
$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t$$

Adam  

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t.$$
  
 $v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2.$   
 $\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t.$ 

http://sebastianruder.com/optimizing-gradient-descent/

### Learning representations

Raw data that lives in some arbitrary (high-dimensional) space



### Representation: smoothness

#### Digits (MNIST)

Ø З з .3 

Embedding learned by variational autoencoder (VAE)

Bedroom (LSUN)



Embedding learned by generative adversarial networks (GAN)

### Representation: smoothness



RNN autoencoder https://arxiv.org/abs/1704.03477

- 字種成東字推
- 符利對亞型斷
- 到用抗語進的
- 字條網言行新
- 符件絡字自方
- 一生對體動法

GAN github.com/kaonashi-tyc/zi2zi

### Representation: linearity



Pretrained word vectors for >70 languages are publicly available

### Representation: linearity





man with glasses

man without glasses

woman without glasses



woman with glasses

### **Representation learning**

This is one of Crichton's best books. The characters of Karen Ross, Peter Elliot, Munro, and Amy are beautifully developed and their interactions are exciting, complex, and fast-paced throughout this impressive novel. And about 99.8 percent of that got lost in the film. Seriously, the screenplay AND the directing were horrendous and clearly done by people who could not fathom what was good about the novel. I can't fault the actors because frankly, they never had a chance to make this turkey live up to Crichton's original work. I know good novels, especially those with a science fiction edge, are hard to bring to the screen in a way that lives up to the original. But this may be the absolute worst disparity in quality between novel and screen adaptation ever. The book is really, really good. The movie is just dreadful.

"Sentiment neuron"

Trained on 82 million Amazon reviews to predict the next character

method: multiplicative LSTM

Reference: https://blog.openai.com/unsupervised-sentiment-neuron/

### Practical neural network structure

#### **Convolutional network**

- Image recognition
- Image super-resolution / in-painting
- Play Go

#### Long Short Term Memory

- Machine translation
- Basic question answering
- Text understanding e.g. sentiment analysis

### Convolutional network / Convnet



#### Human level image recognition

https://www.clarifai.com/

### Convolutional network / Convnet

### Traditional architecture

### The essential part



Convolution (Spatial weight sharing)

<b>1</b> _×1	<b>1</b> _×0	<b>1</b> _×1	0	0
<b>0</b> _×0	<b>1</b> _×1	<b>1</b> _×0	1	0
<b>0</b> <sub>×1</sub>	<b>0</b> <sub>×0</sub>	<b>1</b> _×1	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved Feature

The recipe that won ImageNet 2012 competition (1.3M images in 1000 categories):

```
ConvNet + ReLU + dropout + GPU
```

16.4% error vs 26.1% 2nd place

### Neural style transfer

Pretrained convnet (e.g. Imagenet classifier)

#### Objective 1:

Minimize distance of high level representations to the content image

#### Objective 2:

Minimize distance of correlation matrix of low level representations to the style image



### Photo style transfer: Improving photorealism

Deep Photo Style Transfer: arXiv:1703.07511



additional regularization term to encourage transformation to be locally affine

# State-of-the-art Computer Vision ConvNet at 2017



Huang et al. Densely Connected Convolutional Networks

#### **Batch normalization layer:**

standardization of representations (Loffe and Szegedy, 2015)

'Shortcut' connections:

address vanishing gradients

#### **Reduce number of weights:**

Use exclusively 1x1 convolution and 3x3 convolution. More layers. No fully connected layer

#### Convolution + Pooling is a general technique for enforcing invariance in representations

Can be extended to introduce translation, rotation, or scaling invariance etc.

#### Mathematical perspective: invariant transformations as symmetry groups

Cohen and Welling, 2016Group Equivariant Convolutional NetworksMallat, 2012Group Invariant Scattering

#### Computational challenge: how to compute efficiently?

Possible transformations grow multiplicatively if we stack invariances Stochastic approximation (one random transformation at a time)?

### Convnet can be easily tricked



#### So do most high dimensional models

### LSTM - recurrent neural networks



### Long Short-term Memory (LSTM) network - Gate

#### Standard structure



Forget gate  $f_t = \sigma_g (W_f x_t + U_f h_{t-1} + b_f)$ Input gate  $i_t = \sigma_g(W_i x_t + U_i h_{t-1} + b_i)$ Output gate  $o_t = \sigma_g (W_o x_t + U_o h_{t-1} + b_o)$ Memory cell  $c_t = f_t \circ c_{t-1} + i_t \circ \sigma_c (W_c x_t + U_c h_{t-1} + b_c)$ Hidden units  $h_t = o_t \circ \sigma_h(c_t)$ from: wikipedia

#### The essential part

### Forget gate

Greff et al., 2015, LSTM: A Search Space Odyssey

#### Sequence to sequence (Seq2Seq)



http://suriyadeepan.github.io/2016-12-31-practical-seq2seq/

### Seq2Seq + Attention



Attention - dynamic weighting of the input

### Neural Turing machine / Differentiable Neural Computer

Memory is an array of vectors.



http://distill.pub/2016/augmented-rnns/

Graves et al. 2014, Neural Turing machines

### Reinforcement learning





Given state, choose action, get reward

Image credit: daily.doodl @ instagram

### Deep Q learning: Predict future rewards with deep networks

Q Learning

Q(state, action) = maximal future rewards (with the optimal actions)

Bellman equation

 $Q(s,a) = r + \gamma max_{a\prime}Q(s',a')$ 

Training: minimize MSE



Minh et al, 2013 Playing Atari with Deep Reinforcement Learning

### AlphaGo - surpass human-level game playing in Go (the nature publication version)





SL policy network: predict expert human moves convnet / GLM

RL policy network: optimized by self-play convnet

REINFORCE algorithm (Williams, 1992)

Value network: predict outcome of self-play convnet

Silver et al., 2016, Mastering the game of Go with deep neural networks and tree search

### AlphaGo - Monte carlo tree search



**b** Tree evaluation from value net

С

Tree evaluation from rollouts





### Deep learning for probabilistic models

#### Why

• Toward tractable inference for more expressive probabilistic models

intractable distributions (unnormalized distributions)

Posterior distribution  $p(\theta \mid \mathbf{X}, \alpha) = \frac{p(\mathbf{X} \mid \theta)p(\theta \mid \alpha)}{p(\mathbf{X} \mid \alpha)} \propto p(\mathbf{X} \mid \theta)p(\theta \mid \alpha)$ Energy-based models  $P(x) = \frac{1}{Z} \exp f(x)$ 

Similar to deep learning, inference method are often gradient based

- Variational inference
- MCMC (e.g. Hamiltonian Monte Carlo uses gradient to speed up sampling)

http://arogozhnikov.github.io/2016/12/19/markov\_chain\_monte\_carlo.html

### Deep learning for probabilistic models

Why

• Toward tractable inference for more expressive probabilistic models

intractable distributions (unnormalized distributions)

Posterior distribution  $p(\theta \mid \mathbf{X}, \alpha) = \frac{p(\mathbf{X} \mid \theta)p(\theta \mid \alpha)}{p(\mathbf{X} \mid \alpha)} \propto p(\mathbf{X} \mid \theta)p(\theta \mid \alpha)$ Energy-based models  $P(x) = \frac{1}{Z} \exp f(x)$ 

Potential approaches for NN-assisted inference

- Neural variational inference
- Neural sampler

### Probabilistic modeling with neural networks: **NVI** Neural variational inference (Kingma and Welling, 2014)

Use neural network for describing P(X|Z) or Q(Z|X)



### Probabilistic modeling with neural networks: NVI Neural variational inference

(Kingma and Welling, 2014 Auto-Encoding Variational Bayes)

The variational objective

Use neural network for describing P(X|Z) or Q(Z|X)



### **Backpropagation over stochastic units:** Reparametrization trick

How to compute good gradient estimate of

$$-D_{KL}(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\boldsymbol{\theta}}(\mathbf{z})) + \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}^{(i)})} \left[\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}|\mathbf{z})\right]$$

Gradient of expectation -> expectation of stochastic gradient



$$\nabla_{\mu,\sigma} E_{z \sim \mu,\sigma}[f(z)] = E_{z \sim \mu,\sigma}[f(z)\nabla_{\mu,\sigma}\log\left(p(z|\mu,\sigma)\right)]$$

 $\nabla_{\mu,\sigma} E_{\epsilon \sim p(\epsilon)}[f(z)] = E_{\epsilon \sim p(\epsilon)}[\nabla_{\mu,\sigma} f\left(g(\mu,\sigma,\epsilon)\right)]$ 

### **Backpropagation over stochastic units:** Reparametrization trick for discrete variables

The Gumbel trick for sampling from discrete distributions  $P(X=k) \propto lpha_k$ 

 $G = -\log(-\log(U))$  with  $U \sim \mathrm{Unif}[0,1]$ 

$$X = rg\max_k \left(\log lpha_k + G_k
ight).$$

Softmax function for approximating the max operation with a differentiable function



Discrete variables can always be represented by binary vectors

Toward flexible and normalized probabilistic models

$$-D_{KL}(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\boldsymbol{\theta}}(\mathbf{z})) + \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}^{(i)})} \left[\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}|\mathbf{z})\right]$$

#### 1. Fully factorized models



Neural autoregressive models

Probability function is fully factorial However, it has to commit to a sequence Toward flexible and normalized probabilistic models

2. Invertible transformations



$$p_X(x) = p_H(f(x)) |\det rac{\partial f(x)}{\partial x}|.$$

#### Examples:



Normalizing flow

## Invertible autoregressive flow



determinant fixed

$$f(\mathbf{z}) = \mathbf{z} + \mathbf{u}h(\mathbf{w}^{\top}\mathbf{z} + b)$$



determinant O(D) time

Kingma, 2017 Invertible autoregressive flow

Rezende 2016. Variational Inference with Normalizing Flows

determinant O(D) time

Hidden variables are equal in dimensionality.

Dinh 2015, NICE: NON-LINEAR INDEPENDENT COMPONENTS ESTIMATION Probabilistic modeling with neural networks: Learn to sample

### Generative adversarial networks



<u>Noise</u>

Probabilistic modeling with neural networks: Learn to sample

Generative adversarial networks



#### Formulating Generative adversarial networks as a proper probabilistic model



Generator network: use x~Generator instead of x~model Discriminator network: f(x)

Question: can we do exact sampling instead of approximate sampling?

#### Formulating Generative adversarial networks as a proper probabilistic model



Two proposals :

- 1. Design Generator that generate reversible sequence of samples (allows Metropolis-Hasting sampling)
- 2. Allows computing (unnormalized) probability for generator networks (allows Metropolis sampling)

Question: can we do exact sampling instead of approximate sampling?

### Coding: choice of packages

	Pros	Cons
pytorch	Fastest automatic differentiation engine	in beta
torch7	Fast and low level Good code base	written in lua
Keras	High-level abstraction	hard for implementing new techniques
Tensorflow	Distributed computation	Poor support for interactive mode



https://github.com/pytorch/tutorials

#### Training algorithm



Conditional distribution

$$p(x|z) = \mathcal{N}(f(z), \sigma(z))$$

Marginal distribution

$$p(x) = \int p(x|z)p(z)dz = ?$$

Proposal\* 2:

$$p(x) = \frac{p(x,z)}{p(z|x)} = \frac{p(z)p(x|z)}{p(z|x)}$$

Problem: P(z|x) is generally as intractable as p(x)

#### Training algorithm



Conditional distribution

$$p(x|z) = \mathcal{N}(f(z), \sigma(z))$$

Marginal distribution

$$p(x) = \int p(x|z)p(z)dz = ?$$

Proposal\* 2:

$$p(x) = \frac{p(x,z)}{p(z|x)} = \frac{p(z)p(x|z)}{p(z|x)}$$

Problem: P(z|x) is generally as intractable as p(x)

#### Training algorithm

Proposal 3: (Variational Autoencoder)



Conditional distribution

??

Ζ

Х

Gaussian random variable

Nonlinear dependency

Gaussian random variable

(with diagonal covariance matrix)

θ

with fixed mean e.g. 0 and e.g. 1

$$p(x|z) = \mathcal{N}(f(z), \sigma(z))$$

Ν

Marginal distribution

$$p(x) = \int p(x|z)p(z)dz = ?$$

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \ge \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}^{(i)}) = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})} \left[ -\log q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}) + \log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}) \right]$$
$$= \left[ -D_{KL}(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\boldsymbol{\theta}}(\mathbf{z})) + \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}^{(i)})} \left[ \log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}|\mathbf{z}) \right] \right]$$

The variational objective



 $-D_{KL}(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\boldsymbol{\theta}}(\mathbf{z})) + \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}^{(i)})} \left[\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}|\mathbf{z})\right]$ 

KL divergence between two Gaussians

Use stochastic gradient estimator

Stochastic gradient for  $p_{-}\theta$  is straightforward

Stochastic gradient for  $q_{\Phi}$ :



$$\nabla_{\mu,\sigma} E_{\epsilon \sim p(\epsilon)}[f(z)] = E_{\epsilon \sim p(\epsilon)}[\nabla_{\mu,\sigma} f\left(g(\mu,\sigma,\epsilon)\right)]$$

Kingma , et al. 2014. Auto-Encoding Variational Bayes







Limitation: Q(Z|X) is still limited to independent Gaussians

#### Reparametrized



P(X|Z)

P(X|Z) being independent Gaussians does not limit its flexibility

Prototype for stochastic neural network



P(f(x)) is in general intractable for stochastic neural networks

$$p_X(x) = p_H(f(x)) |\det rac{\partial f(x)}{\partial x}|.$$

is applicable when f is invertible

Invertible models





Prototype for stochastic neural network



P(f(x)) is in general intractable for stochastic neural networks

$$p_X(x) = p_H(f(x)) |\det rac{\partial f(x)}{\partial x}|.$$

is applicable when f is invertible

Invertible models



#### Fully factorized model



Any other family of tractable models existing?

Route 1: Easy sampling, hard probability evaluation

Prototype for stochastic neural network



Fully factorized model



P(f(x)) is in general intractable for stochastic neural networks

$$p_X(x) = p_H(f(x)) |\det rac{\partial f(x)}{\partial x}|.$$

is applicable when f is invertible

Invertible models



Any other more flexible tractable families existing?

#### Route 2: Easy (unnormalized) probability calculation, hard sampling

Prototype for neural network-based energy model

Use stochastic neural network as proposal distribution to assist sampling



Direct gradient-based MCMC sampling like HMC is feasible, but generally too slow

Is it possible to design trainable reversible conditional stochastic neural network ? Q(x'|x)=Q(x|x') Metropolis sampling

Acceptance probability min

$$\operatorname{n}\left(1,\frac{P(x')}{P(x)}\frac{Q(x)}{Q(x')}\right)$$

Metropolis-Hastings sampling

Acceptance probability m

$$\min\left(1,\frac{P(x')}{P(x)}\right)$$