

# Strong Data Processing Inequalities and $\Phi$ -Sobolev Inequalities for Discrete Channels

Maxim Raginsky\*<sup>†</sup>

March 30, 2016

## Abstract

The noisiness of a channel can be measured by comparing suitable functionals of the input and output distributions. For instance, the worst-case ratio of output relative entropy to input relative entropy for all possible pairs of input distributions is bounded from above by unity, by the data processing theorem. However, for a fixed reference input distribution, this quantity may be strictly smaller than one, giving so-called strong data processing inequalities (SDPIs). The same considerations apply to an arbitrary  $\Phi$ -divergence. This paper presents a systematic study of optimal constants in SDPIs for discrete channels, including their variational characterizations, upper and lower bounds, structural results for channels on product probability spaces, and the relationship between SDPIs and so-called  $\Phi$ -Sobolev inequalities (another class of inequalities that can be used to quantify the noisiness of a channel by controlling entropy-like functionals of the input distribution by suitable measures of input-output correlation). Several applications to information theory, discrete probability, and statistical physics are discussed.

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Notation . . . . .	3
<b>2</b>	<b>Background on <math>\Phi</math>-entropies and <math>\Phi</math>-divergences</b>	<b>4</b>
2.1	Subadditivity of $\Phi$ -entropies . . . . .	8
<b>3</b>	<b>Strong data processing inequalities</b>	<b>10</b>
3.1	A universal upper bound via Markov contraction . . . . .	12
3.2	Bounds via maximal correlation . . . . .	13
3.3	Upper bounds for operator convex $\Phi$ . . . . .	17
3.4	Upper bounds via subgaussian concentration and information-transportation inequalities . . . . .	21
3.5	Tensorization . . . . .	28

\*The author is with the Department of Electrical and Computer Engineering and the Coordinated Science Laboratory, University of Illinois, Urbana, IL 61801, USA. E-mail: maxim@illinois.edu.

<sup>†</sup>This work was supported in part by the NSF under CAREER award no. CCF-1254041 and by the Center for Science of Information (CSoI), an NSF Science and Technology Center, under grant agreement CCF-0939370. The material in this paper was presented in part at the 2013 IEEE International Symposium on Information Theory.

3.6	Mixtures of local channels . . . . .	30
3.7	Comparison of SDPI constants . . . . .	33
3.8	Extremal functions . . . . .	34
<b>4</b>	<b>Connections with <math>\Phi</math>-Sobolev inequalities</b>	<b>37</b>
4.1	General framework . . . . .	37
4.2	Logarithmic Sobolev and Poincaré inequalities . . . . .	40
4.3	The gap between SDPI and $\Phi$ -Sobolev . . . . .	49
<b>5</b>	<b>Some applications</b>	<b>50</b>
5.1	Concentration inequalities . . . . .	50
5.2	Contraction of mutual information in a Markov chain . . . . .	52
5.3	Fastest mixing Markov chain on a graph . . . . .	55
5.4	Mixing times of Swendsen-Wang and heat-bath dynamics . . . . .	58
5.5	Reconstruction in graphical models . . . . .	62
<b>6</b>	<b>Summary of contributions and concluding remarks</b>	<b>64</b>
<b>A</b>	<b>Miscellaneous lemmas</b>	<b>66</b>
<b>B</b>	<b>Proof of Proposition 4.1</b>	<b>69</b>

## 1 Introduction

The well-known data processing inequality for the relative entropy states that, for any two probability distributions  $\mu, \nu$  over an alphabet  $\mathsf{X}$  and for any stochastic transformation (channel)  $K$  with input alphabet  $\mathsf{X}$  and output alphabet  $\mathsf{Y}$ ,

$$D(\nu K \parallel \mu K) \leq D(\nu \parallel \mu),$$

where  $\mu K$  denotes the distribution at the output of  $K$  when the input has distribution  $\mu$  (and similarly for  $\nu K$ ). However, if we fix the *reference distribution*  $\mu$  and vary only  $\nu$ , then in many cases it is possible to show that  $D(\nu K \parallel \mu K)$  is *strictly* smaller than  $D(\nu \parallel \mu)$  unless  $\nu \equiv \mu$ . To capture this effect, we define the quantity

$$\eta(\mu, K) \triangleq \sup_{\nu \neq \mu} \frac{D(\nu K \parallel \mu K)}{D(\nu \parallel \mu)},$$

and we say that the channel  $K$  satisfies a *strong data processing inequality* (SDPI) at input distribution  $\mu$  if  $\eta(\mu, K) < 1$ . In a remarkable paper [1], Ahlswede and Gács have uncovered deep relationships between  $\eta(\mu, K)$  and several other quantities, such as the maximal correlation (see [2] and references therein) and so-called hypercontractivity constants of a certain Markov operator associated to the pair  $(\mu, K)$ . For example, they have shown that if  $\mathsf{X} = \mathsf{Y} = \{0, 1\}$ ,  $\mu = \text{Bern}(1/2)$ , and  $K = \text{BSC}(\varepsilon)$ , then  $\eta(\mu, K) = (1 - 2\varepsilon)^2$ , which is also equal to the squared maximal correlation in the joint distribution  $P_{XY}$  with  $P_X = P_Y = \text{Bern}(1/2)$  and  $P_{Y|X} = \text{BSC}(\varepsilon)$ , the so-called doubly symmetric binary source (DSBS) with parameter  $\varepsilon$  [3].

After the pioneering work of Ahlswede and Gács, the contraction properties of relative entropy (and other  $\Phi$ -divergences [4, 5]) under the action of stochastic transformations have been studied by

several other authors [6–10]. In particular, Cohen et al. [6], who were the first ones to take up this subject after [1], showed that the SDPI constant of any channel  $K$  with respect to any  $\Phi$ -divergence is always upper-bounded by the so-called *Dobrushin contraction coefficient* of  $K$  [11, 12], another well-known numerical measure of the amount of noise introduced by a channel. (This result of Cohen et al. was rediscovered five years later in the machine learning community [13].) In the last couple of years, strong data processing inequalities became the subject of intense interest in the information theory community [14–22] due to their apparent usefulness for establishing various converse results.

In this paper, we revisit the problem of characterizing the strong data processing constant  $\eta(\mu, K)$  [and its generalizations for arbitrary  $\Phi$ -divergence] and establish a number of new upper and lower bounds, as well as new structural results on SDPI constants in product probability spaces. We also address the relationship between strong data processing inequalities and so-called  *$\Phi$ -Sobolev inequalities* [23]. These inequalities also quantify the noisiness of a Markov operator (probability transition kernel) by relating certain “entropy-like” functionals of the input to the rate of increase of suitable “energy-like” quantities from the input to the output. (Logarithmic Sobolev inequalities, widely studied in the theory of probability and Markov chains [8, 24–27], are a special case.) In particular, we show that the optimal constants in  $\Phi$ -Sobolev inequalities for a reversible Markov chain can be related to SDPI constants of certain factorizations of the transition kernel of the chain as a product of a forward channel and a backward channel. Such factorizations correspond to all possible realizations of the one-step transition of the chain as a two-component Gibbs sampler [28], which is a standard technique in Markov chain Monte Carlo [29, 30]. Conversely, for a fixed input distribution  $\mu$  on  $\mathsf{X}$ , the SDPI constants of a given channel  $K$  with input in  $\mathsf{X}$  and output in  $\mathsf{Y}$  are related to  $\Phi$ -Sobolev constants of the reversible Markov chain on  $\mathsf{X}$  obtained by composing the forward channel  $K$  with the backward channel  $K^*$  determined via Bayes’ rule. To keep things simple, we focus on the discrete case, when both  $\mathsf{X}$  and  $\mathsf{Y}$  are finite, although some of our results generalize easily to the case of arbitrary Polish alphabets (see, e.g., [21]).

The remainder of the paper is organized as follows. After giving some necessary background on  $\Phi$ -entropies and  $\Phi$ -divergences in Section 2, we proceed to the study of strong data processing inequalities in Section 3. Next, in Section 4, we define the  $\Phi$ -Sobolev inequalities and characterize their relation with SDPIs. Several examples of applications are given in Section 5. Section 6 provides a summary of key contributions. A number of auxiliary technical results are stated and proved in the Appendices.

## 1.1 Notation

We will denote by  $\mathcal{P}(\mathsf{X})$  the set of all probability distributions on an alphabet  $\mathsf{X}$  and by  $\mathcal{P}_*(\mathsf{X})$  the subset of  $\mathcal{P}(\mathsf{X})$  consisting of all strictly positive distributions. The set of all real-valued functions on  $\mathsf{X}$  is denoted by  $\mathcal{F}(\mathsf{X})$ ;  $\mathcal{F}_*(\mathsf{X})$  and  $\mathcal{F}_*^0(\mathsf{X})$  are the subsets of  $\mathcal{F}(\mathsf{X})$  consisting of all strictly positive and nonnegative functions, respectively. Any channel<sup>1</sup> with input alphabet  $\mathsf{X}$ , output alphabet  $\mathsf{Y}$ , and transition probabilities  $\{K(y|x) : x \in \mathsf{X}, y \in \mathsf{Y}\}$  acts on probability distributions  $\mu \in \mathcal{P}(\mathsf{X})$  from the right by

$$\mu K(y) = \sum_{x \in \mathsf{X}} \mu(x) K(y|x), \quad y \in \mathsf{Y}$$

---

<sup>1</sup>We will also use the terms “stochastic transformation” or “Markov kernel.”

or on functions  $f \in \mathcal{F}(\mathsf{Y})$  from the left by

$$Kf(x) = \sum_{y \in \mathsf{Y}} K(y|x)f(y), \quad x \in \mathsf{X}.$$

The set of all such channels will be denoted by  $\mathcal{M}(\mathsf{Y}|\mathsf{X})$ . The affine map  $\mu \mapsto \mu K$  naturally extends to a linear map on the signed measures on  $\mathsf{X}$ , since any such measure  $\nu$  can be uniquely represented as  $\alpha_1\mu_1 - \alpha_2\mu_2$  for some constants  $\alpha_1, \alpha_2 \geq 0$  and some  $\mu_1, \mu_2 \in \mathcal{P}(\mathsf{X})$ ; thus, we set  $\nu K = \alpha_1\mu_1 K - \alpha_2\mu_2 K$ . The linear map  $f \mapsto Kf$  is positive [i.e.,  $K(\mathcal{F}_*^0(\mathsf{Y})) \subseteq \mathcal{F}_*^0(\mathsf{X})$ ], and unital [i.e.,  $K1 = 1$ , where 1 denotes the constant function that takes the value 1 everywhere on its domain]. If  $\mu \otimes K \in \mathcal{P}(\mathsf{X} \times \mathsf{Y})$  denotes the distribution of a random pair  $(X, Y) \in \mathsf{X} \times \mathsf{Y}$  with  $P_X = \mu$  and  $P_{Y|X} = K$ , then  $Kf(x) = \mathbb{E}[f(Y)|X = x]$  for any  $f \in \mathcal{F}(\mathsf{Y})$  and  $x \in \mathsf{X}$ .

We will say that a pair  $(\mu, K) \in \mathcal{P}(\mathsf{X}) \times \mathcal{M}(\mathsf{Y}|\mathsf{X})$  is *admissible* if  $\mu \in \mathcal{P}_*(\mathsf{X})$  and  $\mu K \in \mathcal{P}_*(\mathsf{Y})$ . For any such pair, there exists a unique channel  $K^* \in \mathcal{M}(\mathsf{X}|\mathsf{Y})$  with the property that

$$\mathbb{E}[g(Y)Kf(Y)] = \mathbb{E}[K^*g(X)f(X)] \quad (1.1)$$

for all  $g \in \mathcal{F}(\mathsf{Y}), f \in \mathcal{F}(\mathsf{X})$ . This *backward* or *adjoint* channel can be specified explicitly via the transition probabilities

$$K^*(x|y) = \frac{K(y|x)\mu(x)}{\mu K(y)}, \quad (x, y) \in \mathsf{X} \times \mathsf{Y} \quad (1.2)$$

(this is simply an application of Bayes' rule). If  $(X, Y) \sim \mu \otimes K$ , then  $K^* = P_{X|Y}$ , so in particular  $K^*f(y) = \mathbb{E}[f(X)|Y = y]$  for any  $f \in \mathcal{F}(\mathsf{X})$  and  $y \in \mathsf{Y}$ . Strictly speaking,  $K^*$  depends on both  $\mu$  and  $K$ , and we may occasionally indicate this fact by writing  $K_\mu^*$  instead of  $K^*$ .

Given a number  $p \in [0, 1]$ , we will often write  $\bar{p}$  for  $1 - p$ . For  $p, q \in [0, 1]$ , we let  $p \star q \triangleq p\bar{q} + \bar{p}q$ . Thus, if  $X \sim \text{Bern}(p)$  and  $Z \sim \text{Bern}(q)$  are independent random variables, then  $Y = X \oplus Z$  has distribution  $\text{Bern}(p \star q)$ . For  $a, b \in \mathbb{R}$ , we let  $a \vee b \triangleq \max\{a, b\}$  and  $a \wedge b \triangleq \min\{a, b\}$ . Other notation and definitions will be introduced in the sequel as needed.

## 2 Background on $\Phi$ -entropies and $\Phi$ -divergences

Let  $\mathcal{F}$  denote the set of all convex functions  $\Phi: \mathbb{R}^+ \rightarrow \mathbb{R}$ . For any  $\Phi \in \mathcal{F}$ , the  $\Phi$ -*entropy* of a nonnegative real-valued random variable  $U$  is defined by

$$\text{Ent}_\Phi[U] \triangleq \mathbb{E}[\Phi(U)] - \Phi(\mathbb{E}U), \quad (2.1)$$

provided  $\mathbb{E}[\Phi(U)] < \infty$  (see [23] and [31, Chap. 14]). For example, if  $\Phi(u) = u^2$ , then  $\text{Ent}_\Phi[U] = \text{Var}[U]$ ; if  $\Phi(u) = u \log u$ , then

$$\text{Ent}_\Phi[U] = \mathbb{E}[U \log U] - \mathbb{E}[U] \log \mathbb{E}[U].$$

The  $\Phi$ -entropy is nonnegative by Jensen's inequality.

The  $\Phi$ -divergences<sup>2</sup> between probability distributions [4, 5] arise as a special case of the above definition. Fix some  $\mu \in \mathcal{P}_*(\mathsf{X})$  (this restriction is sufficient for our purposes, and helps avoid

---

<sup>2</sup>We use the term “ $\Phi$ -divergence” instead of the more common “ $f$ -divergence” because we reserve  $f$  for real-valued functions on  $\mathsf{X}$ .

certain technicalities involving division by zero). Then, for any  $\Phi \in \mathcal{F}$ , the  $\Phi$ -divergence between an arbitrary probability distribution  $\nu \in \mathcal{P}(\mathsf{X})$  and  $\mu$  is defined as

$$D_{\Phi}(\nu||\mu) \triangleq \mathbb{E}_{\mu} \left[ \Phi \left( \frac{d\nu}{d\mu} \right) \right] - \Phi(1).$$

Note that this differs from the usual definition by the subtraction of  $\Phi(1)$ . There are two reasons behind this modification: (a)  $D_{\Phi}(\mu||\mu) = 0$  for any  $\mu$ ,<sup>3</sup> and (b) any two  $\Phi, \Phi'$  such that  $\Phi - \Phi'$  is affine determine the same divergence. If we now consider a random variable  $X \in \mathsf{X}$  with distribution  $\mu$  and let  $f = d\nu/d\mu$ , then

$$D_{\Phi}(\nu||\mu) = \text{Ent}_{\Phi} [f(X)].$$

Moreover, if  $\Phi(1) = 0$ , we can write  $D_{\Phi}(\nu||\mu) = \mathbb{E}_{\mu}[\Phi \circ f]$  since  $\mathbb{E}[f(X)] = 1$ . Here are some important examples of  $\Phi$ -divergences [5]:

1. The relative entropy

$$D(\nu||\mu) = \mathbb{E}_{\nu} \left[ \log \frac{d\nu}{d\mu} \right] = \mathbb{E}_{\mu} \left[ \frac{d\nu}{d\mu} \log \frac{d\nu}{d\mu} \right]$$

is a  $\Phi$ -divergence with  $\Phi(u) = u \log u$ .

2. The total variation distance

$$\|\nu - \mu\|_{\text{TV}} = \frac{1}{2} \mathbb{E}_{\mu} \left| \frac{d\nu}{d\mu} - 1 \right|$$

is a  $\Phi$ -divergence with  $\Phi(u) = \frac{1}{2}|u - 1|$ .

3. The  $\chi^2$ -divergence

$$\chi^2(\nu||\mu) = \mathbb{E}_{\mu} \left[ \left( \frac{d\nu}{d\mu} - 1 \right)^2 \right]$$

is a  $\Phi$ -divergence with  $\Phi(u) = (u - 1)^2$  or  $\Phi(u) = u^2 - 1$ . This is a particular instance of the fact that any two  $\Phi, \Phi' \in \mathcal{F}$  that differ by an affine function determine the same divergence.

4. The squared Hellinger distance

$$H^2(\nu, \mu) = \mathbb{E}_{\mu} \left[ \left( \sqrt{\frac{d\nu}{d\mu}} - 1 \right)^2 \right]$$

is a  $\Phi$ -divergence with  $\Phi(u) = (\sqrt{u} - 1)^2$  or  $\Phi(u) = 2 - 2\sqrt{u}$ .

---

<sup>3</sup>However, unless  $u \mapsto \Phi(u)$  is strictly convex at 1,  $D_{\Phi}(\nu||\mu) = 0$  does not necessarily imply that  $\nu = \mu$ .

An important class of  $\Phi$ -divergences arises in the context of Bayesian estimation. Given a parameter  $\lambda \in (0, 1)$ , consider a random pair  $(\Theta, X)$  with

$$\Theta \sim \text{Bern}(\lambda) \quad \text{and} \quad P_{X|\Theta=\theta} = \begin{cases} \mu, & \text{if } \theta = 0 \\ \nu, & \text{if } \theta = 1 \end{cases}.$$

Fix an action space  $\mathbf{A}$  and a loss function  $\ell : \{0, 1\} \times \mathbf{A} \rightarrow \mathbb{R}$  — in other words, if  $\Theta = \theta$  and an action  $a \in \mathbf{A}$  is selected, then we incur the loss of  $\ell(\theta, a)$ . Consider the problem of selecting an action in  $\mathbf{A}$  based on some observation  $Z$  related to  $(\Theta, X)$  via the Markov chain  $\Theta \rightarrow X \rightarrow Z$  — i.e.,  $Z$  and  $\Theta$  are conditionally independent given  $X$ . If  $A = \gamma(Z)$  for some function  $\gamma$ , then we incur the average loss

$$\mathbb{E}[\ell(\Theta, \gamma(Z))] = \bar{\lambda} \mathbb{E}[\ell(0, \gamma(Z))] + \lambda \mathbb{E}[\ell(1, \gamma(Z))].$$

The goal is to pick  $\gamma$  to minimize this expected loss for a given observation channel  $P_{Z|X}$ . In the extreme case when  $Z$  is independent of  $X$ , the best we can do is to take

$$a^* = \arg \min_{a \in \mathbf{A}} [\bar{\lambda} \ell(0, a) + \lambda \ell(1, a)],$$

giving us the average loss of

$$L_\lambda^* \triangleq \inf_{a \in \mathbf{A}} [\bar{\lambda} \ell(0, a) + \lambda \ell(1, a)].$$

On the other hand, if  $Z = X$ , then we can attain the *minimum Bayes risk*

$$\begin{aligned} L_\lambda^*(\nu, \mu) &\triangleq \inf_{\gamma} \mathbb{E}[\ell(\Theta, \gamma(X))] \\ &= \inf_{\gamma} \left\{ \bar{\lambda} \int_{\mathbf{X}} \ell(0, \gamma(x)) \nu(\mathrm{d}x) + \lambda \int_{\mathbf{X}} \ell(1, \gamma(x)) \mu(\mathrm{d}x) \right\}, \end{aligned}$$

where the infimum is over all measurable functions  $\gamma : \mathbf{X} \rightarrow \mathbf{A}$ . The following result is well-known (see, e.g., [32, p. 882]), but the proof is so simple that we give it here:

**Proposition 2.1.** *The quantity*

$$D_{\ell, \lambda}(\nu \| \mu) \triangleq L_\lambda^* - L_\lambda^*(\nu, \mu)$$

*is a  $\Phi$ -divergence.*

*Proof.* Define the function

$$\Phi_{\ell, \lambda}(u) \triangleq \sup_{a \in \mathbf{A}} [L_\lambda^* - \bar{\lambda} \ell(0, a) - \lambda \ell(1, a)u], \quad u \geq 0.$$

Being a pointwise supremum of affine functions of  $u$ , it is convex. Moreover,  $\Phi_{\ell, \lambda}(1) = 0$ . With this, we can write

$$\begin{aligned} L_\lambda^* - L_\lambda^*(\nu, \mu) &= \sup_{\gamma} \left( L_\lambda^* - \int_{\mathbf{X}} \mu(\mathrm{d}x) \left[ \bar{\lambda} \ell(0, \gamma(x)) + \lambda \frac{\mathrm{d}\nu}{\mathrm{d}\mu}(x) \ell(1, \gamma(x)) \right] \right) \\ &= \int_{\mathbf{X}} \mu(\mathrm{d}x) \sup_{a \in \mathbf{A}} \left[ L_\lambda^* - \bar{\lambda} \ell(0, a) - \lambda \ell(1, a) \frac{\mathrm{d}\nu}{\mathrm{d}\mu}(x) \right] \\ &= \mathbb{E}_\mu \left[ \Phi_{\ell, \lambda} \left( \frac{\mathrm{d}\nu}{\mathrm{d}\mu} \right) \right]. \end{aligned}$$

□

We consider two particular cases:

- $A = \{0, 1\}$ ,  $\ell(\theta, a) = \mathbf{1}_{\{\theta \neq a\}}$ . An easy calculation shows that  $L_\lambda^* = \lambda \wedge \bar{\lambda}$  and

$$\begin{aligned}\Phi_{\ell, \lambda}(u) &= [\lambda \wedge \bar{\lambda} - \bar{\lambda}u] \vee [\lambda \wedge \bar{\lambda} - \bar{\lambda}] \\ &= \lambda \wedge \bar{\lambda} - (\lambda u) \wedge \bar{\lambda}.\end{aligned}$$

Alternatively, we can write

$$L_\lambda^* = \frac{1}{2} - \frac{1}{2} \|\text{Bern}(\lambda) - \text{Bern}(\bar{\lambda})\|_{\text{TV}} = \frac{1}{2} - \frac{1}{2} |1 - 2\lambda|$$

and

$$L_\lambda^*(\nu, \mu) = \frac{1}{2} - \frac{1}{2} \|\lambda\nu - \bar{\lambda}\mu\|_{\text{TV}},$$

where the total variation norm  $\|\nu\|_{\text{TV}}$  of a signed measure  $\nu$  on  $\mathbf{X}$  is given by

$$\|\nu\|_{\text{TV}} = \frac{1}{2} \sum_{x \in \mathbf{X}} |\nu(x)|.$$

The optimal decision function is

$$\gamma^*(x) = \mathbf{1}_{\{\lambda \frac{d\nu}{d\mu}(x) \leq \bar{\lambda}\}}.$$

The resulting divergence is known as the *Bayes* or *statistical information* [33]

$$B_\lambda(\nu \|\mu) = \frac{1}{2} \|\lambda\nu - \bar{\lambda}\mu\|_{\text{TV}} - \frac{1}{2} |1 - 2\lambda|.$$

In fact, any  $\Phi$ -divergence can be expressed as an integral of statistical informations [5, Thm. 11]: for any  $\Phi \in \mathcal{F}$ , there exists a unique Borel measure  $\mathbf{M}_\Phi$  on  $[0, 1]$ , such that

$$D_\Phi(\nu \|\mu) = \int_{[0,1]} B_\lambda(\nu \|\mu) \mathbf{M}_\Phi(d\lambda). \quad (2.2)$$

- $A = \mathbb{R}$ ,  $\ell(\theta, a) = (a - \theta)^2$ . Then  $L_\lambda^* = \lambda\bar{\lambda}$  and

$$\Phi_{\ell, \lambda}(u) = \lambda\bar{\lambda} \left( 1 - \frac{u}{\lambda u + \bar{\lambda}} \right),$$

which gives

$$L_\lambda^*(\nu, \mu) = \lambda\bar{\lambda} \mathbb{E}_\mu \left[ \frac{d\nu/d\mu}{\lambda d\nu/d\mu + \bar{\lambda}} \right],$$

with the optimum decision function

$$\gamma^*(x) = \frac{\lambda \frac{d\nu}{d\mu}(x)}{\lambda \frac{d\nu}{d\mu}(x) + \bar{\lambda}}.$$

The corresponding divergence is then given by

$$\begin{aligned} D_{\ell,\lambda}(\nu\|\mu) &= \lambda\bar{\lambda} \left( 1 - \mathbb{E}_\mu \left[ \frac{d\nu/d\mu}{\lambda d\nu/d\mu + \bar{\lambda}} \right] \right) \\ &= (\lambda\bar{\lambda})^2 \mathbb{E}_\mu \left[ \frac{(d\nu/d\mu - 1)^2}{\lambda d\nu/d\mu + \bar{\lambda}} \right], \end{aligned}$$

where the second expression follows after some algebraic manipulations. Note that the functions  $u \mapsto \frac{(u-1)^2}{\lambda u + \bar{\lambda}}$  for  $0 < \lambda < 1$  also belong to  $\mathcal{F}$ . The divergences generated by these functions (modulo multiplicative constants) have appeared throughout the statistical literature [34, 35]. In particular, Le Cam [34] considers the case  $\lambda = 1/2$  with the above Bayesian hypothesis testing interpretation, while Györfi and Vajda [35] look at arbitrary  $\lambda$  (including the endpoints 0 and 1). For our purposes, it will be convenient to work with the function  $u \mapsto \lambda\bar{\lambda} \frac{(u-1)^2}{\lambda u + \bar{\lambda}}$ , which gives the *Le Cam divergence* with parameter  $\lambda \in (0, 1)$ :

$$\text{LC}_\lambda(\nu\|\mu) \triangleq \lambda\bar{\lambda} \mathbb{E}_\mu \left[ \frac{(d\nu/d\mu - 1)^2}{\lambda d\nu/d\mu + \bar{\lambda}} \right] \equiv \frac{1}{\lambda\bar{\lambda}} D_{\ell,\lambda}(\nu\|\mu). \quad (2.3)$$

The Le Cam divergences  $\text{LC}_0(\cdot\|\cdot)$  and  $\text{LC}_1(\cdot\|\cdot)$  are also well-defined and are identically zero.

More examples of  $\Phi$ -divergences, as well as a wide variety of inequalities between them, can be found in [36].

From now on, when dealing with quantities indexed by  $\Phi$ , we will often substitute  $\Phi$  with some mnemonic notation related to the corresponding  $\Phi$ -divergence, e.g.,  $\text{TV}$ ,  $\chi^2$ , etc. Moreover, for the case of the relative entropy we will often omit the index  $\Phi$  altogether and write  $\text{Ent}(\cdot)$ ,  $D(\cdot\|\cdot)$ , etc.

## 2.1 Subadditivity of $\Phi$ -entropies

Let  $U$  and  $Y$  be jointly distributed random variables, where  $U$  takes nonnegative real values and  $Y$  is arbitrary. Given a function  $\Phi \in \mathcal{F}$ , define the *conditional  $\Phi$ -entropy* of  $U$  given  $Y$ :

$$\text{Ent}_\Phi[U|Y] \triangleq \mathbb{E}[\Phi(U)|Y] - \Phi(\mathbb{E}[U|Y]). \quad (2.4)$$

This is a random variable, since it depends on  $Y$ . Combining (2.4) with (2.1) gives the following generalization of the law of total variance:

$$\text{Ent}_\Phi[U] = \mathbb{E}[\text{Ent}_\Phi[U|Y]] + \text{Ent}_\Phi[\mathbb{E}[U|Y]] \quad (2.5)$$

(see [23, pp. 351–352]).

**Remark 2.1.** We may think of

$$J_\Phi(U|Y) \triangleq \mathbb{E}[\text{Ent}_\Phi[U|Y]]$$

as a kind of “Fisher  $\Phi$ -information” about  $U$  contained in  $Y$ .<sup>4</sup> Indeed, let us consider the following special case: let  $(Y, Y')$  be an exchangeable pair on some space  $\mathcal{Y}$  (i.e.,  $P_{Y,Y'}(y, y') = P_{Y,Y'}(y', y)$ ) for

---

<sup>4</sup>We are grateful to P. Tetali for suggesting this interpretation.

all  $y, y'$ ), and let  $U = f(Y)$  for some  $f : \mathsf{Y} \rightarrow \mathbb{R}^+$ . Let  $K$  be the stochastic transformation  $P_{Y'|Y}$ . Then  $\mathbb{E}[U|Y'] = \mathbb{E}[f(Y)|Y'] = K^*f(Y')$  has the same distribution as  $\mathbb{E}[f(Y')|Y] = K^*f(Y)$ , and

$$\begin{aligned} J_\Phi(U|Y) &= \text{Ent}_\Phi[U] - \text{Ent}_\Phi[\mathbb{E}[U|Y]] \\ &= \text{Ent}_\Phi[f(Y)] - \text{Ent}_\Phi[K^*f(Y)]. \end{aligned}$$

By convexity of  $\Phi$ ,

$$\Phi(u + v) \geq \Phi(u) + v\Phi'(u).$$

If we write  $K^* = \text{id} + L$ , where  $\text{id}$  is the identity operator on  $\mathcal{F}(\mathsf{Y})$ , then

$$\begin{aligned} J_\Phi(U|Y) &= \text{Ent}_\Phi[f(Y)] - \text{Ent}_\Phi[f(Y) + Lf(Y)] \\ &\leq -\mathbb{E}[\Phi'(f(Y))Lf(Y)]. \end{aligned}$$

Moreover, if we have a continuous-time reversible Markov chain on  $\mathsf{Y}$  with stationary distribution  $P_Y$  and with infinitesimal generator  $L$ , then  $(Y_0, Y_t)$  is an exchangeable pair for each  $t$ , and

$$\begin{aligned} J_\Phi(f(Y_0)|Y_t) &= \text{Ent}_\Phi[f(Y_0)] - \text{Ent}_\Phi[K_t^*f(Y_0)] \\ &= -t \mathbb{E}[\Phi'(f(Y_0))Lf(Y_0)] + o(t) \end{aligned}$$

Dividing both sides by  $t$  and taking the limit as  $t \rightarrow 0$ , we get

$$\left. \frac{d}{dt} J_\Phi(f(Y_0)|Y_t) \right|_{t=0} = \lim_{t \rightarrow 0} \frac{J_\Phi(f(Y_0)|Y_t)}{t} = -\mathbb{E}[\Phi'(f(Y_0))Lf(Y_0)],$$

which coincides with the  $\Phi$ -Fisher information functional of Chafaï [23, Eq. (1.14)].  $\diamond$

We say that the  $\Phi$ -entropy is *subadditive* if the inequality

$$\text{Ent}_\Phi[f(X^n)] \leq \sum_{i=1}^n \mathbb{E} \left[ \text{Ent}_\Phi[f(X^n)|X^{\setminus i}] \right] \quad (2.6)$$

holds for any tuple  $X^n = (X_1, \dots, X_n)$  of independent random variables taking values in some spaces  $\mathsf{X}_1, \dots, \mathsf{X}_n$  and for any function  $f : \mathsf{X}_1 \times \dots \times \mathsf{X}_n \rightarrow \mathbb{R}^+$ , such that  $\text{Ent}_\Phi[f(X^n)] < +\infty$ . Here,  $X^{\setminus i}$  denotes the  $(n-1)$ -tuple  $(X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n)$  obtained by deleting  $X_i$  from  $X^n$ . We are interested in the following question: what conditions on  $\Phi$  ensure that this subadditivity property holds?

For example, if  $\Phi(u) = u^2$ , then  $\text{Ent}_\Phi[U] = \text{Var}[U]$ , and in this case the subadditivity property (2.6) is the well-known *Efron–Stein–Steele inequality* [37, 38]

$$\text{Var}[U] \leq \sum_{i=1}^n \mathbb{E} \left[ \text{Var}[U|X^{\setminus i}] \right], \quad U = f(X^n).$$

It is also not hard to show that the “ordinary” entropy  $\text{Ent}[U]$  [i.e., the  $\Phi$ -entropy with  $\Phi(u) = u \log u$ ] is subadditive. In general, an induction argument can be used to show that subadditivity is equivalent to the following convexity property [39]: for any two probability spaces  $(\mathsf{X}_1, \nu_1)$  and  $(\mathsf{X}_2, \nu_2)$  and any function  $f : \mathsf{X}_1 \times \mathsf{X}_2 \rightarrow \mathbb{R}^+$ ,

$$\text{Ent}_\Phi \left[ \int_{\mathsf{X}_2} f(X_1, x_2) \nu_2(dx_2) \right] \leq \int_{\mathsf{X}_2} \text{Ent}_\Phi[f(X_1, x_2)] \nu_2(dx_2), \quad (2.7)$$

where  $X_1 \sim \nu_1$ . The following criterion for subadditivity is useful [39, 40]:

**Proposition 2.2.** *Let  $\mathcal{C}$  be the class of all convex functions  $\Phi: \mathbb{R}^+ \rightarrow \mathbb{R}$  that are twice differentiable on  $(0, \infty)$ , and such that either  $\Phi$  is affine or  $\Phi'' > 0$  and  $1/\Phi''$  is concave. Then the  $\Phi$ -entropy is subadditive for all  $\Phi \in \mathcal{C}$ . Conversely, if  $\Phi$  is twice differentiable with  $\Phi'' > 0$  and the  $\Phi$ -entropy is subadditive, then  $1/\Phi''$  is concave.*

### 3 Strong data processing inequalities

We now turn to the main subject of the paper: strong data processing inequalities.

**Definition 3.1.** *Given an admissible pair  $(\mu, K) \in \mathcal{P}_*(\mathbf{X}) \times \mathcal{M}(\mathbf{Y}|\mathbf{X})$  and a function  $\Phi \in \mathcal{F}$ , we say that  $K$  satisfies a  $\Phi$ -type strong data processing inequality (SDPI) at  $\mu$  with constant  $c \in [0, 1)$ , or  $\text{SDPI}_\Phi(\mu, c)$  for short, if*

$$D_\Phi(\nu K \| \mu K) \leq c D_\Phi(\nu \| \mu) \quad (3.1)$$

for all  $\nu \in \mathcal{P}(\mathbf{X})$ . We say that  $K$  satisfies  $\text{SDPI}_\Phi(c)$  if it satisfies  $\text{SDPI}_\Phi(\mu, c)$  for all  $\mu \in \mathcal{P}_*(\mathbf{X})$ .

We are interested in the tightest constants in SDPIs; with that in mind, we define

$$\begin{aligned} \eta_\Phi(\mu, K) &\triangleq \sup_{\nu \neq \mu} \frac{D_\Phi(\nu K \| \mu K)}{D_\Phi(\nu \| \mu)}, \\ \eta_\Phi(K) &\triangleq \sup_{\mu \in \mathcal{P}_*(\mathbf{X})} \eta_\Phi(\mu, K). \end{aligned}$$

For future reference, we record the following straightforward results:

**Proposition 3.1** (Functional form of SDPI). *Fix an admissible pair  $(\mu, K)$  and let  $(X, Y)$  be a random pair with probability law  $\mu \otimes K$ . Then  $\eta_\Phi(\mu, K) \leq c$  if and only if the inequality*

$$\text{Ent}_\Phi[f(X)] \leq \frac{1}{1-c} \mathbb{E}[\text{Ent}_\Phi[f(X)|Y]] \quad (3.2)$$

holds for all nonconstant  $f \in \mathcal{F}_*^0(\mathbf{X})$  with  $\mathbb{E}[f(X)] = 1$ . Consequently,

$$\eta_\Phi(\mu, K) = \sup \left\{ \frac{\text{Ent}_\Phi[K^* f(Y)]}{\text{Ent}_\Phi[f(X)]} : f \in \mathcal{F}_*^0(\mathbf{X}), f \neq \text{const}, \mathbb{E}[f(X)] = 1 \right\} \quad (3.3)$$

$$= 1 - \inf \left\{ \frac{\mathbb{E}[\text{Ent}_\Phi[f(X)|Y]]}{\text{Ent}_\Phi[f(X)]} : f \in \mathcal{F}_*^0(\mathbf{X}), f \neq \text{const}, \mathbb{E}[f(X)] = 1 \right\}. \quad (3.4)$$

*Proof.* Fix a probability distribution  $\nu \neq \mu$  and let  $f = d\nu/d\mu$ . Then  $f \neq \text{const}$ ,  $\mathbb{E}[f(X)] = 1$ , and

$$\frac{d(\nu K)}{d(\mu K)} = K^* f$$

by Lemma A.1 in the Appendix. Therefore,

$$D_\Phi(\nu \| \mu) = \text{Ent}_\Phi \left[ \frac{d\nu}{d\mu}(X) \right] \quad \text{and} \quad D_\Phi(\nu K \| \mu K) = \text{Ent}_\Phi \left[ \frac{d(\nu K)}{d(\mu K)}(Y) \right].$$

Conversely, for any nonconstant  $f \in \mathcal{F}_*^0(\mathbf{X})$  with  $\mathbb{E}[f(X)] = 1$  there exists a probability distribution  $\nu \in \mathcal{P}(\mathbf{X})$  such that  $\nu \neq \mu$  and  $f = d\nu/d\mu$ . In that case, the above formulas for the  $\Phi$ -entropies hold as well.

Now, if  $c = 1$ , then (3.2) holds trivially, so assume  $c < 1$ . In that case, the result follows from Eq. (3.2) and the law of total  $\Phi$ -entropy, Eq. (2.5).  $\square$

**Definition 3.2.** We say that the  $\Phi$ -entropy  $\text{Ent}_\Phi[\cdot]$  is homogeneous if there exists some function  $\kappa: (0, \infty) \rightarrow (0, \infty)$ , such that the equality

$$\text{Ent}_\Phi[cU] = \kappa(c) \text{Ent}_\Phi[U] \quad (3.5)$$

holds for any nonnegative random variable  $U$  such that  $\text{Ent}_\Phi[U] < +\infty$  and for any positive real number  $c$ .

For example,  $\Phi(u) = u \log u$  satisfies (3.5) with  $\kappa(c) = c$ , while  $\Phi(u) = \frac{u^\alpha - 1}{\alpha - 1}$ ,  $\alpha > 1$ , satisfies (3.5) with  $\kappa(c) = c^\alpha$ .

**Proposition 3.2.** Suppose that (3.5) holds. Then

$$\begin{aligned} \eta_\Phi(\mu, K) &= \sup \left\{ \frac{\text{Ent}_\Phi[K^*f(Y)]}{\text{Ent}_\Phi[f(X)]} : f \in \mathcal{F}_*^0(\mathsf{X}), f \neq \text{const} \right\} \\ &= 1 - \inf \left\{ \frac{\mathbb{E}[\text{Ent}_\Phi[f(X)|Y]]}{\text{Ent}_\Phi[f(X)]} : f \in \mathcal{F}_*^0(\mathsf{X}), f \neq \text{const} \right\}. \end{aligned} \quad (3.6)$$

Moreover, if  $\kappa$  is an invertible function, then

$$\eta_\Phi(\mu, K) = \eta_\Phi(\mu, K, t) \triangleq \sup \left\{ \frac{\text{Ent}_\Phi[K^*f(Y)]}{\text{Ent}_\Phi[f(X)]} : f \in \mathcal{F}_*^0(\mathsf{X}), \text{Ent}_\Phi[f(X)] \leq t \right\}, \quad \forall t > 0. \quad (3.7)$$

Again,  $(X, Y)$  is a random pair with law  $\mu \otimes K$ .

*Proof.* Eq. (3.6) is obvious from homogeneity. To prove (3.7), pick an arbitrary nonconstant  $f \in \mathcal{F}_*^0(\mathsf{X})$  and let

$$c = \kappa^{-1} \left( \frac{t}{\text{Ent}_\Phi[f(X)]} \right).$$

Let  $g = cf$ . Then  $\text{Ent}_\Phi[g(X)] = \text{Ent}_\Phi[cf(X)] = \kappa(c) \text{Ent}_\Phi[f(X)] = t$ . Therefore,

$$\text{Ent}_\Phi[K^*g(Y)] \leq \eta_\Phi(\mu, K, t) \text{Ent}_\Phi[g(X)].$$

Since  $\text{Ent}_\Phi[K^*g(Y)] = \text{Ent}_\Phi[cK^*g(Y)] = \kappa(c) \text{Ent}_\Phi[K^*g(Y)]$ , and since  $c > 0$  by the properties of  $\kappa$ , we conclude that  $\text{Ent}_\Phi[K^*f(Y)] \leq \eta_\Phi(\mu, K, t) \text{Ent}_\Phi[f(X)]$ , which implies that  $\eta_\Phi(\mu, K) \leq \eta_\Phi(\mu, K, t)$ . The reverse inequality,  $\eta_\Phi(\mu, K) \geq \eta_\Phi(\mu, K, t)$ , is obvious.  $\square$

**Proposition 3.3** (Convexity in the kernel). For a given choice of  $\mathsf{X}, \mathsf{Y}$ , and  $\mu \in \mathcal{P}_*(\mathsf{X})$ , the SDPI constants  $\eta_\Phi(\mu, K)$  and  $\eta_\Phi(K)$  are convex in  $K \in \mathcal{M}(\mathsf{Y}|\mathsf{X})$ .

*Proof.* For fixed  $\nu, \mu \in \mathcal{P}(\mathsf{X})$ , the functional  $K \mapsto \frac{D_\Phi(\nu K \| \mu K)}{D_\Phi(\nu \| \mu)}$  is convex because of the joint convexity of  $D_\Phi(\cdot \| \cdot)$  [41, Lemma 4.1].<sup>5</sup> Now,

$$\eta_\Phi(\mu, K) = \sup_\nu \frac{D_\Phi(\nu K \| \mu K)}{D_\Phi(\nu \| \mu)} \quad \text{and} \quad \eta_\Phi(K) = \sup_\mu \sup_\nu \frac{D_\Phi(\nu K \| \mu K)}{D_\Phi(\nu \| \mu)}$$

are pointwise suprema of convex functionals of  $K$ , and therefore are convex in  $K$ .  $\square$

<sup>5</sup>Joint convexity of  $D_\Phi(\cdot \| \cdot)$  follows from the fact that, for any convex function  $\Phi: \mathbb{R}^+ \rightarrow \mathbb{R}$ , the *perspective function*  $(p, q) \mapsto q\Phi(p/q)$  is jointly convex in  $(p, q) \in \mathbb{R}^+ \times \mathbb{R}^+$  [42, Prop. 2.2.1].

### 3.1 A universal upper bound via Markov contraction

A universal upper bound on  $\eta_\Phi(K)$  was originally obtained by Cohen et al. [6] in the discrete case and subsequently extended by Del Moral et al. [10] to the general case. We state this bound and give a proof which is more information-theoretic in nature:

**Theorem 3.1.** *Define the Dobrushin contraction coefficient [11, 12] of a channel  $K \in \mathcal{M}(\mathsf{Y}|\mathsf{X})$  by*

$$\vartheta(K) \triangleq \max_{x, x' \in \mathsf{X}} \|K(\cdot|x) - K(\cdot|x')\|_{\text{TV}}. \quad (3.8)$$

Then for any  $\Phi \in \mathcal{F}$  we have

$$\eta_\Phi(K) \leq \vartheta(K). \quad (3.9)$$

Moreover,  $\eta_{\text{TV}}(K) \equiv \vartheta(K)$ .

*Proof.* By the integral representation (2.2), it suffices to show that (3.9) holds for the statistical informations  $B_\lambda(\cdot|\cdot)$ ,  $0 \leq \lambda \leq 1$ . For that, we need the following *strong Markov contraction lemma* [6, Lemma 3.2]: for any signed measure  $\tilde{\nu}$  on  $\mathsf{X}$  and any Markov kernel  $K \in \mathcal{M}(\mathsf{Y}|\mathsf{X})$ ,

$$\|\tilde{\nu}K\|_{\text{TV}} \leq \vartheta(K)\|\tilde{\nu}\|_{\text{TV}} + \frac{1 - \vartheta(K)}{2}|\tilde{\nu}(\mathsf{X})|. \quad (3.10)$$

Let  $\tilde{\nu} = \lambda\nu - \bar{\lambda}\mu$ . Then  $\tilde{\nu}K = \lambda\nu K - \bar{\lambda}\mu K$  and  $\tilde{\nu}(\mathsf{X}) = 2\lambda - 1$ . Thus, using (3.10), we get

$$\|\lambda\nu K - \bar{\lambda}\mu K\|_{\text{TV}} \leq \vartheta(K)\|\lambda\nu - \bar{\lambda}\mu\|_{\text{TV}} + \frac{1 - \vartheta(K)}{2}|1 - 2\lambda|.$$

Therefore,

$$\begin{aligned} B_\lambda(\nu K \|\mu K) &= \frac{1}{2}\|\lambda\nu K - \bar{\lambda}\mu K\|_{\text{TV}} - \frac{1}{2}|1 - 2\lambda| \\ &\leq \vartheta(K) \cdot \left( \frac{1}{2}\|\lambda\nu - \bar{\lambda}\mu\|_{\text{TV}} - \frac{1}{2}|1 - 2\lambda| \right) \\ &= \vartheta(K) \cdot B_\lambda(\nu \|\mu). \end{aligned}$$

This establishes the bound (3.9). It remains to show that this bound is achieved for  $\|\cdot\|_{\text{TV}}$ .

To that end, let us first assume that  $|\mathsf{X}| > 2$ . Let  $x_0, x_1 \in \mathsf{X}$  achieve the maximum in (3.8), pick some  $\varepsilon_1, \varepsilon_2, \varepsilon \in (0, 1)$  such that  $\varepsilon_1 \neq \varepsilon_2$ ,  $\varepsilon_1 + \varepsilon < 1$ ,  $\varepsilon_2 + \varepsilon < 1$ , and consider the following probability distributions:

- $\nu$  that puts the mass  $1 - \varepsilon_1 - \varepsilon$  on  $x_0$ ,  $\varepsilon_1$  on  $x_1$ , and distributes the remaining mass of  $\varepsilon$  evenly among the set  $\mathsf{X} \setminus \{x_0, x_1\}$ ;
- $\mu$  that puts the mass  $1 - \varepsilon_2 - \varepsilon$  on  $x_0$ ,  $\varepsilon_2$  on  $x_1$ , and distributes the remaining mass of  $\varepsilon$  evenly among the set  $\mathsf{X} \setminus \{x_0, x_1\}$ .

Then a simple calculation gives

$$\begin{aligned} \|\nu - \mu\|_{\text{TV}} &= |\varepsilon_1 - \varepsilon_2| \\ \|\nu K - \mu K\|_{\text{TV}} &= |\varepsilon_1 - \varepsilon_2| \cdot \|K(\cdot|x_0) - K(\cdot|x_1)\|_{\text{TV}} \\ &= \vartheta(K) \cdot \|\nu - \mu\|_{\text{TV}}. \end{aligned}$$

For  $|\mathsf{X}| = 2$ , the idea is the same, except that there is no need for the extra slack  $\varepsilon$ . □

**Remark 3.1.** Theorem 3.1 says that any channel  $K$  with  $\vartheta(K) < 1$  satisfies an SDPI for any  $\Phi \in \mathcal{F}$  at any reference input distribution  $\mu \in \mathcal{P}(\mathsf{X})$ . However, the bounds it gives are generally loose. For example, for  $K = \text{BSC}(\varepsilon)$  with  $\varepsilon \in (0, 1)$ , we have  $\vartheta(K) = |1 - 2\varepsilon| < 1$ , so by Theorem 3.1

$$\eta_\Phi(\text{Bern}(p), \text{BSC}(\varepsilon)) \leq |1 - 2\varepsilon| < 1$$

for all  $\Phi \in \mathcal{F}$  and all  $p \in [0, 1]$ . However, as we know from [1],

$$\eta(\text{Bern}(1/2), \text{BSC}(\varepsilon)) = (1 - 2\varepsilon)^2 < |1 - 2\varepsilon|.$$

Later on, we will develop tighter bounds on SDPI constants for a broad class of  $\Phi$ -entropies.  $\diamond$

**Remark 3.2.** Suppose that the channel  $K \in \mathcal{M}(\mathsf{Y}|\mathsf{X})$  has the following property: There exist a constant  $0 < \alpha \leq 1$  and a probability distribution  $\tilde{\mu} \in \mathcal{P}(\mathsf{Y})$ , such that

$$K(y|x) \geq \alpha \tilde{\mu}(y) \tag{3.11}$$

for all  $x \in \mathsf{X}$  and  $y \in \mathsf{Y}$  (in Markov chain theory, this is known as a Doeblin minorization condition [43, Sec. 4.3.3]). Then  $\eta_\Phi(K) \leq 1 - \alpha$ . This bound can be proved using a nice operational argument. Indeed, if (3.11) holds, then

$$\tilde{K}(y|x) \triangleq \frac{K(y|x) - \alpha \tilde{\mu}(y)}{1 - \alpha}, \quad (x, y) \in \mathsf{X} \times \mathsf{Y}$$

defines a channel from  $\mathsf{X}$  to  $\mathsf{Y}$ . Let  $\mathbf{e}$  be a special erasure symbol, and let  $E_\alpha \in \mathcal{M}(\mathsf{X} \cup \{\mathbf{e}\}|\mathsf{X})$  denote the symmetric erasure channel on  $\mathsf{X}$  with erasure probability  $\alpha$ : any input symbol  $x \in \mathsf{X}$  is erased with probability  $\alpha$  and reproduced exactly with probability  $\bar{\alpha}$ . Then a simple calculation shows that  $K = T \circ E_\alpha$ , where the channel  $T \in \mathcal{M}(\mathsf{Y}|\mathsf{X} \cup \{\mathbf{e}\})$  is defined by

$$\begin{aligned} T(\cdot|x) &= \tilde{K}(\cdot|x), & x \in \mathsf{X} \\ T(\cdot|\mathbf{e}) &= \tilde{\mu}(\cdot). \end{aligned}$$

In that case, for any  $\mu, \nu \in \mathcal{P}(\mathsf{X})$ ,

$$\begin{aligned} D_\Phi(\nu K \| \mu K) &= D_\Phi((\nu E_\alpha)T \| (\mu E_\alpha)T) \\ &\leq D_\Phi(\nu E_\alpha \| \mu E_\alpha) \\ &= D_\Phi(\bar{\alpha}\nu + \alpha\delta_{\mathbf{e}} \| \bar{\alpha}\mu + \alpha\delta_{\mathbf{e}}) \\ &\leq \bar{\alpha}D_\Phi(\nu \| \mu), \end{aligned}$$

where the first inequality is by the usual data processing inequality, while the second inequality is by convexity. It is not hard to show that if (3.11) holds, then  $\vartheta(K) \leq 1 - \alpha$ .  $\diamond$

### 3.2 Bounds via maximal correlation

For any pair  $(\mu, K) \in \mathcal{P}(\mathsf{X}) \times \mathcal{M}(\mathsf{Y}|\mathsf{X})$ , the *maximal correlation* is defined as

$$S(\mu, K) \triangleq \sup_{f, g} \mathbb{E}[f(X)g(Y)],$$

where  $(X, Y) \sim \mu \otimes K$ , and the supremum is over all  $f \in \mathcal{F}(\mathsf{X})$ ,  $g \in \mathcal{F}(\mathsf{Y})$  satisfying  $\mathbb{E}[f(X)] = \mathbb{E}[g(Y)] = 0$  and  $\mathbb{E}[f^2(X)] = \mathbb{E}[g^2(Y)] = 1$  (see [2] and the references therein). The square of  $S(\mu, K)$  is the SDPI constant of the pair  $(\mu, K)$  for the  $\chi^2$ -divergence:

**Theorem 3.2.** Consider the  $\chi^2$ -divergence

$$\chi^2(\nu\|\mu) = \mathbb{E}_\mu \left[ \left( \frac{d\nu}{d\mu} - 1 \right)^2 \right] \equiv \text{Var}_\mu \left[ \frac{d\nu}{d\mu} \right].$$

Then, for any  $(\mu, K) \in \mathcal{P}(\mathsf{X}) \times \mathcal{M}(\mathsf{Y}|\mathsf{X})$ ,

$$\eta_{\chi^2}(\mu, K) = S^2(\mu, K).$$

**Remark 3.3.** This result has appeared in the literature in different forms (see, e.g., [10]). We give a short proof for completeness.  $\diamond$

*Proof.* For this proof, it is convenient to use operator-theoretic ideas, following Witsenhausen [2] (see also [44]). If we equip the space  $\mathcal{F}(\mathsf{X})$  with the inner product

$$\langle f, g \rangle_\mu \triangleq \mathbb{E}[f(X)g(X)], \quad \text{where } X \sim \mu$$

then it becomes the Hilbert space  $L^2(\mathsf{X}, \mu)$ ; the Hilbert space  $L^2(\mathsf{Y}, \mu K)$  is constructed in the same way. Moreover, the channels  $K$  and  $K^*$  become mutually adjoint linear operators  $K : L^2(\mathsf{Y}, \mu K) \rightarrow L^2(\mathsf{X}, \mu)$  and  $K^* : L^2(\mathsf{X}, \mu) \rightarrow L^2(\mathsf{Y}, \mu K)$ , i.e.,

$$\langle f, Kg \rangle_\mu = \langle K^*f, g \rangle_{\mu K}, \quad \forall f \in L^2(\mathsf{X}, \mu), g \in L^2(\mathsf{Y}, \mu K).$$

For  $f = d\nu/d\mu$ , we have  $\chi^2(\nu\|\mu) = \text{Var}[f(X)]$  and  $\chi^2(\nu K\|\mu K) = \text{Var}[K^*f(Y)]$ . Using this together with the fact that  $\text{Var}[U + c] = \text{Var}[U]$  for any  $c \in \mathbb{R}$  and that  $\mathbb{E}[K^*f(Y)] = \mathbb{E}[f(X)]$ , we can write

$$\begin{aligned} \eta_{\chi^2}(\nu, K) &= \sup_{\nu \neq \mu} \frac{\chi^2(\nu K\|\mu K)}{\chi^2(\nu\|\mu)} \\ &= \sup_{f \in \mathcal{H}_0(\mathsf{X})} \frac{\text{Var}[K^*f(Y)]}{\text{Var}[f(X)]}, \end{aligned}$$

where  $\mathcal{H}_0(\mathsf{X})$  is the closed linear subspace of  $L^2(\mathsf{X}, \mu)$  consisting of all  $f$  satisfying  $\langle f, 1 \rangle_\mu = 0$ , i.e.,  $\mathbb{E}[f(X)] = 0$ . For any  $f \in \mathcal{H}_0(\mathsf{X})$ ,

$$\text{Var}[f(X)] = \|f\|_\mu^2, \quad \text{Var}[K^*f(Y)] = \|K^*f\|_{\mu K}^2.$$

Since  $K$  and  $K^*$  are adjoint operators, we have

$$\|K^*f\|_{\mu K}^2 = \langle K^*f, K^*f \rangle_{\mu K} = \langle f, KK^*f \rangle_\mu,$$

which gives

$$\eta_{\chi^2}(\mu, K) = \sup_{f \in \mathcal{H}_0(\mathsf{X})} \frac{\langle f, KK^*f \rangle_\mu}{\langle f, f \rangle_\mu}.$$

Moreover,  $K^*$  maps  $\mathcal{H}_0(\mathsf{X})$  into  $\mathcal{H}_0(\mathsf{Y})$ , and  $K$  maps  $\mathcal{H}_0(\mathsf{Y})$  into  $\mathcal{H}_0(\mathsf{X})$ . Thus, by the Courant–Fischer–Weyl minimax principle [45],  $\eta_{\chi^2}(\nu, \mu)$  is the largest eigenvalue of the operator  $KK^*$  :

$\mathcal{H}_0(\mathsf{X}) \rightarrow \mathcal{H}_0(\mathsf{X})$ . The square root of this largest eigenvalue is the largest singular value of the operator  $K^* : \mathcal{H}_0(\mathsf{X}) \rightarrow \mathcal{H}_0(\mathsf{Y})$ , so, by definition,

$$\begin{aligned} \sqrt{\eta_{\chi^2}(\nu, \mu)} &= \sup_{f,g} \langle K^* f, g \rangle \\ &= \sup_{f,g} \mathbb{E}[\mathbb{E}[f(X)|Y]g(Y)] \\ &= \sup_{f,g} \mathbb{E}[f(X)g(Y)], \end{aligned}$$

where the supremum is over all  $f \in \mathcal{H}_0(\mathsf{X})$  and  $g \in \mathcal{H}_0(\mathsf{Y})$  with  $\|f\|_\mu = \|g\|_{\mu K} = 1$ . This is precisely the maximal correlation  $S(\mu, K)$ .  $\square$

**Remark 3.4** (Maximum correlation and the spectral gap). In the literature on Markov chains (see, e.g., [25, 46, 47]), one often sees the following definition: given a pair  $(\mu, K) \in \mathcal{P}(\mathsf{X}) \times \mathcal{M}(\mathsf{X}|\mathsf{X})$  such that  $\mu$  is invariant w.r.t.  $K$ , i.e.,  $\mu = \mu K$ , the (*absolute*) *spectral gap* of  $K$  is equal to

$$\gamma_*(\mu, K) \triangleq 1 - \sup_{f \in \mathcal{H}_0(\mathsf{X})} \frac{\|K^* f\|_\mu}{\|f\|_\mu}$$

(here we are using the Hilbert space notation from the proof above). Thus, the spectral gap and the maximal correlation are related by  $\gamma_*(\mu, K) = 1 - S(\mu, K) = 1 - \sqrt{\eta_{\chi^2}(\mu, K)}$ .  $\diamond$

The maximal correlation  $S^2(\mu, K)$  also provides a *lower bound* on the SDPI constants  $\eta_\Phi(\mu, K)$  for a certain subset of  $\mathcal{F}$ :

**Theorem 3.3.** *For any  $\Phi \in \mathcal{F}$  which is three times differentiable and has  $\Phi''(1) > 0$ , we have*

$$\eta_\Phi(\mu, K) \geq S^2(\mu, K), \tag{3.12}$$

$$\eta_\Phi(K) \geq S^2(K), \tag{3.13}$$

where  $S^2(K) \triangleq \sup_{\mu \in \mathcal{P}(\mathsf{X})} S^2(\mu, K)$ .

**Remark 3.5.** The second bound, Eq. (3.13), was proved by Cohen et al. [6], generalizing the results of Ahlswede and Gács [1] for  $\Phi(u) = u \log u$ . However, more or less the same proof technique also gives the distribution-dependent bound (3.12). A recent paper of Polyanskiy and Wu [21] presents an extension of Theorem 3.3 to abstract alphabets.

*Proof.* Without loss of generality, we assume that  $\Phi(1) = 0$ . Let us expand  $\Phi$  in a Taylor series around  $u = 1$ :

$$\begin{aligned} \Phi(u) &= \Phi(1) + \Phi'(1)(u-1) + \frac{1}{2}\Phi''(1)(u-1)^2 + o((u-1)^2) \\ &= \Phi'(1)(u-1) + \frac{1}{2}\Phi''(1)(u-1)^2 + o((u-1)^2), \end{aligned}$$

where the second step uses the fact that  $\Phi(1) = 0$ . Therefore, for any bounded real-valued random variable  $U$  and any  $\varepsilon > 0$  such that  $1 + \varepsilon U \geq 0$  a.s., we have

$$\text{Ent}_\Phi[1 + \varepsilon U] = \frac{\Phi''(1)}{2} \varepsilon^2 \text{Var}[U] + O(\varepsilon^3).$$

Now, fix an admissible pair  $(\mu, K)$ . For any  $\nu \neq \mu$ , consider the mixture  $\nu_\varepsilon \triangleq \bar{\varepsilon}\mu + \varepsilon\nu$ . Let  $f = d\nu/d\mu - 1$ . Then

$$\begin{aligned} D_\Phi(\nu_\varepsilon\|\mu) &= \text{Ent}_\Phi[1 + \varepsilon f(X)] \\ &= \frac{\Phi''(1)}{2}\varepsilon^2 \text{Var}[f(X)] + o(\varepsilon^2) \\ &= \frac{\Phi''(1)}{2}\varepsilon^2 \chi^2(\nu\|\mu) + o(\varepsilon^2) \end{aligned}$$

and

$$\begin{aligned} D_\Phi(\nu_\varepsilon K\|\mu K) &= \text{Ent}_\Phi[1 + \varepsilon K^* f(Y)] \\ &= \frac{\Phi''(1)}{2}\varepsilon^2 \text{Var}[K^* f(Y)] + o(\varepsilon^2) \\ &= \frac{\Phi''(1)}{2}\varepsilon^2 \chi^2(\nu K\|\mu K) + o(\varepsilon^2), \end{aligned}$$

where in the first step we have used Lemma A.1 in the Appendix and the linearity of  $K^*$ . Using the fact that  $\Phi''(1) > 0$ , for any  $\varepsilon > 0$  we have

$$\begin{aligned} \eta_\Phi(\mu, K) &\geq \sup_{\nu \neq \mu} \frac{D_\Phi(\nu_\varepsilon K\|\mu K)}{D_\Phi(\nu_\varepsilon\|\mu)} \\ &= \sup_{\nu \neq \mu} \frac{\chi^2(\nu K\|\mu K) + o(\varepsilon)}{\chi^2(\nu\|\mu) + o(\varepsilon)}. \end{aligned}$$

Taking the limit as  $\varepsilon \searrow 0$ , we get

$$\eta_\Phi(\mu, K) \geq \sup_{\nu \neq \mu} \frac{\chi^2(\nu K\|\mu K)}{\chi^2(\nu\|\mu)} = \eta_{\chi^2}(\mu, K).$$

This proves (3.12), and (3.13) follows after taking the supremum over all  $\mu$ .  $\square$

For example, the function  $\Phi(u) = u \log u$  that induces the usual relative entropy satisfies the conditions of Theorem 3.3, as does the function  $\Phi(u) = (\sqrt{u} - 1)^2$  that gives rise to the squared Hellinger distance.

Under additional regularity conditions on  $\Phi$ , we can obtain an upper bound on  $\eta_\Phi$  which is proportional to the maximal correlation  $S^2(\mu, K)$ :

**Theorem 3.4.** *Suppose that  $\Phi \in \mathcal{F}$  is twice differentiable, strictly convex, has a nonincreasing second derivative, and the function*

$$\Psi(u) \triangleq \frac{\Phi(u) - \Phi(0)}{u} \tag{3.14}$$

*is concave. Then, for any admissible pair  $(\mu, K)$ ,*

$$\eta_\Phi(\mu, K) \leq \frac{2\Psi'(1)}{\Phi''(1/\mu_*)} S^2(\mu, K), \tag{3.15}$$

*where  $\mu_* \triangleq \min_{x \in X} \mu(x)$  is the smallest mass of  $\mu$ .*

**Remark 3.6.** It can be shown (see, e.g., [31, Lm. 14.5]) that if  $\Phi \in \mathcal{F} \cap \mathcal{C}$ , where the function class  $\mathcal{C}$  is defined in Proposition 2.2, then the function  $\Psi$  defined in (3.14) is concave. For example, if  $\Phi(u) = u \log u$ , then  $\Psi(u) = \log u$ .  $\diamond$

*Proof.* Let  $(X, Y)$  be a random pair with law  $\mu \otimes K$ . Fix any probability distribution  $\nu \neq \mu$  and let  $f = d\nu/d\mu$ . Then we have the following chain of estimates:

$$\text{Ent}_\Phi[K^*f(Y)] \leq \Psi'(1) \text{Var}[K^*f(Y)] \quad (3.16)$$

$$\leq \Psi'(1) S^2(\mu, K) \text{Var}[f(X)] \quad (3.17)$$

$$\leq \frac{2\Psi'(1)}{\Phi''(\|f(X)\|_\infty)} S^2(\mu, K) \text{Ent}_\Phi[f(X)], \quad (3.18)$$

where (3.16) is by Lemma A.2 in Appendix A, (3.17) is by Theorem 3.2, and (3.18) is by Lemma A.3 in Appendix A. Now, since

$$\|f(X)\|_\infty = \left\| \frac{d\nu}{d\mu} \right\|_\infty \leq \frac{1}{\mu_*},$$

we have  $\Phi''(\|f(X)\|_\infty) \geq \Phi''(1/\mu_*)$ . By the arbitrariness of  $\nu$  (and hence  $f$ ), we obtain (3.15).  $\square$

For example, functions of the form  $\Phi_p(u) = \frac{u^p-1}{p-1}$  for  $1 < p \leq 2$  satisfy the conditions of the theorem with  $\Psi_p(u) = \frac{u^{p-1}}{p-1}$  and  $\Phi_p''(u) = pu^{p-2}$ . This gives the bound

$$\eta_{\Phi_p}(\mu, K) \leq \frac{2\mu_*^{p-2}}{p} S^2(\mu, K). \quad (3.19)$$

Note that  $\Phi_2(u) = u^2 - 1$  induces the  $\chi^2$ -divergence, so  $\eta_{\Phi_2}(\mu, K) = \eta_{\chi^2}(\mu, K) = S^2(\mu, K)$ , and in that case the bound (3.19) holds with equality. Moreover, as  $p \searrow 1$ , we have  $\text{Ent}_{\Phi_p}[U] \rightarrow \text{Ent}[U]$ , and in that limit (3.19) becomes

$$\eta(\mu, K) \leq \frac{2}{\mu_*} S^2(\mu, K). \quad (3.20)$$

Of course, the bound (3.19) is nontrivial only if  $S^2(\mu, K) < \frac{p}{2\mu_*^{p-2}}$ ; similarly, the bound (3.20) is nontrivial only if  $S^2(\mu, K) < \frac{\mu_*}{2}$ . As recently shown by Makur and Zheng [22], the constant 2 in (3.20) can be reduced to 1, but it is not clear how to extend their techniques to  $\Phi(u) \neq u \log u$ .

### 3.3 Upper bounds for operator convex $\Phi$

Theorem 3.4 gives an upper bound on the SDPI constant  $\eta_\Phi(\mu, K)$  in terms of the squared maximal correlation  $S^2(\mu, K)$ , but this bound has a multiplicative constant that depends on  $\mu$ . Given the lower bound of Theorem 3.3, it is natural to ask whether there is a matching upper bound without such a multiplicative constant. A partial result in this direction was obtained by Choi et al. [7], who showed that the equality  $\eta_\Phi(K) = \eta_{\chi^2}(K) = S^2(K)$  holds for all functions  $\Phi \in \mathcal{F}$  that are *operator convex* (see below for definitions). In this section, we will derive a *distribution-dependent* upper bound on  $\eta_\Phi(\mu, K)$  that implies the result of Choi et al.

In preparation for this result, we first need some facts from matrix analysis [45]. Let  $H_n$  denote the space of all  $n \times n$  Hermitian matrices, and let  $H_n(I)$  denote the subset of  $H_n$  consisting of all matrices whose eigenvalues lie in a given finite or infinite interval  $I$  of the real line. Any function  $\Phi : I \rightarrow \mathbb{R}$  can be extended to a matrix-valued function  $\Phi : H_n(I) \rightarrow H_n$  as follows:

- if  $A \in H_n(I)$  is diagonal, i.e.,  $A = \text{diag}(a_1, \dots, a_n)$  for some  $a_1, \dots, a_n \in I$ , then we let

$$\Phi(A) \triangleq \text{diag}(\Phi(a_1), \dots, \Phi(a_n)).$$

- if  $A \in H_n(I)$  can be diagonalized as  $A = U\Lambda U^*$ , where  $U$  is a unitary  $n \times n$  matrix and  $\Lambda \in H_n(I)$  is diagonal, then we let

$$\Phi(A) \triangleq U\Phi(\Lambda)U^*.$$

We introduce the following partial order on  $H_n$ : given any two  $A, B \in H_n$ , we write  $A \preceq B$  if  $B - A$  is positive semidefinite. We say that a function  $\Phi : I \rightarrow \mathbb{R}$  is *n-convex* if

$$\Phi(\lambda A + (1 - \lambda)B) \preceq \lambda\Phi(A) + (1 - \lambda)\Phi(B)$$

for all  $A, B \in H_n(I)$  and all  $\lambda \in [0, 1]$ . If  $\Phi$  is *n-convex* for all  $n \in \mathbb{N}$ , then we say that it is *operator convex*. By definition, any operator convex function is *a fortiori* convex in the ordinary sense, but the converse is generally not true. We are particularly interested in functions  $\Phi : \mathbb{R}^+ \rightarrow \mathbb{R}$  that are operator convex; here are some examples and counterexamples [45, Ch. V]:

- $\Phi(u) = u \log u$  is operator convex;
- $\Phi(u) = u^p$  is operator convex if and only if  $p \in [-1, 0] \cup [1, 2]$ .
- $\Phi(u) = -u^p$  is operator convex for  $0 \leq p \leq 1$ .

In general, it is not easy to determine whether a given function is operator convex. However, there is a deep result known as *Loewner's theorem* [48], which shows that operator convex functions possess very special integral representations:

**Theorem 3.5.** *A function  $\Phi : \mathbb{R}^+ \rightarrow \mathbb{R}$  with  $\Phi(0) = 0$  is operator convex if and only if there exist some constants  $\alpha \in \mathbb{R}, \beta \geq 0$  and a positive measure  $\nu$  on  $\mathbb{R}^+$  satisfying  $\int_0^\infty (1 + t^2)^{-1} \nu(dt) < \infty$ , such that*

$$\Phi(u) = \alpha u + \beta u^2 + \int_0^\infty \left( \frac{tu}{1+t^2} - \frac{u}{u+t} \right) \nu(dt). \quad (3.21)$$

For example, the operator convex function  $\Phi(u) = u \log u$  can be represented in the form (3.21) with  $\alpha = \beta = 0$  and with  $\nu$  given by the restriction of the Lebesgue measure to  $\mathbb{R}^+$  [45, Example V.4.18]; the operator convex function  $\Phi(u) = u^p$ ,  $1 < p < 2$ , can be represented in the form (3.21) with

$$\alpha = \cos \frac{\pi p}{2}, \quad \beta = 0, \quad \nu(dt) = \frac{\sin(\pi p)}{\pi} t^p dt$$

[45, Example V.4.19].

We also recall the definition of the Le Cam divergence with parameter  $\lambda \in (0, 1)$ , cf. Eq. (2.3):

$$\text{LC}_\lambda(\nu \parallel \mu) = \lambda \bar{\lambda} \mathbb{E}_\mu \left[ \frac{(\text{d}\nu/\text{d}\mu - 1)^2}{\lambda \text{d}\nu/\text{d}\mu + \bar{\lambda}} \right] = 1 - \mathbb{E}_\mu \left[ \frac{\text{d}\nu/\text{d}\mu}{\lambda \text{d}\nu/\text{d}\mu + \bar{\lambda}} \right],$$

which is a  $\Phi$ -divergence with

$$\Phi(u) = 1 - \frac{u}{\lambda u + \bar{\lambda}}.$$

Note that  $\text{LC}_0(\|\cdot\|) = \text{LC}_1(\|\cdot\|) = 0$ . For  $\lambda \in (0, 1)$ , consider the SDPI constant

$$\eta_{\text{LC}_\lambda}(\mu, K) = \sup_{\nu \neq \mu} \frac{\text{LC}_\lambda(\nu K \|\mu K)}{\text{LC}_\lambda(\nu \|\mu)}.$$

Now we are in a position to state our result:

**Theorem 3.6.** *Suppose that  $\Phi \in \mathcal{F}$  is operator convex. Then*

$$S^2(\mu, K) \leq \eta_\Phi(\mu, K) \leq \max \left( S^2(\mu, K), \sup_{0 < \lambda < 1} \eta_{\text{LC}_\lambda}(\mu, K) \right). \quad (3.22)$$

**Remark 3.7.** Since all explicit examples of functions in  $\mathcal{C}$  seem to be operator convex, it is tempting to think that all operator convex  $\Phi$  are elements of the function class  $\mathcal{C}$  (cf. Proposition 2.2). However, this is not the case. For example, the function  $\Phi(u) = (\sqrt{u} - 1)^2$ , which generates the Hellinger divergence, is operator convex. However,  $1/\Phi''(u) = 2u^{3/2}$  is not concave, so  $\Phi \notin \mathcal{C}$ .  $\diamond$

*Proof.* By Loewner's theorem (Theorem 3.5),  $\Phi$  admits the integral representation (3.21). Any  $\Phi$  that can be represented in this form is infinitely differentiable and strictly convex at  $u = 1$ . Therefore,  $\eta_\Phi(\mu, K) \geq S^2(\mu, K)$  by Theorem 3.3. This establishes first inequality in Eq. (3.22).

Now we prove the second inequality in (3.22). First, let us rewrite (3.21) as

$$\Phi(u) = \beta u^2 - \int_0^\infty \frac{u}{u+t} v(dt) + A(u),$$

where  $A(u)$  is an affine function. A change of variables  $\lambda = \frac{1}{t+1}$  gives

$$\Phi(u) = \beta u^2 - \int_0^1 \frac{\lambda u}{\lambda u + \bar{\lambda}} \Upsilon(d\lambda) + A(u), \quad (3.23)$$

where  $\Upsilon$  is some positive measure on  $[0, 1]$ . Since any two elements of  $\mathcal{F}$  that differ by an affine function determine the same divergence, Eq. (3.23) allows us to express  $D_\Phi(\nu \|\mu)$  as

$$D_\Phi(\nu \|\mu) = \beta \chi^2(\nu \|\mu) + \int_0^1 \lambda \text{LC}_\lambda(\nu \|\mu) \Upsilon(d\lambda)$$

The same holds for  $\nu K$  and  $\mu K$ , so

$$\begin{aligned} D_\Phi(\nu K \|\mu K) &= \beta \chi^2(\nu K \|\mu K) + \int_0^1 \lambda \text{LC}_\lambda(\nu K \|\mu K) \Upsilon(d\lambda) \\ &\leq \beta S^2(\mu, K) \chi^2(\nu \|\mu) + \int_0^1 \lambda \eta_{\text{LC}_\lambda}(\mu, K) \text{LC}_\lambda(\nu \|\mu) \Upsilon(d\lambda) \\ &\leq \max \left( S^2(\mu, K), \sup_{0 < \lambda < 1} \eta_{\text{LC}_\lambda}(\mu, K) \right) \cdot D_\Phi(\nu \|\mu). \end{aligned}$$

□

We can now recover the result of Choi et al. [7] as a corollary:

**Corollary 3.1.** *Suppose that  $\Phi \in \mathcal{F}$  is operator convex. Then*

$$\eta_{\Phi}(K) = S^2(K)$$

for any discrete channel  $K$ .

**Remark 3.8.** Since  $\Phi(u) = u \log u$  is operator convex, this is a broad generalization of a result of Ahlswede and Gács [1, Thm. 8]. It should be emphasized that Corollary 3.1 does not mean that  $\eta_{\Phi}(\mu, K) = S^2(\mu, K)$  for a given input distribution  $\mu \in \mathcal{P}_*(\mathcal{X})$ ; however, this may be the case for specific choices of  $\mu$  and  $K$ , as we show in the example after the proof.  $\diamond$

*Proof.* It suffices to show that

$$\sup_{0 < \lambda < 1} \eta_{\text{LC}_{\lambda}}(\mu, K) \leq S^2(K).$$

To that end, we first note that the Le Cam divergence  $\text{LC}_{\lambda}(\nu \parallel \mu)$  can be written as a convex combination of two  $\chi^2$ -divergences:

$$\text{LC}_{\lambda}(\nu \parallel \mu) = \lambda \chi^2(\nu \parallel \lambda \nu + \bar{\lambda} \mu) + \bar{\lambda} \chi^2(\mu \parallel \lambda \nu + \bar{\lambda} \mu).$$

From this, it follows that

$$\begin{aligned} \text{LC}_{\lambda}(\nu K \parallel \mu K) &= \lambda \chi^2(\nu K \parallel \lambda \nu K + \bar{\lambda} \mu K) + \bar{\lambda} \chi^2(\mu K \parallel \lambda \nu K + \bar{\lambda} \mu K) \\ &\leq S^2(K) [\lambda \chi^2(\nu \parallel \lambda \nu + \bar{\lambda} \mu) + \bar{\lambda} \chi^2(\mu \parallel \lambda \nu + \bar{\lambda} \mu)] \\ &= S^2(K) \text{LC}_{\lambda}(\nu \parallel \mu). \end{aligned}$$

□

**Example 3.1.** Let  $\mu = \text{Bern}(1/2)$  and  $K = \text{BSC}(\varepsilon)$ . For any  $q \neq 1/2$  and  $\nu = \text{Bern}(q)$ , we have  $\mu K = \text{Bern}(1/2)$  and  $\nu K = \text{Bern}(q \star \varepsilon)$ . Moreover,

$$\chi^2(\nu \parallel \mu) = \chi^2(\text{Bern}(q) \parallel \text{Bern}(1/2)) = (1 - 2q)^2$$

and

$$\chi^2(\nu K \parallel \mu K) = \chi^2(\text{Bern}(q \star \varepsilon) \parallel \text{Bern}(1/2)) = (1 - 2(q \star \varepsilon))^2 = (1 - 2\varepsilon)^2 (1 - 2q)^2.$$

Therefore,

$$\eta_{\chi^2}(\mu, K) \equiv S^2(\mu, K) = (1 - 2\varepsilon)^2. \quad (3.24)$$

Moreover, for any  $\lambda \in (0, 1)$ ,

$$\begin{aligned} \text{LC}_{\lambda}(\nu \parallel \mu) &= \text{LC}_{\lambda}(\text{Bern}(q) \parallel \text{Bern}(1/2)) \\ &= 1 - \frac{1}{2} \left( \frac{2q}{2\lambda q + \bar{\lambda}} + \frac{2\bar{q}}{2\lambda \bar{q} + \bar{\lambda}} \right) \\ &= \frac{\lambda \bar{\lambda} (1 - 2q)^2}{1 - \lambda^2 (1 - 2q)^2} \end{aligned}$$

and

$$\begin{aligned} \text{LC}_\lambda(\nu K \|\mu K) &= \text{LC}_\lambda(\text{Bern}(q \star \varepsilon) \|\text{Bern}(1/2)) \\ &= \frac{\lambda \bar{\lambda} (1 - 2(q \star \varepsilon))^2}{1 - \lambda^2 (1 - 2(q \star \varepsilon))^2} \\ &= \frac{\lambda \bar{\lambda} (1 - 2\varepsilon)^2 (1 - 2q)^2}{1 - \lambda^2 (1 - 2\varepsilon)^2 (1 - 2q)^2}. \end{aligned}$$

Both of these divergences are invariant with respect to the transformation  $q \mapsto 1 - q$ , so

$$\eta_{\text{LC}_\lambda}(\mu, K) = \sup_{0 \leq q < 1/2} \frac{(1 - 2\varepsilon)^2 (1 - \lambda^2 (1 - 2q)^2)}{1 - (1 - 2\varepsilon)^2 \lambda^2 (1 - 2q)^2} = (1 - 2\varepsilon)^2, \quad (3.25)$$

where the supremum is achieved at  $q = 1/2$  (but not at any  $q \neq 1/2$ ). (As an aside, it is not hard to show that the expression under the supremum in (3.25) is a concave function of  $q$ .) Comparing Eqs. (3.24) and (3.25), we see that

$$\eta_{\chi^2}(\mu, K) = \sup_{0 < \lambda < 1} \eta_{\text{LC}_\lambda}(\mu, K) = (1 - 2\varepsilon)^2.$$

Therefore, by Theorem 3.6,

$$\eta_\Phi(\text{Bern}(1/2), \text{BSC}(\varepsilon)) = (1 - 2\varepsilon)^2$$

for all operator convex  $\Phi \in \mathcal{F}$ .

### 3.4 Upper bounds via subgaussian concentration and information-transportation inequalities

Fix an admissible pair  $(\mu, K) \in \mathcal{P}_*(\mathsf{X}) \times \mathcal{M}(\mathsf{Y}|\mathsf{X})$ , and let  $(X, Y)$  be a random pair with probability law  $\mu \otimes K$ . We expect the SDPI constant  $\eta(\mu, K)$  to be small if the channel output  $Y$  of  $K$  is nearly independent of the channel input  $X \sim \mu$ . In this section, we present upper bounds on  $\eta(\mu, K)$  that capture this intuition in terms of the properties of the posterior likelihood ratio

$$a(x, y) \triangleq \frac{dP_{X|Y=y}}{dP_X}(x) = \frac{K^*(x|y)}{\mu(x)}. \quad (3.26)$$

Theorems 3.7 and 3.8 quantify near-independence by looking at how tightly the random variable  $a(X, y)$  concentrates around its expected value 1 for each fixed  $y$ . Moreover, Theorem 3.8 shows a connection between SDPI for the relative entropy and *information-transportation inequalities* introduced in the pioneering work of Marton [49, 50].

First, we collect some preliminaries. A real-valued random variable  $U$  is called *subgaussian with parameter  $v$*  (or  $v$ -subgaussian) if  $\mathbb{E}[e^{t(U - \mathbb{E}U)}] \leq e^{vt^2/2}$  for all  $t \in \mathbb{R}$  [31, Sec. 2.3]. For any  $v$ -subgaussian random variable  $U$  we have the tail estimate

$$\mathbb{P}(|U - \mathbb{E}U| \geq t) \leq 2e^{-t^2/2v}, \quad \forall t \in \mathbb{R}.$$

To get the tightest such bound, we define the *subgaussian constant*

$$\sigma^2(U) \triangleq \inf \left\{ v \geq 0 : \mathbb{E}[e^{t(U - \mathbb{E}U)}] \leq e^{vt^2/2}, t \in \mathbb{R} \right\}.$$

With these definitions in place, we have the following theorem:

**Theorem 3.7.** For each  $y \in \mathsf{Y}$ , let  $\sigma^2(y) \triangleq \sigma^2(a(X, y))$ . Then

$$\eta(\mu, K) \leq 2 \mathbb{E}[\sigma^2(Y)]. \quad (3.27)$$

*Proof.* Fix any  $\nu \in \mathcal{P}(\mathsf{X})$  and let  $f = d\nu/d\mu$ . Observe that  $\mathbb{E}[a(X, y)] = \mathbb{E}[f(X)] = 1$ . Then

$$\begin{aligned} D(\nu K \| \mu K) &= \text{Ent}[K^* f(Y)] \\ &\leq \text{Var}[K^* f(Y)] \\ &= \sum_{y \in \mathsf{Y}} \mu K(y) (K^* f(y) - 1)^2 \\ &= \sum_{y \in \mathsf{Y}} \mu K(y) \left| \sum_{x \in \mathsf{X}} \mu(x) [a(x, y) f(x) - 1] \right|^2 \\ &= \sum_{y \in \mathsf{Y}} \mu K(y) |\text{Cov}(a(X, y), f(X))|^2, \end{aligned} \quad (3.28)$$

where the inequality is by Lemma A.2 in Appendix A. Next, we make use of the fact that

$$\text{Ent}[U] \geq \mathbb{E}[UZ] - \mathbb{E}[U] \log \mathbb{E}[e^Z] \quad (3.29)$$

for any random variable  $Z$  jointly distributed with  $U$  and satisfying  $\mathbb{E}[e^Z] < \infty$  (see, e.g., [31, Thm. 4.13]; in fact, this bound holds with equality for  $Z = \log U$ ). If we fix an arbitrary  $y \in \mathsf{Y}$  and then use (3.29) with  $U = f(X)$  and  $Z = \pm t(a(X, y) - 1)$  for some  $t > 0$ , we get

$$\begin{aligned} |\text{Cov}(a(X, y), f(X))| &\leq \frac{1}{t} \left( \log \mathbb{E}[e^{t(a(X, y) - 1)}] + \text{Ent}[f(X)] \right) \\ &\leq \frac{\sigma^2(y)t}{2} + \frac{\text{Ent}[f(X)]}{t}. \end{aligned}$$

Since this holds for an arbitrary  $t$ , we have

$$|\text{Cov}(a(X, y), f(X))| \leq \inf_{t > 0} \left\{ \frac{\sigma^2(y)t}{2} + \frac{\text{Ent}[f(X)]}{t} \right\} = \sqrt{2\sigma^2(y) \text{Ent}[f(X)]}.$$

Using this estimate in (3.28), we get (3.27).  $\square$

In order to apply Theorem 3.7, we need to compute or upper-bound the subgaussian constant  $\sigma^2(a(X, y))$  for each  $y \in \mathsf{Y}$ . In some situations, it is possible to derive exact expressions for subgaussian constants (as we show in the examples below); when the function  $x \mapsto a(x, y)$  is Lipschitz for each  $y \in \mathsf{Y}$ , one can derive upper bounds using *information-transportation inequalities* introduced in the pioneering work of Marton [49, 50] (see, e.g., the text of Villani [51]). If we endow the input alphabet  $\mathsf{X}$  with a metric  $d$ , then we can define the  $L^1$  Wasserstein distance (or *optimal transportation distance*) on  $\mathcal{P}(\mathsf{X})$  by

$$W_1(\mu, \nu) \triangleq \inf \{ \mathbb{E}[d(X, \bar{X})] : P_{X\bar{X}} \in \mathcal{P}(\mathsf{X} \times \mathsf{X}), P_X = \mu, P_{\bar{X}} = \nu \}$$

For example, for the trivial metric  $d(x, x') = \mathbf{1}\{x \neq x'\}$  we recover the total variation distance:  $W_1(\mu, \nu) = \|\mu - \nu\|_{\text{TV}}$ . Given a function  $f : \mathsf{X} \rightarrow \mathbb{R}$ , denote by

$$\delta(f) \triangleq \sup_{\substack{x, x' \in \mathsf{X} \\ x \neq x'}} \frac{|f(x) - f(x')|}{d(x, x')}$$

the *oscillation* (or the *Lipschitz norm*) of  $f$  w.r.t. the metric  $d$ .

**Theorem 3.8.** Fix an admissible pair  $(\mu, K) \in \mathcal{P}_*(\mathsf{X}) \times \mathcal{M}(\mathsf{Y}|\mathsf{X})$ . Suppose that  $\mu$  satisfies an information-transportation inequality with constant  $c > 0$ , i.e.,

$$W_1(\nu, \mu) \leq \sqrt{2c D(\nu\|\mu)}, \quad \forall \nu \neq \mu. \quad (3.30)$$

Then

$$\eta(\mu, K) \leq 2c \mathbb{E} [\delta^2(a(\cdot, Y))]. \quad (3.31)$$

*Proof.* By a result of Bobkov and Götze [52], a probability measure  $\mu \in \mathcal{P}(\mathsf{X})$  satisfies (3.30) if and only if

$$\mathbb{E}_\mu \left[ e^{t(f(X) - \mathbb{E}f(X))} \right] \leq e^{\frac{ct^2}{2}}, \quad t \in \mathbb{R} \quad (3.32)$$

for every  $f \in \mathcal{F}(\mathsf{X})$  with  $\delta(f) \leq 1$ . In particular, if (3.30) holds, then, by rescaling (3.32), we get

$$\mathbb{E}_\mu \left[ e^{t(a(X, y) - 1)} \right] \leq \exp \left( \frac{c\delta^2(a(\cdot, y)) t^2}{2} \right).$$

This implies that  $\sigma^2(a(X, y)) \leq c\delta^2(a(\cdot, y))$ . Substituting this into (3.27), we get (3.31).  $\square$

**Example 3.2** (Binary symmetric channels with asymmetric inputs). Let  $\mathsf{X} = \mathsf{Y} = \{0, 1\}$ ,  $\mu = \text{Bern}(p)$ ,  $K = \text{BSC}(\varepsilon)$ . We take the trivial metric  $d(x, x') = \mathbf{1}\{x \neq x'\}$ . In this case, Theorems 3.7 and 3.8 give the same bound. Indeed, by a result of Ordentlich and Weinberger [53],  $\mu = \text{Bern}(p)$  satisfies an information-transportation inequality

$$\|\nu - \mu\|_{\text{TV}} \leq \sqrt{2c(p) D(\nu\|\mu)}, \quad c(p) \triangleq \frac{p - \bar{p}}{2(\log p - \log \bar{p})}, \quad (3.33)$$

and the constant in front of the relative entropy is optimal, i.e.,

$$\inf_{\nu} \frac{D(\nu\|\mu)}{\|\nu - \mu\|_{\text{TV}}^2} = \frac{1}{2c(p)}.$$

[The inequality (3.33) is a distribution-dependent refinement of Pinsker's inequality, where we fix  $\mu$  and vary only  $\nu$ .] A simple calculation gives

$$\delta(a(\cdot, 0)) = \left| \frac{1 - 2\varepsilon}{1 - \varepsilon \star p} \right|, \quad \delta(a(\cdot, 1)) = \left| \frac{1 - 2\varepsilon}{\varepsilon \star p} \right|,$$

where  $\varepsilon \star p = \varepsilon \bar{p} + \bar{\varepsilon} p$ . Therefore, applying Theorem 3.8, we get the bound

$$\eta(\text{Bern}(p), \text{BSC}(\varepsilon)) \leq \frac{2c(p)(1 - 2\varepsilon)^2}{(1 - \varepsilon \star p)(\varepsilon \star p)}, \quad (3.34)$$

This bound is, unfortunately, loose. Indeed, if we take the limit  $p \searrow 1/2$ , then we get

$$\eta(\text{Bern}(1/2), \text{BSC}(\varepsilon)) \leq 2(1 - 2\varepsilon)^2. \quad (3.35)$$

which is off by a factor of 2, but still tighter than the Dobrushin contraction bound  $|1 - 2\varepsilon|$  (Theorem 3.1) in the range  $1/4 < \varepsilon < 3/4$ . Figure 1 shows a plot of the maximum value of the

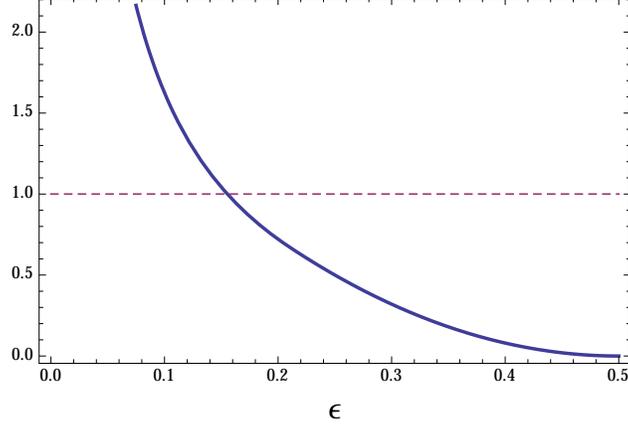


Figure 1: Maximum value of the right-hand side of (3.34) over  $p \in [0, 1]$  for each fixed  $\varepsilon$ .

right-hand side of (3.34) over  $p$  for each fixed value of the crossover probability  $\varepsilon$ ; from this, we see that the bound is nontrivial (i.e., takes values strictly smaller than 1) for  $\varepsilon \gtrsim 0.156$ .

In order to apply Theorem 3.7, we need to know the subgaussian constants of  $a(X, y)$ ,  $y \in \{0, 1\}$ . By a result of Bobkov et al. [54], for any function  $f : \{0, 1\} \rightarrow \mathbb{R}$  and for  $X \sim \text{Bern}(p)$  we have

$$2\sigma^2(f(X)) = 2c(p) |f(0) - f(1)|^2. \quad (3.36)$$

Applying (3.36) to  $f = a(\cdot, 0)$  and  $a(\cdot, 1)$ , we get

$$2\sigma^2(0) = 2c(p) \left| \frac{1 - 2\varepsilon}{1 - \varepsilon \star p} \right|^2, \quad 2\sigma^2(1) = 2c(p) \left| \frac{1 - 2\varepsilon}{\varepsilon \star p} \right|^2,$$

and indeed Theorem 3.7 gives the same bound (3.34).

**Example 3.3** (Binary input channels). Let  $\mathsf{X} = \{0, 1\}$  with  $\mu = \text{Bern}(p)$ , and consider an arbitrary channel  $K \in \mathcal{M}(\mathsf{Y}|\mathsf{X})$  with a finite (not necessarily binary) output alphabet  $\mathsf{Y}$ . Then

$$a(x, y) = \frac{K^*(x|y)}{\mu(x)} = \frac{K(y|x)}{\mu K(y)},$$

where  $\mu K(y) = \bar{p}K(y|0) + pK(y|1)$ . If we again take  $d$  to be the trivial metric, then the same analysis as in the previous example can be used to show that

$$2\sigma^2(y) = 2c(p) \frac{|K(y|0) - K(y|1)|^2}{\mu K(y)^2},$$

and Theorem 3.7 gives the bound

$$\eta(\text{Bern}(p), K) \leq 2c(p) \sum_{y \in \mathsf{Y}} \frac{|K(y|0) - K(y|1)|^2}{\bar{p}K(y|0) + pK(y|1)}.$$

**Example 3.4** (Random walk on a graph). Consider a connected undirected graph  $G = (\mathsf{V}, \mathsf{E})$  without self-loops or multiple edges, and let  $\mathsf{X} = \mathsf{Y} = \mathsf{V}$ . If the vertices  $x$  and  $y$  are connected by

an edge, we shall write  $x \leftrightarrow y$ ; the degree of a vertex  $x$  is defined as  $\deg_G(x) \triangleq |\{y \in \mathbf{V} : x \leftrightarrow y\}|$ . Define a probability measure  $\mu = \mu_G \in \mathcal{P}(\mathbf{V})$  by

$$\mu_G(x) \triangleq \frac{\deg_G(x)}{2|\mathbf{E}|}, \quad x \in \mathbf{V}.$$

Fix a parameter  $\varepsilon \in (0, 1)$ , and consider a channel  $K_G^{(\varepsilon)}$  with

$$K_G^{(\varepsilon)}(y|x) = \begin{cases} \bar{\varepsilon}, & \text{if } x = y \\ \frac{\varepsilon}{\deg_G(x)}, & \text{if } x \leftrightarrow y \\ 0, & \text{otherwise} \end{cases}. \quad (3.37)$$

Again, let  $d$  be the trivial metric,  $d(x, x') = \mathbf{1}\{x \neq x'\}$ . Then  $W_1(\nu, \mu) = \|\nu - \mu\|_{\text{TV}}$ , and we can take  $c = 1/4$  in (3.30), which is then just Pinsker's inequality. It is not hard to show that  $K_G^{(\varepsilon)}$  is *reversible* w.r.t.  $\mu_G$ , i.e.,

$$\mu_G(x)K_G^{(\varepsilon)}(y|x) = \mu_G(y)K_G^{(\varepsilon)}(x|y), \quad \forall x, y \in \mathbf{V}.$$

Therefore,  $\mu_G K_G^{(\varepsilon)} = \mu_G$ , so the posterior likelihood ratio is given by

$$a(x, y) = \frac{K_G^{(\varepsilon)}(y|x)}{\mu_G(y)} = \frac{2|\mathbf{E}|}{\deg_G(y)} K_G^{(\varepsilon)}(y|x).$$

Now, from the definition (3.37) of  $K_G^{(\varepsilon)}$  it follows that

$$\left| K_G^{(\varepsilon)}(y|x) - K_G^{(\varepsilon)}(y|x') \right| = \begin{cases} \left| \bar{\varepsilon} - \frac{\varepsilon}{\deg_G(x')} \mathbf{1}\{x' \leftrightarrow y\} \right|, & \text{if } x = y \\ \left| \bar{\varepsilon} - \frac{\varepsilon}{\deg_G(x)} \mathbf{1}\{x \leftrightarrow y\} \right|, & \text{if } x' = y \\ \frac{\varepsilon}{\deg_G(x)}, & \text{if } x \leftrightarrow y, x' \not\leftrightarrow y \\ \frac{\varepsilon}{\deg_G(x')}, & \text{if } x \not\leftrightarrow y, x' \leftrightarrow y \\ \varepsilon \left| \frac{1}{\deg_G(x)} - \frac{1}{\deg_G(x')} \right|, & \text{if } x \leftrightarrow y, x' \leftrightarrow y \\ 0, & \text{if } x \not\leftrightarrow y, x' \not\leftrightarrow y \end{cases}$$

where  $x \not\leftrightarrow y$  means that  $x$  and  $y$  are not connected by an edge and that  $x \neq y$ . Therefore,

$$\begin{aligned} \delta^2(a(\cdot, y)) &= \frac{4|\mathbf{E}|^2}{\deg_G(y)} \max_{x, x' \in \mathbf{V}} \left| K_G^{(\varepsilon)}(y|x) - K_G^{(\varepsilon)}(y|x') \right|^2 \\ &= \frac{4|\mathbf{E}|^2}{\deg_G(y)^2} \left( \Delta_0(y, \varepsilon) \vee \Delta_1(y, \varepsilon) \vee \Delta_2(y, \varepsilon) \right), \end{aligned}$$

where

$$\Delta_0(y, \varepsilon) \triangleq \max_{x \in \mathbf{V} \setminus \{y\}} \left( \frac{\varepsilon}{\deg_G(x)} \right)^2 \mathbf{1}\{\deg_G(y) < |\mathbf{V}| - 1\} \quad (3.38a)$$

$$\Delta_1(y, \varepsilon) \triangleq \max_{x \in \mathbf{V} \setminus \{y\}} \left| \bar{\varepsilon} - \frac{\varepsilon}{\deg_G(x)} \mathbf{1}\{x \leftrightarrow y\} \right|^2 \quad (3.38b)$$

$$\Delta_2(y, \varepsilon) \triangleq \max_{x, x' \in \mathbf{V} \setminus \{y\}} \varepsilon^2 \left| \frac{1}{\deg_G(x)} - \frac{1}{\deg_G(x')} \right|^2 \mathbf{1}\{x \leftrightarrow y, x' \leftrightarrow y\}. \quad (3.38c)$$

Theorem 3.8 then gives the bound

$$\eta\left(\mu_G, K_G^{(\varepsilon)}\right) \leq |\mathbb{E}| \sum_{y \in \mathbb{V}} \frac{\Delta_0(y, \varepsilon) \vee \Delta_1(y, \varepsilon) \vee \Delta_2(y, \varepsilon)}{\deg_G(y)} \quad (3.39)$$

(note that  $\deg_G(y) > 0$  for each  $y$ , since  $G$  is connected).

For example, if  $G$  is a complete graph, then  $\Delta_0(y, \varepsilon) = \Delta_2(y, \varepsilon) = 0$  for all  $y$ , while

$$\Delta_1(y, \varepsilon) = \left(1 - \frac{|\mathbb{V}|}{|\mathbb{V}| - 1} \varepsilon\right)^2, \quad y \in \mathbb{V}$$

so we get the bound

$$\eta\left(\mu_G, K_G^{(\varepsilon)}\right) \leq \frac{|\mathbb{V}|^2}{2} \left(1 - \frac{|\mathbb{V}|}{|\mathbb{V}| - 1} \varepsilon\right)^2, \quad (3.40)$$

which is nontrivial (i.e., strictly smaller than unity) in the range

$$\frac{(|\mathbb{V}| - 1)(|\mathbb{V}| - \sqrt{2})}{|\mathbb{V}|^2} < \varepsilon < \frac{(|\mathbb{V}| - 1)(|\mathbb{V}| + \sqrt{2})}{|\mathbb{V}|^2}.$$

For the complete graph on the two-point set  $\mathbb{V} = \{0, 1\}$ , the channel  $K_G^{(\varepsilon)}$  is just BSC( $\varepsilon$ ), and the bound (3.40) reduces to (3.35).

As another example, let  $G$  be the path graph on the ternary vertex set  $\mathbb{V} = \{0, 1, 2\}$ , i.e.,  $\mathbb{E} = \{\{0, 1\}, \{1, 2\}\}$ . Then  $\mu_G(0) = \mu_G(2) = 1/4$  and  $\mu_G(1) = 1/2$ . From (3.38), we get

$$\begin{aligned} \Delta_0(y, \varepsilon) &= \begin{cases} \varepsilon^2, & y \in \{0, 2\} \\ 0, & y = 1 \end{cases} \\ \Delta_1(y, \varepsilon) &= \begin{cases} \left(1 - \frac{3\varepsilon}{2}\right)^2, & y \in \{0, 2\} \\ (1 - 2\varepsilon)^2, & y = 1 \end{cases} \\ \Delta_2(y, \varepsilon) &= 0, \quad y \in \{0, 1, 2\}. \end{aligned}$$

Substituting this into (3.39), we get

$$\eta\left(\mu_G, K_G^{(\varepsilon)}\right) \leq \begin{cases} 13\varepsilon^2 - 16\varepsilon + 5, & 0 \leq \varepsilon \leq 0.4 \\ 8\varepsilon^2 - 4\varepsilon + 1, & 0.4 \leq \varepsilon \leq 1 \end{cases}. \quad (3.41)$$

This bound, plotted in Figure 2, is nontrivial only in the range  $\frac{8-2\sqrt{3}}{13} < \varepsilon < \frac{1}{2}$ .

**Example 3.5** (General discrete channel). Consider arbitrary finite alphabets  $\mathbb{X}$  and  $\mathbb{Y}$ , together with an admissible pair  $(\mu, K) \in \mathcal{P}_*(\mathbb{X}) \times \mathcal{M}(\mathbb{Y}|\mathbb{X})$ . If we endow  $\mathbb{X}$  with the trivial metric  $d(x, x') = \mathbf{1}\{x \neq x'\}$ , then  $\mu$  will satisfy the information-transportation inequality (3.30) for  $W_1(\nu, \mu) = \|\nu - \mu\|_{\text{TV}}$  with optimal ( $\mu$ -dependent) constant  $c(\beta_\mu)$ , where the function  $c(\cdot)$  is defined in (3.33), and

$$\beta_\mu \triangleq \min \{\mu(\mathbb{A}) : \mathbb{A} \subseteq \mathbb{X}, \mu(\mathbb{A}) \geq 1/2\}$$

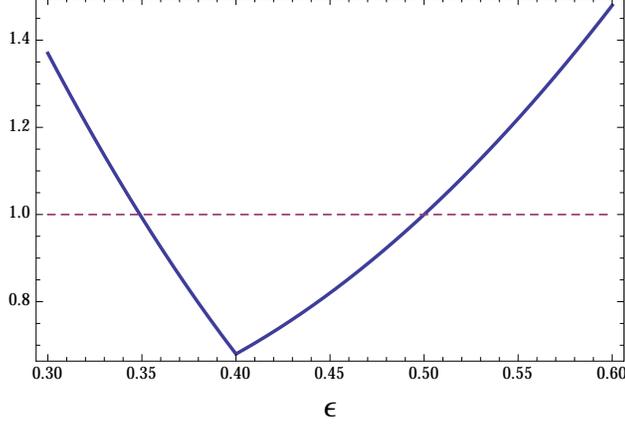


Figure 2: The bound of Eq. (3.41) as a function of the noise parameter  $\varepsilon$ .

is a measure of “imbalance” of  $\mu$  — in particular, when  $\mu$  is the uniform distribution on  $\mathsf{X}$  and  $|\mathsf{X}|$  is even,  $\beta_\mu = 1/2$ . Again, this is just the distribution-dependent refinement of Pinsker’s inequality [53]. Then

$$\begin{aligned} \delta^2(a(\cdot, y)) &= \max_{x, x' \in \mathsf{X}} \left| \frac{K^*(x|y)}{\mu(x)} - \frac{K^*(x'|y)}{\mu(x')} \right|^2 \\ &= \frac{1}{\mu K(y)^2} \max_{x, x' \in \mathsf{X}} |K(y|x) - K(y|x')|^2 \\ &= \frac{1}{\mu K(y)^2} \delta^2(K(y|\cdot)), \end{aligned}$$

so Theorem 3.8 gives the bound

$$\eta(\mu, K) \leq 2c(\beta_\mu) \sum_{y \in \mathsf{Y}} \frac{\delta^2(K(y|\cdot))}{\mu K(y)}. \quad (3.42)$$

In general, the bounds of Theorems 3.7 and 3.8 are nontrivial only for channels that are “sufficiently noisy,” in the sense that the posterior likelihood ratio (3.26) is nearly constant as a function of the input symbol  $x$  for any fixed output symbol  $y$ . In particular, the function  $x \mapsto a(x, y)$  is constant for each  $y \in \mathsf{Y}$  if and only if the output of  $K$  is independent of the input, i.e., if  $\eta(\mu, K) = 0$ . However, these bounds may be useful for capturing the *scaling* of the SDPI constant  $\eta(\mu, K)$  with various parameters of the problem. To the best of our knowledge, the first bound on  $\eta(\mu, K)$  in terms of a certain concentration property of the posterior likelihood ratio is due to Weitz [55] (see also [56]), and can be stated in our notation as follows:

$$\eta(\mu, K) \leq \left( \frac{c}{(\mu K)_*} \right)^2 \mathbb{E}[\tau(Y)], \quad (3.43)$$

where  $c > 0$  is some numerical constant,  $(\mu K)_* = \min_{y \in \mathsf{Y}} \mu K(y)$ , and

$$\tau(y) \triangleq \inf \left\{ t \geq 0 : \mathbb{P}(|a(X, y) - 1| > t) \leq e^{-2/t} \right\}.$$

Since the function  $t \mapsto e^{-2/t}$  is increasing, converges to 1 as  $t \rightarrow \infty$ , and to 0 as  $t \searrow 0$ , the quantity  $\tau(y)$  should be very close to zero for the bound (3.43) to be nontrivial. In contrast to the bounds of Theorems 3.7 and 3.8, which involve quantities pertaining to *large deviations* of  $a(X, y)$  from its mean, Weitz's bound is in terms of a quantity that has to do with *small deviations* of  $a(X, y)$  from its mean.

### 3.5 Tensorization

So far, we have considered the case of a single channel. However, many problems in information theory involve multiple uses of the same channel (or, more generally, transmission of correlated data over a memoryless channel with time-varying transition probabilities). In this context, it is of interest to determine whether the resulting “super-channel” inherits any SDPI-type behavior from the constituent channels.

In precise terms, let  $(\mu_1, K_1), \dots, (\mu_n, K_n)$  be  $n$  admissible pairs, where, for each  $i$ ,  $\mu_i \in \mathcal{P}_*(\mathsf{X}_i)$  and  $K_i \in \mathcal{M}(\mathsf{Y}_i|\mathsf{X}_i)$  for some alphabets  $\mathsf{X}_i, \mathsf{Y}_i$ . Fix some  $\Phi \in \mathcal{F}$ , a product distribution  $\mu = \mu_1 \otimes \dots \otimes \mu_n \in \mathcal{P}_*(\mathsf{X}_1 \times \dots \times \mathsf{X}_n)$ , and a product channel  $K = K_1 \otimes \dots \otimes K_n \in \mathcal{M}(\mathsf{Y}_1 \times \dots \times \mathsf{Y}_n|\mathsf{X}_1 \times \dots \times \mathsf{X}_n)$ . We say that the SDPI constant  $\eta_\Phi(\mu, K)$  *tensorizes* if

$$\eta_\Phi(\mu, K) = \max_{1 \leq i \leq n} \eta_\Phi(\mu_i, K_i).$$

For instance, Witsenhausen [2] showed that  $\eta_{\chi^2}(\mu, K)$  tensorizes, while a recent paper by Anantharam et al. [15] presents two different proofs of the tensorization property of  $\eta(\mu, K)$ . In each case, the proof relies on specific properties of the underlying  $\Phi$  — Witsenhausen exploits the connection between  $\eta_{\chi^2}(\mu, K)$  and the eigenvalues of the linear operator  $KK^* : L^2(\mathsf{X}, \mu) \rightarrow L^2(\mathsf{X}, \mu)$ , whereas Anantharam et al. use the chain rule for the relative entropy. The question is, can one give a unified proof of tensorization for a broader class of functions  $\Phi \in \mathcal{F}$  that contains both  $\Phi(u) = (u - 1)^2$  and  $\Phi(u) = u \log u$ ? As we show next, the answer is ‘yes’ for all functions  $\Phi$  whose  $\Phi$ -entropies are subadditive and homogeneous in the sense of Definition 3.2.

**Theorem 3.9** (Tensorization). *Suppose that  $\Phi \in \mathcal{F}$  induces a subadditive and homogeneous  $\Phi$ -entropy. Consider any  $n$  admissible pairs  $(\mu_i, K_i) \in \mathcal{P}_*(\mathsf{X}_i) \times \mathcal{M}(\mathsf{Y}_i|\mathsf{X}_i)$ . Then*

$$\eta_\Phi(\mu_1 \otimes \dots \otimes \mu_n, K_1 \otimes \dots \otimes K_n) = \max_{1 \leq i \leq n} \eta_\Phi(\mu_i, K_i).$$

*Proof.* For the sake of brevity, let  $\eta = \eta_\Phi(\mu_1 \otimes \dots \otimes \mu_n, K_1 \otimes \dots \otimes K_n)$ ,  $\eta_i = \eta_\Phi(\mu_i, K_i)$ ,  $\mu = \mu_1 \otimes \dots \otimes \mu_n$ , and  $K = K_1 \otimes \dots \otimes K_n$ .

To show that  $\eta \geq \eta_i$  for all  $i$ , take  $\nu \in \mathcal{P}(\mathsf{X}_1 \times \dots \times \mathsf{X}_n)$  of the form  $\mu_1 \otimes \dots \otimes \mu_{i-1} \otimes \nu_i \otimes \mu_{i+1} \otimes \dots \otimes \mu_n$  for some  $\nu_i \in \mathcal{P}(\mathsf{X}_i) \setminus \{\mu_i\}$ . Then

$$\begin{aligned} D_\Phi(\nu|\mu) &= D_\Phi(\nu_i|\mu_i), \\ D_\Phi(\nu K|\mu K) &= D_\Phi(\nu_i K_i|\mu_i K_i). \end{aligned}$$

Taking the supremum of  $\frac{D_\Phi(\nu K|\mu K)}{D_\Phi(\nu|\mu)}$  over all such  $\nu$ , we conclude that  $\eta \geq \eta_i$ .

For the reverse inequality  $\eta \leq \max_{1 \leq i \leq n} \eta_i$ , it suffices to consider the case  $n = 2$ ; the general case will follow by induction. Thus, let us fix two admissible pairs  $(\nu_i, K_i) \in \mathcal{P}_*(\mathsf{X}_i) \times \mathcal{M}(\mathsf{Y}_i|\mathsf{X}_i)$ ,  $i = 1, 2$ ,

and an arbitrary nonconstant function  $f \in \mathcal{F}_*^0(\mathbf{X}_1 \times \mathbf{X}_2)$ . Let  $(X_1, X_2, Y_1, Y_2) \in \mathbf{X}_1 \times \mathbf{X}_2 \times \mathbf{Y}_1 \times \mathbf{Y}_2$  be a random tuple, such that

$$P_{X_1 X_2} = \mu_1 \otimes \mu_2, \quad P_{Y_1 Y_2 | X_1 X_2} = K_1 \otimes K_2.$$

Then, from (2.5),

$$\text{Ent}_\Phi [K^* f(Y_1, Y_2)] = \mathbb{E} [\text{Ent}_\Phi [K^* f(Y_1, Y_2) | Y_1]] + \text{Ent}_\Phi [\mathbb{E}[K^* f(Y_1, Y_2) | Y_1]].$$

Define the functions  $f_1 \in \mathcal{F}_*^0(\mathbf{Y}_1 \times \mathbf{X}_2)$  and  $f_2 \in \mathcal{F}_*^0(\mathbf{X}_1 \times \mathbf{Y}_2)$  by

$$\begin{aligned} f_1(y_1, x_2) &= \sum_{x_1 \in \mathbf{X}_1} P_{X_1 | Y_1}(x_1 | y_1) f(x_1, x_2) \\ f_2(x_1, y_2) &= \sum_{x_2 \in \mathbf{X}_2} P_{X_2 | Y_2}(x_2 | y_2) f(x_1, x_2), \end{aligned}$$

which can be written more succinctly as  $f_1 = (K_1^* \otimes \text{id}_2)f$  and  $f_2 = (\text{id}_1 \otimes K_2^*)f$ , where  $\text{id}_1$  and  $\text{id}_2$  are the identity mappings on  $\mathcal{F}(\mathbf{X}_1)$  and  $\mathcal{F}(\mathbf{X}_2)$ . Since  $Y_1$  and  $Y_2$  are independent, we can write

$$\begin{aligned} \text{Ent}_\Phi [K^* f(Y_1, Y_2) | Y_1 = y_1] &= \mathbb{E}[\Phi(K_2^* f_1(y_1, Y_2))] - \Phi(\mathbb{E}[K_2^* f_1(y_1, Y_2)]) \\ &= \text{Ent}_\Phi [K_2^* f_1(y_1, Y_2)] \\ &\leq \eta_2 \text{Ent}_\Phi [f_1(y_1, X_2)] \\ &\leq \eta_2 \sum_{x_1 \in \mathbf{X}_1} P_{X_1 | Y_1}(x_1 | y_1) \text{Ent}_\Phi [f(x_1, X_2)], \end{aligned}$$

where the first inequality uses (3.6), while the second inequality follows from the definition of  $f_1$  and from the convexity property (2.7), which is equivalent to the assumed subadditivity of  $\text{Ent}_\Phi[\cdot]$ . Therefore,

$$\begin{aligned} \mathbb{E} [\text{Ent}_\Phi [K^* f(Y_1, Y_2) | Y_1]] &= \sum_{y_1 \in \mathbf{Y}_1} P_{Y_1}(y_1) \text{Ent}_\Phi [K^* f(Y_1, Y_2) | Y_1 = y_1] \\ &\leq \eta_2 \sum_{y_1 \in \mathbf{Y}_1} P_{Y_1}(y_1) \sum_{x_1 \in \mathbf{X}_1} P_{X_1 | Y_1}(x_1 | y_1) \text{Ent}_\Phi [f(x_1, X_2)] \\ &= \eta_2 \sum_{x_1 \in \mathbf{X}_1} P_{X_1}(x_1) \text{Ent}_\Phi [f(x_1, X_2)] \\ &= \eta_2 \mathbb{E} [\text{Ent}_\Phi [f(X_1, X_2) | X_1]]. \end{aligned}$$

Next, let  $g_2(x_1) = \mathbb{E}[f_2(x_1, Y_2)] = \mathbb{E}[f_2(X_1, Y_2) | X_1 = x_1]$ . Then

$$\begin{aligned} \text{Ent}_\Phi [\mathbb{E}[K^* f(Y_1, Y_2) | Y_1]] &= \text{Ent}_\Phi [K_1^* g_2(Y_1)] \\ &\leq \eta_1 \text{Ent}_\Phi [g_2(X_1)] \\ &= \eta_1 \text{Ent}_\Phi [\mathbb{E}[f(X_1, X_2) | X_1]], \end{aligned}$$

where the first line follows from the fact that  $(X_1, Y_1)$  and  $(X_2, Y_2)$  are independent and from

definitions, whereas in the last line we have used the fact that

$$\begin{aligned}
g_2(x_1) &= \mathbb{E}[f_2(x_1, Y_2)] \\
&= \sum_{y_2 \in \mathcal{Y}_2} P_{Y_2}(y_2) f_2(x_1, y_2) \\
&= \sum_{y_2 \in \mathcal{Y}_2} P_{Y_2}(y_2) \sum_{x_2 \in \mathcal{X}_2} P_{X_2|Y_2}(x_2|y_2) f(x_1, x_2) \\
&= \sum_{x_2 \in \mathcal{X}_2} P_{X_2}(x_2) f(x_1, x_2) \\
&= \mathbb{E}[f(X_1, X_2)|X_1 = x_1].
\end{aligned}$$

Combining everything, we can write

$$\begin{aligned}
\text{Ent}_\Phi [K^* f(Y_1, Y_2)] &\leq \eta_2 \cdot \mathbb{E} [\text{Ent}_\Phi [f(X_1, X_2)|X_1]] + \eta_1 \cdot \text{Ent}_\Phi [\mathbb{E}[f(X_1, X_2)|X_1]] \\
&\leq \max_{i=1,2} \eta_i \cdot \left\{ \mathbb{E} [\text{Ent}_\Phi [f(X_1, X_2)|X_1]] + \text{Ent}_\Phi [\mathbb{E}[f(X_1, X_2)|X_1]] \right\} \\
&= \max_{i=1,2} \eta_i \cdot \text{Ent}_\Phi [f(X_1, X_2)],
\end{aligned}$$

where in the last step we have used the law of total entropy (2.5). Since  $f$  was arbitrary, we obtain the bound  $\eta \leq \max(\eta_1, \eta_2)$ .  $\square$

### 3.6 Mixtures of local channels

Another situation that often arises in stochastic simulation and machine learning is as follows: Fix  $n$  channels  $K_i \in \mathcal{M}(\mathcal{X}_i|\mathcal{X}_i)$ ,  $1 \leq i \leq n$ , and a probability distribution  $p = (p_i)_{i=1}^n$  on the set  $\{1, \dots, n\}$ . Given an input block  $x^n = (x_1, \dots, x_n) \in \mathcal{X}_1 \times \dots \times \mathcal{X}_n$ , a random output block  $Y^n = (Y_1, \dots, Y_n) \in \mathcal{X}_1 \times \dots \times \mathcal{X}_n$  is generated as follows:

1. a random index  $J \in \{1, \dots, n\}$  is drawn according to  $p$ ;
2.  $Y_J$  is drawn according to  $K_J(\cdot|x_J)$ ;
3.  $Y^{\setminus J} = x^{\setminus J}$ .

The overall stochastic transformation is described by the Markov kernel

$$K \triangleq \sum_{i=1}^n p_i (\text{id}_1 \otimes \dots \otimes \text{id}_{i-1} \otimes K_i \otimes \text{id}_{i+1} \otimes \dots \otimes \text{id}_n),$$

where, for each  $i$ ,  $\text{id}_i$  is the identity mapping on  $\mathcal{F}(\mathcal{X}_i)$ . Now let us also fix  $n$  probability distributions  $\mu_i \in \mathcal{P}(\mathcal{X}_i)$ ,  $1 \leq i \leq n$ . The question is: how does the SDPI constant  $\eta_\Phi(\mu_1 \otimes \dots \otimes \mu_n, K)$  for some  $\Phi \in \mathcal{F}$  depend on  $p$  and on the individual SDPI constants  $\eta_\Phi(\mu_i, K_i)$ ?

**Theorem 3.10.** *Under the same conditions as in Theorem 3.9,*

$$1 - \eta_\Phi(\mu_1 \otimes \dots \otimes \mu_n, K) \geq \min_{1 \leq i \leq n} p_i (1 - \eta_\Phi(\mu_i, K_i)). \quad (3.44)$$

*Proof.* Once again, it suffices to consider the case  $n = 2$ . Thus, we fix two admissible pairs  $(\mu_i, K_i) \in \mathcal{P}(\mathcal{X}_i) \times \mathcal{M}(\mathcal{X}_i|\mathcal{X}_i)$ ,  $i \in \{1, 2\}$  and a parameter  $p \in [0, 1]$ , and consider the channel

$$K = p(K_1 \otimes \text{id}_2) + \bar{p}(\text{id}_1 \otimes K_2).$$

Let  $\mu = \mu_1 \otimes \mu_2$  denote the reference input distribution. We need to show that

$$1 - \eta_\Phi(\mu, K) \geq \min \left( p(1 - \eta_\Phi(\mu_1, K_1)), \bar{p}(1 - \eta_\Phi(\mu_2, K_2)) \right). \quad (3.45)$$

As in the proof of Theorem 3.9, we adopt the shorthand notation  $\eta_i = \eta_\Phi(\mu_i, K_i)$  and

$$\eta = \eta_\Phi(\mu, K) = \eta_\Phi(\mu_1 \otimes \mu_2, K).$$

Let  $(X_1, Y_1, X_2, Y_2)$  be a random tuple with  $(X_1, X_2) \sim \mu$  and  $P_{Y_1, Y_2|X_1, X_2} = K$ . Also, define the Radon–Nikodym derivatives

$$g_1(y_1, y_2) \triangleq \frac{d(\mu_1 K_1 \otimes \mu_2)}{d(\mu K)}(y_1, y_2) = \frac{\mu_1 K_1(y_1) \mu_2(y_2)}{p \mu_1 K_1(y_1) \mu_2(y_2) + \bar{p} \mu_1(y_1) \mu_2 K_2(y_2)} \quad (3.46)$$

and

$$g_2(y_1, y_2) \triangleq \frac{d(\mu_1 \otimes \mu_2 K_2)}{d(\mu K)}(y_1, y_2) = \frac{\mu_1(y_1) \mu_2 K_2(y_2)}{p \mu_1(y_1) \mu_2 K_2(y_2) + \bar{p} \mu_1(y_1) \mu_2 K_2(y_2)}. \quad (3.47)$$

A simple calculation shows that

$$\begin{aligned} P_{X_1, X_2|Y_1, Y_2}(\cdot|y_1, y_2) &= K^*(\cdot|y_1, y_2) \\ &= p g_1(y_1, y_2) (K_1^* \otimes \text{id}_2)(\cdot|y_1, y_2) + \bar{p} g_2(y_1, y_2) (\text{id}_1 \otimes K_2^*)(\cdot|y_1, y_2). \end{aligned}$$

Now consider an arbitrary nonconstant function  $f \in \mathcal{F}_*^0(\mathcal{X}_1 \times \mathcal{X}_2)$ . Then

$$\begin{aligned} &\text{Ent}_\Phi [f(X_1, X_2)] - \text{Ent}_\Phi [K^* f(Y_1, Y_2)] \\ &= \mathbb{E} [\text{Ent}_\Phi [f(X_1, X_2)|Y_1, Y_2]] \\ &= \sum_{y_1 \in \mathcal{Y}_1} \sum_{y_2 \in \mathcal{Y}_2} P_{Y_1, Y_2}(y_1, y_2) \left[ \sum_{x_1 \in \mathcal{X}_1} \sum_{x_2 \in \mathcal{X}_2} P_{X_1, X_2|Y_1, Y_2}(x_1, x_2|y_1, y_2) \Phi(f(x_1, x_2)) \right. \\ &\quad \left. - \Phi \left( \sum_{x_1 \in \mathcal{X}_1} \sum_{x_2 \in \mathcal{X}_2} P_{X_1, X_2|Y_1, Y_2}(x_1, x_2|y_1, y_2) f(x_1, x_2) \right) \right] \\ &= \sum_{y_1 \in \mathcal{Y}_1} \sum_{y_2 \in \mathcal{Y}_2} \mu K(y_1, y_2) \left[ p g_1(y_1, y_2) \sum_{x_1 \in \mathcal{X}_1} \sum_{x_2 \in \mathcal{X}_2} K_1^* \otimes \text{id}_2(x_1, x_2|y_1, y_2) \Phi(f(x_1, x_2)) \right. \\ &\quad \left. + \bar{p} g_2(y_1, y_2) \sum_{x_1 \in \mathcal{X}_1} \sum_{x_2 \in \mathcal{X}_2} \text{id}_1 \otimes K_2^*(x_1, x_2|y_1, y_2) \Phi(f(x_1, x_2)) \right] \\ &\quad - \sum_{y_1 \in \mathcal{Y}_1} \sum_{y_2 \in \mathcal{Y}_2} \mu K(y_1, y_2) \Phi \left( p g_1(y_1, y_2) \sum_{x_1 \in \mathcal{X}_1} \sum_{x_2 \in \mathcal{X}_2} K_1^* \otimes \text{id}_2(x_1, x_2|y_1, y_2) f(x_1, x_2) \right. \\ &\quad \left. + \bar{p} g_2(y_1, y_2) \sum_{x_1 \in \mathcal{X}_1} \sum_{x_2 \in \mathcal{X}_2} \text{id}_1 \otimes K_2^*(x_1, x_2|y_1, y_2) f(x_1, x_2) \right). \end{aligned}$$

From this, using the fact that  $\Phi$  is convex and that  $pg_1 + \bar{p}g_2 = 1$ , we get

$$\begin{aligned}
& \text{Ent}_\Phi [f(X_1, X_2)] - \text{Ent}_\Phi [K^* f(Y_1, Y_2)] \\
& \geq p \sum_{y_1 \in Y_1} \sum_{y_2 \in Y_2} \mu K(y_1, y_2) g_1(y_1, y_2) \left[ \sum_{x_1 \in X_1} \sum_{x_2 \in X_2} (K_1^* \otimes \text{id}_2)(x_1, x_2 | y_1, y_2) \Phi(f(x_1, x_2)) \right. \\
& \quad \left. - \Phi \left( \sum_{x_1 \in X_1} \sum_{x_2 \in X_2} (K_1^* \otimes \text{id}_2)(x_1, x_2 | y_1, y_2) f(x_1, x_2) \right) \right] \\
& \quad + \bar{p} \sum_{y_1 \in Y_1} \sum_{y_2 \in Y_2} \mu K(y_1, y_2) g_2(y_1, y_2) \left[ \sum_{x_1 \in X_1} \sum_{x_2 \in X_2} (\text{id}_1 \otimes K_2^*)(x_1, x_2 | y_1, y_2) \Phi(f(x_1, x_2)) \right. \\
& \quad \left. - \Phi \left( \sum_{x_1 \in X_1} \sum_{x_2 \in X_2} (\text{id}_1 \otimes K_2^*)(x_1, x_2 | y_1, y_2) f(x_1, x_2) \right) \right] \\
& = p \sum_{y_2 \in Y_2} \mu_2(y_2) \sum_{y_1 \in Y_1} \mu_1 K_1(y_1) \left[ \sum_{x_1 \in X_1} K_1^*(x_1 | y_1) \Phi(f(x_1, y_2)) - \Phi \left( \sum_{x_1 \in X_1} K_1^*(x_1 | y_1) f(x_1, y_2) \right) \right] \\
& \quad + \bar{p} \sum_{y_1 \in Y_1} \mu_1(y_1) \sum_{y_2 \in Y_2} \mu_2 K_2(y_2) \left[ \sum_{x_2 \in X_2} K_2^*(x_2 | y_2) \Phi(f(y_1, x_2)) - \Phi \left( \sum_{x_2 \in X_2} K_2^*(x_2 | y_2) f(y_1, x_2) \right) \right], \tag{3.48}
\end{aligned}$$

where in the last step we have used the definitions (3.46) and (3.47) of  $g_1$  and  $g_2$ . Now consider a random tuple  $(X_1, X_2, U, V)$ , such that

1.  $U \longrightarrow X_1 \longrightarrow X_2 \longrightarrow V$  is a Markov chain;
2.  $P_{X_1 X_2} = \mu = \mu_1 \otimes \mu_2$ ;
3.  $P_{U|X_1} = K_1$ ;
4.  $P_{V|X_2} = K_2$ .

Using these definitions in (3.48) gives

$$\begin{aligned}
& \text{Ent}_\Phi [f(X_1, X_2)] - \text{Ent}_\Phi [K^* f(Y_1, Y_2)] \\
& \geq p \sum P_{X_2}(x_2) \mathbb{E} [\text{Ent}_\Phi [f(X_1, x_2) | U]] + \bar{p} \sum P_{X_1}(x_1) \mathbb{E} [\text{Ent}_\Phi [f(x_1, X_2) | V]] \\
& \geq p(1 - \eta_1) \sum_{x_2 \in X_2} P_{X_2}(x_2) \text{Ent}_\Phi (f(X_1, x_2)) + \bar{p}(1 - \eta_2) \sum_{x_1 \in X_1} P_{X_1}(x_1) \text{Ent}_\Phi [f(x_1, X_2)] \\
& \geq \min \left( p(1 - \eta_1), \bar{p}(1 - \eta_2) \right) \left\{ \mathbb{E} [\text{Ent}_\Phi [f(X_1, X_2) | X_2]] + \mathbb{E} [\text{Ent}_\Phi [f(X_1, X_2) | X_1]] \right\} \tag{3.49} \\
& \geq \min \left( p(1 - \eta_1), \bar{p}(1 - \eta_2) \right) \text{Ent}_\Phi [f(X_1, X_2)], \tag{3.50}
\end{aligned}$$

where (3.49) is by the independence of  $X_1$  and  $X_2$ , while (3.50) is by the assumed subadditivity of  $\text{Ent}_\Phi[\cdot]$ . Since  $f$  was arbitrary, we see that the inequality (3.45) indeed holds.  $\square$

**Example 3.6.** Let  $X_1 = \dots = X_n = \{0, 1\}$ ,  $\mu_1 = \dots = \mu_n = \text{Bern}(1/2)$ , and  $K_1 = \dots = K_n = \text{BSC}(\varepsilon)$ . Take  $p$  to be the uniform distribution on  $\{1, \dots, n\}$ . Then  $K$  acts as follows: Given an  $n$ -bit input string  $x^n = (x_1, \dots, x_n)$ , we pick one of the bits uniformly at random and flip it with probability  $\varepsilon$ ; the remaining bits stay the same. Then

$$\eta(\text{Bern}(1/2)^{\otimes n}, K) \leq 1 - \frac{1 - (1 - 2\varepsilon)^2}{n} = 1 - \frac{4\varepsilon\bar{\varepsilon}}{n}.$$

In particular, when  $\varepsilon = 1/2$ , we get the upper bound of  $1 - 1/n$ .

We can also consider flipping bits in blocks: Let  $\mathcal{B} = \{B_m\}_{m=1}^k$  be a disjoint partition of the set  $\{1, \dots, n\}$  into  $k$  blocks. We pick a block uniformly at random, and then independently flip each bit in that block with probability  $\varepsilon$ . Denoting the resulting channel by  $K_{\mathcal{B}}$ , we have

$$\eta(\text{Bern}(1/2)^{\otimes n}, K_{\mathcal{B}}) \leq 1 - \frac{4\varepsilon\bar{\varepsilon}}{k}. \quad (3.51)$$

To prove this, let  $\mu^{(m)} = \bigotimes_{i \in B_m} \mu_i$  and  $K^{(m)} = \bigotimes_{i \in B_m} K_i$ . Then  $\mu = \mu_1 \otimes \dots \otimes \mu_n = \mu^{(1)} \otimes \dots \otimes \mu^{(k)}$ , and by Theorem 3.10 we have

$$\eta(\text{Bern}(1/2)^{\otimes n}, K_{\mathcal{B}}) \leq 1 - \frac{1}{k} \min_{1 \leq m \leq k} \left(1 - \eta(\mu^{(m)}, K^{(m)})\right). \quad (3.52)$$

Since each  $\mu^{(m)}$  is a product measure and each  $K^{(m)}$  is a tensor product of BSCs, Theorem 3.9 gives

$$\eta(\mu^{(m)}, K^{(m)}) = \max_{i \in B_m} \eta(\mu_i, K_i) = \eta(\text{Bern}(1/2), \text{BSC}(\varepsilon)) = (1 - 2\varepsilon)^2.$$

Substituting this into (3.52), we get (3.51). For  $k = 1$ ,  $K_{\mathcal{B}} \equiv \text{BSC}(\varepsilon)^{\otimes n}$ , which has  $\eta = (1 - 2\varepsilon)^2$  by Theorem 3.9. The bound of Eq. (3.51) is then achieved with equality.

### 3.7 Comparison of SDPI constants

The following theorem shows that an upper bound on an SDPI constant for one source-channel pair can be converted into an upper bound for another such pair via a change-of-measure argument:

**Theorem 3.11.** *Let  $(\mu, K), (\bar{\mu}, \bar{K}) \in \mathcal{P}_*(\mathsf{X}) \times \mathcal{M}(\mathsf{Y}|\mathsf{X})$  be two admissible pairs. Then, for any  $\Phi \in \mathcal{F}$  that satisfies the homogeneity condition (3.5),*

$$\eta_{\Phi}(\mu, K) \leq 1 - \frac{a}{A} (1 - \eta_{\Phi}(\bar{\mu}, \bar{K})), \quad (3.53)$$

where

$$A \triangleq \max_{(x,y) \in \mathsf{X} \times \mathsf{Y}} \frac{\bar{\mu} \otimes \bar{K}(x, y)}{\mu \otimes K(x, y)} \quad \text{and} \quad a \triangleq \min_{x \in \mathsf{X}} \frac{\bar{\mu}(x)}{\mu(x)}.$$

**Remark 3.9.** It is easy to see that  $0 < a \leq A$ . Indeed, the first inequality holds since  $\mu, \bar{\mu} \in \mathcal{P}_*(\mathsf{X})$ . For the second, by definition of  $a$  and  $A$ , for every  $x \in \mathsf{X}$  we have

$$a\mu(x) \leq \bar{\mu}(x) = \sum_{y \in \mathsf{Y}} \bar{\mu} \otimes \bar{K}(x, y) \leq A \sum_{y \in \mathsf{Y}} \mu \otimes K(x, y) = A\mu(x).$$

*Proof.* Consider random pairs  $(X, Y)$  and  $(\bar{X}, \bar{Y})$  with respective probability laws  $\mu \otimes K$  and  $\bar{\mu} \otimes \bar{K}$ . Using Eq. (3.6) and the law of total entropy Eq. (2.5),

$$\eta_{\Phi}(\mu, K) = \sup \left\{ \frac{\text{Ent}_{\Phi} [\mathbb{E}[f(X)|Y]]}{\text{Ent}_{\Phi} [f(X)]} : f \in \mathcal{F}_*^0(\mathbf{X}) \text{ } f \neq \text{const} \right\} \quad (3.54)$$

$$\begin{aligned} &= \sup \left\{ \frac{\text{Ent}_{\Phi} [f(X)] - \mathbb{E} [\text{Ent}_{\Phi} [f(X)|Y]]}{\text{Ent}_{\Phi} [f(X)]} : f \in \mathcal{F}_*^0(\mathbf{X}), f \neq \text{const} \right\} \\ &= 1 - \inf \left\{ \frac{\mathbb{E} [\text{Ent}_{\Phi} [f(X)|Y]]}{\text{Ent}_{\Phi} [f(X)]} : f \in \mathcal{F}_*^0(\mathbf{X}), f \neq \text{const} \right\}. \end{aligned} \quad (3.55)$$

Using Lemma A.4 in Appendix A, we can write

$$\begin{aligned} \mathbb{E} [\text{Ent}_{\Phi} [f(X)|Y]] &= \inf_{\xi \in \mathcal{F}_*^0(\mathbf{Y})} \mathbb{E} [\Phi(f(X)) - \Phi(\xi(Y)) - (f(X) - \xi(Y))\Phi'(\xi(Y))] \\ &= \inf_{\xi \in \mathcal{F}_*^0(\mathbf{Y})} \mathbb{E} \left[ \frac{d(\mu \otimes K)}{d(\bar{\mu} \otimes \bar{K})}(\bar{X}, \bar{Y}) (\Phi(f(\bar{X})) - \Phi(\xi(\bar{Y})) - (f(\bar{X}) - \xi(\bar{Y}))\Phi'(\xi(\bar{Y}))) \right] \\ &\geq \frac{1}{A} \inf_{\xi \in \mathcal{F}_*^0(\mathbf{Y})} \mathbb{E} [\Phi(f(\bar{X})) - \Phi(\xi(\bar{Y})) - (f(\bar{X}) - \xi(\bar{Y}))\Phi'(\xi(\bar{Y}))] \\ &= \frac{1}{A} \mathbb{E} [\text{Ent}_{\Phi} [f(\bar{X})|\bar{Y}]], \end{aligned}$$

where the inequality follows from the definition of  $A$  and from the convexity of  $\Phi$ . An analogous argument gives the inequality

$$\text{Ent}_{\Phi} [f(X)] \leq \frac{1}{a} \text{Ent}_{\Phi} [f(\bar{X})].$$

Using these estimates in (3.55), we get

$$\begin{aligned} \eta_{\Phi}(\mu, K) &\leq 1 - \frac{a}{A} \inf \left\{ \frac{\mathbb{E} [\text{Ent}_{\Phi} [f(\bar{X})|\bar{Y}]]}{\text{Ent}_{\Phi} [f(\bar{X})]} : f \in \mathcal{F}_*^0(\mathbf{X}), f \neq \text{const} \right\} \\ &= 1 - \frac{a}{A} (1 - \eta_{\Phi}(\bar{\mu}, \bar{K})). \end{aligned}$$

□

**Corollary 3.2.** *If two channels  $K, \bar{K} \in \mathcal{M}(\mathbf{Y}|\mathbf{X})$  are such that  $\bar{K}(y|x) \leq AK(y|x)$  for all  $(x, y) \in \mathbf{X} \times \mathbf{Y}$ , then*

$$\eta_{\Phi}(\mu, K) \leq 1 - \frac{1}{A} (1 - \eta_{\Phi}(\mu, \bar{K}))$$

for any  $\Phi \in \mathcal{F}$  satisfying (3.5) and any  $\mu \in \mathcal{P}_*(\mathbf{X})$ .

### 3.8 Extremal functions

In this section, we will characterize the extremal functions  $f \in \mathcal{F}_*^0(\mathbf{X})$  that attain the infimum in (3.3). In particular, we will prove that, for any sufficiently smooth  $\Phi$ , these functions are solutions of the variational equation

$$\mathbb{E} \left[ \Phi' (\mathbb{E}[f(\bar{X})|Y]) \Big| X \right] - \Phi'(1) = \eta (\Phi' (f(X)) - \Phi'(1)), \quad (3.56)$$

with  $\eta = \eta_\Phi(\mu, K)$  under the constraint  $f \in \mathcal{F}_*^0(\mathbf{X})$  and  $\mathbb{E}[f(X)] = 1$ . Here, the random triple  $(X, \bar{X}, Y)$  is such that  $X \rightarrow Y \rightarrow \bar{X}$  is a Markov chain, and both  $(X, Y)$  and  $(\bar{X}, Y)$  have law  $\mu \otimes K$ . Written more compactly, (3.56) takes the form

$$K(\Phi' \circ K^* f) - \Phi'(1) = \eta (\Phi' \circ f - \Phi'(1)). \quad (3.57)$$

**Theorem 3.12.** *Suppose  $\Phi \in \mathcal{F}$  has the following properties:*

- (a) *It is three times differentiable with  $\Phi''(1) > 0$ .*
- (b) *The  $\Phi$ -entropy functional is homogeneous in the sense of Definition 3.2.*
- (c) *There exists a constant  $c > 0$ , such that*

$$\text{Ent}_\Phi[U] = c \mathbb{E} [U (\Phi'(U) - \Phi'(1))] \quad (3.58)$$

*for any nonnegative-valued random variable  $U$  with  $\mathbb{E}U = 1$ .*

*Then either  $\eta_\Phi(\mu, K) = S^2(\mu, K)$ , or there exists a nonconstant function  $f \in \mathcal{F}_*^0(\mathbf{X})$ , such that (3.57) holds with  $\eta = \eta_\Phi(\mu, K)$ . Moreover,  $\eta_\Phi(\mu, K)$  is the smallest constant  $\eta > 0$ , for which (3.57) has a solution among nonconstant functions in  $\mathcal{F}_*^0(\mathbf{X})$ .*

**Remark 3.10.** The functions  $\Phi(u) = u \log u$  and  $\Phi(u) = \frac{u^p - 1}{p - 1}$ ,  $1 < p \leq 2$ , satisfy the condition (3.58) (details are provided in the examples after the proof). On the other hand, the function  $\Phi(u) = -\log u$  cannot satisfy (3.58) for any choice of  $c$ , since  $\text{Ent}_\Phi[U] = -\mathbb{E} \log U$ , while  $\mathbb{E}[U(\Phi'(U) - \Phi'(1))] = \mathbb{E}U - 1 = 0$  for any nonnegative-valued  $U$  with  $\mathbb{E}U = 1$ .  $\diamond$

**Remark 3.11.** We emphasize that the variational equation (3.57) may have multiple solutions, not all of which are actually extremal. In general, it is not easy to obtain explicit closed-form expressions for the extremal solutions of (3.57).  $\diamond$

*Proof.* Suppose that  $\eta_\Phi(\mu, K) > S^2(\mu, K)$ , for otherwise there is nothing to prove. We seek to minimize the functional

$$W(f) \triangleq \frac{\text{Ent}_\Phi[f(X)] - \text{Ent}_\Phi[K^* f(Y)]}{\text{Ent}_\Phi[f(X)]} = \frac{\mathbb{E}[\text{Ent}_\Phi[f(X)|Y]]}{\text{Ent}_\Phi[f(X)]}$$

over all  $f \in \mathcal{F}_*^0(\mathbf{X})$ . By homogeneity,  $W(cf) = W(f)$  for all  $c > 0$ , so without loss of generality we can restrict the minimization to  $f \in \mathcal{M} \triangleq \{f \in \mathcal{F}_*^0(\mathbf{X}) : \mathbb{E}[f(X)] = 1\}$ . For  $\varepsilon > 0$ , define the set

$$\mathcal{M}_\varepsilon \triangleq \{f \in \mathcal{F}_*^0(\mathbf{X}) : \mathbb{E}[f(X)] = 1 \text{ and } \|f - 1\|_\infty < \varepsilon\},$$

so  $\mathcal{M} = \mathcal{M}_\infty$ . From the Taylor expansion

$$\Phi(1 + u) = \Phi(1) + \Phi'(1)u + \frac{\Phi''(1)}{2}u^2 + O(u^3),$$

we have, for every  $f \in \mathcal{M}_\varepsilon$ ,

$$\text{Ent}_\Phi[f(X)] = \frac{\Phi''(1)}{2} \text{Var}[f(X)] + O(\varepsilon^3)$$

and

$$\mathbb{E} [\text{Ent}_\Phi[f(X)|Y]] = \frac{\Phi''(1)}{2} (\text{Var}[f(X)] - \text{Var}[\mathbb{E}[f(X)|Y]]) + O(\varepsilon^3).$$

Therefore,

$$\inf_{f \in \mathcal{M}_\varepsilon} W(f) = \inf_{f \in \mathcal{M}_\varepsilon} \frac{\text{Var}[f(X)] - \text{Var}[\mathbb{E}[f(X)|Y]] + O(\varepsilon^3)}{\text{Var}[f(X)] + O(\varepsilon^3)},$$

which implies that, for any  $\delta \in (0, 1)$  there exists some  $\varepsilon_0 = \varepsilon_0(\delta)$ , such that

$$\begin{aligned} \inf_{f \in \mathcal{M}_{\varepsilon_0}} W(f) &\geq \inf_{g \in \mathcal{M}} \frac{\text{Var}[g(X)] - \text{Var}[\mathbb{E}[g(X)|Y]]}{\text{Var}[g(X)]} + \delta \\ &= 1 - S^2(\mu, K) + \delta \\ &> 1 - \eta_\Phi(\mu, K) + \delta. \end{aligned}$$

On the other hand, since  $\inf_{f \in \mathcal{M}} W(f) = 1 - \eta_\Phi(\mu, K)$ , any  $f$  that minimizes  $W$ , if it exists, must lie in  $\mathcal{M} \setminus \mathcal{M}_{\varepsilon_0}$ , i.e., it must be nonconstant. It remains to show the existence of such a minimizing  $f$ . Since any  $f \in \mathcal{M}$  satisfies  $\|f\|_\infty \leq 1/\mu_*$ , where  $\mu_*$  is the smallest (positive) mass of  $\mu$ , the set  $\mathcal{M} \setminus \mathcal{M}_{\varepsilon_0}$  is a closed and bounded subset of a finite-dimensional linear space, hence compact. The denominator of  $W(f)$  is positive for all  $f \in \mathcal{M} \setminus \mathcal{M}_{\varepsilon_0}$ , so  $W$  is a continuous functional on the compact set  $\mathcal{M} \setminus \mathcal{M}_{\varepsilon_0}$  and thus attains its infimum on some nonconstant  $f \in \mathcal{M}$ .

Now, let  $f$  be such a minimizing function. We use a variational argument following Bobkov and Tetali [26, Sec. 6]. Given an arbitrary  $g \in \mathcal{F}(\mathbf{X})$ , the perturbed function  $f + \varepsilon g$  is nonnegative for all sufficiently small  $\varepsilon > 0$ . Consequently, by definition of  $\eta_\Phi$ ,

$$(1 - \eta_\Phi) \text{Ent}_\Phi[f(X) + \varepsilon g(X)] \leq \mathbb{E} [\text{Ent}_\Phi[f(X) + \varepsilon g(X)|Y]]. \quad (3.59)$$

Applying the Taylor expansion

$$\begin{aligned} \text{Ent}_\Phi[U + \varepsilon V] &= \mathbb{E}[\Phi(U + \varepsilon V)] - \Phi(\mathbb{E}U + \varepsilon \mathbb{E}V) \\ &= \text{Ent}_\Phi[U] + \varepsilon \mathbb{E}[(\Phi'(U) - \Phi'(\mathbb{E}U))V] + O(\varepsilon^2) \end{aligned}$$

to  $U = f(X)$ ,  $V = g(X)$  first for  $X \sim \mu$  and then for  $X \sim K^*(\cdot|y)$ ,  $y \in \mathbf{Y}$ , and then using the fact that  $(1 - \eta_\Phi) \text{Ent}_\Phi[f(X)] = \mathbb{E} [\text{Ent}_\Phi[f(X)|Y]]$  by the extremality of  $f$ , we have

$$\begin{aligned} &(1 - \eta_\Phi) \text{Ent}_\Phi[f(X) + \varepsilon g(X)] - \mathbb{E} [\text{Ent}_\Phi[f(X) + \varepsilon g(X)|Y]] \\ &= \varepsilon \mathbb{E} \left[ (1 - \eta_\Phi)(\Phi'(f(X)) - \Phi'(1))g(X) - \mathbb{E}[(\Phi'(f(X)) - \Phi'(\mathbb{E}[f(X)|Y]))g(X)|Y] \right] + o(\varepsilon) \\ &= \varepsilon \mathbb{E} \left[ \Phi'(\mathbb{E}[f(X)|Y])\mathbb{E}[g(X)|Y] - \eta_\Phi(\Phi'(f(X)) - \Phi'(1))g(X) \right] + o(\varepsilon) \\ &= \varepsilon \mathbb{E} \left[ (\mathbb{E}[\Phi'(\mathbb{E}[f(\bar{X})|Y])|X] - \eta_\Phi \Phi'(f(X)) - (1 - \eta_\Phi)\Phi'(1))g(X) \right] + o(\varepsilon), \end{aligned} \quad (3.60)$$

where in the last line  $\bar{X}$  is an independent and identically distributed copy of  $X$  given  $Y$ , and we have used the fact that  $\mathbb{E}[\xi(Y)\mathbb{E}[\gamma(X)|Y]] = \mathbb{E}[\mathbb{E}[\xi(Y)|X]\gamma(X)]$  for any pair  $\gamma \in \mathcal{F}(\mathbf{X})$ ,  $\xi \in \mathcal{F}(\mathbf{Y})$ . Now, by (3.59), the leftmost quantity in (3.60) is nonpositive, whereas the rightmost quantity will be nonpositive for all sufficiently small  $\varepsilon > 0$  if and only if

$$\mathbb{E} \left[ (\mathbb{E}[\Phi'(\mathbb{E}[f(\bar{X})|Y])|X] - \eta_\Phi \Phi'(f(X)) - (1 - \eta_\Phi)\Phi'(1))g(X) \right] = 0.$$

Since  $g$  is arbitrary and  $\mu \in \mathcal{P}_*(\mathsf{X})$ , the minimizing function  $f \in \mathcal{F}_*^0(\mathsf{X})$  with  $\mathbb{E}[f(X)] = 1$  must satisfy

$$\mathbb{E}[\Phi'(\mathbb{E}[f(\bar{X})|Y])|X] - \Phi'(1) = \eta_\Phi(\Phi'(f(X)) - \Phi'(1)),$$

which is precisely (3.57).

It remains to show minimality. To that end, let  $\tilde{\eta} > 0$  be another constant such that there exists some function  $\tilde{f} \in \mathcal{F}_*^0(\mathsf{X})$  with  $\mathbb{E}[\tilde{f}(X)] = 1$  satisfying

$$\mathbb{E}[\Phi'(\mathbb{E}[\tilde{f}(\bar{X})|Y])|X] - \Phi'(1) = \tilde{\eta}(\Phi'(\tilde{f}(X)) - \Phi'(1)) \quad (3.61)$$

Multiplying both sides of (3.61) by  $\tilde{f}(X)$ , taking expectations, and using (3.58), we get

$$\text{Ent}_\Phi[\mathbb{E}[\tilde{f}(X)|Y]] = \tilde{\eta} \text{Ent}_\Phi[\tilde{f}(X)].$$

By definition of  $\eta_\Phi(\mu, K)$ , we must have  $\tilde{\eta} \leq \eta_\Phi(\mu, K)$ .  $\square$

The proof of the theorem shows that if  $\eta_\Phi(\mu, K) > S^2(\mu, K)$ , then Eq. (3.57) admits a nontrivial (i.e., nonconstant) solution. The contrapositive of this statement gives:

**Corollary 3.3.** *If the infimum in (3.3) is not achieved, i.e., if the SDPI (3.1) is strict unless  $f \equiv 1$ , then  $\eta_\Phi(\mu, K) = S^2(\mu, K)$ .*

**Remark 3.12.** Equivalently,  $\eta_\Phi(\mu, K) = S^2(\mu, K)$  if for an arbitrary  $\gamma > 0$  the only solution to Eq. (3.57) among  $f \in \mathcal{F}_*^0(\mathsf{X})$  with  $\mathbb{E}[f(X)] = 1$  is the trivial solution  $f \equiv 1$ .  $\diamond$

Here are a couple of specific examples:

- For  $\Phi(u) = u \log u$ , we have  $\Phi'(u) = \log u + 1$ , and  $\Phi$  satisfies the conditions (a)–(c) of Theorem 3.12. In particular, Eq. (3.58) holds with  $c = 1$ . The variational equation (3.57) becomes

$$K(\log K^* f) = \eta \log f, \quad \eta = \eta(\mu, K).$$

- For  $\Phi(u) = \frac{u^p - 1}{p - 1}$ ,  $1 < p \leq 2$ , we have  $\Phi'(u) = \frac{pu^{p-1}}{p-1}$ , and  $\Phi$  satisfies the conditions (a)–(c). In this case, (3.58) holds with  $c = p$ . The variational equation takes the form

$$K((K^* f)^{p-1}) = \eta_p f^{p-1} + 1 - \eta_p, \quad \eta_p = \eta_{\Phi_p}(\mu, K).$$

## 4 Connections with $\Phi$ -Sobolev inequalities

### 4.1 General framework

Strong data processing inequalities for a pair  $(\mu, K)$  can be interpreted in terms of the effect of the adjoint channel  $K^*$  on the  $\Phi$ -entropies of suitably normalized nonnegative functions of the input, see Proposition 3.1. In this section, we show that there is a close relationship between SDPIs and another class of functional inequalities — the so-called  *$\Phi$ -Sobolev inequalities* [23, 31] that relate the  $\Phi$ -entropy  $\text{Ent}_\Phi[f(X)]$  of an arbitrary function of the input  $X \sim \mu$  to some measure of correlation between  $f(X)$  and the output  $Y \sim \mu K$ .

We will measure correlation in the following way. For any triple  $(U, V, Z)$  of jointly distributed random variables, where  $U, V$  are real-valued, we define

$$\mathcal{E}(U, V|Z) \triangleq \mathbb{E} [(U - \mathbb{E}[U|Z])(V - \mathbb{E}[V|Z])]. \quad (4.1)$$

This quantity has an estimation-theoretic interpretation: since  $e(U|Z) \triangleq U - \mathbb{E}[U|Z]$  is the error of a minimum mean-square error (MMSE) estimator of  $U$  given  $Z$ , and  $\mathbb{E}[e(U|Z)] = 0$ ,  $\mathcal{E}(U, V|Z)$  is the covariance of  $e(U|Z)$  and  $e(V|Z)$ :

$$\mathcal{E}(U, V|Z) = \text{Cov}[e(U|Z), e(V|Z)].$$

In particular,

$$\mathcal{E}(U, U|Z) = \mathbb{E} \left[ (U - \mathbb{E}[U|Z])^2 \right] \equiv \text{MMSE}(U|Z),$$

the MMSE achievable in estimating  $U$  from  $Z$ . We pause to record a few key properties of  $\mathcal{E}$  (see Appendix B for the proof):

**Proposition 4.1.** *The functional  $\mathcal{E}$  defined in (4.1) has the following properties:*

1. *Symmetry* –  $\mathcal{E}(U, V|Y) = \mathcal{E}(V, U|Y)$ .
2. *Linearity* –  $\mathcal{E}(aU + bU', V|Y) = a\mathcal{E}(U, V|Y) + b\mathcal{E}(U', V|Y)$  for any constants  $a, b \in \mathbb{R}$ .
3. *Degeneracy* – If  $U$  is constant a.s., then  $\mathcal{E}(U, V|Y) = 0$ .
4. *Representation in terms of an exchangeable pair* – Let  $(X, Y) \in \mathsf{X} \times \mathsf{Y}$  be a random pair with  $P_X = \mu \in \mathcal{P}(\mathsf{X})$  and  $P_{Y|X} = K \in \mathcal{M}(\mathsf{Y}|\mathsf{X})$ , where  $(\mu, K)$  is an admissible pair. Then for any two functions  $f, g \in \mathcal{F}(\mathsf{X})$ ,

$$\mathcal{E}(f(X), g(X)|Y) = \frac{1}{2} \mathbb{E} [(f(X) - f(X'))(g(X) - g(X'))] \quad (4.2)$$

$$= \mathbb{E} \left[ (f(X) - f(X'))_+ (g(X) - g(X')) \right], \quad (4.3)$$

where  $(u)_+ \triangleq u \vee 0$ , and  $(X, X')$  is a pair of  $\mathsf{X}$ -valued random variables with  $P_X = \mu$  and  $P_{X'|X} = K^*K$ .

**Remark 4.1.** The terminology in Item 4 merits some discussion. It is not hard to show (and, in fact, we do show it in the proof of the proposition) that the joint distribution  $P_{XX'}$  has the following symmetry property:

$$P_{XX'}(x, x') = P_{XX'}(x', x), \quad \forall x, x' \in \mathsf{X}. \quad (4.4)$$

In other words, the random variables  $X$  and  $X'$  form an *exchangeable pair*.  $\diamond$

Generalizing the definition due to Chafaï [23], we now introduce  $\Phi$ -Sobolev inequalities:

**Definition 4.1.** *Consider an admissible pair  $(\mu, K)$  and a random pair  $(X, Y)$  with probability distribution  $\mu \otimes K$ . Fix a function  $\Phi \in \mathcal{F}$ . We say that  $(\mu, K)$  satisfies a  $\Phi$ -Sobolev inequality with constant  $\alpha \geq 0$  if there exists some function  $\Psi : \mathbb{R}^+ \rightarrow \mathbb{R}$ , such that the inequality*

$$\text{Ent}_\Phi[f(X)] \leq \alpha \mathcal{E}(f(X), \Psi \circ f(X)|Y) \quad (4.5)$$

holds for all  $f \in \mathcal{F}_*^0(\mathsf{X})$ .

Now we are ready to state our main result that relates SDPIs to  $\Phi$ -Sobolev inequalities:

**Theorem 4.1.** *Suppose that  $\Phi \in \mathcal{F}$  is such that  $\Phi(0) < \infty$ , the function  $\Psi$  defined in (3.14) is concave, and the corresponding  $\Phi$ -entropy is homogeneous. Then  $\eta_\Phi(\mu, K) \leq c$  implies that the pair  $(\mu, K)$  satisfies the  $\Phi$ -Sobolev inequality of the form (4.5) with constant  $\alpha = (1 - c)^{-1}$ .*

*Proof.* For any  $u > 0$  we can write  $\Phi(u) = u\Psi(u) + \Phi(0)$ . Thus, for any real-valued random variable  $U$  which is a.s. strictly positive and a jointly distributed random variable  $Y$ , we have

$$\begin{aligned} \text{Ent}_\Phi[U|Y] &= \mathbb{E}[U\Psi(U)|Y] - \mathbb{E}[U|Y]\Psi(\mathbb{E}[U|Y]) \\ &\leq \mathbb{E}[U\Psi(U)|Y] - \mathbb{E}[U|Y]\mathbb{E}[\Psi(U)|Y], \end{aligned} \quad (4.6)$$

where the second line is by the concavity of  $\Psi$ . Now let  $U = f(X)$  for some  $f \in \mathcal{F}_*(\mathbf{X})$ . Using (4.6) and Proposition 3.1, we get

$$\begin{aligned} \text{Ent}_\Phi[f(X)] &\leq \frac{1}{1-c} \mathbb{E}[f(X)\Psi(f(X)) - \mathbb{E}[f(X)|Y]\mathbb{E}[\Psi(f(X)|Y)]] \\ &= \frac{1}{1-c} \mathcal{E}(f(X), \Psi \circ f(X)|Y), \end{aligned}$$

where the second line follows from the easily verified identity  $\mathcal{E}(U, V|Z) = \mathbb{E}[UV - \mathbb{E}[U|Z]\mathbb{E}[V|Z]] = \mathbb{E}[U(V - \mathbb{E}[V|Z])]$ .  $\square$

Theorem 4.1 provides a route to  $\Phi$ -Sobolev inequalities via SDPIs — any good upper bound on  $\eta_\Phi(\mu, K)$  would automatically translate into a bound on the constant in the corresponding  $\Phi$ -Sobolev inequality. Such functional inequalities are a powerful tool in applied probability (for example, in the context of quantifying the convergence of Markov chains to equilibrium); in the next section, we will illustrate this on the particular case of Poincaré inequalities (corresponding to  $\Phi(u) = u^2 - 1$  and log-Sobolev inequalities (corresponding to  $\Phi(u) = u \log u$ ).

It is often useful to estimate the  $\Phi$ -entropy of a composite function  $F \circ f$  (we will see examples of this later on). The following result contains Theorem 5 of [40] as a special case:

**Theorem 4.2.** *Suppose that the assumptions of Theorem 4.1 hold, and that the function  $\Psi$  is differentiable. Let  $F : \mathbb{R} \rightarrow \mathbb{R}^+$  be a convex, differentiable, nondecreasing function, such that  $\Psi \circ F$  is convex. Then, for any  $f \in \mathcal{F}(\mathbf{X})$ ,*

$$\text{Ent}_\Phi[F \circ f(X)] \leq \frac{1}{1-c} \mathbb{E} \left[ F'^2(f(X)) \Psi'(F \circ f(X)) (f(X) - f(X'))_+^2 \right],$$

where  $(X, X')$  is an exchangeable pair of random variables with  $P_X = \mu$  and  $P_{X'|X} = K^*K$ , and  $\Psi'$  denotes the right derivative of  $\Psi$ . Similarly, if  $F$  is nonincreasing, then

$$\text{Ent}_\Phi[F \circ f(X)] \leq \frac{1}{1-c} \mathbb{E} \left[ F'^2(f(X)) \Psi'(F \circ f(X)) (f(X') - f(X))_+^2 \right],$$

*Proof.* We only consider the case when  $F$  is nondecreasing, since the other case is handled similarly. Suppose that  $u > v$ . Then, by monotonicity and convexity of  $F$ ,

$$0 \leq F(u) - F(v) \leq F'(u)(u - v).$$

Moreover, because  $f$  is convex, the function  $\Psi$  defined in (3.14) is nondecreasing. Using this together with the assumed convexity of  $\Psi \circ F$ , we have

$$0 \leq \Psi(F(u)) - \Psi(F(v)) \leq \Psi'(F(u))F'(u)(u - v).$$

Thus, when  $u > v$ ,

$$(F(u) - F(v))(\Psi(F(u)) - \Psi(F(v))) \leq F'^2(u)\Psi'(F(u))(u - v)^2. \quad (4.7)$$

Therefore, using Theorem 4.1, we can write

$$\begin{aligned} \text{Ent}_\Phi[F \circ f(X)] &\leq \frac{1}{1-c} \mathcal{E}(F \circ f(X), \Psi \circ F \circ f(X)|Y) \\ &= \frac{1}{1-c} \mathbb{E} \left[ (F(f(X)) - F(f(X')))_{+} \cdot (\Psi(F(f(X))) - \Psi(F(f(X')))) \right] \\ &\leq \frac{1}{1-c} \mathbb{E} \left[ F'^2(f(X)) \Psi'(F(f(X))) (f(X) - f(X'))_{+}^2 \right], \end{aligned}$$

where the second step is by (4.3), while the last step is by (4.7).  $\square$

## 4.2 Logarithmic Sobolev and Poincaré inequalities

We now particularize the above general results to two specific types of functional inequalities:

- logarithmic Sobolev inequalities, with  $\Phi(u) = u \log u$ ;
- Poincaré inequalities, with  $\Phi(u) = u^2 - 1$ .

These inequalities are well-known in functional analysis and probability theory (see, e.g., [8, 24–27]). We will first introduce our definitions of these inequalities following the ideas laid down in the preceding section, and then show how these definitions are related to the “standard” ones.

We start with Poincaré inequalities:

**Definition 4.2.** *We say that an admissible pair  $(\mu, K) \in \mathcal{P}(\mathbf{X}) \times \mathcal{M}(\mathbf{Y}|\mathbf{X})$  satisfies a Poincaré inequality with constant  $\alpha \geq 0$  if*

$$\text{Var} [f(X)] \leq \alpha \mathcal{E}(f(X), f(X)|Y)$$

for all  $f \in \mathcal{F}_*^0(\mathbf{X})$ , where  $(X, Y)$  is a random pair with probability law  $\mu \otimes K$ . The Poincaré constant of  $(\mu, K)$  is given by

$$\lambda(\mu, K) \triangleq \inf_{f \in \mathcal{F}_*^0(\mathbf{X})} \frac{\mathcal{E}(f(X), f(X)|Y)}{\text{Var} [f(X)]},$$

where we adopt the convention that  $\frac{0}{0} = +\infty$ .

According to the above definition,  $\alpha^* = \frac{1}{\lambda(\mu, K)}$  is the smallest value of  $\alpha$  for which the pair  $(\mu, K)$  will satisfy a Poincaré inequality. Moreover, we have the following:

**Proposition 4.2.** For any admissible pair  $(\mu, K)$ ,

$$\lambda(\mu, K) = 1 - S^2(\mu, K).$$

That is,  $(\mu, K)$  satisfies a Poincaré inequality with constant  $\alpha$  if and only if  $\eta_{\chi^2}(\mu, K) \leq 1 - 1/\alpha$ .

*Proof.* The function  $\Phi(u) = u^2 - 1$  satisfies the conditions of Theorem 4.1 with  $\Psi(u) = u$ , and  $\text{Ent}_\Phi[U] = \text{Var}[U]$ . Therefore, if  $\eta_{\chi^2}(\mu, K) \leq c$ , then for any  $f \in \mathcal{F}_*^0(\mathsf{X})$  we have

$$\begin{aligned} \text{Var}[f(X)] &\leq \frac{1}{1-c} \mathcal{E}(f(X), \Psi \circ f(X) | Y) \\ &= \frac{1}{1-c} \mathcal{E}(f(X), f(X) | Y), \end{aligned}$$

which implies that the pair  $(\mu, K)$  satisfies Poincaré with constant  $\alpha = \frac{1}{1-c}$ . Therefore,

$$\begin{aligned} \lambda(\mu, K) &\geq \sup \{1 - c : \eta_{\chi^2}(\mu, K) \leq c\} \\ &= 1 - \eta_{\chi^2}(\mu, K) \\ &= 1 - S^2(\mu, K), \end{aligned}$$

where the last step is by Theorem 3.2.

Conversely, suppose that  $(\mu, K)$  satisfies Poincaré with constant  $\alpha$ . A simple computation shows

$$\mathcal{E}(f(X), f(X) | Y) = \text{Var}[f(X)] - \text{Var}[K^* f(Y)].$$

Therefore,

$$\text{Var}[K^* f(Y)] \leq \left(1 - \frac{1}{\alpha}\right) \text{Var}[f(X)]$$

for any  $f \in \mathcal{F}_*^0(\mathsf{X})$ . This, in turn, implies that

$$\begin{aligned} S^2(\mu, K) &= \sup_{f \in \mathcal{F}_*^0(\mathsf{X})} \frac{\text{Var}[K^* f(Y)]}{\text{Var}[f(X)]} \\ &\leq \inf \left\{ 1 - \frac{1}{\alpha} : \frac{1}{\alpha} \leq \lambda(\mu, K) \right\} \\ &= 1 - \lambda(\mu, K). \end{aligned}$$

□

Now let us consider log-Sobolev inequalities:

**Definition 4.3.** We say that an admissible pair  $(\mu, K) \in \mathcal{P}(\mathsf{X}) \times \mathcal{M}(\mathsf{Y}|\mathsf{X})$  satisfies a logarithmic Sobolev inequality with constant  $\alpha \geq 0$  if

$$\text{Ent}[f(X)] \leq \alpha \mathcal{E}(f(X), \log f(X) | Y)$$

for all  $f \in \mathcal{F}_*(\mathsf{X})$ , where  $(X, Y)$  is a random pair with probability law  $\mu \otimes K$ . The log-Sobolev constant of  $(\mu, K)$  is given by

$$\rho_1(\mu, K) \triangleq \inf_{f \in \mathcal{F}_*(\mathsf{X})} \frac{\mathcal{E}(f(X), \log f(X) | Y)}{\text{Ent}[f(X)]}, \quad (4.8)$$

again with the convention that  $\frac{0}{0} = +\infty$ .

The following is an extension of Prop. 5.1 in [10] to the case  $X \neq Y$ , and with explicit constants:

**Proposition 4.3.** *For any admissible pair  $(\mu, K)$ ,*

$$1 - \eta(\mu, K) \leq \rho_1(\mu, K) \leq 1 - \frac{(1 - \log 2) \log 2}{2} \eta(\mu, K) \leq 1 - \frac{1}{10} \eta(\mu, K). \quad (4.9)$$

*That is, if  $\eta(\mu, K) \leq c$ , then  $(\mu, K)$  satisfies a log-Sobolev inequality with constant  $\alpha = \frac{1}{1-c}$ . Conversely, if  $(\mu, K)$  satisfies log-Sobolev with constant  $\alpha$ , then*

$$\eta(\mu, K) \leq \frac{2}{(1 - \log 2) \log 2} \left(1 - \frac{1}{\alpha}\right) \leq 10 \left(1 - \frac{1}{\alpha}\right).$$

*Proof.* The first inequality in (4.9) follows from Theorem 4.1 with  $\Phi(u) = u \log u$  and  $\Psi(u) = \log u$ . To prove the second inequality, we borrow (and slightly streamline) an ingenious idea from [10]. Let us fix an arbitrary  $f \in \mathcal{F}_*(X)$ . Then

$$\begin{aligned} \mathbb{E} [\text{Ent} [f(X)|Y]] &= \sum_{y \in Y} \mu K(y) \sum_{x \in X} K^*(x|y) f(x) \log \frac{f(x)}{K^* f(y)} \\ &= \sum_{y \in Y} \mu K(y) \text{Ent} [f(X)|Y = y]. \end{aligned} \quad (4.10)$$

By [10, Lm. 5.2], the entropy  $\text{Ent}[U]$  of any nonnegative real-valued random variable  $U$  with  $\mathbb{E}[U \log U] < \infty$  admits the integral representation

$$\text{Ent}[U] = \frac{1}{2} \int_0^\infty e^{-t} \mathbb{E} \left[ (U - \bar{U}) \log \frac{e^{-t}U + 1 - e^{-t}}{e^{-t}\bar{U} + 1 - e^{-t}} \right] dt, \quad (4.11)$$

where  $\bar{U}$  is an independent copy of  $U$ . Applying (4.11) to each term in (4.10), we obtain

$$\text{Ent}[f(X)|Y = y] = \frac{1}{2} \sum_{x, \bar{x} \in X} K^*(x|y) K^*(\bar{x}|y) \left[ \int_0^\infty (f_t(x) - f_t(\bar{x})) (\log f_t(x) - \log f_t(\bar{x})) dt \right],$$

where  $f_t(x) \triangleq e^{-t}f(x) + 1 - e^{-t}$ . Averaging this w.r.t.  $Y \sim \mu K$  gives

$$\begin{aligned} &\mathbb{E} [\text{Ent} [f(X)|Y]] \\ &= \frac{1}{2} \sum_{y \in Y} \sum_{x, \bar{x} \in X} \mu K(y) K^*(x|y) K^*(\bar{x}|y) \left[ \int_0^\infty (f_t(x) - f_t(\bar{x})) (\log f_t(x) - \log f_t(\bar{x})) dt \right] \\ &= \frac{1}{2} \sum_{x, \bar{x} \in X} \mu(x) \sum_{y \in Y} K(y|x) K^*(\bar{x}|y) \left[ \int_0^\infty (f_t(x) - f_t(\bar{x})) (\log f_t(x) - \log f_t(\bar{x})) dt \right] \\ &= \frac{1}{2} \int_0^\infty \mathbb{E} [(f_t(X) - f_t(X')) (\log f_t(X) - \log f_t(X'))] dt \\ &= \int_0^\infty \mathcal{E}(f_t(X), \log f_t(X)|Y) dt, \end{aligned} \quad (4.12)$$

where  $(X, X')$  is an exchangeable pair with joint law  $P_{XX'} = \mu \otimes K^*K$ , and in the last step we have used Eq. (4.2). From (4.12) and the definition of the log-Sobolev constant, it follows that

$$\mathbb{E} [\text{Ent} [f(X)|Y]] \geq \rho_1(\mu, K) \int_0^\infty \text{Ent} [f_t(X)] dt. \quad (4.13)$$

Now consider the function  $\xi(u) \triangleq (u+1)\log(u+1) - u$ ,  $u \geq -1$ . This function is nonnegative, nonincreasing on  $[-1, 0]$ , nondecreasing on  $\mathbb{R}^+$ , and

$$\inf_{u \geq -1} \frac{\xi(u/2)}{\xi(u)} = \frac{1 - \log 2}{2}.$$

By monotonicity,  $\xi(cu) \geq \xi(u/2) \geq \frac{1-\log 2}{2}\xi(u)$  for all  $u \geq -1$  and for any  $1/2 \leq c \leq 1$ . Therefore,

$$\begin{aligned} \int_0^\infty \text{Ent}[f_t(X)] dt &= \int_0^\infty \mathbb{E}[f_t(X) \log f_t(X) - f_t(X) + 1] dt \\ &= \int_0^\infty \mathbb{E}[\xi(f_t(X) - 1)] dt \\ &= \int_0^\infty \mathbb{E}[\xi(e^{-t}(f(X) - 1))] dt \\ &\geq \int_0^{\log 2} \mathbb{E}[\xi(e^{-t}(f(X) - 1))] dt \\ &\geq \frac{1 - \log 2}{2} \int_0^{\log 2} \mathbb{E}[\xi(f(X) - 1)] dt \\ &= \frac{(1 - \log 2) \log 2}{2} \text{Ent}[f(X)]. \end{aligned} \tag{4.14}$$

Using (4.14) in (4.12), we obtain

$$\mathbb{E}[\text{Ent}[f(X)|Y]] \geq \frac{(1 - \log 2) \log 2}{2} \rho_1(\mu, K) \text{Ent}[f(X)].$$

Since  $f$  was arbitrary, this implies the second inequality in (4.9).  $\square$

Now let us see how these results are related to the standard formulation of log-Sobolev inequalities in a discrete setting (see, e.g., [25–27]). Given a finite set  $\mathsf{X}$ , we fix an admissible pair  $(\mu, M) \in \mathcal{P}(\mathsf{X}) \times \mathcal{M}(\mathsf{X}|\mathsf{X})$ , such that the Markov kernel  $M$  is *reversible* w.r.t.  $\mu$ :

$$\mu(x)M(x'|x) = \mu(x')M(x|x'), \quad \forall x, x' \in \mathsf{X} \tag{4.15}$$

(nonreversible kernels can be handled as well, but we will not need this generalization here). From (4.15), it follows that  $M$  leaves  $\mu$  invariant:  $\mu M = \mu$ . Define the *Dirichlet form*  $\mathcal{E}_{\mu, M} : \mathcal{F}(\mathsf{X}) \times \mathcal{F}(\mathsf{X}) \rightarrow \mathbb{R}$  by

$$\begin{aligned} \mathcal{E}_{\mu, M}(f, g) &\triangleq \frac{1}{2} \sum_{x, x' \in \mathsf{X}} (f(x) - f(x')) (g(x) - g(x')) \mu(x)M(x'|x) \\ &\equiv \frac{1}{2} \mathbb{E}[(f(X) - f(X')) (g(X) - g(X'))], \end{aligned} \tag{4.16}$$

where  $(X, X') \in \mathsf{X} \times \mathsf{X}$  is a random pair with probability law  $\mu \otimes M$ . Our “overloading” of the notation  $\mathcal{E}(\cdot, \cdot)$  [compare with Eq. (4.1)] is not accidental. To see this, we first need a definition:

**Definition 4.4.** Fix some alphabet  $\mathsf{Y}$  and a channel  $K \in \mathcal{M}(\mathsf{Y}|\mathsf{X})$ . We say that the pair  $(\mu, M)$  factors through  $K$  if  $M = K_\mu^* \circ K$ , i.e., if

$$M(x'|x) = \sum_{y \in \mathsf{Y}} K_\mu^*(x'|y)K(y|x), \quad \forall (x, x') \in \mathsf{X} \times \mathsf{X}.$$

In other words,  $(\mu, M)$  factors through  $K$  if we can generate a copy of  $(X, X')$  according to the following two-stage procedure, starting with a draw  $X \sim \mu$ :

1. Pass  $X$  through the channel  $K$  to get  $Y$ .
2. Pass  $Y$  through the adjoint channel  $K_\mu^*$  to get  $X'$ .

This is nothing but the well-known *two-stage* (or *two-component*) *Gibbs sampler* [28, 30].

**Proposition 4.4.** *The random variables  $X$  and  $X'$  form an exchangeable pair. Moreover, if  $(\mu, M)$  factors through some channel  $K \in \mathcal{M}(\mathsf{Y}|\mathsf{X})$ , then*

$$\mathcal{E}_{\mu, M}(f, g) = \mathcal{E}(f(X), g(X)|Y), \quad \forall f, g \in \mathcal{F}(\mathsf{X}) \quad (4.17)$$

where  $(X, Y) \in \mathsf{X} \times \mathsf{Y}$  is a random pair with law  $\mu \otimes K$ .

*Proof.* Exchangeability of  $(X, X')$  follows from the reversibility condition (4.15). The identity (4.17) is a consequence of (4.16) and Prop. 4.1, Part 4).  $\square$

With these definitions out of the way, we can introduce the hierarchy of log-Sobolev inequalities following Mossel et al. [27]:

**Definition 4.5.** *The pair  $(\mu, M)$  satisfies log-Sobolev inequality of order  $p \in \mathbb{R}^+ \setminus \{0, 1\}$  with constant  $c$ , or  $\text{LSI}_p(c)$ , if*

$$\text{Ent}[f^p(X)] \leq \frac{cp^2}{4(p-1)} \mathcal{E}_{\mu, M}(f^{p-1}, f), \quad \forall f \in \mathcal{F}_*^0(\mathsf{X});$$

$\text{LSI}_1(c)$  if

$$\text{Ent}[f(X)] \leq \frac{c}{4} \mathcal{E}_{\mu, M}(f, \log f), \quad \forall f \in \mathcal{F}_*(\mathsf{X});$$

and  $\text{LSI}_0(c)$  if

$$\text{Var}[\log f(X)] \leq -\frac{c}{2} \mathcal{E}_{\mu, M}(f, 1/f), \quad \forall f \in \mathcal{F}_*(\mathsf{X}).$$

Another important functional inequality relates the variance to the Dirichlet form  $\mathcal{E}_{\mu, M}$ :

**Definition 4.6.**  *$(\mu, M)$  satisfies a Poincaré inequality with constant  $c \geq 0$ , or  $\text{PI}(c)$ , if*

$$\text{Var}[f(X)] \leq c \mathcal{E}_{\mu, M}(f, f), \quad \forall f \in \mathcal{F}_*^0(\mathsf{X}).$$

We are interested in the tightest constants in log-Sobolev inequalities for  $p \in [0, 2]$ . With that in mind, we define

$$\tilde{\rho}_p(\mu, M) \triangleq \frac{p^2}{4(p-1)} \inf_{f \in \mathcal{F}_*(\mathsf{X})} \frac{\mathcal{E}_{\mu, M}(f^{p-1}, f)}{\text{Ent}[f^p(X)]}$$

for  $p \notin \{0, 1\}$ , with the convention  $\frac{0}{0} = \infty$ . The constants  $\tilde{\rho}_0, \tilde{\rho}_1$  are defined analogously. The *Poincaré constant* is

$$\tilde{\lambda}(\mu, M) \triangleq \inf_{f \in \mathcal{F}_*(\mathsf{X})} \frac{\mathcal{E}_{\mu, M}(f, f)}{\text{Var}[f(X)]}.$$

Mossel et al. [27] proved that the function  $p \mapsto \tilde{\rho}_p(\mu, M)$  is nonincreasing:

$$\tilde{\rho}_0(\mu, M) \geq \tilde{\rho}_p(\mu, M) \geq \tilde{\rho}_q(\mu, M), \quad 0 \leq p \leq q \leq 2 \quad (4.18)$$

and moreover  $\tilde{\rho}_0(\mu, M) = \frac{1}{2}\tilde{\lambda}(\mu, M)$ . Log-Sobolev and Poincaré inequalities arise naturally in the study of the continuous-time random walk on  $\mathsf{X}$  with infinitesimal generator  $L = M - I$ . This is a pure-jump Markov process with state space  $\mathsf{X}$  that jumps from state  $x$  to another state  $x'$  with probability  $M(x'|x)$ , and the times between successive jumps are i.i.d.  $\text{Exp}(1)$  random variables. Let  $\{X_t\}_{t \geq 0}$  denote this process with  $X_0 \sim \mu$ , where  $X_t \sim \mu$  for all  $t$  by stationarity. For each  $t \geq 0$ , define the mapping  $P_t : \mathcal{F}(\mathsf{X}) \rightarrow \mathcal{F}(\mathsf{X})$  by

$$P_t f(x) \triangleq \mathbb{E}[f(X_t) | X_0 = x].$$

Then one can prove the following (see, e.g., [46, Prop. 1.7]):

1.  $\text{Var}[P_t f(X_0)] \leq e^{-t/c} \text{Var}[f(X_0)]$  for all  $f \in \mathcal{F}(\mathsf{X})$  and all  $t \geq 0$  if and only if the pair  $(\mu, M)$  satisfies  $\text{PI}(c)$ .
2.  $\text{Ent}[P_t f(X_0)] \leq e^{-t/c} \text{Ent}[f(X_0)]$  for all  $f \in \mathcal{F}_*^0(\mathsf{X})$  and all  $t \geq 0$  if and only if the pair  $(\mu, M)$  satisfies  $\text{LSI}_1(4c)$ .

In other words, the Poincaré inequality and the log-Sobolev inequality for  $p = 1$  completely characterize the exponential rate of decay of variance and entropy, respectively, along the trajectory of  $\{X_t\}$  with  $X_0 \sim \mu$ . In particular, if for each  $t \geq 0$  we consider the channel  $M_t \in \mathcal{M}(\mathsf{X}|\mathsf{X})$  with transition probabilities  $M_t(x'|x) = \mathbb{P}(X_t = x' | X_0 = x)$ , then

$$\eta_{\chi^2}(\mu, M_t) \leq e^{-\tilde{\lambda}(\mu, M)t} \quad \text{and} \quad \eta(\mu, M_t) \leq e^{-4\tilde{\rho}_1(\mu, M)t}.$$

The main utility of the log-Sobolev inequality for  $p = 2$  is that the Dirichlet form  $\mathcal{E}_{\mu, M}(f, f)$  is much easier to deal with than  $\mathcal{E}_{\mu, M}(f, \log f)$ ; by monotonicity property of the log-Sobolev constants [cf. (4.18)], we end up with the handy estimate

$$\eta(\mu, M_t) \leq e^{-4\tilde{\rho}_2(\mu, M)t}.$$

Thus, it is important to obtain tight upper and lower bounds on the Poincaré and the log-Sobolev constants of the pair  $(\mu, M)$ . We now show that such bounds can be given in terms of the SDPI constant  $\eta(\mu, K)$  of any channel  $K$  that the pair  $(\mu, M)$  factors through; conversely, we can obtain bounds on  $\eta(\mu, K)$  in terms of log-Sobolev and Poincaré constants of the pair  $(\mu, K_\mu^* \circ K)$ . We start with the Poincaré constant, in which case we have the following exact characterization:

**Theorem 4.3.** *The functional  $K \mapsto S^2(\mu, K)$  is constant on the collection*

$$\mathcal{M}(\mu, M) \triangleq \{K : (\mu, M) \text{ factors through } K\},$$

*and its value there is equal to  $1 - \tilde{\lambda}(\mu, M)$ . Equivalently, if  $K \in \mathcal{M}(\mu, M)$ , then*

$$\eta_{\chi^2}(\mu, K) = 1 - \tilde{\lambda}(\mu, M).$$

*Proof.* We need to show the following: if  $(\mu, M)$  factors through  $K$ , then

$$\mathcal{E}_{\mu, M}(f, f) = \text{Var}[f(X)] - \text{Var}[K_{\mu}^* f(Y)], \quad (4.19)$$

where  $(X, Y) \in \mathsf{X} \times \mathsf{Y}$  is a random pair with law  $\mu \otimes K$ . Assuming this is true, we then have

$$\begin{aligned} \tilde{\lambda}(\mu, M) &= \inf_{f \in \mathcal{F}(\mathsf{X})} \frac{\mathcal{E}_{\mu, M}(f, f)}{\text{Var}[f(X)]} \\ &= \inf_{f \in \mathcal{F}(\mathsf{X})} \frac{\text{Var}[f(X)] - \text{Var}[K_{\mu}^* f(Y)]}{\text{Var}[f(X)]} \\ &= 1 - \eta_{\chi^2}(\mu, K) \\ &= 1 - S^2(\mu, K). \end{aligned}$$

Noting that  $\tilde{\lambda}(\mu, M)$  is independent of the choice of  $K$ , we obtain the statement of the theorem.

It remains to prove (4.19). Fixing  $K$ , we have

$$\text{Var}[f(X)] - \text{Var}[K_{\mu}^* f(Y)] = \mathbb{E}[f^2(X)] - \mathbb{E}[(\mathbb{E}[f(X)|Y])^2],$$

where

$$\begin{aligned} \mathbb{E}[(\mathbb{E}[f(X)|Y])^2] &= \sum_{x, x' \in \mathsf{X}} \sum_{y \in \mathsf{Y}} \mu K(y) K_{\mu}^*(x|y) K_{\mu}^*(x'|y) f(x) f(x') \\ &= \sum_{x, x' \in \mathsf{X}} \sum_{y \in \mathsf{Y}} \mu(x) K(y|x) K_{\mu}^*(x'|y) f(x) f(x') \\ &= \sum_{x, x' \in \mathsf{X}} \mu(x) M(x'|x) f(x) f(x') \\ &= \mathbb{E}[f(X) f(X')]. \end{aligned}$$

Therefore,

$$\begin{aligned} \mathcal{E}_{\mu, M}(f, f) &= \frac{1}{2} \mathbb{E}[(f(X) - f(X'))^2] \\ &= \mathbb{E}[f^2(X)] - \mathbb{E}[f(X) f(X')] \\ &= \mathbb{E}[f^2(X)] - \mathbb{E}[(\mathbb{E}[f(X)|Y])^2] \\ &= \text{Var}[f(X)] - \text{Var}[K_{\mu}^* f(Y)]. \end{aligned}$$

□

**Example 4.1** (Doubly symmetric binary source). Consider the case  $\mathsf{X} = \{0, 1\}$ ,  $\mu = \text{Bern}(1/2)$ ,  $M = \text{BSC}(\varepsilon)$  with  $\varepsilon \leq 1/2$ . The resulting exchangeable pair  $(X, X')$  is the *doubly symmetric binary source* (DSBS) with parameter  $\varepsilon$  [3]. It is a matter of simple computation to show that the pair  $(\mu, M)$  factors through  $K = \text{BSC}(\delta(\varepsilon))$  with

$$\delta(\varepsilon) = \frac{1 + \sqrt{1 - 2\varepsilon}}{2}. \quad (4.20)$$

We know that

$$S^2(\text{Bern}(1/2), \text{BSC}(\delta(\varepsilon))) = (1 - 2\delta(\varepsilon))^2 = 1 - 2\varepsilon,$$

which therefore gives

$$\tilde{\lambda}(\text{Bern}(1/2), \text{BSC}(\varepsilon)) = 2\varepsilon.$$

For any  $f \in \mathcal{F}(X)$ , we can compute the Dirichlet form

$$\begin{aligned} \mathcal{E}_{\mu, M}(f, f) &= \frac{1}{2} [\mu(0)M(1|0) + \mu(1)M(0|1)] (f(0) - f(1))^2 \\ &= \frac{\varepsilon}{2} (f(0) - f(1))^2, \end{aligned}$$

which gives us the Poincaré inequality

$$\text{Var}[f(X)] \leq \frac{1}{4} (f(0) - f(1))^2$$

(see, e.g., [26, Ex. 3.9]). Note that this inequality is independent of the crossover probability  $\varepsilon$ .

Next, we consider the case of the log-Sobolev constant  $\tilde{\rho}_1(\mu, M)$ , for which we can only give upper and lower bounds:

**Theorem 4.4.** *The functional  $K \mapsto \rho_1(\mu, K)$  is constant on the collection of all channels  $K$  such that  $M = K_\mu^* \circ K$ , where it takes the value  $4\tilde{\rho}_1(\mu, M)$ . Moreover, if  $(\mu, M)$  factors through  $K$ , the log-Sobolev constant  $\tilde{\rho}_1(\mu, M)$  satisfies*

$$1 - \eta(\mu, K) \leq 4\tilde{\rho}_1(\mu, M) \leq 1 - \frac{(1 - \log 2) \log 2}{2} \eta(\mu, K). \quad (4.21)$$

*Proof.* Choose any channel  $K$ , such that  $M = K_\mu^* \circ K$ . Then, with  $(X, Y) \sim \mu \otimes K$ ,

$$1 - \eta(\mu, K) \leq \rho_1(\mu, K) \quad (4.22)$$

$$= \inf_{f \in \mathcal{F}_*(X)} \frac{\mathcal{E}(f(X), \log f(X)|Y)}{\text{Ent}[f(X)]} \quad (4.23)$$

$$= \inf_{f \in \mathcal{F}_*(X)} \frac{\mathcal{E}_{\mu, M}(f, \log f)}{\text{Ent}[f(X)]} \quad (4.24)$$

$$= 4\tilde{\rho}_1(\mu, M), \quad (4.25)$$

where (4.22) is by Proposition 4.3; (4.23) is by (4.8); (4.24) is by Proposition 4.4; (4.25) is by definition of  $\tilde{\rho}_1(\mu, M)$ .

This proves the first inequality in (4.9). The second inequality follows from the upper bound on  $\rho_1(\mu, K)$  in Proposition 4.3, as well as from the just established fact that  $\rho_1(\mu, K) = 4\tilde{\rho}_1(\mu, M)$ .  $\square$

**Example 4.2** (Doubly symmetric binary source, continued). Consider again the case of the DSBS with parameter  $\varepsilon \leq 1/2$ . From the previous example, we know that  $(\mu, M)$  factors through  $K = \text{BSC}(\delta(\varepsilon))$  with crossover probability  $\delta(\varepsilon)$  given by (4.20). For this channel, we have

$$\eta(\text{Bern}(1/2), \text{BSC}(\delta(\varepsilon))) = 1 - 2\varepsilon.$$

Applying this and Theorem 4.4, we get the following upper and lower bounds on the log-Sobolev constant  $\tilde{\rho}_1$ :

$$\frac{\varepsilon}{2} \leq \tilde{\rho}_1(\text{Bern}(1/2), \text{BSC}(\varepsilon)) \leq \frac{1}{4} \left[ 1 - \frac{(1 - \log 2) \log 2}{2} (1 - 2\varepsilon) \right].$$

Unfortunately, neither of the bounds is tight, since the log-Sobolev constant in this case is known exactly:  $\tilde{\rho}_1(\text{Bern}(1/2), \text{BSC}(\varepsilon)) = \varepsilon$  [26, Ex. 3.9]. A sharp bound can be obtained from the monotonicity property (4.18) of the log-Sobolev constants:

$$\begin{aligned} \tilde{\rho}_1(\text{Bern}(1/2), \text{BSC}(\varepsilon)) &\leq \tilde{\rho}_0(\text{Bern}(1/2), \text{BSC}(\varepsilon)) \\ &= \frac{1}{2} \tilde{\lambda}(\text{Bern}(1/2), \text{BSC}(\varepsilon)) \\ &= \varepsilon. \end{aligned}$$

Finally, we consider the log-Sobolev constant  $\tilde{\rho}_2(\mu, M)$ :

**Theorem 4.5.** *For any channel  $K \in \mathcal{M}(Y|X)$  such that  $M = K_\mu^* \circ K$ ,*

$$\eta(\mu, K) \leq 1 - \tilde{\rho}_2(\mu, M). \quad (4.26)$$

*Proof.* We use the following delicate convexity bound for the function  $\Phi(u) = u \log u$  [8]:

$$\Phi(u) \geq \Phi(v) + (1 + \log v)(u - v) + (\sqrt{u} - \sqrt{v})^2, \quad \forall u, v \geq 0. \quad (4.27)$$

Let  $(X, Y)$  be a random pair with law  $\mu \otimes K$ . Fix any function  $f \in \mathcal{F}_*^0(X)$  with  $\mathbb{E}[f(X)] = 1$  and use the bound (4.27) to get

$$\Phi(f(x)) \geq \Phi(K_\mu^* f(y)) + (1 + \log K_\mu^* f(y)) (f(x) - K_\mu^* f(y)) + (\sqrt{f(x)} - \sqrt{K_\mu^* f(y)})^2 \quad (4.28)$$

Taking conditional expectation  $\mathbb{E}[\cdot|Y]$  of both sides of (4.28), we obtain

$$\begin{aligned} \mathbb{E}[\Phi(f(X))|Y] &\geq \Phi(\mathbb{E}[f(X)|Y]) + \mathbb{E}\left[ (\sqrt{f(X)} - \sqrt{\mathbb{E}[f(X)|Y]})^2 | Y \right] \\ &\geq \Phi(\mathbb{E}[f(X)|Y]) + \mathbb{E}\left[ (\sqrt{f(X)} - \mathbb{E}[\sqrt{f(X)}|Y])^2 | Y \right], \end{aligned}$$

where we have used the fact that

$$\mathbb{E}[(U - \mathbb{E}[U|Y])^2 | Y] = \inf_{f \in \mathcal{F}(Y)} \mathbb{E}[(U - f(Y))^2 | Y].$$

for any real-valued random variable  $U$  jointly distributed with  $Y$ . Next we take the expectation w.r.t.  $Y$  to get

$$\begin{aligned} \text{Ent}[f(X)] &\geq \text{Ent}[\mathbb{E}[f(X)|Y]] + \mathcal{E}(\sqrt{f(X)}, \sqrt{f(X)}|Y) \\ &= \text{Ent}[\mathbb{E}[f(X)|Y]] + \mathcal{E}_{\mu, M}(\sqrt{f}, \sqrt{f}) \end{aligned}$$

where we have used the fact that  $\text{Ent}[U] = \mathbb{E}[\Phi(U)]$  for all nonnegative random variables  $U$  with  $\mathbb{E}U = 1$ , as well as Proposition 4.4. Using this and the definition of  $\tilde{\rho}_2(\mu, M)$ , we get

$$\text{Ent}[\mathbb{E}[f(X)|Y]] \leq (1 - \tilde{\rho}_2(\mu, M)) \text{Ent}[f(X)].$$

Since  $f$  was arbitrary, we get the bound (4.26).  $\square$

### 4.3 The gap between SDPI and $\Phi$ -Sobolev

As evident from the proof of Theorem 4.1, we need to invoke Jensen's inequality in order to pass from a  $\Phi$ -entropy SDPI to a  $\Phi$ -Sobolev inequality. This observation prompts us to investigate the gap between these two inequalities:

**Theorem 4.6.** *If  $\Phi$  satisfies the conditions of Theorem 4.1, then for any  $f \in \mathcal{F}_*^0(\mathbf{X})$*

$$\mathcal{E}(f(X), \Psi \circ f(X)|Y) = \mathbb{E}[\text{Ent}_\Phi[f(X)|Y]] + \mathbb{E}[f(X) \text{Ent}_{-\Psi}[f(X)|Y]] \quad (4.29)$$

*Therefore, if  $\eta_\Phi(P_X, P_{Y|X}) \leq c$  for some  $0 \leq c < 1$ , and if the function  $u \mapsto -\Psi(u)$  is strictly convex at  $u = 1$ , then the  $\Phi$ -Sobolev inequality*

$$\text{Ent}_\Phi[f(X)] \leq \frac{1}{1-c} \mathcal{E}(f(X), \Psi \circ f(X)|Y) \quad (4.30)$$

*is strict for any nonconstant  $f$ . If  $\Psi$  is affine, then  $\mathcal{E}(f(X), \Psi \circ f(X)|Y) = \mathbb{E}[\text{Ent}_\Phi[f(X)|Y]]$ , and in that case  $\eta_\Phi(P_X, P_{Y|X}) \leq c$  is equivalent to (4.30).*

**Remark 4.2.** When  $\Psi$  is affine,  $\Phi$  is of the form  $\Phi(u) = au^2 + bu + c$  for some  $a \geq 0, b, c \in \mathbb{R}$ . Thus, the SDPI for  $\chi^2$ -divergence is equivalent to the corresponding  $\Phi$ -Sobolev inequality (which in this case is precisely the Poincaré inequality).  $\diamond$

*Proof.* By definition of  $\mathcal{E}$ ,

$$\begin{aligned} & \mathcal{E}(f(X), \Psi \circ f(X)|Y) \\ &= \mathbb{E}\left[f(X)\Psi(f(X)) - \mathbb{E}[f(X)|Y]\mathbb{E}[\Psi(f(X)|Y)]\right] \\ &= \mathbb{E}[f(X)\Psi(f(X))] - \mathbb{E}[\mathbb{E}[f(X)|Y]\Psi(\mathbb{E}[f(X)|Y])] \\ &\quad + \mathbb{E}[\mathbb{E}[f(X)|Y]\Psi(\mathbb{E}[f(X)|Y])] - \mathbb{E}[\mathbb{E}[f(X)|Y]\mathbb{E}[\Psi(f(X)|Y)]] \\ &= \mathbb{E}[\text{Ent}_\Phi[f(X)|Y]] + \mathbb{E}\left[f(X)\{\mathbb{E}[-\Psi(f(X)|Y)] - (-\Psi(\mathbb{E}[f(X)|Y]))\}\right]. \end{aligned} \quad (4.31)$$

Since  $-\Psi$  is convex, we recognize the quantity in the curly braces in (4.31) as the conditional entropy  $\text{Ent}_{-\Psi}[f(X)|Y]$ . This proves (4.29).

From (4.29) we see that  $\mathcal{E}(f(X), \Psi \circ f(X)|Y) = \mathbb{E}[\text{Ent}_\Phi[f(X)|Y]]$  for a given nonconstant  $f \in \mathcal{F}_*^0(\mathbf{X})$  if and only if  $\text{Ent}_{-\Psi}[f(X)|Y] = 0$  a.s. If  $-\Psi$  is strictly convex at 1, then  $\text{Ent}_{-\Psi}[U] = 0$  if and only if  $U$  is a.s. constant; thus, in this case, the inequality (4.30) is strict for any nonconstant  $f$ . If  $\Psi$  is affine, then  $\text{Ent}_{-\Psi}[U] = 0$  for all  $U$ , so in that case  $\mathcal{E}(f(X), \Psi \circ f(X)|Y) = \mathbb{E}[\text{Ent}_\Phi[f(X)|Y]]$ .  $\square$

As a corollary, we obtain the following useful formula that expresses the covariance between  $f(X)$  and  $\Psi \circ f(X)$  in terms of entropies:

**Corollary 4.1.**

$$\text{Cov}[f(X), \Psi \circ f(X)] = \text{Ent}_\Phi[f(X)] + \mathbb{E}[f(X) \text{Ent}_{-\Psi}[f(X)]]. \quad (4.32)$$

*Proof.* Consider any pair  $(X, Y)$ , where  $Y$  is independent of  $X$ . In that case,  $\mathcal{E}(f(X), g(X)|Y) = \text{Cov}[f(X), g(X)]$  for any pair  $f, g \in \mathcal{F}(\mathbf{X})$ , whereas  $\text{Ent}_\Phi[f(X)|Y] = \text{Ent}_\Phi[f(X)]$  for any  $\Phi \in \mathcal{F}$ . The formula (4.32) follows from these observations.  $\square$

## 5 Some applications

### 5.1 Concentration inequalities

One of the main uses of logarithmic Sobolev inequalities is in the context of concentration inequalities: Given a probability space  $(X, \mu)$  and a function  $f \in \mathcal{F}(X)$ , the objective is to obtain tight upper bounds on the deviation probabilities  $\mathbb{P}[f(X) - \mathbb{E}f(X) \geq t]$  for  $t \geq 0$ , where  $X \sim \mu$ . A general procedure that allows one to pass from a suitable log-Sobolev inequality to a Gaussian tail bound of the form

$$\mathbb{P}[f(X) - \mathbb{E}f(X) \geq t] \leq e^{-\kappa t^2}, \quad t \geq 0 \quad (5.1)$$

for some  $\kappa > 0$  and for all  $f$  in a suitable subset of  $\mathcal{F}(X)$  is called the *Herbst argument* [31, 57, 58], and can be summarized as follows (see, e.g., [58, Chap. 3]):

We start with a pair  $(\mathcal{A}, \Gamma)$ , where:

1.  $\mathcal{A} \subseteq \mathcal{F}(X)$  is a class of real-valued functions on  $X$ , such that  $af + b \in \mathcal{A}$  for all  $f \in \mathcal{A}$  and all  $a \geq 0, b \in \mathbb{R}$ .
2.  $\Gamma : \mathcal{A} \rightarrow \mathcal{F}_*^0(X)$  is an operator with the property that  $\Gamma(af + b) = a\Gamma f$  for all  $f \in \mathcal{A}$  and all  $a \geq 0, b \in \mathbb{R}$ .

We then say that  $\mu$  satisfies a modified log-Sobolev inequality with constant  $c > 0$  on  $(\mathcal{A}, \Gamma)$  if

$$\text{Ent}[e^{f(X)}] \leq \frac{c}{2} \mathbb{E} \left[ e^{f(X)} |\Gamma f(X)|^2 \right], \quad \forall f \in \mathcal{A}. \quad (5.2)$$

Here is how we pass from (5.2) to a Gaussian tail bound of the form (5.1). Without loss of generality, we may assume that  $\mathbb{E}[f(X)] = 0$ . For any  $\lambda \geq 0$ ,  $\lambda f \in \mathcal{A}$  and  $\Gamma(\lambda f) = \lambda \Gamma f$ . Therefore, replacing  $f$  with  $\lambda f$  in (5.2), we arrive at

$$\begin{aligned} \text{Ent} \left[ e^{\lambda f(X)} \right] &\leq \frac{c\lambda^2}{2} \mathbb{E} \left[ e^{\lambda f(X)} |\Gamma f(X)|^2 \right] \\ &\leq \frac{c\|\Gamma f\|_\infty^2 \lambda^2}{2} \mathbb{E} \left[ e^{\lambda f(X)} \right], \end{aligned} \quad (5.3)$$

where  $\|\Gamma f\|_\infty \triangleq \sup_{x \in X} |\Gamma f(x)|$ . If we define the tilted distribution  $d\mu^{(\lambda)} \triangleq e^{\lambda f} d\mu / \mathbb{E}[e^{\lambda f(X)}]$ , then

$$D(\lambda) \triangleq D(\mu^{(\lambda)} \| \mu) = \frac{\text{Ent} \left[ e^{\lambda f(X)} \right]}{\mathbb{E} \left[ e^{\lambda f(X)} \right]}.$$

Therefore, from (5.3) we get

$$D(\lambda) \leq \frac{c\|\Gamma f\|_\infty^2 \lambda^2}{2}, \quad \forall \lambda \geq 0.$$

On the other hand, if we define the logarithmic moment-generating function  $\Lambda(\lambda) \triangleq \log \mathbb{E}[e^{\lambda f(X)}]$ , then it is a matter of simple calculus to show that

$$D(\lambda) = \lambda^2 \frac{d}{d\lambda} \left( \frac{\Lambda(\lambda)}{\lambda} \right). \quad (5.4)$$

Combining (5.3) and (5.4), we get the differential inequality

$$\frac{d}{d\lambda} \left( \frac{\Lambda(\lambda)}{\lambda} \right) \leq \frac{c \|\Gamma f\|_\infty^2}{2},$$

which can be integrated to give  $\Lambda(\lambda) \leq \frac{c \|\Gamma f\|_\infty^2 \lambda^2}{2}$ . This shows that  $f$  is  $v$ -subgaussian with  $v = c \|\Gamma f\|_\infty^2$ , and therefore it satisfies (5.1) with  $\kappa = 1/2v = 1/(2c \|\Gamma f\|_\infty^2)$  (cf. Section 3.4). Effectively,  $\|\Gamma f\|_\infty$  is a measure of the ‘‘variability’’ of  $f$ .

We now show that we can use any reversible Markov kernel  $M$  on  $\mathsf{X}$  that leaves  $\mu$  invariant as a yardstick for measuring the variability of functions in  $\mathcal{F}(\mathsf{X})$ , and that the constant  $c$  in the log-Sobolev inequality (5.2) can be expressed in terms of the relative-entropy SDPI constants  $\eta(\mu, K)$ , where  $K$  runs over all factorizations  $M = K_\mu^* K$ . Following Houdré and Tetali [59], let us define the  $\ell_2$  positive discrete gradient operator  $\nabla_2^+ : \mathcal{F}(\mathsf{X}) \rightarrow \mathcal{F}_*^0(\mathsf{X})$  via

$$\nabla_2^+ f(x) \triangleq \left( \sum_{x' \in \mathsf{X}} M(x'|x) (f(x) - f(x'))_+^2 \right)^{1/2}.$$

It is easy to see that the pair  $(\mathcal{A}, \Gamma) = (\mathcal{F}(\mathsf{X}), \nabla_2^+)$  satisfies the requirements 1 and 2 listed in the preceding paragraph.

**Theorem 5.1.** *Consider a pair  $(\mu, M) \in \mathcal{P}_*(\mathsf{X}) \times \mathcal{M}(\mathsf{Y}|\mathsf{X})$ , where  $M$  is reversible w.r.t.  $\mu$ . Then the following modified log-Sobolev inequality holds for all  $f \in \mathcal{F}(\mathsf{X})$ :*

$$\text{Ent} \left[ e^{f(X)} \right] \leq \frac{c}{2} \mathbb{E} \left[ e^{f(X)} |\nabla_2^+ f(X)|^2 \right],$$

where  $X \sim \mu$ , and

$$c = \inf \left\{ \frac{2}{1 - \eta(\mu, K)} : M = K_\mu^* K \right\}. \quad (5.5)$$

*Proof.* Suppose that  $M$  factors through some channel  $K \in \mathcal{M}(\mathsf{Y}|\mathsf{X})$ . As before, let  $(X, X')$  be an exchangeable pair with law  $\mu \otimes K^* K \equiv \mu \otimes M$ . Applying Theorem 4.2 to  $\Phi(u) = u \log u$  and  $F(u) = e^u$ , we conclude that any  $f \in \mathcal{F}(\mathsf{X})$  satisfies

$$\begin{aligned} \text{Ent} \left[ e^{f(X)} \right] &\leq \frac{1}{1 - \eta(\mu, K)} \mathbb{E} \left[ e^{f(X)} (f(X) - f(X'))_+^2 \right] \\ &= \frac{1}{1 - \eta(\mu, K)} \sum_{x \in \mathsf{X}} \mu(x) e^{f(x)} \sum_{x' \in \mathsf{X}} K^* K(x'|x) (f(x) - f(x'))_+^2 \\ &= \frac{1}{1 - \eta(\mu, K)} \sum_{x \in \mathsf{X}} \mu(x) e^{f(x)} \sum_{x' \in \mathsf{X}} M(x'|x) (f(x) - f(x'))_+^2 \\ &= \frac{1}{1 - \eta(\mu, K)} \sum_{x \in \mathsf{X}} \mu(x) e^{f(x)} |\nabla_2^+ f(x)|^2 \\ &= \frac{1}{1 - \eta(\mu, K)} \mathbb{E} \left[ e^{f(X)} |\nabla_2^+ f(X)|^2 \right]. \end{aligned}$$

Optimizing over the choice of  $K$ , we see that (5.2) holds with  $c$  given by (5.5).  $\square$

## 5.2 Contraction of mutual information in a Markov chain

Consider a Markov chain  $U \rightarrow X \rightarrow Y$ , where the joint law  $P_{XY}$  is fixed, while the alphabet  $\mathsf{U}$  of  $U$  and the conditional distribution  $P_{U|X}$  are allowed to vary arbitrarily. By the data processing inequality for the mutual information,  $I(U; Y) \leq I(U; X)$  for any choice of  $P_{U|X}$ . The question is: what is the maximum value of the ratio  $\frac{I(U; Y)}{I(U; X)}$  that can be achieved by any choice of  $P_{U|X}$ ? The following claim was made by Erkip and Cover [60]:

$$\sup_{P_{U|X}} \frac{I(U; Y)}{I(U; X)} = S^2(P_X, P_{Y|X}), \quad (5.6)$$

where  $S^2$  is the squared maximal correlation (see Section 3.2). However, Anantharam et al. in a recent preprint [15] pointed out a flaw in the proof of (5.6), and showed instead that

$$\sup_{P_{U|X}} \frac{I(U; Y)}{I(U; X)} = \eta(P_X, P_{Y|X}), \quad (5.7)$$

where  $\eta$  is the relative-entropy SDPI constant. Moreover, they provided an explicit example of a source-channel pair  $(P_X, P_{Y|X})$ , for which the mutual-information ratio on the left-hand sides of Eqs. (5.6) and (5.7) is strictly larger than  $S^2(P_X, P_{Y|X})$ .

We will now present a generalization of the result of Anantharam et al., and show, as a consequence, that  $S^2(P_X, P_{Y|X})$  can indeed be expressed as a supremum of the ratio of two information-like quantities pertaining to the Markov chain  $U \rightarrow X \rightarrow Y$  with an arbitrary choice of  $P_{U|X}$ . Fix a function  $\Phi \in \mathcal{F}$ . Given a random pair  $(U, V)$ , we define the *mutual  $\Phi$ -information* [9]<sup>6</sup> as

$$\begin{aligned} I_\Phi(U; V) &\triangleq D_\Phi(P_{UV} \| P_U \otimes P_V) \\ &= \int P_U(du) \int P_V(dv) \Phi \left( \frac{dP_{V|U=u}}{dP_V}(v) \right) \\ &= \int P_U(du) D_\Phi(P_{V|U=u} \| P_V). \end{aligned} \quad (5.8)$$

If  $U$  and  $V$  are related via a Markov kernel  $K$  (i.e.,  $P_{UV} = P_U \otimes K$ ), we may also use the notation  $I_\Phi(P_U, K)$  to indicate the fact that the  $\Phi$ -information is a functional of the source distribution and the kernel that generates the random output given the input.

**Theorem 5.2.** *If  $\Phi \in \mathcal{F}$  is differentiable, and its derivative is uniformly bounded in some neighborhood of 1, then*

$$\sup_{P_{U|X}} \frac{I_\Phi(U; Y)}{I_\Phi(U; X)} = \eta_\Phi(P_X, P_{Y|X}).$$

*Proof.* Define a probability measure  $Q \in \mathcal{P}(\mathsf{U})$  by

$$Q(u) \triangleq \frac{P_U(u) D_\Phi(P_{X|U=u} \| P_X)}{\sum_{u \in \mathsf{U}} P_U(u) D_\Phi(P_{X|U=u} \| P_X)}.$$

<sup>6</sup>Palomar and Verdú [61] define  $\Phi$ -information between  $U$  and  $V$  as  $D_\Phi(P_U \otimes P_V \| P_{UV})$ . Their definition is equivalent to Eq. (5.8) if we replace  $\Phi$  with its *Csiszár conjugate*  $\Phi^*(u) \triangleq u\Phi(1/u)$  [5].

This measure is supported on the set  $\tilde{\mathcal{U}} \triangleq \{u \in \mathcal{U} : D_{\Phi}(P_{X|U=u} \| P_X) > 0\}$ . From data processing, we have the inclusion  $\{u \in \mathcal{U} : D_{\Phi}(P_{Y|U=u} \| P_Y) > 0\} \subseteq \tilde{\mathcal{U}}$ . Taking all of this into account, we can write

$$\begin{aligned} \frac{I_{\Phi}(U; Y)}{I_{\Phi}(U; X)} &= \frac{\sum_{u \in \tilde{\mathcal{U}}} P_U(u) D_{\Phi}(P_{Y|U=u} \| P_Y)}{\sum_{u \in \tilde{\mathcal{U}}} P_U(u) D_{\Phi}(P_{X|U=u} \| P_X)} \\ &= \sum_{u \in \tilde{\mathcal{U}}} Q(u) \frac{D_{\Phi}(P_{Y|U=u} \| P_Y)}{D_{\Phi}(P_{X|U=u} \| P_X)} \\ &\leq \max_{u \in \tilde{\mathcal{U}}} \frac{D_{\Phi}(P_{Y|U=u} \| P_Y)}{D_{\Phi}(P_{X|U=u} \| P_X)} \\ &= \max_{u \in \tilde{\mathcal{U}}} \frac{D_{\Phi}(P_{X|U=u} P_{Y|X} \| P_X P_{Y|X})}{D_{\Phi}(P_{X|U=u} \| P_X)} \\ &\leq \eta_{\Phi}(P_X, P_{Y|X}). \end{aligned}$$

To prove the reverse inequality, we adopt the construction from [15]. Fix an arbitrary  $Q_X \in \mathcal{P}(\mathcal{X})$ . For any  $\varepsilon \in (0, 1)$  small enough so that  $\nu \triangleq P_X - \varepsilon Q_X$  is a nonnegative measure, let  $P_U^{(\varepsilon)} = \text{Bern}(\varepsilon)$  and define  $P_{X|U}^{(\varepsilon)}$  by

$$P_{X|U=0}^{(\varepsilon)} = Q_X, \quad P_{X|U=1}^{(\varepsilon)} = \frac{\nu}{\bar{\varepsilon}}.$$

With these choices,  $P_U^{(\varepsilon)} P_{X|U}^{(\varepsilon)} = \varepsilon P_{X|U=0}^{(\varepsilon)} + \bar{\varepsilon} P_{X|U=1}^{(\varepsilon)} = \varepsilon Q_X + P_X - \varepsilon P_X = P_X$ . For any  $\eta > 0$ , define the function

$$L_{\eta}(\varepsilon) \triangleq I_{\Phi}\left(P_U^{(\varepsilon)}, P_{Y|X} \circ P_{X|U}^{(\varepsilon)}\right) - \eta I_{\Phi}\left(P_U^{(\varepsilon)}, P_{X|U}^{(\varepsilon)}\right).$$

A simple calculation gives

$$\begin{aligned} I_{\Phi}\left(P_U^{(\varepsilon)}, P_{Y|X} \circ P_{X|U}^{(\varepsilon)}\right) &= \varepsilon D_{\Phi}(P_{X|U=0}^{(\varepsilon)} P_{Y|X} \| P_Y) + \bar{\varepsilon} D_{\Phi}(P_{X|U=1}^{(\varepsilon)} P_{Y|X} \| P_Y) \\ &= \varepsilon D_{\Phi}(Q_X P_{Y|X} \| P_Y) + \bar{\varepsilon} D_{\Phi}\left(\frac{P_X - \varepsilon Q_X}{\bar{\varepsilon}} P_{Y|X} \| P_Y\right) \\ &= \varepsilon D_{\Phi}(Q_X P_{Y|X} \| P_Y) + \bar{\varepsilon} D_{\Phi}\left(\frac{P_Y - \varepsilon Q_X P_{Y|X}}{\bar{\varepsilon}} \| P_Y\right), \end{aligned}$$

where in the last line we have used the fact that any Markov kernel extends to a linear map on signed measures. Similarly,

$$\begin{aligned} I_{\Phi}\left(P_U^{(\varepsilon)}, P_{X|U}^{(\varepsilon)}\right) &= \varepsilon D_{\Phi}(P_{X|U=0}^{(\varepsilon)} \| P_X) + \bar{\varepsilon} D_{\Phi}(P_{X|U=1}^{(\varepsilon)} \| P_X) \\ &= \varepsilon D_{\Phi}(Q_X \| P_X) + \bar{\varepsilon} D_{\Phi}\left(\frac{P_X - \varepsilon Q_X}{\bar{\varepsilon}} \| P_X\right). \end{aligned}$$

Let  $f = \frac{dQ_X}{dP_X}$  and  $g^{(\varepsilon)} = \frac{1-\varepsilon f}{\bar{\varepsilon}}$ . Then, by virtue of our choice of  $\varepsilon$ ,  $g^{(\varepsilon)} \in \mathcal{F}_*^0(\mathcal{X})$ , and  $\mathbb{E}[g^{(\varepsilon)}(X)] = 1$ . With these definitions, we can rewrite the above expressions as

$$I_{\Phi}\left(P_U^{(\varepsilon)}, P_{Y|X} \circ P_{X|U}^{(\varepsilon)}\right) = \varepsilon \text{Ent}_{\Phi}\left[P_{Y|X}^* f(Y)\right] + \bar{\varepsilon} \text{Ent}_{\Phi}\left[P_{Y|X}^* g^{(\varepsilon)}(Y)\right]$$

and

$$I_{\Phi} \left( P_U^{(\varepsilon)}, P_{X|U}^{(\varepsilon)} \right) = \varepsilon \text{Ent}_{\Phi} [f(X)] + \bar{\varepsilon} \text{Ent}_{\Phi} [g^{(\varepsilon)}(X)].$$

Consequently,

$$\begin{aligned} \frac{d}{d\varepsilon} L_{\eta}(\varepsilon) &= \text{Ent}_{\Phi} \left[ P_{Y|X}^* f(Y) \right] - \eta \text{Ent}_{\Phi} [f(X)] + \frac{d}{d\varepsilon} \left\{ \bar{\varepsilon} \left( \text{Ent}_{\Phi} \left[ P_{Y|X}^* g^{(\varepsilon)}(Y) \right] - \eta \text{Ent}_{\Phi} [g^{(\varepsilon)}(X)] \right) \right\} \\ &= \text{Ent}_{\Phi} \left[ P_{Y|X}^* f(Y) \right] - \eta \text{Ent}_{\Phi} [f(X)] + \eta \text{Ent}_{\Phi} [g^{(\varepsilon)}(X)] - \text{Ent}_{\Phi} \left[ P_{Y|X}^* g^{(\varepsilon)}(Y) \right] \\ &\quad + \bar{\varepsilon} \frac{d}{d\varepsilon} \left\{ \text{Ent}_{\Phi} \left[ P_{Y|X}^* g^{(\varepsilon)}(Y) \right] - \eta \text{Ent}_{\Phi} [g^{(\varepsilon)}(X)] \right\}. \end{aligned}$$

Now let us choose  $Q_X$  so that  $\frac{D_{\Phi}(Q_Y \| P_Y)}{D_{\Phi}(Q_X \| P_X)} > \eta_{\Phi}(P_X, P_{Y|X}) - \delta$  for some small  $\delta > 0$ . Then, for any  $\eta < \eta_{\Phi}(P_X, P_{Y|X}) - \delta$  we have

$$\left. \frac{d}{d\varepsilon} L_{\eta}(\varepsilon) \right|_{\varepsilon=0} = \text{Ent}_{\Phi} \left[ P_{Y|X}^* f(Y) \right] - \eta \text{Ent}_{\Phi} [f(X)] > 0,$$

where we have used Lemma A.5 in Appendix A, and where the strict inequality holds due to our choice of  $\eta$ . Thus, the function  $\varepsilon \mapsto L_{\eta}(\varepsilon)$  is strictly increasing in some neighborhood of 0. Since  $L_{\eta}(0) = 0$ , there exists some value  $\varepsilon_0 > 0$ , such that  $L_{\eta}(\varepsilon_0) > 0$ , i.e.,

$$\sup_{P_{U|X}} \frac{I_{\Phi}(U; Y)}{I_{\Phi}(U; X)} \geq \frac{I_{\Phi} \left( P_U^{(\varepsilon_0)}, P_{Y|X} \circ P_{X|U}^{(\varepsilon_0)} \right)}{I_{\Phi} \left( P_U^{(\varepsilon_0)}, P_{X|U}^{(\varepsilon_0)} \right)} > \eta.$$

Since this holds for all  $0 < \eta < \eta_{\Phi}(P_X, P_{Y|X}) - \delta$ , and  $\delta > 0$  was arbitrary, we conclude, upon taking  $\delta \searrow 0$ , that

$$\sup_{P_{U|X}} \frac{I_{\Phi}(U; Y)}{I_{\Phi}(U; X)} \geq \eta_{\Phi}(P_X, P_{Y|X}).$$

Since we already established the reverse inequality, the theorem is proved.  $\square$

Thus, if  $\Phi(u) = u \log u$ , we recover the result of Anantharam et al. [15]; on the other hand, choosing  $\Phi(u) = (u - 1)^2$ , we can express the squared maximal correlation  $S^2(P_X, P_{Y|X})$  as

$$S^2(P_X, P_{Y|X}) = \sup_{P_{U|X}} \frac{I_{\chi^2}(U; Y)}{I_{\chi^2}(U; X)},$$

where the  $\chi^2$ -information  $I_{\chi^2}(U; V)$  is the variance of the Radon–Nikodym derivative  $\frac{dP_{UV}}{d(P_U \otimes P_V)}$  w.r.t. the product distribution  $P_U \otimes P_V$ . We also have the following result:

**Corollary 5.1.** *Let  $(X, Y)$  be a random pair taking values in a finite product space  $\mathsf{X} \times \mathsf{Y}$ , such that  $P_X \in \mathcal{P}_*(\mathsf{X})$  and  $P_Y \in \mathcal{P}_*(\mathsf{Y})$ . Then for any  $\Phi \in \mathcal{F}$  satisfying the conditions of Theorem 5.2,*

$$\eta_{\Phi}(P_X, P_{Y|X}) \eta_{\Phi}(P_Y, P_{X|Y}) \geq \frac{I_{\Phi}(X; X')}{I_{\Phi}(X; X')} \vee \frac{I_{\Phi}(Y; Y')}{I_{\Phi}(Y; Y')}, \quad (5.9)$$

where  $(X, X')$  is an exchangeable pair generated according to the Markov chain

$$X \xrightarrow{P_{Y|X}} Y \xrightarrow{P_{X|Y}} X' \quad (5.10)$$

and  $(Y, Y')$  is an exchangeable pair generated according to the Markov chain

$$Y \xrightarrow{P_{X|Y}} X \xrightarrow{P_{Y|X'}} Y'.$$

*Proof.* Applying Theorem 5.2 to the Markov chain (5.10) gives

$$\eta_{\Phi}(P_Y, P_{X|Y}) \geq \frac{I_{\Phi}(X; X')}{I_{\Phi}(X; Y)}.$$

On the other hand,

$$\begin{aligned} I_{\Phi}(X; Y) &= \sum_{x \in \mathbf{X}} P_X(x) D_{\Phi}(P_{Y|X=x} \| P_Y) \\ &= \sum_{x \in \mathbf{X}} P_X(x) D_{\Phi}(\delta_x P_{Y|X} \| P_X P_{Y|X}) \\ &\leq \eta_{\Phi}(P_X, P_{Y|X}) \sum_{x \in \mathbf{X}} P_X(x) D_{\Phi}(\delta_x \| P_X) \\ &= \eta_{\Phi}(P_X, P_{Y|X}) I_{\Phi}(X; X), \end{aligned}$$

where  $\delta_x$  denotes the Dirac measure located at  $x$ . Combining these estimates gives (5.9). Interchanging the roles of  $X$  and  $Y$ , we obtain an analogous bound involving  $I_{\Phi}(Y; Y')$  and  $I_{\Phi}(Y; Y)$ .  $\square$

For example, if  $\Phi(u) = u \log u$ , the bound (5.9) becomes

$$\eta(P_X, P_{Y|X}) \eta(P_Y, P_{X|Y}) \geq \frac{I(X; X')}{H(X)} \sqrt{\frac{I(Y; Y')}{H(Y)}},$$

where  $H(X)$  is the usual Shannon entropy of  $X$ . If  $\Phi(u) = (u - 1)^2$ , then we have

$$S^2(P_X, P_{Y|X}) S^2(P_Y, P_{X|Y}) \geq \frac{I_{\chi^2}(X; X')}{|\mathbf{X}| - 1} \sqrt{\frac{I_{\chi^2}(Y; Y')}{|\mathbf{Y}| - 1}}.$$

Corollary 5.1 may be useful for obtaining lower bounds on the mixing time of Gibbs samplers. It also shows that the modified log-Sobolev constant  $c$  defined in (5.5) is bounded from below as

$$c \geq 2H(X)H(X|X'),$$

where  $(X, X')$  is an exchangeable pair with  $P_X = \mu$  and  $P_{X'|X} = M$ .

### 5.3 Fastest mixing Markov chain on a graph

Let  $G = (\mathbf{V}, \mathbf{E})$  be a connected undirected graph with vertex set  $\mathbf{V}$  and edge set  $\mathbf{E} \subseteq \mathbf{V} \times \mathbf{V}$ . Since  $G$  is undirected,  $(x, x') \in \mathbf{E} \Rightarrow (x', x) \in \mathbf{E}$ . We assume that each vertex has a self-loop, i.e.,  $(x, x) \in \mathbf{E}$  for all  $x \in \mathbf{V}$ . Consider a (discrete-time) Markov chain  $\{X_t\}_{t=0,1,\dots}$  with states in  $\mathbf{V}$ , whose one-step transition probability matrix  $K$  has the following properties:

1. It is symmetric, i.e.,  $K(x'|x) = K(x|x')$  for all  $x, x' \in \mathbb{V}$ .
2. It respects the graph structure, i.e.,  $K(x'|x) \neq 0$  only if  $(x, x') \in \mathbb{E}$ .

Let  $\mu$  be the uniform distribution on  $\mathbb{V}$ . The first property of  $K$  implies that it is reversible with respect to  $\mu$ , so that  $\mu = \mu K$ . Let  $\nu$  be the distribution of the initial state  $X_0$ , and let  $\nu_t$  denote the distribution of  $X_t$ , the state at time  $t$ , so that  $\nu_t = \nu K^t$ . If the Markov chain is irreducible and aperiodic (which will be the case if  $K(x|x) > 0$  for all  $x \in \mathbb{V}$ ), then  $\nu_t$  will converge to  $\mu$ . There are multiple ways of quantifying the rate of convergence; we introduce the following definition:

**Definition 5.1.** *Given a convex function  $\Phi \in \mathcal{F}$ , the  $\Phi$ -mixing time of  $K$  is the function  $\tau_\Phi(K, \cdot) : \mathbb{R}^+ \rightarrow \mathbb{N}$ , defined by*

$$\tau_\Phi(K, \varepsilon) \triangleq \min \left\{ t \in \mathbb{N} : \sup_{\nu \in \mathcal{P}(\mathbb{V})} D_\Phi(\nu K^t \| \mu) \leq \varepsilon \right\}$$

Unsurprisingly, the mixing time is controlled by the SDPI constant  $\eta_\Phi(\mu, K)$ :

**Proposition 5.1.** *Suppose  $\Phi(0) < \infty$ , and let  $n = |\mathbb{V}|$ . Then*

$$\tau_\Phi(K, \varepsilon) \leq \frac{\log(D_{\Phi, n}^*/\varepsilon)}{\log(1/\eta_\Phi(\mu, K))}, \quad (5.11)$$

where  $D_{\Phi, n}^* \triangleq \frac{\Phi(n)}{n} + (1 - \frac{1}{n})\Phi(0)$ .

*Proof.* For any  $t \geq 0$  and any  $\nu \in \mathcal{P}(\mathbb{V})$ ,

$$D_\Phi(\nu K^t \| \mu) = D_\Phi(\nu K^t \| \mu K^t) \leq \left( \eta_t(\mu, K) \right)^t D_\Phi(\nu \| \mu).$$

where we have used the fact that  $\mu$  is  $K$ -invariant. Since  $\Phi$ -divergences are convex, and since a convex function on a compact convex set attains its maximum on an extreme point, we have

$$D_\Phi(\nu \| \mu) \leq \max_{x \in \mathbb{V}} D_\Phi(\delta_x \| \mu),$$

where  $\delta_x$  is the Dirac measure located at  $x$ . Moreover, for any  $x \in \mathbb{V}$ ,

$$\begin{aligned} D_\Phi(\delta_x \| \mu) &= \frac{1}{n} \sum_{x' \in \mathbb{V}} \Phi \left( \frac{\delta_x(x')}{1/n} \right) \\ &= \frac{\Phi(n)}{n} + \left( 1 - \frac{1}{n} \right) \Phi(0) \\ &\equiv D_{\Phi, n}^*. \end{aligned}$$

Since  $\nu$  was arbitrary, we have

$$\sup_{\nu \in \mathcal{P}(\mathbb{V})} D_\Phi(\nu K^t \| \mu) \leq D_{\Phi, n}^* \left( \eta_t(\mu, K) \right)^t.$$

Solving for the smallest  $t$  that would make the right-hand side smaller than  $\varepsilon$ , we obtain (5.11).  $\square$

It is customary to fix some value of  $\varepsilon$  (for discrete-time chains, a common choice is  $1/2$ ), and to speak about the scaling of the mixing times in terms of the parameters of the graph and the Markov chain. For example, if  $\Phi(u) = \frac{1}{2}|u - 1|$ , then the chain with one-step transition kernel  $K$  mixes in  $O\left(\frac{1}{\log \vartheta(K)^{-1}}\right)$  steps (TV), where  $\vartheta(K)$  is the Dobrushin coefficient of  $K$ ; for  $\Phi(u) = (u - 1)^2$ , the chain mixes in  $O\left(\frac{\log n}{\log[S^2(\mu, K)]^{-1}}\right)$  steps ( $\chi^2$ ), where  $S^2(\mu, K)$  is the maximal correlation; and for  $\Phi(u) = u \log u$ , the chain mixes in  $O\left(\frac{\log \log n}{\log \eta(\mu, K)^{-1}}\right)$  steps (relative entropy). Thus, if  $\eta_\Phi(\mu, K)$  is small, the corresponding Markov chain will mix faster in the sense that it will take fewer steps for the  $\Phi$ -divergence between the current state distribution and the uniform distribution on  $\mathbf{V}$  to fall below a given value. This motivates the following

**Fastest mixing Markov chain (FMMC) problem:** Let  $\mathcal{M}(G) \subset \mathcal{M}(\mathbf{V}|\mathbf{V})$  be the set of all Markov kernels  $K \in \mathcal{M}(\mathbf{V}|\mathbf{V})$  satisfying the conditions listed in the beginning of this section. For a fixed convex function  $\Phi \in \mathcal{F}$ ,

$$\begin{aligned} & \text{minimize} && \eta_\Phi(\mu, K) \\ & \text{subject to} && K \in \mathcal{M}(G) \end{aligned}$$

**Proposition 5.2.** *For any  $\Phi \in \mathcal{F}$ , the FMMC problem is a convex program.*

*Proof.* The constraint set  $\mathcal{M}(G)$  is convex. To see this, consider any two  $K_1, K_2 \in \mathcal{M}(G)$ , and let  $K = \lambda K_1 + \bar{\lambda} K_2$  for some  $\lambda \in (0, 1)$ . Since both  $K_1$  and  $K_2$  are symmetric, for any pair  $x, x' \in \mathbf{X}$  we have

$$K(x'|x) = \lambda K_1(x'|x) + \bar{\lambda} K_2(x'|x) = \lambda K_1(x|x') + \bar{\lambda} K_2(x|x') = K(x|x').$$

Similarly, suppose that  $(x, x') \notin \mathbf{E}$ . Then  $K_1(x'|x) = K_2(x'|x) = 0$ , so  $K(x'|x) = 0$  as well. Thus,  $K \in \mathcal{M}(G)$ . The objective function  $K \mapsto \eta_\Phi(\mu, K)$  is likewise convex, by Proposition 3.3.  $\square$

For  $\Phi(u) = (u - 1)^2$ , the FMMC problem was studied by Boyd et al. [62], who showed that it can be equivalently represented by a semidefinite program (SDP), for which efficient solvers are available. For a general  $\Phi$ , there is not much one can say without exploiting specific properties of that  $\Phi$  or any symmetries of the graph  $G$ ; however, we can provide bounds on the values of the FMMC problems for different choices of  $\Phi$ . With that in mind, let  $\eta_\Phi^*(G)$  denote the minimum value of the FMMC objective a given choice of  $\Phi$  and  $G$ :

$$\eta_\Phi^*(G) \triangleq \inf_{K \in \mathcal{M}(G)} \eta_\Phi(\mu, K).$$

Then we observe the following:

- $\eta_\Phi^*(G) \leq \inf_{K \in \mathcal{M}(G)} \vartheta(K)$  for any  $\Phi \in \mathcal{F}$ . This follows from the fact that  $\eta_\Phi(\mu, K) \leq \eta_\Phi(K) \leq \vartheta(K)$ , by Theorem 3.1.
- If  $\Phi$  is three times differentiable and  $\Phi''(1) > 0$ , then  $\eta_\Phi^*(G) \geq \eta_{\chi^2}^*(G)$ . This follows from the fact that, for such  $\Phi$ ,  $\eta_\Phi(\mu, K) \geq \eta_{\chi^2}(\mu, K)$  [cf. Theorem 3.3]. The quantity  $\eta_{\chi^2}^*(G)$  and the corresponding convex program were studied extensively by Boyd et al. [62].

The above definition of mixing time can be generalized to any other invariant distribution  $\mu$  on  $\mathbb{V}$ : Let  $\mathcal{M}_\mu(G) \subset \mathcal{M}(\mathbb{V}|\mathbb{V})$  be the set of all Markov kernels  $K$ , such that:

1.  $\mu(x)K(x'|x) = \mu(x')K(x|x')$  for all  $x, x' \in \mathbb{V}$ .
2.  $K(x'|x) \neq 0$  only if  $(x, x') \in \mathbb{E}$ .

Then the same definition of the mixing time applies, and we have the bound

$$\tau_\Phi(K, \varepsilon) \leq \frac{\log(D_{\Phi, \mu}^*/\varepsilon)}{\log(1/\eta_\Phi(\mu, K))},$$

where

$$D_{\Phi, \mu}^* \triangleq \max_{x \in \mathbb{V}} \left\{ \mu(x) \Phi\left(\frac{1}{\mu(x)}\right) + (1 - \mu(x)) \Phi(0) \right\}.$$

We can then consider the appropriate modification of the FMMC problem, and the same arguments as before can be used to show that it is given by a convex program.

#### 5.4 Mixing times of Swendsen-Wang and heat-bath dynamics

Let  $G = (\mathbb{V}, \mathbb{E})$  be an undirected graph without self-loops. In this case, we can identify the edge set of  $G$  with a subset of  $\binom{\mathbb{V}}{2}$ , the set of all two-element subsets of  $\mathbb{V}$ . If two vertices  $u, v \in \mathbb{V}$  are connected by an edge, we will write  $u \leftrightarrow v$ . Fix an integer  $q \geq 2$ , and consider the set  $\mathbb{X} = \mathbb{X}_q = \{1, \dots, q\}^{\mathbb{V}}$  of tuples  $x = (x_v : v \in \mathbb{V})$  with coordinates in  $\{1, \dots, q\}$ . The elements of  $\mathbb{X}$  are called *q-colorings of G*, and we say that  $x \in \mathbb{X}$  is a *proper q-coloring* if  $x_u \neq x_v$  whenever  $u \leftrightarrow v$ .

The problem of computing the number  $P_G(q)$  of proper  $q$ -colorings of an arbitrary  $G$  (or even deciding whether it is nonzero) is intractable, although it is known that  $P_G(q)$  is polynomial in  $q$ . A related problem of drawing a  $q$ -coloring of  $G$  uniformly at random (assuming  $P_G(q) > 0$ ) is also intractable [63]. However, it turns out that the problem of computing (or approximating)  $P_G(q)$  is closely related to the problem of sampling from the so-called *q-state Potts model*, described by the Gibbs distribution

$$\mathbb{P}_{\beta, q}(x) \triangleq \frac{1}{Z(\beta, q)} \exp\left(\beta \sum_{\substack{u, v \in \mathbb{V} \\ u \leftrightarrow v}} \mathbf{1}\{x_u = x_v\}\right), \quad (5.12)$$

where the parameter  $\beta \geq 0$  is called the inverse temperature, and  $Z(\beta, q)$  is the normalization constant known as the partition function. In particular,  $P_G(q) = \lim_{\beta \rightarrow \infty} Z(\beta, q)$ . Direct sampling from  $\mathbb{P}_{\beta, q}$  is also intractable, so one resorts to Markov Chain Monte Carlo (MCMC) methods: Pick a Markov kernel  $K \in \mathcal{M}(\mathbb{X}|\mathbb{X})$  that leaves the Gibbs distribution (5.12) invariant, pick an arbitrary initial configuration  $x_0 \in \mathbb{X}$ , and for each  $t = 0, 1, \dots$  generate a random configuration  $X_{t+1}$  according to  $K(\cdot|X_t)$ . With a good choice of  $K$ , the distribution of  $X_t$  will rapidly converge to  $\mathbb{P}_{\beta, q}$ . Two popular choices of  $K$  are the *heat-bath* (or *Glauber*) *dynamics* and the *Swendsen-Wang dynamics* [64, 65]. They are defined as follows:

**Heat-bath dynamics.** At each time step  $t$ , given the current configuration  $x_t = (x_{v,t})_{v \in \mathbf{V}}$ , we pick a vertex  $v \in \mathbf{V}$  uniformly at random, assign it a new random color  $X_{v,t+1} \in \{1, \dots, q\}$  according to the conditional distribution  $\mathbb{P}_{\beta,q}(X_{v,t+1} | X_t \setminus v = x_t \setminus v)$ , and set  $X_{t+1} \setminus v = x_t \setminus v$ . Here,  $x_t \setminus v = (x_{u,t})_{u \in \mathbf{V} \setminus \{v\}}$  is the time- $t$  configuration of all the vertices except  $v$ . Thus, the transition probabilities of the heat-bath Markov chain are given by the Markov kernel

$$K_{\beta,q}^{\text{HB}}(x'|x) = \frac{1}{|\mathbf{V}|} \sum_{v \in \mathbf{V}} \mathbb{P}_{\beta,q}(x'_v | x \setminus v) \mathbf{1}\{(x') \setminus v = x \setminus v\}. \quad (5.13)$$

**Swendsen-Wang dynamics.** This construction is based on a coupling of the  $q$ -Potts model and the so-called *random-cluster* (or *Fortuin-Kasteleyn*) model on  $G$ . The latter is defined as follows [65]. Let  $\mathbf{Y} = 2^{\mathbf{E}} = \{\mathbf{A} : \mathbf{A} \subseteq \mathbf{E}\}$  and fix a parameter  $p \in (0, 1)$ . Then the random-cluster model is described by the following probability measure on  $\mathbf{Y}$ :

$$\mathbb{Q}_{p,q}(\mathbf{A}) \triangleq \frac{1}{\tilde{Z}(p,q)} \left(\frac{p}{\bar{p}}\right)^{|\mathbf{A}|} q^{C(\mathbf{A})}, \quad \forall \mathbf{A} \in \mathbf{Y} \quad (5.14)$$

where  $\tilde{Z}(\cdot, \cdot)$  is the partition function,  $\bar{p} = 1 - p$ , and  $C(\mathbf{A})$  is the number of connected components of the induced graph  $(\mathbf{V}, \mathbf{A})$ . It can be shown that

$$\tilde{Z}(p, q) = Z(\log(1/\bar{p}), q),$$

where  $Z(\cdot, \cdot)$  is the partition function for the  $q$ -Potts model. Now let  $p = 1 - e^{-\beta}$ , and consider the following probability measure on the Cartesian product  $\mathbf{X} \times \mathbf{Y}$ :

$$\mathbb{M}(x, \mathbf{A}) = \frac{1}{\tilde{Z}(p,q)} \left(\frac{p}{\bar{p}}\right)^{|\mathbf{A}|} \mathbf{1}\{\mathbf{A} \subset \mathbf{E}(x)\} \quad (5.15)$$

$$= \frac{1}{Z(\beta, q)} (e^{\beta} - 1)^{|\mathbf{A}|} \mathbf{1}\{\mathbf{A} \subset \mathbf{E}(x)\}. \quad (5.16)$$

where  $\mathbf{E}(x) \triangleq \{\{u, v\} \in \mathbf{E} : x_u \neq x_v\}$  is the set of edges on which  $x$  violates the proper  $q$ -coloring constraint. It can be shown that  $\mathbb{M}$  is a coupling of  $\mathbb{P}_{\beta,q}$  and  $\mathbb{Q}_{p,q}$  with  $p = 1 - e^{-\beta}$ , i.e., if  $(X, Y)$  is a random pair with law  $\mathbb{M}$ , then  $P_X = \mathbb{P}_{\beta,q}$  and  $P_Y = \mathbb{Q}_{1-e^{-\beta},q}$ .

With these definitions at hand, we can describe the Swendsen-Wang algorithm:

- Start with an arbitrary initial configuration  $x_0 \in \mathbf{X}$
- For each  $t = 0, 1, 2, \dots$ 
  - Draw a random set  $\mathbf{A}_t \in \mathbf{Y}$  according to the conditional distribution  $\mathbb{M}(Y = \cdot | X = x_t)$ .
  - Draw  $X_{t+1}$  from the conditional distribution  $\mathbb{M}(X = \cdot | Y = \mathbf{A}_t)$ .

In words, given  $x_t$ , we draw  $\mathbf{A}_t$  by deleting each edge of  $\mathbf{E}(x_t)$  independently with probability  $p = 1 - e^{-\beta}$ ; given  $\mathbf{A}_t$ , we draw  $X_{t+1}$  by assigning a random color independently to each connected component of  $(\mathbf{V}, \mathbf{A}_t)$  and coloring all vertices in the same component with the same color. Thus, the Swendsen-Wang dynamics is a two-stage Gibbs sampler that generates a trajectory  $\{(X_t, Y_t)\}_{t \geq 0}$  according to

$$\dots \longrightarrow X_t \xrightarrow{\mathbb{M}_{Y|X}} Y_t \xrightarrow{\mathbb{M}_{X|Y}} X_{t+1} \longrightarrow \dots$$

The discrete-time process  $\{X_t\}_{t \geq 0}$  is a Markov chain with one-step transition kernel

$$\begin{aligned} K_{\beta,q}^{\text{SW}}(x'|x) &= \mathbb{M}_{X|Y} \circ \mathbb{M}_{Y|X}(x'|x) \\ &= \sum_{\mathbf{A} \in \mathcal{Y}} \mathbb{M}_{X|Y}(x'|\mathbf{A}) \mathbb{M}_{Y|X}(\mathbf{A}|x). \end{aligned}$$

By construction, the Markov kernel  $K_{\beta,q}^{\text{SW}}$  is reversible w.r.t. the Gibbs measure  $\mathbb{P}_{\beta,q}$ .

With each of these two algorithms, the hope is that the corresponding Markov chain mixes rapidly, i.e., the distribution of the state  $X_t$  converges quickly to  $\mathbb{P}_{\beta,q}$  as  $t \rightarrow \infty$ . Just as in the previous section, for a given divergence-generating function  $\Phi \in \mathcal{F}$ , the rate at which  $D_\Phi(P_{X_t} \| \mathbb{P}_{\beta,q})$  converges to zero is controlled by the SDPI constant  $\eta_\Phi(\mathbb{P}_{\beta,q}, K_{\beta,q}^\bullet)$ , where  $\bullet$  is either HB or SW. The heat-bath algorithm is widely used because it is easy to implement. On the other hand, the popularity of the Swendsen-Wang algorithm is due to the fact that, empirically, it tends to mix rapidly for a wide variety of graphs and small values of  $q$  (however, see [66] for examples of slow mixing of Swendsen-Wang). In a recent paper, Ullrich [67] showed that the spectral gap of Swendsen-Wang is lower-bounded by a constant multiple of the spectral gap of the heat-bath kernel, where the constant depends on the number of colors  $q$ , the inverse temperature  $\beta$ , and the maximum degree  $\Delta$  of  $G$ . Now, the spectral gap can be related to the SDPI constant for the  $\chi^2$ -divergence (see Remark 3.4), so Ullrich's result can immediately be converted into a statement about the  $\chi^2$  SDPI constants of Swendsen-Wang and heat-bath kernels. The theorem below sharpens and extends the bound of Ullrich to other  $\Phi$ -divergences; just like in [67], the theorem allows us to convert any available upper bound for the heat-bath kernel into an upper bound for the Swendsen-Wang kernel (or, conversely, any lower bound for Swendsen-Wang into a lower bound for heat-bath).

**Theorem 5.3.** *For any  $\Phi \in \mathcal{F}$  that satisfies the generalized homogeneity condition (3.5),*

$$\eta_\Phi(\mathbb{P}_{\beta,q}, K_{\beta,q}^{\text{SW}}) \leq \frac{q^{2\Delta+1} e^{4\beta\Delta} - 1}{q^{2\Delta+1} e^{4\beta\Delta} - \left[ \eta_\Phi(\mathbb{P}_{\beta,q}, K_{\beta,q}^{\text{HB}}) \right]^2}, \quad (5.17)$$

where  $\Delta = \max_{v \in V} \deg_G(v)$  is the maximum degree of  $G$ .

**Remark 5.1.** In the notation of this paper, the main result of [67] can be written as

$$\sqrt{\eta_{\chi^2}(\mathbb{P}_{\beta,q}, K_{\beta,q}^{\text{SW}})} \leq \frac{2q^{4\Delta+2} e^{8\beta\Delta} - 1 + \sqrt{\eta_{\chi^2}(\mathbb{P}_{\beta,q}, K_{\beta,q}^{\text{HB}})}}{2q^{4\Delta+2} e^{8\beta\Delta}}. \quad (5.18)$$

Particularizing our bound (5.17) to the case  $\Phi(u) = (u-1)^2$ , we see that it is tighter than (5.18). A plot of the two bounds as a function of the  $\chi^2$  SDPI constant of the heat-bath dynamics is shown in Figure 5.4 for  $q = 2$ ,  $\Delta = 3$ , and  $\beta = 0.001$ . (Admittedly, both bounds are fairly crude even for small values of  $q$  and  $\Delta$ , due to the presence of  $O(q^\Delta)$  terms.)  $\diamond$

*Proof.* We borrow a clever trick of Ullrich [67] and compare the Swendsen-Wang kernel  $K_{\beta,q}^{\text{SW}}$  to  $K = K_{\beta,q}^{\text{HB}} \circ K_{\beta,q}^{\text{SW}} \circ K_{\beta,q}^{\text{HB}}$ . Since the Gibbs distribution  $\mathbb{P}_{\beta,q}$  is invariant under both the SW and the

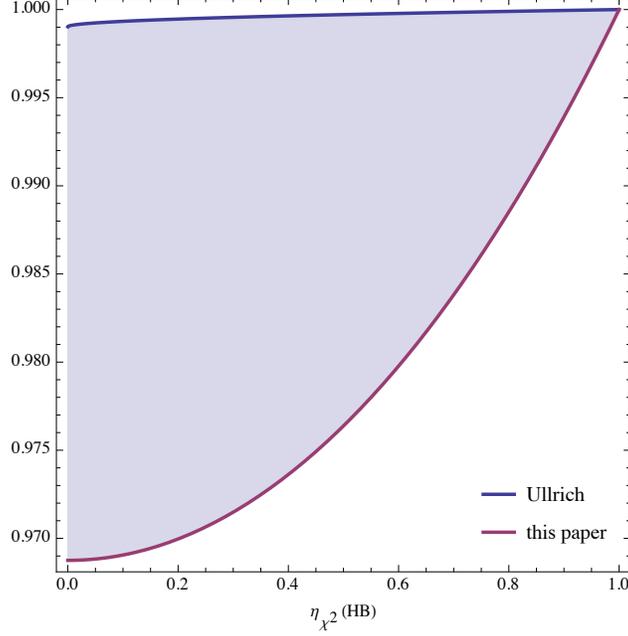


Figure 3: Upper bounds of Ullrich and Theorem 5.3 as a function of the  $\chi^2$  SDPI constant of the heat-bath dynamics, for  $q = 2$ ,  $\Delta = 3$ ,  $\beta = 0.001$ .

HB kernels, it is also the invariant distribution of  $K$ . Moreover, for any  $\nu \neq \mathbb{P}_{\beta,q}$ , we have

$$\begin{aligned}
D_{\Phi}(\nu K \parallel \mathbb{P}_{\beta,q}) &= D_{\Phi}(\nu K \parallel \mathbb{P}_{\beta,q} K) \\
&= D_{\Phi}\left(\nu(K_{\beta,q}^{\text{HB}} \circ K_{\beta,q}^{\text{SW}} \circ K_{\beta,q}^{\text{HB}}) \parallel \mathbb{P}_{\beta,q}(K_{\beta,q}^{\text{HB}} \circ K_{\beta,q}^{\text{SW}} \circ K_{\beta,q}^{\text{HB}})\right) \\
&\leq [\eta_{\Phi}(\mathbb{P}_{\beta,q}, K_{\beta,q}^{\text{HB}})]^2 \eta_{\Phi}(\mathbb{P}_{\beta,q}, K_{\beta,q}^{\text{SW}}) D_{\Phi}(\nu \parallel \mathbb{P}_{\beta,q}),
\end{aligned}$$

where we have repeatedly exploited the invariance of  $\mathbb{P}_{\beta,q}$  w.r.t. the SW and the HB kernels. Since  $\nu$  was arbitrary, we conclude that

$$\eta_{\Phi}(\mathbb{P}_{\beta,q}, K) \leq [\eta_{\Phi}(\mathbb{P}_{\beta,q}, K_{\beta,q}^{\text{HB}})]^2 \eta_{\Phi}(\mathbb{P}_{\beta,q}, K_{\beta,q}^{\text{SW}}). \quad (5.19)$$

On the other hand, Ullrich also proved that

$$\max_{x, x' \in \mathcal{X}} \frac{K(x'|x)}{K_{\beta,q}^{\text{SW}}(x'|x)} \leq q^{2\Delta+1} e^{4\beta\Delta}. \quad (5.20)$$

From Eq. (5.20) and Corollary 3.2, we get the estimate

$$\eta_{\Phi}(\mathbb{P}_{\beta,q}, K_{\beta,q}^{\text{SW}}) \leq 1 - \frac{1}{q^{2\Delta+1} e^{4\beta\Delta}} (1 - \eta_{\Phi}(\mathbb{P}_{\beta,q}, K)). \quad (5.21)$$

Finally, using (5.19) in (5.21) and rearranging, we obtain (5.17).  $\square$

## 5.5 Reconstruction in graphical models

The Potts model described in the preceding section is an example of a *probabilistic graphical model* (or a pairwise Markov random field) [68]. Any such model is specified by a pair  $(G, \mathbb{U})$ , where  $G = (\mathbf{V}, \mathbf{E})$  is an undirected graph and  $\mathbb{U} = \{\mathbb{U}_e\}_{e \in \mathbf{E}}$  is a collection of symmetric *edge potentials*  $\mathbb{U}_e : \Omega \times \Omega \rightarrow \mathbb{R}^+$ . Here,  $\Omega$  is a finite set often referred to as *state* or *spin space*. The configuration space of the graphical model is the set  $\mathbf{X} = \Omega^{\mathbf{V}}$  of all tuples  $x = (x_v)_{v \in \mathbf{V}}$ , where each  $x_v$  takes values in  $\Omega$ . Once  $G$  and  $\mathbb{U}$  are fixed, we consider the following probability measure on  $\mathbf{X}$ :

$$\mathbb{P}_{G, \mathbb{U}}(x) = \frac{1}{Z(G, \mathbb{U})} \prod_{\{u, v\} \in \mathbf{E}} \mathbb{U}_{\{u, v\}}(x_u, x_v),$$

where  $Z$  is the normalization constant. For example, the  $q$ -state Potts model on  $G$  [cf. Eq. (5.12)] is of this form with  $\Omega = \{1, \dots, q\}$  and  $\mathbb{U}_{uv}(x_u, x_v) = \exp(\beta \mathbf{1}\{x_u = x_v\})$ .

The *reconstruction problem* (see, e.g., [69, 70]) for the graphical model  $(G, \mathbb{U})$  can be stated informally as follows: Given two disjoint sets of vertices  $\mathbf{A}$  and  $\mathbf{B}$ , how much can we infer about the configuration  $X_{\mathbf{A}} \triangleq (X_v)_{v \in \mathbf{A}}$  on  $\mathbf{A}$  by observing  $X_{\mathbf{B}}$ ? For a precise definition, let  $d_G$  denote the *graph distance* on  $G$ , i.e.,  $d_G(u, v)$  is the number of edges on the shortest path between  $u$  and  $v$ .

**Definition 5.2.** *Given a function  $\Phi \in \mathcal{F}$ , we say that the probabilistic graphical model  $(G, \mathbb{U})$  is not  $\Phi$ -reconstructible if for any set of vertices  $\mathbf{A}$  there exist some constants  $C_{\mathbf{A}}, c_{\mathbf{A}} > 0$ , such that*

$$I_{\Phi}(X_{\mathbf{A}}; X_{\mathbf{B}}) \leq C_{\mathbf{A}} e^{-c_{\mathbf{A}} d_G(\mathbf{A}, \mathbf{B})}$$

for all sets  $\mathbf{B}$  disjoint from  $\mathbf{A}$ , where  $d_G(\mathbf{A}, \mathbf{B}) \triangleq \min_{u \in \mathbf{A}, v \in \mathbf{B}} d_G(u, v)$ . Here, the  $\Phi$ -information is computed w.r.t. the marginal distribution of  $(X_{\mathbf{A}}, X_{\mathbf{B}})$  induced by  $\mathbb{P}_{G, \mathbb{U}}$ .

Alternatively, we may examine correlations between functions of  $X_{\mathbf{A}}$  and  $X_{\mathbf{B}}$ :

**Definition 5.3.** *The graphical model  $(G, \mathbb{U})$  has exponential decay of correlations if for any  $\mathbf{A} \subset \mathbf{V}$  there exist positive constants  $C_{\mathbf{A}}, c_{\mathbf{A}}$ , such that, for any set of vertices  $\mathbf{B}$  disjoint from  $\mathbf{A}$  and for any two functions  $f \in \mathcal{F}(X_{\mathbf{A}})$  and  $g \in \mathcal{F}(X_{\mathbf{B}})$ ,*

$$\text{Cov}[f(X_{\mathbf{A}}), g(X_{\mathbf{B}})] \leq C_{\mathbf{A}} e^{-c_{\mathbf{A}} d_G(\mathbf{A}, \mathbf{B})} \sqrt{\text{Var}[f(X_{\mathbf{A}})] \text{Var}[g(X_{\mathbf{B}})]}.$$

We can now establish the following result:

**Theorem 5.4.** *Suppose that  $\Phi \in \mathcal{F}$  is twice differentiable and strictly convex, its second derivative is nonincreasing, and the function  $\Psi$  defined in (3.14) is concave. Then  $(G, \mathbb{U})$  is not  $\Phi$ -reconstructible if and only if it has exponential decay of correlations.*

*Proof.* We first show that exponential decay of correlations is equivalent to  $(G, \mathbb{U})$  not being  $\chi^2$ -reconstructible. With a slight abuse of notation, we will denote by  $P_{\mathbf{A}}$  the marginal distribution of  $X_{\mathbf{A}}$ , etc. By definition of maximal correlation,  $(G, \mathbb{U})$  has exponential decay of correlation if and only if for any  $\mathbf{A} \subset \mathbf{V}$  there exist some  $C_{\mathbf{A}}, c_{\mathbf{A}} > 0$ , such that

$$S^2(P_{\mathbf{A}}, P_{\mathbf{B}|\mathbf{A}}) \leq C_{\mathbf{A}} e^{-c_{\mathbf{A}} d_G(\mathbf{A}, \mathbf{B})} \tag{5.22}$$

for all sets of vertices  $\mathbf{B}$  with  $\mathbf{B} \cap \mathbf{A} = \emptyset$ . Now let  $\tilde{X}_{\mathbf{A}}$  denote the support of  $P_{\mathbf{A}}$ . Using the definition of  $\chi^2$ -information and Theorem 3.2, we can write

$$\begin{aligned}
I_{\chi^2}(X_{\mathbf{A}}; X_{\mathbf{B}}) &= \sum_{x_{\mathbf{A}} \in \tilde{X}_{\mathbf{A}}} P_{\mathbf{A}}(x_{\mathbf{A}}) \chi^2(P_{X_{\mathbf{B}}|X_{\mathbf{A}}=x_{\mathbf{A}}} \| P_{\mathbf{B}}) \\
&= \sum_{x_{\mathbf{A}} \in \tilde{X}_{\mathbf{A}}} P_{\mathbf{A}}(x_{\mathbf{A}}) \chi^2(\delta_{x_{\mathbf{A}}} P_{\mathbf{B}|\mathbf{A}} \| P_{\mathbf{A}} P_{\mathbf{B}|\mathbf{A}}) \\
&\leq S^2(P_{\mathbf{A}}, P_{\mathbf{B}|\mathbf{A}}) \sum_{x_{\mathbf{A}}} P_{\mathbf{A}}(x_{\mathbf{A}}) \chi^2(\delta_{x_{\mathbf{A}}} \| P_{\mathbf{A}}) \\
&= S^2(P_{\mathbf{A}}, P_{\mathbf{B}|\mathbf{A}}) I_{\chi^2}(X_{\mathbf{A}}; X_{\mathbf{A}}) \\
&= S^2(P_{\mathbf{A}}, P_{\mathbf{B}|\mathbf{A}}) (|\mathbf{X}_{\mathbf{A}}| - 1).
\end{aligned}$$

From this and from (5.22), we see that exponential decay of correlations implies that  $(G, \mathbb{U})$  is not  $\chi^2$ -reconstructible. The converse statement follows from the inequality  $S^2(P_{\mathbf{A}}, P_{\mathbf{B}|\mathbf{A}}) \leq I_{\chi^2}(X_{\mathbf{A}}; X_{\mathbf{B}})$  [18, Prop. 12].

To complete the proof, let  $(\bar{X}_{\mathbf{A}}, \bar{X}_{\mathbf{B}})$  be a random pair with probability law  $P_{\mathbf{A}} \otimes P_{\mathbf{B}}$ . Then

$$I_{\Phi}(X_{\mathbf{A}}; X_{\mathbf{B}}) = \text{Ent}_{\Phi} [f(\bar{X}_{\mathbf{A}}, \bar{X}_{\mathbf{B}})],$$

where we have defined

$$f(x_{\mathbf{A}}, x_{\mathbf{B}}) \triangleq \frac{P_{\mathbf{A}|\mathbf{B}}(x_{\mathbf{A}}|x_{\mathbf{B}})}{P_{\mathbf{A}}(x_{\mathbf{A}})}.$$

In particular,  $I_{\chi^2}(X_{\mathbf{A}}; X_{\mathbf{B}}) = \text{Var} [f(\bar{X}_{\mathbf{A}}, \bar{X}_{\mathbf{B}})]$ . Therefore, applying Lemmas A.2 and A.3 in Appendix A and using the fact that  $\|f(\bar{X}_{\mathbf{A}}, \bar{X}_{\mathbf{B}})\|_{\infty} \leq 1/p_*^{\mathbf{A}}$ , where  $p_*^{\mathbf{A}} \triangleq \min_{x_{\mathbf{A}} \in \tilde{X}_{\mathbf{A}}} P_{\mathbf{A}}(x_{\mathbf{A}})$  is the minimum nonzero probability of any configuration in  $\mathbf{A}$ , we get

$$\frac{\Phi''(1/p_*^{\mathbf{A}})}{2} I_{\chi^2}(X_{\mathbf{A}}; X_{\mathbf{B}}) \leq I_{\Phi}(X_{\mathbf{A}}; X_{\mathbf{B}}) \leq \Psi'(1) I_{\chi^2}(X_{\mathbf{A}}; X_{\mathbf{B}}).$$

Since  $\Phi$  is strictly convex,  $\Phi''$  is everywhere positive. This inequality shows that the graphical model  $(G, \mathbb{U})$  is not  $\Phi$ -reconstructible if and only if it is not  $\chi^2$ -reconstructible, which in turn is equivalent to exponential decay of correlations.  $\square$

A related notion of correlation decay has to do with the diminishing influence of “far away” spins. A key property of Gibbs measures is the following conditional independence relation: for any  $\mathbf{A} \subset \mathbf{V}$ , the outer boundary of  $\mathbf{A}$ , denoted by  $\partial\mathbf{A}$ , is the set of all  $v \in \mathbf{A}^c$ , such that  $\{u, v\} \in \mathbf{E}$  for some  $u \in \mathbf{A}$ . Then under  $\mathbb{P}_{G, \mathbb{U}}$ ,

$$X_{\mathbf{A}} \longrightarrow X_{\partial\mathbf{A}} \longrightarrow X_{\mathbf{A}^c}$$

is a Markov chain. That is, the configuration of spins in a given set  $\mathbf{A}$  of vertices is conditionally independent of all other spins given the configuration of the neighbors of  $\mathbf{A}$ . The following definition formalizes the notion that the influence of the spins in the boundary of  $\mathbf{A}$  on the spins in any subset of  $\mathbf{A}$  should decay with the distance from that subset to the boundary:

**Definition 5.4.** *The graphical model  $(G, \mathbb{U})$  has the spatial mixing property if there exist positive constants  $C, c$ , such that, for any two sets of vertices  $\mathbf{B} \subset \mathbf{A} \subset \mathbf{V}$  and for any two boundary configurations  $x_{\partial\mathbf{A}}, \bar{x}_{\partial\mathbf{A}}$ ,*

$$\left\| \frac{P_{\mathbf{B}|\partial\mathbf{A}}(\cdot|x_{\partial\mathbf{A}})}{P_{\mathbf{B}|\partial\mathbf{A}}(\cdot|\bar{x}_{\partial\mathbf{A}})} - 1 \right\|_{\infty} \leq C|\mathbf{A}|e^{-cd_G(\mathbf{B}, \partial\mathbf{A})}. \quad (5.23)$$

**Remark 5.2.** This mixing condition is slightly stronger than the condition proposed by Weitz [55], which is in turn stronger (but more generally applicable) than the complete analyticity condition of Dobrushin and Shlosman [71]. The latter is only applicable to the case when the underlying graph  $G$  is the square lattice  $\mathbb{Z}^d$ .  $\diamond$

If  $(G, \mathbb{U})$  has spatial mixing, then one would expect the relative-entropy SDPI constant of the channel  $P_{\mathbf{B}|\partial\mathbf{A}}$  at  $P_{\partial\mathbf{A}}$  to decay exponentially with the distance  $d_G(\mathbf{B}, \partial\mathbf{A})$ . This is indeed the case:

**Theorem 5.5.** *Suppose that  $(G, \mathbb{U})$  has the spatial mixing property. Then*

$$\eta(P_{\partial\mathbf{A}}, P_{\mathbf{B}|\partial\mathbf{A}}) \leq \frac{2C^2|\mathbf{A}|^2}{P_{\mathbf{B}}^*} e^{-2cd_G(\mathbf{B}, \partial\mathbf{A})}. \quad (5.24)$$

*Proof.* Using (3.42), we can upper-bound  $\eta(P_{\partial\mathbf{A}}, P_{\mathbf{B}|\partial\mathbf{A}})$  as follows:

$$\eta(P_{\partial\mathbf{A}}, P_{\mathbf{B}|\partial\mathbf{A}}) \leq \frac{1}{2} \sum_{x_{\mathbf{B}}} \frac{1}{P_{\mathbf{B}}(x_{\mathbf{B}})} \max_{x_{\partial\mathbf{A}}, x'_{\partial\mathbf{A}}} |P_{\mathbf{B}|\partial\mathbf{A}}(x_{\mathbf{B}}|x_{\partial\mathbf{A}}) - P_{\mathbf{B}|\partial\mathbf{A}}(x_{\mathbf{B}}|x'_{\partial\mathbf{A}})|^2.$$

If we now pick an arbitrary boundary configuration  $\bar{x}_{\partial\mathbf{A}}$ , then we can write

$$\begin{aligned} \eta(P_{\partial\mathbf{A}}, P_{\mathbf{B}|\partial\mathbf{A}}) &\leq 2 \sum_{x_{\mathbf{B}}} \frac{1}{P_{\mathbf{B}}(x_{\mathbf{B}})} \max_{x_{\partial\mathbf{A}}} |P_{\mathbf{B}|\partial\mathbf{A}}(x_{\mathbf{B}}|x_{\partial\mathbf{A}}) - P_{\mathbf{B}|\partial\mathbf{A}}(x_{\mathbf{B}}|\bar{x}_{\partial\mathbf{A}})|^2 \\ &\leq 2 \sum_{x_{\mathbf{B}}} \frac{P_{\mathbf{B}|\partial\mathbf{A}}(x_{\mathbf{B}}|\bar{x}_{\partial\mathbf{A}})}{P_{\mathbf{B}}(x_{\mathbf{B}})} \max_{x_{\partial\mathbf{A}}} \left| \frac{P_{\mathbf{B}|\partial\mathbf{A}}(x_{\mathbf{B}}|x_{\partial\mathbf{A}})}{P_{\mathbf{B}|\partial\mathbf{A}}(x_{\mathbf{B}}|\bar{x}_{\partial\mathbf{A}})} - 1 \right|^2. \end{aligned}$$

Using (5.23), we get (5.24).  $\square$

## 6 Summary of contributions and concluding remarks

In this paper, we have attempted to give a systematic and unified presentation of strong data processing inequalities (SDPIs) for discrete channels. As a reminder, given a convex function  $\Phi : \mathbb{R}^+ \rightarrow \mathbb{R}$ , we say that a channel  $K \in \mathcal{M}(\mathbf{Y}|\mathbf{X})$  satisfies an SDPI with constant  $c \in [0, 1)$  at input distribution  $\mu$  if

$$D_{\Phi}(\nu K || \mu K) \leq c D_{\Phi}(\nu || \mu) \quad (6.1)$$

for all  $\nu \neq \mu$ . We denote the best constant in the above inequality by  $\eta_{\Phi}(\mu, K)$ , and let  $\eta_{\Phi}(K) \triangleq \sup_{\mu} \eta_{\Phi}(\mu, K)$ . For the reader's convenience, we summarize the key novel contributions:

- For all sufficiently smooth  $\Phi$ ,  $\eta_{\Phi}(\mu, K)$  is lower-bounded by the squared maximal correlation  $S^2(\mu, K)$ , which is also the SDPI constant of  $K$  at  $\mu$  for the  $\chi^2$ -divergence (Theorem 3.3). This refines the inequality  $\eta(\mu, K) \geq S^2(\mu, K)$  due to Ahlswede and Gács [1], as well as the inequality  $\eta_{\Phi}(K) \geq S^2(K) \equiv \sup_{\mu} S^2(\mu, K)$  due to Cohen et al. [6].

- For all operator convex  $\Phi$  (see Section 3.3 for definitions and examples), we have proved the upper bound

$$\eta_{\Phi}(\mu, K) \leq \max \left( S^2(\mu, K), \sup_{0 < \lambda < 1} \eta_{\text{LC}_{\lambda}}(\mu, K) \right)$$

(Theorem 3.6), where  $\text{LC}_{\lambda}(\cdot \| \cdot)$  denotes the Le Cam divergence with parameter  $\lambda$  (see Section 2). This refines the inequality  $\eta_{\Phi}(K) \leq S^2(K)$  for all operator convex  $\Phi$ , due to Choi et al. [7], and reduces to it upon taking the supremum of both sides w.r.t.  $\mu$ .

- For  $\Phi(u) = u \log u$  (which gives the usual relative entropy), the SDPI constant  $\eta_{\Phi}(\mu, K)$  can be upper-bounded in terms of the subgaussian constant  $\sigma^2(y)$  of the posterior likelihood ratio  $a(X, y) = \frac{K^*(X|y)}{\mu(X)}$  for each  $y \in \mathcal{Y}$ , where  $X \sim \mu$ . Smaller value of  $\sigma^2(y)$  indicates that  $a(X, y) \approx 1$  with high probability, which means that the observation  $Y = y$  is nearly uninformative about the input  $X$ . Theorem 3.7 gives the inequality  $\eta(\mu, K) \leq 2 \mathbb{E}[\sigma^2(Y)]$ , which can be weakened to the bound of Theorem 3.8 using information-transportation inequalities.
- Under mild regularity conditions on  $\Phi$ , the SDPI constants *tensorize*: given a product distribution  $\mu_1 \otimes \dots \otimes \mu_n$  and a product channel  $K_1 \otimes \dots \otimes K_n$ ,

$$\eta_{\Phi}(\mu_1 \otimes \dots \otimes \mu_n, K_1 \otimes \dots \otimes K_n) = \max_{1 \leq i \leq n} \eta_{\Phi}(\mu_i, K_i)$$

(Theorem 3.9). This extends previous tensorization results for  $\Phi(u) = (u - 1)^2$  due to Witsenhausen [2] and for  $\Phi(u) = u \log u$  due to Anantharam et al. [15]. Theorem 3.10 gives a tensorization inequality for *mixtures* of local channels, i.e., when an input block of length  $n$  is transformed to an output block of length  $n$  by drawing a coordinate index  $I$  at random from  $\{1, \dots, n\}$  and then passing the  $I$ th symbol through the channel  $K_I$ .

- Section 4 is dedicated to an exposition of the deep links between SDPIs and  $\Phi$ -Sobolev inequalities [23], which provide a powerful tool for nonasymptotic quantitative analysis of convergence to equilibrium in Markov processes and other random dynamical systems. For the specific case of  $\Phi(u) = u \log u$ , we have obtained a number of inequalities relating the optimal constants in log-Sobolev inequalities for a reversible Markov chain  $M$  with invariant distribution  $\mu$  to relative-entropy SDPI constants  $\eta(\mu, K)$  for any channel  $K$  with the property that  $M = K_{\mu}^* \circ K$ , where  $K_{\mu}^*$  is the adjoint, or backward, channel associated to the pair  $(\mu, K)$  [see Eq. (1.2) for the definition].
- Section 5 presents several applications of the results of preceding sections to information theory, discrete probability, and statistical physics. In particular, we discuss a connection between the strong data processing property and the concentration-of-measure phenomenon; generalize a recent result of Anantharam et al. [15] on the strong contraction of mutual information in discrete Markov chains<sup>7</sup> to a more general notion of  $\Phi$ -information; relate the problem of computing SDPI constants (which is a convex program) to the problem of finding the fastest mixing Markov chain on a graph; sharpen a recent result of Ullrich [67] on the mixing time of two popular MCMC schemes for a certain class of graphical models; and outline an SDPI-based characterization of the decay of correlations in discrete graphical models.

---

<sup>7</sup>See [21] for an extension of this result to abstract alphabets.

After the original breakthrough work of Ahlswede and Gács [1], strong data processing inequalities have received a great deal of attention, with a recent surge of research activity motivated by problems in information theory. Recent work by Polyanskiy and Wu [18] has uncovered certain limitations of SDPIs. For example, in the setting of continuous alphabets and additive-noise channels, they have shown that it is possible for a channel  $K$  to have  $\eta_{\Phi}(K) = 1$  and still satisfy a weaker “nonlinear” strong data processing inequality of the form

$$D_{\Phi}(\nu K \parallel \mu K) \leq F_{\Phi}(D_{\Phi}(\nu \parallel \mu))$$

for some increasing function  $F_{\Phi} : \mathbb{R}^+ \rightarrow \mathbb{R}^+$  with  $F_{\Phi}(0) = 0$ , such that  $F_{\Phi}(u) < u$  for all sufficiently small  $u \neq 0$ . Nevertheless, SDPIs still remain a versatile tool for many problems of current theoretical and practical interest.

## Acknowledgments

The author would like to thank V. Anantharam, S. Kamath, A. Kontorovich, C. Nair, Y. Polyanskiy, I. Sason, P. Tetali, R. van Handel, and Y. Wu for many useful and stimulating discussions, and the two anonymous reviewers and the Associate Editor for their meticulous reading of the manuscript and for numerous useful suggestions and corrections. The author would also like to separately thank one of the anonymous reviewers for a suggestion on how to streamline the proof of Theorem 3.12, as well as for pointing out a subtle issue pertaining to Theorem 4.6.

## A Miscellaneous lemmas

**Lemma A.1.** *Let  $(\mu, K) \in \mathcal{P}_*(\mathsf{X}) \times \mathcal{M}(\mathsf{Y}|\mathsf{X})$  be an admissible pair, and consider any other  $\nu \in \mathcal{P}(\mathsf{X})$ . If  $f = d\nu/d\mu$ , then*

$$K^* f = \frac{d(\nu K)}{d(\mu K)},$$

where  $K^* = K_{\mu}^* \in \mathcal{M}(\mathsf{Y}|\mathsf{X})$  is the backward channel induced by the pair  $(\mu, K)$ , cf. Eqs. (1.1), (1.2).

*Proof.* A direct calculation:

$$\begin{aligned} K^* f(y) &= \sum_{x \in \mathsf{X}} K^*(x|y) f(x) \\ &= \sum_{x \in \mathsf{X}} \frac{K(y|x) \mu(x)}{\mu K(y)} \frac{\nu(x)}{\mu(x)} \\ &= \frac{1}{\mu K(y)} \sum_{x \in \mathsf{X}} \nu(x) K(y|x) \\ &= \frac{\nu K(y)}{\mu K(y)} \\ &= \frac{d(\nu K)(y)}{d(\mu K)} \end{aligned}$$

for any  $y \in \mathsf{Y}$ . □

**Lemma A.2.** Suppose  $\Phi \in \mathcal{F}$  is differentiable, and the function  $\Psi(u) = \frac{\Phi(u) - \Phi(0)}{u}$  is concave. Then for any nonnegative random variable  $U$  with  $\mathbb{E}U = 1$ ,

$$\text{Ent}_\Phi[U] \leq \Psi(1 + \text{Var}[U]) - \Psi(1) \leq \Psi'(1) \text{Var}[U]. \quad (\text{A.1})$$

*Proof.* We can assume that  $\text{Var}[U] < \infty$ , because otherwise there is nothing to prove. Let  $P$  denote the law of  $U$ . Since  $U$  is nonnegative and has unit mean,  $Q(\text{d}u) \triangleq uP(\text{d}u)$  is a probability measure. Therefore,

$$\begin{aligned} \text{Ent}_\Phi[U] &= \mathbb{E}_P[\Phi(U)] - \Phi(1) \\ &= \mathbb{E}_Q[\Psi(U)] - \Psi(1) \\ &\leq \Psi(\mathbb{E}_Q U) - \Psi(1) \\ &= \Psi(\mathbb{E}[U^2]) - \Psi(1) \\ &= \Psi(1 + \text{Var}[U]) - \Psi(1), \end{aligned}$$

where the third step is by Jensen's inequality, and the remaining steps follow from definitions. This proves the first inequality in (A.1). Now, since  $\Psi$  is concave, we have

$$\Psi(1 + \text{Var}[U]) \leq \Psi(1) + \Psi'(1) \text{Var}[U].$$

Using this, we obtain the second inequality.  $\square$

**Lemma A.3.** Suppose  $\Phi \in \mathcal{F}$  is twice differentiable, and  $\Phi''$  is nonincreasing. Then for any nonnegative random variable  $U$  with  $\mathbb{E}U = 1$  and  $\|U\|_\infty < \infty$ ,

$$\text{Ent}_\Phi[U] \geq \frac{\Phi''(\|U\|_\infty)}{2} \text{Var}[U]. \quad (\text{A.2})$$

*Proof.* By Taylor's theorem, for any  $u \geq 0$  we have

$$\Phi(u) - \Phi(1) = \Phi'(1)(u - 1) + \frac{\Phi''(v)}{2}(u - 1)^2$$

for some  $v \in [u \wedge 1, u \vee 1]$ . Since  $\Phi''$  is nonincreasing,  $\Phi''(v) \geq \Phi''(u \vee 1) \geq \Phi''(\|U \vee 1\|_\infty) = \Phi''(\|U\|_\infty)$ , where the equality is a consequence of the assumption that  $\mathbb{E}U = 1$ . Taking expectations w.r.t.  $U$ , we obtain (A.2).  $\square$

**Lemma A.4.** Let  $U$  and  $Z$  be two jointly distributed random variables, where  $U$  is real-valued and nonnegative, and  $Z$  takes values in an arbitrary set  $\mathcal{Z}$ . Then, for any  $\Phi \in \mathcal{F}$ , the expectation of the conditional  $\Phi$ -entropy  $\text{Ent}_\Phi[U|Z]$  admits the following variational representation:

$$\mathbb{E}[\text{Ent}_\Phi[U|Z]] = \inf_{\xi \in \mathcal{F}_*^0(\mathcal{Z})} \mathbb{E}[\Phi(U) - \Phi(\xi(Z)) - (U - \xi(Z))\Phi'(\xi(Z))],$$

where  $\Phi'$  denotes the right derivative of  $\Phi$  (which exists due to convexity).

*Proof.* This lemma is a generalization of Lemma 14.4 in [31]. Fix an arbitrary  $\xi \in \mathcal{F}_*^0(\mathcal{Z})$ . Then, by convexity of  $\Phi$ , for any  $z \in \mathcal{Z}$  we have

$$\Phi(\mathbb{E}[U|Z = z]) \geq \Phi(\xi(z)) + \Phi'(\xi(z))(\mathbb{E}[U|Z = z] - \xi(z)).$$

From this, we get

$$\begin{aligned}\text{Ent}_\Phi [U|Z = z] &= \mathbb{E}[\Phi(U)|Z = z] - \Phi(\mathbb{E}[U|Z = z]) \\ &\leq \mathbb{E}[\Phi(U)|Z = z] - \Phi(\xi(z)) - \Phi'(\xi(z))(\mathbb{E}[U|Z = z] - \xi(z)).\end{aligned}$$

Taking expectations of both sides w.r.t.  $Z$ , we see that

$$\mathbb{E}[\text{Ent}_\Phi[U|Z]] \leq \mathbb{E}[\Phi(U) - \Phi(\xi(Z)) - (U - \xi(Z))\Phi'(\xi(Z))] \quad (\text{A.3})$$

for any  $\xi \in \mathcal{F}_*^0(\mathbf{Z})$ . On the other hand, if we take  $\xi(z) = \mathbb{E}[U|Z = z]$ , then the bound in (A.3) is achieved with equality.  $\square$

**Lemma A.5.** *Let  $\Phi \in \mathcal{F}$  be a differentiable function, such that  $\Phi'(u)$  is uniformly bounded in some neighborhood of  $u = 1$ . Then for any nonnegative real-valued random variable  $U$  with  $\mathbb{E}U = 1$  and  $\|U\|_\infty < \infty$ , we have*

$$\left. \frac{d}{d\varepsilon} \text{Ent}_\Phi \left[ \frac{1 - \varepsilon U}{\bar{\varepsilon}} \right] \right|_{\varepsilon=0} = 0. \quad (\text{A.4})$$

*Proof.* Since  $\mathbb{E}U = 1$ , for all sufficiently small  $\varepsilon > 0$  we have

$$\text{Ent}_\Phi \left[ \frac{1 - \varepsilon U}{\bar{\varepsilon}} \right] = \mathbb{E} \left[ \Phi \left( \frac{1 - \varepsilon U}{\bar{\varepsilon}} \right) \right].$$

By our assumptions on  $\Phi$ , there exists a constant  $C > 0$ , such that

$$\left| \frac{d}{d\varepsilon} \Phi \left( \frac{1 - \varepsilon u}{\bar{\varepsilon}} \right) \right| = \left| \frac{1 - u}{\bar{\varepsilon}^2} \Phi' \left( \frac{1 - \varepsilon u}{\bar{\varepsilon}} \right) \right| \leq C|u - 1|$$

for all sufficiently small  $\varepsilon > 0$ . Therefore, by the dominated convergence theorem, we can interchange expectation and derivative to get

$$\left. \frac{d}{d\varepsilon} \text{Ent}_\Phi \left[ \frac{1 - \varepsilon U}{\bar{\varepsilon}} \right] \right|_{\varepsilon=0} = \Phi'(1)\mathbb{E}[(1 - U)] = 0.$$

$\square$

## B Proof of Proposition 4.1

Items 1)–3) are obvious. We prove 4). To that end, we first analyze the joint distribution of  $X$  and  $X'$ . First of all, for any  $x, x' \in \mathbf{X}$ , using the definition of  $K^*$ , we can write

$$\begin{aligned}
P_{XX'}(x, x') &= \mu(x)K^*K(x'|x) \\
&= \mu(x) \sum_{y \in \mathbf{Y}} K^*(x'|y)K(y|x) \\
&= \mu(x) \sum_{y \in \mathbf{Y}} \frac{K(y|x')\mu(x')}{\mu K(y)} K(y|x) \\
&= \mu(x') \sum_{y \in \mathbf{Y}} \frac{K(y|x)\mu(x)}{\mu K(y)} K(y|x') \\
&= \mu(x') \sum_{y \in \mathbf{Y}} K^*(x|y)K(y|x) \\
&= \mu(x')K^*K(x|x') \\
&= P_{XX'}(x', x).
\end{aligned}$$

In other words, the distribution of  $P_{XX'}$  is *exchangeable* (or  $(X, X')$  is an *exchangeable pair*). This implies, in particular, that the marginal distribution  $P_{X'}$  is the same as  $P_X$ , i.e.,  $\mu$ . Moreover, for any function  $f \in \mathcal{F}(\mathbf{X})$  and any  $x \in \mathbf{X}$ ,

$$\begin{aligned}
\mathbb{E}[f(X')|X = x] &= \sum_{x' \in \mathbf{X}} K^*K(x'|x)f(x') \\
&= \sum_{x' \in \mathbf{X}} \sum_{y \in \mathbf{Y}} K^*(x'|y)K(y|x)f(x') \\
&= \sum_{y \in \mathbf{Y}} K(y|x) \sum_{x' \in \mathbf{X}} K^*(x'|y)f(x') \\
&= \sum_{y \in \mathbf{Y}} K(y|x)K^*f(y) \\
&= KK^*f(x).
\end{aligned}$$

Using these facts, we can write

$$\begin{aligned}
&\mathbb{E}[(f(X) - f(X')) (g(X) - g(X'))] \\
&= \mathbb{E}[f(X)g(X)] + \mathbb{E}[f(X')g(X')] - \left( \mathbb{E}[f(X)g(X')] - \mathbb{E}[f(X')g(X)] \right) \\
&= 2\left\{ \mathbb{E}[f(X)g(X)] - \mathbb{E}[f(X)g(X')] \right\},
\end{aligned}$$

where

$$\begin{aligned}
\mathbb{E}[f(X)g(X')] &= \mathbb{E}[f(X)\mathbb{E}[g(X')|X]] \\
&= \mathbb{E}[f(X)KK^*g(X)] \\
&= \mathbb{E}[K^*f(X)K^*g(X)] \\
&= \mathbb{E}[\mathbb{E}[f(X)|Y]\mathbb{E}[g(X)|Y]] \\
&= \mathbb{E}[f(X)\mathbb{E}[g(X)|Y]].
\end{aligned}$$

Accordingly, we have

$$\begin{aligned}\mathbb{E} [(f(X) - f(X')) (g(X) - g(X'))] &= 2 \left\{ \mathbb{E} [f(X) (g(X) - \mathbb{E}[g(X)|Y])] \right\} \\ &= 2 \mathcal{E}(f(X), g(X)|Y),\end{aligned}$$

where the second step follows from the identity  $\mathcal{E}(U, V|Y) = \mathbb{E}[U(V - \mathbb{E}[V|Y])]$ . This proves (4.2). To prove (4.3), write

$$\begin{aligned}\mathbb{E} [(f(X) - f(X')) (g(X) - g(X'))] &= \mathbb{E} [1_{\{f(X) > f(X')\}} (f(X) - f(X')) (g(X) - g(X'))] \\ &\quad + \mathbb{E} [1_{\{f(X) < f(X')\}} (f(X) - f(X')) (g(X) - g(X'))] \\ &= 2 \mathbb{E} [1_{\{f(X) > f(X')\}} (f(X) - f(X')) (g(X) - g(X'))] \\ &= 2 \mathbb{E} [(f(X) - f(X'))_+ (g(X) - g(X'))],\end{aligned}$$

where the second step is by exchangeability of  $X$  and  $X'$ .

## References

- [1] R. Ahlswede and P. Gács, “Spreading of sets in product spaces and hypercontraction of the Markov operator,” *Ann. Probab.*, vol. 4, no. 6, pp. 925–939, 1976.
- [2] H. S. Witsenhausen, “On sequences of pairs of dependent random variables,” *SIAM J. Appl. Math.*, vol. 28, no. 1, pp. 100–113, January 1975.
- [3] A. D. Wyner, “The common information of two dependent random variables,” *IEEE Trans. Inform. Theory*, vol. 21, no. 2, pp. 163–179, March 1975.
- [4] I. Csiszár, “Information-type measures of difference of probability distributions and indirect observations,” *Stud. Sci. Math. Hung.*, vol. 2, pp. 299–318, 1967.
- [5] F. Liese and I. Vajda, “On divergences and informations in statistics and information theory,” *IEEE Trans. Inform. Theory*, vol. 52, no. 10, pp. 4394–4412, October 2006.
- [6] J. E. Cohen, Y. Iwasa, G. Rautu, M. B. Ruskai, E. Seneta, and G. Zbăganu, “Relative entropy under mappings by stochastic matrices,” *Lin. Algebra Appl.*, vol. 179, pp. 211–235, 1993.
- [7] M. Choi, M. B. Ruskai, and E. Seneta, “Equivalence of certain entropy contraction coefficients,” *Lin. Algebra Appl.*, vol. 208/209, pp. 29–36, 1994.
- [8] L. Miclo, “Remarques sur l’hypercontractivité et l’évolution de l’entropie pour des chaînes de Markov finies,” *Séminaire de probabilités (Strasbourg)*, vol. 31, pp. 136–167, 1997.
- [9] J. E. Cohen, J. H. B. Kemperman, and G. Zbăganu, *Comparisons of Stochastic Matrices, With Applications in Information Theory, Statistics, Economics, and Population Sciences*. Boston: Birkhäuser, 1998.
- [10] P. Del Moral, M. Ledoux, and L. Miclo, “On contraction properties of Markov kernels,” *Prob. Theory Rel. Fields*, vol. 126, pp. 395–420, 2003.

- [11] R. L. Dobrushin, “Central limit theorems for nonstationary Markov chains, I,” *Theory Probab. Appl.*, vol. 1, pp. 65–80, 1956.
- [12] —, “Central limit theorems for nonstationary Markov chains, II,” *Theory Probab. Appl.*, vol. 1, pp. 365–425, 1956.
- [13] X. Boyen and D. Koller, “Tractable inference for complex stochastic processes,” in *Proc. 14th Annual Conf. on Uncertainty in Artif. Intel.*, Madison, WI, July 1998, pp. 33–42.
- [14] S. Kamath and V. Anantharam, “Non-interactive simulation of joint distributions: The Hirschfeld–Gebelein–Rényi maximal correlation and the hypercontractivity ribbon,” in *Proc. 50th Annu. Allerton Conf. on Commun., Control, and Comput.*, Monticello, IL, October 2012.
- [15] V. Anantharam, A. Gohari, S. Kamath, and C. Nair, “On maximal correlation, hypercontractivity, and the data processing inequality studied by Erkip and Cover,” 2013, arXiv preprint. [Online]. Available: <http://arxiv.org/abs/1304.6133>
- [16] T. Courtade, “Outer bounds for multiterminal source coding via a strong data processing inequality,” in *Proc. Int. IEEE Symp. on Inform. Theory*, Istanbul, Turkey, July 2013, pp. 559–563.
- [17] M. Raginsky, “Logarithmic Sobolev inequalities and strong data processing theorems for discrete channels,” in *Proc. Int. IEEE Symp. on Inform. Theory*, Istanbul, Turkey, July 2013, pp. 419–423.
- [18] Y. Polyanskiy and Y. Wu, “Dissipation of information in channels with input constraints,” *IEEE Trans. Inform. Theory*, vol. 62, no. 1, pp. 35–55, January 2016.
- [19] V. Anantharam, A. Gohari, S. Kamath, and C. Nair, “On hypercontractivity and a data processing inequality,” in *Proc. Int. IEEE Symp. on Inform. Theory*, Honolulu, HI, July 2014, pp. 3022–3026.
- [20] J. Liu, P. Cuff, and S. Verdú, “Key capacity with limited one-way communication for product sources,” in *Proc. IEEE Int. Symp. Inform. Theory*, Honolulu, HI, July 2014, pp. 1146–1150.
- [21] Y. Polyanskiy and Y. Wu, “Strong data-processing inequalities for channels and Bayesian networks,” arXiv.org preprint 1508.06025. [Online]. Available: <http://arxiv.org/abs/1508.06025>
- [22] A. Makur and L. Zheng, “Bounds between contraction coefficients,” 2015, arXiv preprint 1510.01844. [Online]. Available: <http://arxiv.org/abs/1510.01844>
- [23] D. Chafaï, “Entropies, convexity, and functional inequalities: on  $\Phi$ -entropies and  $\Phi$ -Sobolev inequalities,” *J. Math. Kyoto Univ.*, vol. 44, no. 2, pp. 325–363, 2004.
- [24] D. Bakry, “L’hypercontractivité et son utilisation en théorie des semigroupes,” in *Lectures on Probability Theory*. Springer, 1994, vol. 1581, pp. 1–114.
- [25] P. Diaconis and L. Saloff-Coste, “Logarithmic Sobolev inequalities for finite Markov chains,” *Ann. Appl. Probab.*, vol. 6, no. 3, pp. 695–750, 1996.

- [26] S. G. Bobkov and P. Tetali, “Modified logarithmic Sobolev inequalities in discrete settings,” *J. Theor. Prob.*, vol. 19, no. 2, pp. 289–336, 2006.
- [27] E. Mossel, K. Oleszkiewicz, and A. Sen, “On reverse hypercontractivity,” *Geom. Funct. Anal.*, vol. 23, no. 3, pp. 1062–1097, 2013.
- [28] P. Diaconis, K. Khare, and L. Saloff-Coste, “Stochastic alternating projections,” *Illinois J. Math.*, vol. 54, no. 3, pp. 963–979, 2010.
- [29] W. R. Gilks, S. Richardson, and D. Spiegelhalter, Eds., *Markov Chain Monte Carlo in Practice*. Chapman & Hall, 1996.
- [30] C. P. Robert and G. Casella, *Monte Carlo Statistical Methods*, 2nd ed. Springer, 2004.
- [31] S. Boucheron, G. Lugosi, and P. Massart, *Concentration Inequalities: A Nonasymptotic Theory of Independence*. Oxford Univ. Press, 2013.
- [32] X. Nguyen, M. J. Wainwright, and M. I. Jordan, “On surrogate loss functions and  $f$ -divergences,” *Ann. Statist.*, vol. 37, no. 2, pp. 876–904, 2009.
- [33] M. H. DeGroot, “Uncertainty, information, and sequential experiments,” *Ann. Math. Statist.*, vol. 33, no. 2, pp. 404–419, 1962.
- [34] L. Le Cam, *Asymptotic Methods in Statistical Decision Theory*. Springer, 1986.
- [35] L. Györfi and I. Vajda, “A class of modified Pearson and Neyman statistics,” *Statistics and Decisions*, vol. 19, no. 3, pp. 239–252, 2001.
- [36] I. Sason and S. Verdú, “ $f$ -divergence inequalities,” 2015, arXiv preprint 1508.00335. [Online]. Available: <http://arxiv.org/abs/1508.00335>
- [37] B. Efron and C. Stein, “The jackknife estimate of variance,” *Ann. Statist.*, vol. 9, pp. 586–596, 1981.
- [38] J. M. Steele, “An Efron–Stein inequality for nonsymmetric statistics,” *Ann. Statist.*, vol. 14, pp. 753–758, 1986.
- [39] R. Latała and K. Oleszkiewicz, “Between Sobolev and Poincaré,” in *Geometric Aspects of Functional Analysis*, ser. Lecture Notes in Mathematics. Springer, 2000, vol. 1745, pp. 147–168.
- [40] S. Boucheron, O. Bousquet, G. Lugosi, and P. Massart, “Moment inequalities for functions of independent random variables,” *Ann. Probab.*, vol. 33, no. 2, pp. 514–560, 2005.
- [41] I. Csiszár and P. C. Shields, “Information theory and statistics: A tutorial,” *Foundations and Trends in Communications and Information Theory*, vol. 1, no. 4, pp. 417–528, 2004.
- [42] J. Hiriart-Urruty and C. Lemaréchal, *Fundamentals of Convex Analysis*. Berlin: Springer, 2001.
- [43] O. Cappé, E. Moulines, and T. Rydén, *Inference in Hidden Markov Models*. Springer, 2005.

- [44] O. V. Sarmanov, “Maximal coefficient of correlation (nonsymmetric case),” *Doklady Akad. Nauk SSSR*, vol. 121, no. 1, pp. 52–55, 1958.
- [45] R. Bhatia, *Matrix Analysis*. New York: Springer, 1997.
- [46] R. Montenegro and P. Tetali, “Mathematical aspects of mixing times in Markov chains,” *Foundations and Trends in Theoretical Computer Science*, vol. 1, no. 3, pp. 237–354, 2006.
- [47] D. A. Levin, Y. Peres, and E. L. Wilmer, *Markov Chains and Mixing Times*. Amer. Math. Soc., 2008.
- [48] F. Hansen, “The fast track to Loewner’s theorem,” *Lin. Algebra Appl.*, vol. 438, pp. 4557–4571, 2013.
- [49] K. Marton, “A simple proof of the blowing up lemma,” *IEEE Trans. Inform. Theory*, vol. 32, no. 3, pp. 445–446, 1986.
- [50] —, “Bounding  $\bar{d}$ -distance by informational divergence: a method to prove measure concentration,” *Ann. Probab.*, vol. 24, no. 2, pp. 857–866, 1996.
- [51] C. Villani, *Topics in Optimal Transportation*, ser. Graduate Studies in Mathematics. Providence, RI: Amer. Math. Soc., 2003, vol. 58.
- [52] S. G. Bobkov and F. Götze, “Exponential integrability and transportation cost related to logarithmic Sobolev inequalities,” *J. Funct. Anal.*, vol. 163, pp. 1–28, 1999.
- [53] E. Ordentlich and M. J. Weinberger, “A distribution dependent refinement of Pinsker’s inequality,” *IEEE Trans. Inform. Theory*, vol. 51, no. 5, pp. 1836–1840, May 2005.
- [54] S. G. Bobkov, C. Houdré, and P. Tetali, “The subgaussian constant and concentration inequalities,” *Israel J. Math.*, vol. 156, no. 1, pp. 255–283, December 2006.
- [55] D. Weitz, “Mixing in time and space for discrete spin systems,” Ph.D. dissertation, University of California, Berkeley, 2004.
- [56] F. Martinelli, A. Sinclair, and D. Weitz, “Glauber dynamics on trees: boundary conditions and mixing time,” *Commun. Math. Phys.*, vol. 250, pp. 301–334, 2004.
- [57] M. Ledoux, *The Concentration of Measure Phenomenon*. Amer. Math. Soc., 2001.
- [58] M. Raginsky and I. Sason, *Concentration of Measure Inequalities in Information Theory, Communications, and Coding*, 2nd ed. Now Publishers, 2014.
- [59] C. Houdré and P. Tetali, “Concentration of measure for products of Markov kernels and graph products via functional inequalities,” *Comb. Probab. Comput.*, vol. 10, pp. 1–28, 2001.
- [60] E. Erkip and T. M. Cover, “The efficiency of investment information,” *IEEE Trans. Inform. Theory*, vol. 44, no. 3, pp. 1026–1040, May 1998.
- [61] D. P. Palomar and S. Verdú, “Lautum information,” *IEEE Trans. Inform. Theory*, vol. 54, no. 3, pp. 964–975, March 2008.

- [62] S. Boyd, P. Diaconis, and L. Xiao, “Fastest mixing Markov chain on a graph,” *SIAM Review*, vol. 46, no. 4, pp. 667–689, 2004.
- [63] M. Jerrum, *Counting, Sampling, and Integrating: Algorithms and Complexity*. Birkhäuser, 2003.
- [64] G. Winkler, *Image Analysis, Random Fields, and Markov Chain Monte Carlo Methods: A Mathematical Introduction*, 2nd ed. Springer, 2003.
- [65] G. Grimmett, *The Random Cluster Model*. Berlin: Springer, 2006.
- [66] C. Borgs, J. T. Chayes, and P. Tetali, “Tight bounds for mixing of the Swendsen–Wang algorithm at the Potts transition point,” *Prob. Theory Rel. Fields*, vol. 152, pp. 509–557, 2012.
- [67] M. Ullrich, “Comparison of Swendsen–Wang and heat-bath dynamics,” *Random Struct. Alg.*, vol. 42, pp. 520–535, 2012.
- [68] M. J. Wainwright and M. I. Jordan, “Graphical models, exponential families, and variational inference,” *Foundations and Trends in Machine Learning*, vol. 1, no. 1-2, pp. 1–305, December 2008.
- [69] A. Montanari and N. Gerschenfeld, “Reconstruction for models on random graphs,” in *Proc. 48th IEEE Symp. on Foundations of Comp. Sci.*, 2007, pp. 194–204.
- [70] N. Bhatnagar, J. Vera, E. Vigoda, and D. Weitz, “Reconstruction for colorings on trees,” *SIAM J. Discrete Math.*, vol. 25, no. 2, pp. 809–826, 2011.
- [71] R. L. Dobrushin and S. B. Shlosman, “Completely analytical interactions: constructive description,” *J. Stat. Phys.*, vol. 46, no. 5/6, pp. 983–1014, 1987.