

# Entropy contractions in Markov chains: half-step, full-step and continuous-time

Pietro Caputo\*    Zongchen Chen†    Yuzhou Gu‡    Yury Polyanskiy§

## Abstract

This paper considers the speed of convergence (mixing) of a finite Markov kernel  $P$  with respect to the Kullback-Leibler divergence (entropy). Given a Markov kernel one defines either a discrete-time Markov chain (with the  $n$ -step transition kernel given by the matrix power  $P^n$ ) or a continuous-time Markov process (with the time- $t$  transition kernel given by  $e^{t(P-\text{Id})}$ ). The contraction of entropy for  $n = 1$  or  $t = 0+$  are given by the famous functional inequalities, the *strong data processing inequality (SDPI)* and the *modified log-Sobolev inequality (MLSI)*, respectively. When  $P = KK^*$  is written as the product of a kernel and its adjoint, one could also consider the “half-step” contraction, which is the SDPI for  $K$ , while the SDPI for  $P$  is called the “full-step” contraction. Del Moral, Ledoux and Miclo (PTRF, 2003) claimed that these contraction coefficients (half-step, full-step, and continuous-time) are generally comparable, that is their ratio is bounded from above and below by two absolute constants. We disprove this and related conjectures by working out a number of different counterexamples. In particular, we construct (a) a continuous-time Markov process that contracts arbitrarily faster than its discrete-time counterpart; and (b) a kernel  $P$  such that  $P^{m+1}$  contracts arbitrarily better than  $P^m$ . Hence, our main conclusion is that the four standard inequalities comparing five known notions of entropy contraction are generally not improvable (even in the subclass of factorizable Markov chains).

In the process of analyzing the counterexamples, we survey and sharpen the tools for bounding the contraction coefficients and characterize properties of extremizers of the respective functional inequalities, showing, for example, that while MLSI extremizer always has full support (unless MLSI constant equals twice the spectral gap), the SDPI extremizers can have partial support. As our examples range from Bernoulli-Laplace, random walks on graphs to birth-death chains, the paper is also intended as a tutorial on computing MLSI, SDPI and other constants for these types of commonly occurring Markov chains.

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Continuous-time contraction notions . . . . .	4
1.2	Discrete-time contraction notions . . . . .	5
1.3	Continuous-time versus discrete-time contraction . . . . .	6

---

\*pietro.caputo@uniroma3.it. Università Roma Tre.

†chenzongchen@gatech.edu. Georgia Institute of Technology.

‡yuzhougu@nyu.edu. New York University.

§yp@mit.edu. Massachusetts Institute of Technology.

<b>2</b>	<b>Factorizable kernels</b>	<b>8</b>
2.1	Important classes of factorizable kernels . . . . .	9
2.2	Non-uniqueness of factorization . . . . .	10
2.3	Characterizations of factorizable kernels . . . . .	11
<b>3</b>	<b>Comparison between constants</b>	<b>12</b>
3.1	Log-Sobolev constant $\rho$ vs half-step entropy contraction $\alpha$ . . . . .	12
3.2	Half-step entropy contraction $\alpha$ vs full-step entropy contraction $\delta$ . . . . .	13
3.3	Full-step entropy contraction $\delta$ vs modified log-Sobolev constant $\rho_0$ . . . . .	18
3.4	Modified log-Sobolev constant $\rho_0$ vs Poincaré constant $\lambda$ . . . . .	21
3.5	Other comparisons . . . . .	21
3.6	Comments on several previous works . . . . .	21
<b>4</b>	<b>Extremal functions</b>	<b>22</b>
4.1	Log-Sobolev constant $\rho$ , modified log-Sobolev constant $\rho_0$ , Poincaré constant $\lambda$ . . . . .	22
4.2	Half-step entropy contraction $\alpha$ and full-step entropy contraction $\delta$ . . . . .	23

# 1 Introduction

Markov chains are widely used in almost all areas of science and engineering, in the form of MCMC averaging in numerical analysis, approximation algorithms in computer science, generative modeling in artificial intelligence, and so on. Perhaps the most important problem in the study of Markov chains is to understand their equilibration properties. One example of such property is the mixing time, characterizing the time it takes for the marginal distribution to come close to the stationary one, as measured by some statistical distance: the total variation (most commonly) or Wasserstein distance,  $\chi^2$  or Kullback-Leibler (KL) divergence.

In this paper, we focus on Markov chains on a finite state space. There are two common ways to define or implement Markov processes. In a discrete-time Markov chain, the state is updated at every discrete time  $t \in \mathbb{Z}_{\geq 1}$ ; meanwhile, in a continuous-time Markov chain, the update times are distributed as a Poisson point process on the positive real axis  $\mathbb{R}_{\geq 0}$ . For example, if the single-step kernel for the former is given by a row-stochastic matrix  $P$  then the corresponding continuous-time chain has kernel  $T_t = e^{t(P-\text{Id})}$ . Due to concentration of the number of updates in the time interval  $[0, t]$ , both versions are known to have almost equal mixing times in *total variation* (e.g., [LP17, Theorem 20.3]). Thus, for the study of total-variation mixing time probabilists are free to switch between the discrete and continuous times as convenient.

One common approach (e.g., [DSC96]) to show rapid mixing of Markov chains is to prove that they “make progress” step-wise in terms of some  $f$ -divergence, most commonly the  $\chi^2$  or the KL divergence. For discrete-time Markov processes, this means that the  $f$ -divergence to the target distribution decreases by a certain factor in every step, which can be described by the contraction coefficient of the associated Markov kernel. For continuous-time chains, this corresponds to the derivative of  $f$ -divergence with respect to time is suitably bounded away from zero from below, and hence the  $f$ -divergence to the target distribution decreases at a certain speed. For the KL divergence, the respective contraction inequalities are known as the *strong data processing inequality (SDPI)* and the *modified log-Sobolev inequality (MLSI)*, respectively for discrete-time and continuous-time. See Eqs. (5) and (13) below for formal definitions.

For a large family of Markov chains, including the Glauber dynamics and random walks on high-dimensional simplices (e.g., [KM17]), the transition kernel  $P$  consists of two stages, for which  $P = KK^*$  is written as the product of a kernel and its adjoint. The first stage ( $K$ ) is often described

as the downward, forward, or noising move, and the second stage ( $K^*$ ) as the upward, backward, or denoising move. Recent success on analyzing such two-stage processes (e.g., [CLV21, BCC<sup>+</sup>22]) considers the decay of  $f$ -divergence in a single stage  $K$  or  $K^*$ , which for many specific problems turns out to be easier to handle. For such chains it is then natural to define a *half-step* contraction coefficient (corresponding to application of  $K$  only) and ask how it compares with full-step, multi-step and continuous-time ones.

Given the equivalence between continuous-time and discrete-time mixing times, one naturally expects similar equivalence between the contraction coefficients. This turns out to be true for the  $\chi^2$  contraction ([Rag16, Remark 4.2]). Thus, it was not surprising when the work [DMLM03] claimed to show such equivalence for the KL divergence contraction as well (that is showing that the ratio of the MLSI and SDPI constants is universally bounded). Specifically, they claimed an equivalence (up to universal multiplicative factors) between the MLSI and the *half-step SDPI*, which implies in particular that the MLSI and the *full-step SDPI* (i.e., continuous-time and discrete-time entropy contraction) are equivalent. The present paper originated from us discovering a gap in their proof (see Section 3.6) and realizing that their claim cannot hold true because the ratio between the MLSI constant and the half-step SDPI constant can be arbitrarily large for the random transposition model (Example 24). However, the question of whether the MLSI and the full-step SDPI are equivalent remained open.

We answer this question by presenting an example (Example 20) separating the MLSI and the full-step SDPI. That is, there exist cases where the contraction rate of the discrete-time chain can be much slower than that of the continuous-time one. The example is adapted from [Mün23]’s counterexample to the Peres-Tetali conjecture, although there seems to be no direct relationship between the properties used here and op. cit.

The separation between continuous-time and discrete-time entropy contraction leads us to consider the related question of comparing the entropy contraction of the  $m$ -step kernel  $P^m$  versus the  $(m + 1)$ -step one  $P^{m+1}$ . If a separation is possible, then it would explain how a continuous-time chain (which averages over many  $P^m$ ’s) could have better convergence properties at finite time than the corresponding discrete-time chain (see discussions at the end of Section 3.2). It turns out that a counterexample is again possible (see Example 18 below).

To avoid trivial counterexamples, we restrict our attention to *factorizable* kernels, which are kernels that can be written in the form  $P = KK^*$ . As discussed above, this is a natural class of Markov kernels to study. All our examples are factorizable.

In all, the main purpose of this paper is to give a self-contained and thorough introduction of all notions of relative entropy decay for finite-state Markov chains, for both continuous-time and discrete-time versions. We summarize known comparisons among these notions (including, the recent half-step contraction) and we give examples, demonstrating, that in all cases where comparisons are not available there exist counterexamples, in the sense that the ratio can be arbitrarily large. Along the way, we correct several misstatements that appeared in previous works, we show how to get sharp upper and lower bounds for these coefficients and study the extremizers in the respective functional inequalities.

**Organization.** This paper’s content is succinctly summarized in three tables: Table 1 gives definitions of 5 main contraction coefficients, Table 2 lists all known comparison inequalities along with (new) counterexamples for the missing comparisons, Table 3 summarizes the list of counterexamples. The rest of the introduction defines all notions rigorously and recalls the standard comparison chain (27):

$$\rho \leq \alpha \leq \delta \leq \rho_0 \leq 2\lambda.$$

Then, after introducing factorizable Markov kernels ( $P = KK^*$ ) in Section 2, Section 3 shows that

every inequality above can have arbitrarily high ratio, even when restricted to factorizable kernels. Section 4 concludes with some results on the properties of extremizers of functional inequalities on complete and complete bipartite graphs.

**Notation.** Throughout the paper, let  $\mathcal{X}$  be a finite set,  $P : \mathcal{X} \rightarrow \mathcal{X}$  be a Markov kernel with an invariant distribution  $\pi$ . Assume that  $(\pi, P)$  is reversible, i.e.,  $\pi(x)P(x, y) = \pi(y)P(y, x)$  for all  $x, y \in \mathcal{X}$ . For a Markov kernel  $K : \mathcal{X} \rightarrow \mathcal{Y}$  and a distribution  $\pi$  on  $\mathcal{X}$ , we define the reverse channel  $K_\pi^* : \mathcal{Y} \rightarrow \mathcal{X}$  as

$$K_\pi^*(y, x) = \begin{cases} \frac{\pi(x)K(x, y)}{(\pi K)(y)}, & \text{if } (\pi K)(y) > 0, \\ \pi(x), & \text{otherwise.} \end{cases} \quad (1)$$

Note that  $\pi$  is an invariant distribution of  $KK_\pi^*$  and  $(\pi, KK_\pi^*)$  is reversible.

### 1.1 Continuous-time contraction notions

For  $f, g : \mathcal{X} \rightarrow \mathbb{R}$ , define the Dirichlet form as

$$\mathcal{E}_{\pi, P}(f, g) = -\pi[f(Lg)] \quad (2)$$

where  $L = P - I$  is the Markov generator.

The *Poincaré constant* (also called the *spectral gap*)  $\lambda = \lambda(\pi, P)$  is the largest number such that

$$\lambda \text{Var}_\pi(f) \leq \mathcal{E}_{\pi, P}(f, f), \quad \forall f : \mathcal{X} \rightarrow \mathbb{R}, \quad (3)$$

where  $\text{Var}_\pi(f) := \pi[f^2] - (\pi[f])^2$  is the variance of  $f$  under  $\pi$ . When  $f = \frac{d\nu}{d\pi}$  for some  $\nu \in \mathcal{P}(\mathcal{X})$  (where  $\mathcal{P}(\mathcal{X})$  denotes the space of distributions on  $\mathcal{X}$ ), we have  $\text{Var}_\pi(f) = \chi^2(\nu||\pi)$  where  $\chi^2(\cdot||\cdot)$  stands for the  $\chi^2$ -divergence. Eq. (3) is called the *Poincaré inequality*.

The *log-Sobolev constant* (LSC)  $\rho = \rho(\pi, P)$  is the largest number such that

$$\rho \text{Ent}_\pi(f) \leq \mathcal{E}_{\pi, P}(\sqrt{f}, \sqrt{f}), \quad \forall f : \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}, \quad (4)$$

where  $\text{Ent}_\pi(f) := \pi \left[ f \log \frac{f}{\pi[f]} \right]$  is the entropy of  $f$ . When  $f = \frac{d\nu}{d\pi}$  for some  $\nu \in \mathcal{P}(\mathcal{X})$ , we have  $\text{Ent}_\pi(f) = D(\nu||\pi)$  where  $D(\cdot||\cdot)$  stands for the Kullback-Leibler (KL) divergence. Eq. (4) is called the *log-Sobolev inequality* (LSI).

The *modified log-Sobolev constant* (MLSC)  $\rho_0 = \rho_0(\pi, P)$  is the largest number such that

$$\rho_0 \text{Ent}_\pi(f) \leq \mathcal{E}_{\pi, P}(f, \log f), \quad \forall f : \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}. \quad (5)$$

Eq. (5) is called the *modified log-Sobolev inequality* (MLSI).

For reversible  $(\pi, P)$ , we always have ([DSC96, BT06])

$$4\rho \leq \rho_0 \leq 2\lambda. \quad (6)$$

The constants  $\rho$ ,  $\rho_0$ , and  $\lambda$  represent the contraction ability of the continuous-time Markov chain (also known as the Markov semigroup)  $(T_t)_{t \geq 0}$ , where  $T_t = \exp(-tL)$ .<sup>1</sup> These constants are properties of the Markov generator  $L$  (which can be arbitrarily scaled), rather than the Markov kernel  $P$ .

<sup>1</sup>In this work we focus on continuous-time Markov chains of this type and do not consider more general continuous-time Markov chains.

Let  $\nu$  be a distribution on  $\mathcal{X}$ . Define  $\nu_t = \nu T_t$  to be the distribution of the Markov chain at time  $t \geq 0$  when initialized at  $\nu$ , and define  $f_t = \frac{d\nu_t}{d\pi}$  to be the relative density. A direct computation yields

$$\frac{d}{dt} \text{Var}_\pi(f_t) = -2\mathcal{E}_{\pi,P}(f_t, f_t), \quad (7)$$

$$\frac{d}{dt} \text{Ent}_\pi(f_t) = -\mathcal{E}_{\pi,P}(f_t, \log f_t). \quad (8)$$

Therefore the Poincaré inequality and the modified log-Sobolev inequality can be equivalently stated as

$$\frac{d}{dt} \text{Var}_\pi(f_t) \leq -2\lambda \text{Var}_\pi(f_t), \quad (9)$$

$$\frac{d}{dt} \text{Ent}_\pi(f_t) \leq -\rho_0 \text{Ent}_\pi(f_t). \quad (10)$$

Eqs. (9) and (10) can also be understood as alternative definitions for  $\lambda$  and  $\rho_0$ , from which one immediately obtains

$$\text{Var}_\pi(f_t) \leq \exp(-2\lambda t) \text{Var}_\pi(f_0), \quad (11)$$

$$\text{Ent}_\pi(f_t) \leq \exp(-\rho_0 t) \text{Ent}_\pi(f_0). \quad (12)$$

We remark that the log-Sobolev inequality is equivalent to hypercontractivity by [DSC96]. Furthermore, [BG99] shows that the log-Sobolev inequality is equivalent (up to a constant factor) to the Poincaré inequality on an Orlicz space. In this work we consider only the Poincaré inequality on the  $L^2$  space.

We refer the reader to [DSC96, BT06] for more discussions on (modified) log-Sobolev constants.

## 1.2 Discrete-time contraction notions

Another class of contraction notions comes from contraction coefficients for  $f$ -divergences. We refer the reader to [PW24, Chapter 7] for an introduction to  $f$ -divergences. Let  $K : \mathcal{X} \rightarrow \mathcal{Y}$  be a Markov kernel and  $\pi$  be a distribution on  $\mathcal{X}$ . For any  $f$ -divergence, we define the (input-restricted<sup>2</sup>)  $f$ -contraction coefficient

$$\eta_f(\pi, K) := \sup_{\substack{\nu \in \mathcal{P}(\mathcal{X}) \\ 0 < D_f(\nu \| \pi) < \infty}} \frac{D_f(\nu K \| \pi K)}{D_f(\nu \| \pi)}. \quad (13)$$

In other words, we have

$$D_f(\nu K \| \pi K) \leq \eta_f(\pi, K) D_f(\nu \| \pi), \quad (14)$$

known as the *strong data processing inequality* (SDPI). By the data processing inequality (DPI), we always have  $0 \leq \eta_f(\pi, K) \leq 1$ , and a smaller value means a stronger contraction ability for  $f$ -divergence. The most commonly used  $f$ -contraction coefficients include the total variation (TV) contraction coefficient  $\eta_{\text{TV}}(\pi, K)$ , the Kullback-Leibler (KL) contraction coefficient  $\eta_{\text{KL}}(\pi, K)$  (the

---

<sup>2</sup>One could also consider the input-unrestricted contraction coefficient defined as  $\eta_f(K) := \sup_{\pi \in \mathcal{P}(\mathcal{X})} \eta_f(\pi, K)$ . In this paper we focus on the input-restricted version.

subscript KL is sometimes omitted), and the  $\chi^2$ -contraction coefficient  $\eta_{\chi^2}(\pi, K)$ . It is known ([AG76, PW17]) that

$$\eta_{\chi^2}(\pi, K) \leq \eta_{\text{KL}}(\pi, K) \leq \eta_{\text{TV}}(\pi, K). \quad (15)$$

The TV contraction coefficient  $\eta_{\text{TV}}(\pi, K)$  is also known as the Dobrushin's coefficient ([Dob56]), and satisfies

$$\eta_{\text{TV}}(\pi, K) = \max_{x, x' \in \mathcal{X}} \text{TV}(K(x, \cdot), K(x', \cdot)). \quad (16)$$

Via Eq. (15), Eq. (16) provides an easy upper bound for  $\eta_{\chi^2}(\pi, K)$  and  $\eta_{\text{KL}}(\pi, K)$ .

We refer the reader to [Rag16, PW24] for more discussions on contraction coefficients.

Let  $P = KK_{\pi}^*$ . We define the half-step entropy contraction coefficient

$$\alpha(\pi, K) = 1 - \eta_{\text{KL}}(\pi, K) = \inf_{\substack{\nu \in \mathcal{P}(\mathcal{X}) \\ 0 < D(\nu \| \pi) < \infty}} \frac{D(\nu \| \pi) - D(\nu K \| \pi K)}{D(\nu \| \pi)} \quad (17)$$

and the full-step entropy contraction coefficient

$$\delta(\pi, P) = 1 - \eta_{\text{KL}}(\pi, P) = \inf_{\substack{\nu \in \mathcal{P}(\mathcal{X}) \\ 0 < D(\nu \| \pi) < \infty}} \frac{D(\nu \| \pi) - D(\nu P \| \pi)}{D(\nu \| \pi)}. \quad (18)$$

By rearranging, we have inequalities

$$\text{Ent}_{\pi K}(K_{\pi}^* f) - \text{Ent}_{\pi}(f) \leq -\alpha \text{Ent}_{\pi}(f), \quad (19)$$

$$\text{Ent}_{\pi}(P f) - \text{Ent}_{\pi}(f) \leq -\delta \text{Ent}_{\pi}(f), \quad (20)$$

for all  $f : \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$ . Eqs. (19) and (20) can be seen as definitions of  $\alpha$  and  $\delta$ , and can be compared with Eqs. (11) and (12). In the next section, we will discuss relationships between the discrete-time contraction notions  $\alpha$ ,  $\delta$  and the continuous-time contraction notions  $\rho$ ,  $\rho_0$ ,  $\lambda$ .

### 1.3 Continuous-time versus discrete-time contraction

As we discussed above, the log-Sobolev constant  $\rho(\pi, P)$ , modified log-Sobolev constant  $\rho_0(\pi, P)$ , and the Poincaré constant  $\lambda(\pi, P)$  represent the contraction ability of the continuous Markov chain, while the contraction coefficients represent the contraction ability of the discrete-time Markov chain. These constants allow one to derive mixing time bounds for associated Markov chains in the continuous-time and discrete-time settings respectively, see e.g. [Cap23] and the references therein for more details. In this paper, we give a full comparison between the discrete-time contraction notions and the continuous-time contraction notions.

The  $\chi^2$ -contraction coefficient has a close relationship with the Poincaré constant  $\lambda$ . We have

$$1 - \eta_{\chi^2}(\pi, K) = \lambda(\pi, KK_{\pi}^*). \quad (21)$$

For  $P = KK_{\pi}^*$ , the two versions of contraction are equivalent up to a constant factor. By Eq. (21) and  $1 - \lambda(\pi, P^2) = (1 - \lambda(\pi, P))^2$ , we have

$$1 - \sqrt{\eta_{\chi^2}(\pi, P)} = \lambda(\pi, P). \quad (22)$$

In particular,

$$\frac{1}{2}(1 - \eta_{\chi^2}(\pi, P)) \leq \lambda(\pi, P) \leq 1 - \eta_{\chi^2}(\pi, P). \quad (23)$$

This shows that the rates of  $\chi^2$ -divergence contraction for continuous-time and discrete-time Markov chains are within a factor of two of each other.

If we consider entropy (KL divergence) rather than  $\chi^2$ -divergence, then previous works have shown a one-sided inequality. [BCP<sup>+</sup>21] shows that

$$\delta(\pi, P) \leq \rho_0(\pi, P). \quad (24)$$

[Mic97] proves that

$$\rho(\pi, P) \leq \alpha(\pi, K). \quad (25)$$

By the data processing inequality, we have

$$\alpha(\pi, K) \leq \delta(\pi, P). \quad (26)$$

Summarizing the above, we have a chain of inequalities

$$\rho \leq \alpha \leq \delta \leq \rho_0 \leq 2\lambda. \quad (27)$$

In other words, we have the following implications:

$$\begin{aligned} \text{Log-Sobolev Inequality} &\implies \text{Half-Step Entropy Contraction} \implies \text{Full-Step Entropy Contraction} \\ &\implies \text{Modified Log-Sobolev Inequality} \implies \text{Poincaré Inequality} \end{aligned} \quad (28)$$

In [Section 3](#), we show that the gap between any two adjacent constants in [Eq. \(27\)](#) can be arbitrarily large. In particular, unlike the  $\chi^2$ -divergence, the discrete-time entropy contraction  $\delta(\pi, P)$  is not equivalent to the continuous-time entropy contraction  $\rho_0(\pi, P)$ .

[Table 1](#) summarizes the main constants discussed in this paper, and [Table 2](#) summarizes relationships between these contraction notions. While these contraction notions are in general non-equivalent (up to a constant factor), we will see in [Section 3.5](#) that under extra conditions, some of them could become equivalent for certain chains.

Symbol	Name	Inequality (Definition)
$\rho(\pi, P)$	Log-Sobolev Constant	$\rho \text{Ent}_\pi(f) \leq \mathcal{E}_{\pi, P}(\sqrt{f}, \sqrt{f})$
$\alpha(\pi, K)$	Half-Step Entropy Contraction	$\alpha \text{Ent}_\pi(f) \leq \text{Ent}_\pi(f) - \text{Ent}_{\pi K}(K_\pi^* f)$
$\delta(\pi, P)$	Full-Step Entropy Contraction	$\delta \text{Ent}_\pi(f) \leq \text{Ent}_\pi(f) - \text{Ent}_\pi(Pf)$
$\rho_0(\pi, P)$	Modified Log-Sobolev Constant	$\rho_0 \text{Ent}_\pi(f) \leq \mathcal{E}_{\pi, P}(f, \log f)$
$\lambda(\pi, P)$	Poincaré Constant	$\lambda \text{Var}_\pi(f) \leq \mathcal{E}_{\pi, P}(f, f)$

Table 1: Contraction notions discussed in this paper. We assume that  $P = KK_\pi^*$ . See [Table 2](#) for relationships between the constants.

Relationship	Explanation	Reference
$\rho(\pi, P) \leq \alpha(\pi, K)$	Log-Sobolev Inequality $\Rightarrow$ Half-Step Entropy Contraction	[Mic97, Prop. 6]
$\rho(\pi, P) \not\lesssim \alpha(\pi, K)$	Half-Step Entropy Contraction $\not\Rightarrow$ Log-Sobolev Inequality	Section 3.1
$\alpha(\pi, K) \leq \delta(\pi, P)$	Half-Step Entropy Contraction $\Rightarrow$ Full-Step Entropy Contraction	Data processing inequality
$\alpha(\pi, K) \not\lesssim \delta(\pi, P)$	Full-Step Entropy Contraction $\not\Rightarrow$ Half-Step Entropy Contraction	Section 3.2
$\delta(\pi, P^m) \not\lesssim \delta(\pi, P^{m+1})$	$(m+1)$ -Step Entropy Contraction $\not\Rightarrow$ $m$ -Step Entropy Contraction	Section 3.2
$\delta(\pi, P) \leq \rho_0(\pi, P)$	Full-Step Entropy Contraction $\Rightarrow$ Modified Log-Sobolev Inequality	[BCP <sup>+</sup> 21, Lemma 2.7]
$\delta(\pi, P) \not\lesssim \rho_0(\pi, P)$	Modified Log-Sobolev Inequality $\not\Rightarrow$ Full-Step Entropy Contraction	Section 3.3
$\rho_0(\pi, P) \leq 2\lambda(\pi, P)$	Modified Log-Sobolev Inequality $\Rightarrow$ Poincaré Inequality	[BT06, Prop. 3.6]
$\rho_0(\pi, P) \not\lesssim \lambda(\pi, P)$	Poincaré Inequality $\not\Rightarrow$ Modified Log-Sobolev Inequality	Section 3.4

Table 2: Relationships between contraction notions. The setting is the same as Table 1.  $a \not\lesssim b$  means  $\frac{b}{a}$  can be arbitrarily large. Inequality A  $\Rightarrow$  Inequality B means Inequality A with constant  $a$  implies Inequality B with constant  $Ca$  for some absolute constant  $C > 0$ .

## 2 Factorizable kernels

**Definition 1** (Factorizable pairs). Let  $\mathcal{X}$  be a finite set and  $P : \mathcal{X} \rightarrow \mathcal{X}$  be a Markov kernel with invariant distribution  $\pi$ . We say  $(\pi, P)$  is factorizable if  $P = KK_\pi^*$  for some finite Markov kernel  $K : \mathcal{X} \rightarrow \mathcal{Y}$ .

**Lemma 2** (Factorizable implies reversible). *A factorizable pair  $(\pi, P)$  is reversible.*

*Proof.* Suppose  $P = KK_\pi^*$  for a finite Markov kernel  $K : \mathcal{X} \rightarrow \mathcal{Y}$ . For  $x, y \in \mathcal{X}$ , we have

$$\begin{aligned} \pi(x)P(x, y) &= \pi(x) \sum_{z \in \mathcal{Y}} K(x, z)K_\pi^*(z, y) \\ &= \pi(x)\pi(y) \sum_{z \in \mathcal{Y}} \frac{1}{(\pi K)(z)} K(x, z)K(y, z) = \pi(y)P(y, x). \end{aligned} \tag{29}$$

So  $(\pi, P)$  is reversible. □

Factorizability is a reasonable assumption for several reasons. Recall that  $\rho$ ,  $\rho_0$ , and  $\lambda$  are properties of the Markov generator  $L$ , while  $\alpha$  and  $\delta$  are properties of the Markov kernel  $P$ , and the two are related by  $L = P - I$ . If we do not make extra assumptions on  $P$ , then it could happen that for two Markov kernels  $P_1$  and  $P_2$ ,  $P_1$  contracts better than  $P_2$ , but the corresponding Markov generators  $L_1$  and  $L_2$  satisfy that  $L_2$  contracts better than  $L_1$ . Consider the example where  $\mathcal{X} = [2]$ ,  $\pi = \text{Unif}(\mathcal{X})$ ,  $P_c = \begin{pmatrix} 1-c & c \\ c & 1-c \end{pmatrix}$  for  $c \in [0, 1]$ . Then the contraction ability (as Markov kernels) is increasing for  $c \in [0, \frac{1}{2}]$  and decreasing for  $c \in [\frac{1}{2}, 1]$ . However, the contraction



ability for the corresponding Markov generators  $L_c = P_c - I$  is increasing on the whole interval  $[0, 1]$ . To avoid such undesirable behavior, we need to impose extra assumptions like factorizability. Furthermore, Eqs. (25) and (26) imply that

$$\rho(\pi, P) \leq \delta(\pi, P) \quad (30)$$

for factorizable  $P$ . While the statement of Eq. (30) does not involve factorizability, it does not hold if  $P$  is not factorizable. Consider the same example with  $c = 1$ , that is,  $P = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$ . Then  $\delta(\pi, P) = 0$  but  $\rho(\pi, P) = 1$ .

## 2.1 Important classes of factorizable kernels

Another reason that the factorizability assumption is reasonable is that many natural Markov chains considered in the literature are factorizable. In this section we discuss two important classes of factorizable kernels: lazy kernels and the Glauber dynamics.

**Lemma 3** (Lazy implies factorizable). *Let  $(\pi, P)$  be a reversible pair. If  $P(x, x) \geq \frac{1}{2}$  for all  $x \in \mathcal{X}$ , then there exists  $\mathcal{Y}$  and  $K : \mathcal{X} \rightarrow \mathcal{Y}$  such that  $P = KK_\pi^*$ .*

*Proof.* Let  $\mathcal{Y} = \binom{\mathcal{X}}{1} \cup \binom{\mathcal{X}}{2}$ . Let  $K : \mathcal{X} \rightarrow \mathcal{Y}$  be the map

$$K(x, e) = \begin{cases} 2P(x, x) - 1, & \text{if } e = \{x\}, \\ 2P(x, y), & \text{if } e = \{x, y\}, \\ 0, & \text{if } x \notin e. \end{cases} \quad (31)$$

Then  $K_\pi^*(e, \cdot) = \text{Unif}(e)$ . For  $y \neq x \in \mathcal{X}$ , we have

$$(KK_\pi^*)(x, y) = K(x, \{x, y\})K_\pi^*(\{x, y\}, y) = P(x, y). \quad (32)$$

Therefore  $P = KK_\pi^*$ .  $\square$

Lazy chains are quite common. In particular, many of our examples are lazy random walks on regular graphs.

**Definition 4** (Lazy random walk Markov chain). Let  $G = (V, E)$  be a  $d$ -regular graph. We associate with it a canonical Markov chain called the (lazy) random walk. Let  $\mathcal{X} = V$ ,  $\mathcal{Y} = E$ , and  $K : \mathcal{X} \rightarrow \mathcal{Y}$  be  $K(x, e) = \frac{1}{d}\mathbb{1}\{x \in e\}$  for  $x \in \mathcal{X}$ ,  $e \in \mathcal{Y}$ . Then  $K_\pi^*(e, \cdot) = \text{Unif}(e)$  and  $P = KK_\pi^*$  satisfies

$$P(x, y) = \begin{cases} \frac{1}{2}, & \text{if } x = y, \\ \frac{1}{2d}, & \text{if } (xy) \in E, \\ 0, & \text{otherwise.} \end{cases} \quad (33)$$

The Glauber dynamics is an important class of Markov chains that have shown huge practical and theoretical success in sampling spin systems. Let us consider a general  $\alpha$ -weighted block dynamics, defined as follows. Let  $\pi$  be a probability measure on  $\Omega = [q]^n$  (e.g., the Ising model on a graph with  $n$  vertices). For any  $\sigma, \eta \in \Omega$  and  $A \subseteq [n]$ , consider the conditional probability of  $\eta$  given the configuration  $\sigma$  on  $A$ :

$$\pi(\eta|\sigma_A) = \frac{\pi(\eta)\mathbb{1}\{\sigma|_A = \eta|_A\}}{\sum_{\eta' \in \Omega} \pi(\eta')\mathbb{1}\{\sigma|_A = \eta'|_A\}}. \quad (34)$$

Let  $\mathcal{S} = 2^{[n]}$  be the set of all subsets of  $[n]$ . For any probability measure  $\alpha = (\alpha_A)_{A \in \mathcal{S}}$  on  $\mathcal{S}$ , the  $\alpha$ -weighted block dynamics is the Markov chain on  $\Omega$  with transition matrix

$$P(\sigma, \sigma') = \sum_{A \in \mathcal{S}} \alpha_A \pi(\sigma' | \sigma_{A^c}) \quad (35)$$

where  $A^c = [n] \setminus A$ . The case  $\alpha_A = \frac{1}{n} \mathbb{1}\{|A| = 1\}$  is known as the Glauber dynamics. The pair  $(\pi, P)$  is reversible. A factorization of the  $\alpha$ -weighted block dynamics can be defined as follows. Let  $\mathcal{X} = \Omega$ ,  $\mathcal{Y} = \mathcal{S} \times \Omega$ , and  $K : \mathcal{X} \rightarrow \mathcal{Y}$  be defined as

$$K(\sigma, (A, \eta)) = \alpha_A \pi(\eta | \sigma_{A^c}). \quad (36)$$

One can compute that  $(\pi K)(A, \eta) = \alpha_A \pi(\eta)$  and that  $K_\pi^*((A, \eta), \sigma) = \pi(\sigma | \eta_{A^c})$  for all  $(A, \eta) \in \mathcal{S} \times \Omega$ ,  $\sigma \in \Omega$ . Therefore, for all  $\sigma, \sigma' \in \Omega$ ,

$$\sum_{(A, \eta) \in \mathcal{S} \times \Omega} K(\sigma, (A, \eta)) K_\pi^*((A, \eta), \sigma') = \sum_{A \in \mathcal{S}} \alpha_A \sum_{\eta \in \Omega} \pi(\eta | \sigma_{A^c}) \pi(\sigma' | \eta_{A^c}) = P(\sigma, \sigma'). \quad (37)$$

Thus  $P = K K_\pi^*$  is a factorization of the  $\alpha$ -weighted block dynamics. Half-step entropy contractions for such Markov chains have been recently investigated in the context of spin systems under the name of block factorizations of entropy (e.g., [BCC<sup>+</sup>22]).

## 2.2 Non-uniqueness of factorization

In general, a factorizable pair  $(\pi, P)$  can be factorized in more than one different ways, and the associated half-step contraction rates may differ considerably. This is illustrated in the following example.

**Example 5** (Complete graph). Let  $n \geq 3$  be an integer. Let  $\mathcal{X} = [n]$ ,  $\pi = \text{Unif}(\mathcal{X})$ . Let  $P : \mathcal{X} \rightarrow \mathcal{X}$  be the lazy random walk on the complete graph. That is,

$$P(x, y) = \begin{cases} \frac{1}{2}, & \text{if } x = y, \\ \frac{1}{2(n-1)}, & \text{if } x \neq y. \end{cases} \quad (38)$$

Given an integer  $2 \leq \ell \leq n$ , let  $\mathcal{Y}$  be the set of all subsets  $A \subseteq [n]$  with either  $|A| = 1$  or  $|A| = \ell$ , and define  $K(\ell) : \mathcal{X} \rightarrow \mathcal{Y}$  as

$$K(\ell)(x, y) = \begin{cases} \frac{\ell-2}{2(\ell-1)}, & \text{if } y = \{x\}, \\ \frac{\ell}{2(\ell-1)\binom{n-1}{\ell-1}}, & \text{if } |y| = \ell, x \in y, \\ 0, & \text{otherwise.} \end{cases} \quad (39)$$

One can check that  $K(\ell)_\pi^*(y, \cdot) = \text{Unif}(y)$  and that  $K(\ell)K(\ell)_\pi^* = P$ . We also note that when  $\ell = 2$  this reduces to the construction in the proof of [Lemma 3](#).

[BC24, Theorem 1.1] computes the half-step entropy contraction coefficient  $\alpha(\pi, K(\ell))$  for every  $\ell$ , and shows that

$$\alpha(\pi, K(\ell)) = \frac{\ell \log \ell}{2(\ell-1) \log n}, \quad (40)$$

achieved at and only at point distributions. From [Eq. \(40\)](#) we see that  $\alpha(\pi, K(\ell))$  increases with  $\ell$  from the minimum value  $\alpha(\pi, K(2)) = \frac{\log 2}{\log n}$  to the maximum value  $\alpha(\pi, K(n)) = \frac{n}{2(n-1)}$ . The former is of the same magnitude of the log-Sobolev constant  $\rho(\pi, P) = \frac{n-2}{2(n-1) \log(n-1)}$  ([DSC96, Corollary A.5]) and the latter matches asymptotically with the full-step contraction rate  $\delta(\pi, P) = \frac{1}{2} \pm o(1)$  ([GP23, Eq. (131)]).

### 2.3 Characterizations of factorizable kernels

An interesting question is to characterize the set of factorizable kernels for a fixed  $\pi$ . [Lemma 3](#) provides a sufficient condition and it is certainly not necessary. For example, any pair  $(\pi, P)$  satisfying  $P(x, \cdot) = \pi$  for all  $x \in \mathcal{X}$  is factorizable (see [Example 10](#)). On the other hand, a necessary condition for factorizability is positive semidefiniteness.

**Lemma 6** (Factorizable implies positive semidefinite). *Let  $(\pi, P)$  be a reversible pair. If  $P$  is factorizable, then the matrix  $\text{Diag}(\pi)P$  is positive semidefinite (PSD), where  $\text{Diag}(\pi)$  denotes the  $\mathcal{X} \times \mathcal{X}$  diagonal matrix with diagonal  $\pi$ .*

*Proof.* Let  $A = \text{Diag}(\pi)P$ . Let  $K : \mathcal{X} \rightarrow \mathcal{Y}$  be a Markov kernel such that  $P = KK_\pi^*$ . WLOG assume that  $\pi K$  has full support. Then

$$A_{x,y} = \pi(x) \sum_{z \in \mathcal{Y}} K(x,z)K_\pi^*(z,y) = \pi(x)\pi(y) \sum_{z \in \mathcal{Y}} \frac{1}{(\pi K)(z)} K(x,z)K(y,z). \quad (41)$$

So  $A = MM^\top$  where  $M = \text{Diag}(\pi)K \text{Diag}(\pi K)^{-1/2}$ . This finishes the proof.  $\square$

In particular, while a factorizable kernel is not necessarily lazy, when  $\pi$  has full support,  $P$  must have strictly positive diagonal entries.

For a distribution  $\pi$  on  $\mathcal{X}$ , let  $\mathcal{F}_\pi$  denote the set of Markov kernels  $P : \mathcal{X} \rightarrow \mathcal{X}$  such that  $(\pi, P)$  is factorizable.

**Lemma 7** ( $\mathcal{F}_\pi$  is convex). *For any distribution  $\pi$ , the set  $\mathcal{F}_\pi$  is convex.*

*Proof.* Let  $P_0, P_1 \in \mathcal{F}_\pi$  and  $K_0 : \mathcal{X} \rightarrow \mathcal{Y}_0, K_1 : \mathcal{X} \rightarrow \mathcal{Y}_1$  be the corresponding factors. Let  $t \in [0, 1]$ . We prove that  $P_t := (1-t)P_0 + tP_1$  is in  $\mathcal{F}_\pi$ . Let  $K_t : \mathcal{X} \rightarrow \mathcal{Y}_0 \sqcup \mathcal{Y}_1$  be defined as

$$K_t(x, y) = \begin{cases} (1-t)K_0(x, y), & \text{if } y \in \mathcal{Y}_0, \\ tK_1(x, y), & \text{if } y \in \mathcal{Y}_1. \end{cases} \quad (42)$$

Then we can verify that  $K_t$  is a Markov kernel, and

$$K_t(K_t)_\pi^* = (1-t)K_0(K_0)_\pi^* + tK_1(K_1)_\pi^* = P_t. \quad (43)$$

$\square$

For a distribution  $\pi$  on  $\mathcal{X}$ , let  $\mathcal{P}_\pi$  denote the set of Markov kernels  $P : \mathcal{X} \rightarrow \mathcal{X}$  such that  $\text{Diag}(\pi)P$  is PSD. By [Lemmas 6](#) and [7](#),  $\mathcal{F}_\pi$  is a convex subset of  $\mathcal{P}_\pi$ . When the state space is binary, the two sets are equal, but this is not true in general.

**Lemma 8.** *If  $|\mathcal{X}| = 2$ , then for any distribution  $\pi$  on  $\mathcal{X}$ , we have  $\mathcal{F}_\pi = \mathcal{P}_\pi$ .*

*Proof.* Direct calculation shows that  $\mathcal{P}_\pi$  is the convex hull of  $\{J \text{Diag}(\pi), I\}$ , where  $J$  is the all-ones matrix. Both extreme points are in  $\mathcal{F}_\pi$ .  $\square$

**Lemma 9.** *For  $n \geq 5$ ,  $\mathcal{X} = [n]$ , and  $\pi = \text{Unif}(\mathcal{X})$ , the set  $\mathcal{F}_\pi$  is strictly smaller than  $\mathcal{P}_\pi$ .*

*Proof.* An  $n \times n$  matrix  $A$  is called completely positive if it can be written as  $A = MM^\top$  for some (not necessarily square) matrix  $M$  with non-negative entries. Clearly, all completely positive matrices are PSD. It is known ([\[MM62, BSM03\]](#)) that for  $n \leq 4$ , a PSD matrix is completely

positive if and only if all its entries are non-negative, while for  $n \geq 5$ , there exist PSD matrices with strictly positive entries that are not completely positive.

Fix  $n \geq 5$ ,  $\mathcal{X} = [n]$ , and  $\pi = \text{Unif}(\mathcal{X})$ . Then  $\mathcal{P}_\pi$  is exactly the set of doubly stochastic PSD matrices. Let  $A$  be an  $n \times n$  PSD matrix with strictly positive entries that is not completely positive. By Sinkhorn's theorem ([MN61, Sin64]), there is a diagonal matrix  $D$  with strictly positive entries such that  $DAD$  is doubly stochastic. Let  $P = DAD$ . Clearly  $P$  is in  $\mathcal{P}_\pi$ . We claim that  $P$  is not in  $\mathcal{F}_\pi$ . Suppose for the sake of contradiction that  $P = KK_\pi^*$  for some  $K$ . Then

$$A = D^{-1}PD^{-1} = \frac{1}{n}D^{-1}K \text{Diag}(\pi K)^{-1}K^\top D^{-1} = MM^\top \quad (44)$$

where  $M = \frac{1}{\sqrt{n}}D^{-1}K \text{Diag}(\pi K)^{-1/2}$  is non-negative. This contradicts with the assumption that  $A$  is not completely positive.  $\square$

It remains an interesting open problem to characterize  $\mathcal{F}_\pi$ , even for uniform  $\pi$ .

### 3 Comparison between constants

In this section we compare constants in Table 1, showing that there is a superconstant separation between any two of them. Table 3 summarizes examples in this section.

Example	Description	Separation
Example 10	One-step chain	Log-Sobolev Constant $\rho(\pi, P)$ vs Half-Step Entropy Contraction $\alpha(\pi, K)$
Example 11	1-to- $k$ chain	Half-Step Entropy Contraction $\alpha(\pi, K)$ vs Full-Step Entropy Contraction $\delta(\pi, P)$
Example 13	Bernoulli-Laplace model	Half-Step Entropy Contraction $\alpha(\pi, K)$ vs Full-Step Entropy Contraction $\delta(\pi, P)$
Example 16	Three-state chain	One-Step Entropy Contraction $\delta(\pi, P)$ vs Two-Step Entropy Contraction $\delta(\pi, P^2)$
Example 18	Birth-death chain	$m$ -Step Entropy Contraction $\delta(\pi, P^m)$ vs $(m+1)$ -Step Entropy Contraction $\delta(\pi, P^{m+1})$
Example 20	Three-state chain	Full-Step Entropy Contraction $\delta(\pi, P)$ vs Modified Log-Sobolev Constant $\rho_0(\pi, P)$
Example 23	Expander graph	Modified Log-Sobolev Constant $\rho_0(\pi, P)$ vs Poincaré Constant $\lambda(\pi, P)$
Example 24	Random transposition model	Half-Step Entropy Contraction $\alpha(\pi, K)$ vs Modified Log-Sobolev Constant $\rho_0(\pi, P)$

Table 3: Examples in Section 3 and the separations they witness.

#### 3.1 Log-Sobolev constant $\rho$ vs half-step entropy contraction $\alpha$

By [Mic97], we always have  $\rho(\pi, P) \leq \alpha(\pi, K)$  for  $P = KK_\pi^*$ . The following example shows that the gap can be arbitrarily large.

**Example 10** (A one-step Markov chain). Let  $\mathcal{X}$  be a finite set,  $\pi$  be a distribution on  $\mathcal{X}$  with full support. Let  $\mathcal{Y} = \{*\}$  and  $K : \mathcal{X} \rightarrow \mathcal{Y}$  be the unique Markov kernel from  $\mathcal{X}$  to  $\mathcal{Y}$ . Then  $P = KK_\pi^*$  satisfies  $P(x, y) = \pi(y)$  for all  $x, y \in \mathcal{X}$ . By [DSC96, Theorem A.1],

$$\rho(\pi, P) = \frac{1 - 2\pi_*}{\log(1/\pi_* - 1)} \quad (45)$$

where  $\pi_* = \min_{x \in \mathcal{X}} \pi(x)$ . On the other hand,  $\alpha(\pi, K) = 1 - \eta_{\text{KL}}(\pi, K) = 1$ . As  $\pi_* \rightarrow 0$ , we have  $\frac{\alpha(\pi, K)}{\rho(\pi, P)} \rightarrow \infty$ .

### 3.2 Half-step entropy contraction $\alpha$ vs full-step entropy contraction $\delta$

By the data processing inequality, we always have  $\alpha(\pi, K) \leq \delta(\pi, P)$  for  $P = KK_\pi^*$ . Example 5 with  $\ell = 2$  already shows that the gap can be arbitrarily large. Below we present a few different examples.

**Example 11** (A 1-to- $k$  Markov chain). Let  $\mathcal{X} = [n]$ ,  $\mathcal{Y} = [n]^k$ ,  $K(x, y) = \frac{1}{kn^{k-1}} \sum_{j \in [k]} \mathbb{1}\{y_j = x\}$ ,  $\pi = \text{Unif}(\mathcal{X})$ . In other words, on input  $x \in \mathcal{X}$ , this chain generates a uniform length- $k$  output string and then randomly overwrite one of the  $k$  positions with  $x$ . Then  $K_\pi^* : \mathcal{Y} \rightarrow \mathcal{X}$  satisfies  $K_\pi^*(y, x) = \frac{1}{k} \sum_{j \in [k]} \mathbb{1}\{y_j = x\}$ . Let  $P = KK_\pi^*$ . The motivation for this chain comes from analysis of the Glauber dynamics on  $\text{Unif}(\mathcal{X}^k)$ . The Markov kernel  $K_\pi^*$  is the  $k$ -to-1 walk, and its entropy contraction is called entropic independence in [AJK<sup>+</sup>22]. Proposition 12 shows that for constant  $k \geq 2$ , as  $n \rightarrow \infty$ , we have  $\frac{\delta(\pi, P)}{\alpha(\pi, K)} \rightarrow \infty$ .

**Proposition 12** (A 1-to- $k$  Markov chain). Let  $\pi, K, P$  be as in Example 11. Then we have  $\alpha(\pi, K) = O\left(\frac{1}{\log n}\right)$  and  $\delta(\pi, P) \geq 1 - \frac{1}{k}$ . In particular, for fixed  $k \geq 2$ ,  $\frac{\delta(\pi, P)}{\alpha(\pi, K)} = \Omega(\log n)$ .

*Proof. Upper bound on  $\alpha(\pi, K)$ .* Let  $\nu$  be the point distribution at  $1 \in \mathcal{X}$ . Then  $D(\nu \| \pi) = \log n$ . Consider the distribution  $\nu K$ . For  $y \in \mathcal{Y}$ , if  $y$  contains  $i$  copies of 1, then  $(\nu K)(y) = \frac{i}{kn^{k-1}}$ . So

$$\begin{aligned} D(\nu K \| \pi K) &= \sum_{1 \leq i \leq k} \binom{k}{i} (n-1)^{k-i} \cdot \frac{i}{kn^{k-1}} \log \frac{ni}{k} \\ &= \log n - \sum_{1 \leq i \leq k} \binom{k-1}{i-1} \frac{(n-1)^{k-i}}{n^{k-1}} \log \frac{k}{i}. \end{aligned} \quad (46)$$

For constant  $k \geq 2$ , as  $n \rightarrow \infty$ , we have  $D(\nu K \| \pi K) = \log n - \Theta(1)$ . Therefore

$$\alpha(\pi, K) \leq 1 - \frac{D(\nu K \| \pi K)}{D(\nu \| \pi)} = \Theta\left(\frac{1}{\log n}\right). \quad (47)$$

**Lower bound on  $\delta(\pi, P)$ .** Let  $M : \mathcal{Y} \rightarrow [k] \times [n]$  be the Markov kernel defined as  $M(y, \cdot) = \text{Unif}(\{(i, y_i) : i \in [k]\})$ . By the data processing inequality,

$$1 - \delta(\pi, P) = \eta_{\text{KL}}(\pi, P) \leq \eta_{\text{KL}}(\pi K, K_\pi^*) \leq \eta_{\text{KL}}(\pi K, M). \quad (48)$$

Let  $\nu$  be any distribution on  $\mathcal{Y}$ , and  $\nu_i$  ( $i \in [k]$ ) be the  $i$ -th marginal of  $\nu$ . Then

$$D(\nu \| \pi K) = D(\nu \| \pi^{\times k}) = D(\nu \| \nu_1 \times \cdots \times \nu_k) + \sum_{i \in [k]} D(\nu_i \| \pi) \geq \sum_{i \in [k]} D(\nu_i \| \pi). \quad (49)$$

On the other hand,

$$D(\nu M \parallel \pi K M) = D(\nu M \parallel \text{Unif}([k]) \times \pi) = \frac{1}{k} \sum_{i \in [k]} D(\nu_i \parallel \pi). \quad (50)$$

Therefore

$$\eta_{\text{KL}}(\pi, M) \leq \frac{1}{k}. \quad (51)$$

So

$$\delta(\pi, P) \geq 1 - \frac{1}{k}. \quad (52)$$

□

**Example 13** (Bernoulli-Laplace model). Let  $n$  be a positive integer and  $1 \leq k \leq n - 1$ . We define a graph  $G = (V, E)$ . Let  $V = \binom{[n]}{k}$  (i.e., size- $k$  subsets of  $[n]$ ). Equivalently,  $V$  is the set of length- $n$  bit strings with Hamming weight  $k$ . There is an edge  $(x, y)$  for  $x, y \in V$  if and only if  $\|x - y\|_1 = 2$  (considered as elements in  $\{0, 1\}^n$ ). The Bernoulli-Laplace model is the lazy random walk on  $G$  (Definition 4). That is,  $\mathcal{X} = V$ ,  $\pi = \text{Unif}(\mathcal{X})$ ,  $\mathcal{Y} = E$ ,  $K : \mathcal{X} \rightarrow \mathcal{Y}$  is defined as  $K(x, e) = \frac{1}{k(n-k)} \mathbb{1}\{x \in e\}$ . For  $(ij) \in \binom{[n]}{2}$ , define map  $\sigma_{ij} : \mathcal{X} \rightarrow \mathcal{X}$  by swapping the  $i$ -th coordinate and the  $j$ -th coordinate (under the bit string interpretation). Then

$$P(x, \cdot) = \frac{1}{2} \mathbb{1}_x + \frac{1}{2k(n-k)} \sum_{\substack{i \in x \\ j \in [n] \setminus x}} \mathbb{1}_{\sigma_{ij}(x)}. \quad (53)$$

Proposition 14 shows that for constant  $k \geq 1$ , as  $n \rightarrow \infty$ , we have  $\frac{\delta(\pi, P)}{\alpha(\pi, K)} \rightarrow \infty$ .

**Proposition 14** (Bernoulli-Laplace model). *Let  $\pi, K, P$  be as in Example 13 with  $1 \leq k \leq n - 1$ . Then  $\alpha(\pi, K) = O\left(\frac{1}{\log \binom{n}{k}}\right)$  and  $\delta(\pi, P) \geq \frac{n}{2k(n-k)}$ . In particular, for constant  $k \geq 1$ , we have  $\frac{\delta(\pi, P)}{\alpha(\pi, K)} = \Omega(\log n)$ .*

*Proof. Upper bound on  $\alpha(\pi, K)$ .* Let  $\nu$  be the point distribution on any  $x \in \mathcal{X}$ . Then  $D(\nu \parallel \pi) = \log |\mathcal{X}| = \log \binom{n}{k}$ ,

$$\nu K = \frac{1}{k(n-k)} \sum_{\substack{i \in x \\ j \in [n] \setminus x}} \mathbb{1}_{\{x, \sigma_{ij}(x)\}}. \quad (54)$$

Because  $\pi K = \text{Unif}(\mathcal{Y})$ ,

$$D(\nu K \parallel \pi K) = \log |\mathcal{Y}| - \log(k(n-k)) = \log \binom{n}{k} - \log 2. \quad (55)$$

Therefore

$$\alpha(\pi, K) \leq 1 - \frac{D(\nu K \parallel \pi K)}{D(\nu \parallel \pi)} = \frac{\log 2}{\log \binom{n}{k}}. \quad (56)$$

**Lower bound on  $\delta(\pi, P)$ .** We make use of the following useful result from [CMS24].

**Lemma 15** ([CMS24, Theorem 1]). *Let  $d$  be a metric on  $\mathcal{X}$ . Let  $W_p$  denote the Wasserstein  $p$ -distance on  $\mathcal{P}(\mathcal{X})$  If*

$$W_\infty(P(x, \cdot), P(y, \cdot)) \leq d(x, y), \quad \forall x, y \in \mathcal{X}, \quad (57)$$

$$W_1(P(x, \cdot), P(y, \cdot)) \leq (1 - \kappa)d(x, y), \quad \forall x, y \in \mathcal{X}, \quad (58)$$

then

$$\delta(\pi, P) \geq \kappa. \quad (59)$$

For the Bernoulli-Laplace model, we let  $d$  be the graph distance on  $V = \binom{[n]}{k}$ . To prove Eqs. (57) and (58), it suffices to prove the result for adjacent  $x$  and  $y$ . By symmetry, WLOG assume that  $x = \{1, 3, 4, \dots, k+1\}$ ,  $y = \{2, 3, \dots, k+1\}$ . We define a coupling between  $P(x, \cdot)$  and  $P(y, \cdot)$  such that Eqs. (57) and (58) are both satisfied.

- (1) For  $3 \leq i \leq k+1$ ,  $k+2 \leq j \leq n$ , couple  $\sigma_{ij}(x)$  with  $\sigma_{ij}(y)$ . This happens with probability  $\frac{(k-1)(n-k-1)}{2k(n-k)}$  and incurs distance 1.
- (2) For  $k+2 \leq j \leq n$ , couple  $\sigma_{1j}(x)$  with  $\sigma_{2j}(y)$ . This happens with probability  $\frac{n-k-1}{2k(n-k)}$  and incurs distance 0.
- (3) For  $3 \leq i \leq k+1$ , couple  $\sigma_{2i}(x)$  with  $\sigma_{1i}(y)$ . This happens with probabilities  $\frac{k-1}{2k(n-k)}$  and incurs distance 0.
- (4) Couple  $\sigma_{12}(x)$  with  $y$ , and  $x$  with  $\sigma_{12}(y)$  each with weight  $\frac{1}{2k(n-k)}$ . This happens with probability  $\frac{1}{k(n-k)}$  and incurs distance 0.
- (5) At this point, all remaining mass in  $P(x, \cdot)$  (resp.  $P(y, \cdot)$ ) is at  $x$  (resp.  $y$ ). Couple them directly. This happens with probability  $\frac{1}{2} - \frac{1}{2k(n-k)}$  and incurs distance 1.

To summarize, the coupling has distance at most one and expected distance  $1 - \frac{n}{2k(n-k)}$ . By Lemma 15, we have

$$\delta(\pi, P) \geq \frac{n}{2k(n-k)}. \quad (60)$$

□

We remark that for the Bernoulli-Laplace model, the exact value of  $\alpha(\pi, K)$  has been determined in [BC24, Theorem 1.12], where it is shown that Eq. (56) is tight. For  $\delta(\pi, P)$ , by considering a point distribution at any  $x \in \mathcal{X}$ , we have

$$\delta(\pi, P) \leq \frac{\log(2n(n-k))}{2 \log \binom{n}{k}}. \quad (61)$$

Therefore, for  $n, k$  satisfying  $\log k = (1 - \Omega(1)) \log n$ , we have  $\delta(\pi, P) = \Theta\left(\frac{1}{k}\right)$ .

Examples 11 and 13 show that there is a separation between  $\alpha(\pi, K)$  and  $\delta(\pi, P = KK^*)$ . In these examples,  $K_\pi^*$  can be quite different from  $K$ . One natural question is whether there is a separation between  $\delta(\pi, P)$  and  $\delta(\pi, P^2)$ . That is, can running the same Markov kernel twice result in much better contraction than running only once? The following example shows that indeed such a separation exists. This example is adapted from [Mün23]'s counterexample to the Peres-Tetali conjecture. We note, however, that the properties of this chain we use here and in Example 20 are different from those used op. cit.

**Example 16** (A three-state Markov chain). Let  $M$  be a positive real number. Let  $\mathcal{X} = [3]$ ,  $\pi = \left(\frac{M}{M+2}, \frac{1}{M+2}, \frac{1}{M+2}\right)$ . Let  $P : \mathcal{X} \rightarrow \mathcal{X}$  be defined as

$$P = \begin{pmatrix} 1 - \frac{1}{4M} & \frac{1}{4M} & 0 \\ \frac{1}{4} & \frac{1}{2} & \frac{1}{4} \\ 0 & \frac{1}{4} & \frac{3}{4} \end{pmatrix}. \quad (62)$$

It is easy to see that  $\pi$  is the invariant distribution of  $P$  and  $(\pi, P)$  is reversible. Note that by [Lemma 3](#),  $P$  is factorizable. By [Proposition 17](#),  $\frac{\delta(\pi, P^2)}{\delta(\pi, P)} \rightarrow \infty$  as  $M \rightarrow \infty$ .

**Proposition 17** (A three-state Markov chain). *Let  $\pi, P$  be as in [Example 16](#). Then  $\delta(\pi, P) = O\left(\frac{1}{\log M}\right)$  and  $\delta(\pi, P^2) = \Omega(1)$ . In particular,  $\frac{\delta(\pi, P^2)}{\delta(\pi, P)} = \Omega(\log M)$ .*

*Proof. Upper bound on  $\delta(\pi, P)$ .* Let  $\nu$  be the point distribution at  $3 \in \mathcal{X}$ . Then  $D(\nu||\pi) = \log(M+2)$ ,

$$D(\nu P||\pi) = \frac{1}{4} \log \frac{M+2}{4} + \frac{3}{4} \log \frac{3(M+2)}{4} = \log(M+2) - h(1/4) \quad (63)$$

where  $h(x) : [0, 1] \rightarrow [0, \log 2]$  is the binary entropy function

$$h(x) = -x \log x - (1-x) \log(1-x). \quad (64)$$

Therefore

$$\delta(\pi, P) \leq 1 - \frac{D(\nu P||\pi)}{D(\nu||\pi)} = O\left(\frac{1}{\log M}\right). \quad (65)$$

**Lower bound on  $\delta(\pi, P^2)$ .** We prove that for large  $M$ , we have  $P^2(x, 1) = \Omega(1)$  for all  $x \in \mathcal{X}$ . In fact,  $P^2(1, 1) \geq P(1, 1)^2 = \Omega(1)$ ,  $P^2(2, 1) \geq P(2, 1)P(1, 1) = \Omega(1)$ , and  $P^2(3, 1) \geq P(3, 2)P(2, 1) = \Omega(1)$ . So  $\text{TV}(P^2(x, \cdot), P^2(x', \cdot)) = 1 - \Omega(1)$  for all  $x, x' \in \mathcal{X}$ . By [Eq. \(16\)](#), the Dobrushin's coefficient satisfies  $\eta_{\text{TV}}(\pi, P^2) = 1 - \Omega(1)$ . By [Eq. \(15\)](#), we have

$$\delta(\pi, P^2) \geq 1 - \eta_{\text{TV}}(\pi, P^2) = \Omega(1). \quad (66)$$

□

We generalize [Example 16](#) as follows, showing that for any positive integer  $m$ , there exists a Markov kernel such that running  $(m+1)$  steps results in entropy contraction much better than running  $m$  steps. We note that there is a relatively simple characterization of the LSI for birth-death chains ([\[Che05, Che03\]](#)), but for MLSI or SDPI no such characterizations are known, except for partial progress in [\[Rob01, CDPP09\]](#).

**Example 18** (A birth-death Markov chain). We fix a positive integer  $m$  and let  $M$  be a large positive real number. Let  $\mathcal{X} = [m+2]$ ,  $\pi(x) = \frac{1}{M+m+1} + \mathbb{1}\{x=1\} \frac{M-1}{M+m+1}$ . Let  $P : \mathcal{X} \rightarrow \mathcal{X}$  be a birth-death Markov chain, where  $P(x, y) = 0$  for  $|x-y| \geq 2$ ,  $P(x, x-1) = \frac{1}{4}$  for  $2 \leq x \leq m+2$ ,  $P(x, x+1) = \frac{1}{4}$  for  $2 \leq x \leq m+1$ ,  $P(1, 2) = \frac{1}{4M}$ , and  $P(x, x) = 1 - P(x, x-1) - P(x, x+1)$ . It is easy to verify that  $(\pi, P)$  is a reversible pair and  $P(x, x) \geq \frac{1}{2}$  for all  $x \in [m+2]$ . By [Lemma 3](#),  $(\pi, P)$  is factorizable. [Proposition 19](#) shows that as  $M \rightarrow \infty$ , we have  $\frac{\delta(\pi, P^{m+1})}{\delta(\pi, P^m)} \rightarrow \infty$ .



**Proposition 19** (A birth-death Markov chain). *Let  $\pi, P$  be as in [Example 18](#). Then  $\delta(\pi, P^m) = O\left(\frac{1}{\log M}\right)$  and  $\delta(\pi, P^{m+1}) = \Omega(1)$ . In particular,  $\frac{\delta(\pi, P^{m+1})}{\delta(\pi, P^m)} = \Omega(\log M)$ .*

*Proof. Upper bound on  $\delta(\pi, P^m)$ .* Let  $\nu$  be the point distribution at  $m+2 \in \mathcal{X}$ . Then  $D(\nu||\pi) = \log(M+m+1)$ . Note that  $(\nu P^m)(1) = 0$ ,  $(\nu P^m)(x) = c_x$  for some  $c_x = \Theta_m(1)$  for  $2 \leq x \leq m+2$ , where  $\Theta_m$  hides a constant factor depending on  $m$ . Furthermore,  $(c_2, \dots, c_{m+2})$  is a distribution on  $\{2, \dots, m+2\}$ . Then

$$D(\nu P||\pi) = \log(M+m+1) - H(c_2, \dots, c_{m+2}) \quad (67)$$

where  $H$  is the entropy function

$$H(c_2, \dots, c_{m+2}) = - \sum_{2 \leq i \leq m+2} c_i \log c_i. \quad (68)$$

Because  $c_i = \Theta_m(1)$  for all  $2 \leq i \leq m+2$ , we have  $H(c_2, \dots, c_{m+2}) = \Theta_m(1)$ . Therefore

$$\delta(\pi, P^m) \leq 1 - \frac{D(\nu P||\pi)}{D(\nu||\pi)} = \Theta_m\left(\frac{1}{\log M}\right). \quad (69)$$

**Lower bound on  $\delta(\pi, P^{m+1})$ .** Note that for any  $x \in [m+2]$ , we have

$$P^{m+1}(x, 1) \geq P(x, x-1) \cdots P(2, 1) \cdot P(1, 1)^{m+1-x} = \Omega_m(1) \quad (70)$$

where  $\Omega_m$  hides a constant factor depending on  $m$ . By [Eqs. \(15\) and \(16\)](#),

$$\delta(\pi, P^{m+1}) \geq 1 - \eta_{\text{TV}}(\pi, P^{m+1}) = \Omega_m(1). \quad (71)$$

□

If a Markov kernel  $(\pi, P)$  separates  $\delta(\pi, P^m)$  and  $\delta(\pi, P^{m+1})$ , then it also separates  $\delta(\pi, T_t)$  (where  $T_t = e^{t(P-\text{Id})}$ ) and  $\delta(\pi, P^m)$  for finite  $t$ . In other words, the continuous-time chain contracts entropy at finite time better than the discrete-time counterpart. To see this, we note that for any function  $f$  on  $\mathcal{X}$ , we have

$$\begin{aligned} \text{Ent}_\pi(T_t f) &\leq \mathbb{E}_{n \sim \text{Pois}(t)} \text{Ent}_\pi(P^n f) \\ &\leq \mathbb{E}_{n \sim \text{Pois}(t)} [\mathbb{1}\{n \leq m\} \text{Ent}_\pi(f) + \mathbb{1}\{n \geq m+1\} \text{Ent}_\pi(P^{m+1} f)] \\ &\leq \mathbb{P}[\text{Pois}(t) \leq m] \text{Ent}_\pi(f) + \mathbb{P}[\text{Pois}(t) \geq m+1] (1 - \delta(\pi, P^{m+1})) \\ &\leq \text{Ent}_\pi(f) (1 - \mathbb{P}[\text{Pois}(t) \geq m+1] \delta(\pi, P^{m+1})), \end{aligned} \quad (72)$$

where the first step is by convexity of  $\text{Ent}_\pi$ , and the second step is by the data processing inequality. Therefore

$$\delta(\pi, T_t) \geq \mathbb{P}[\text{Pois}(t) \geq m+1] \delta(\pi, P^{m+1}), \quad (73)$$

which is separated from  $\delta(\pi, P^m)$  for finite  $t$ , assuming that  $\delta(\pi, P^{m+1})$  is separated from  $\delta(\pi, P^m)$ .

### 3.3 Full-step entropy contraction $\delta$ vs modified log-Sobolev constant $\rho_0$

By [BCP<sup>+</sup>21, Lemma 2.7], we always have  $\delta(\pi, P) \leq \rho_0(\pi, P)$  for any reversible  $(\pi, P)$ .<sup>3</sup> If we allow non-factorizable  $(\pi, P)$ , then it is easy to give an example where the gap is infinite: take  $\mathcal{X} = [2]$ ,  $\pi = \text{Unif}(\mathcal{X})$ , and  $P = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$ . The following example shows that the gap can be arbitrarily large even if we restrict to factorizable  $(\pi, P)$ .

**Example 20** (A three-state Markov chain). This example is the same as [Example 16](#). [Proposition 21](#) shows that we have  $\frac{\rho_0(\pi, P)}{\delta(\pi, P)} \rightarrow \infty$  as  $M \rightarrow \infty$ .

**Proposition 21** (A three-state Markov chain). *Let  $(\pi, P)$  be as in [Example 16](#). Then  $\delta(\pi, P) = O\left(\frac{1}{\log M}\right)$  and  $\rho_0(\pi, P) = \Theta\left(\frac{\log \log M}{\log M}\right)$ . In particular,  $\frac{\rho_0(\pi, P)}{\delta(\pi, P)} = \Omega(\log \log M)$ .*

*Proof. Upper bound on  $\delta(\pi, P)$ .* We have established in [Proposition 17](#) that  $\delta(\pi, P) = O\left(\frac{1}{\log M}\right)$ .

*Lower bound on  $\rho_0(\pi, P)$ .* We will show that there exists a sufficiently large universal constant  $C > 0$  such that, for any positive function  $f : [3] \rightarrow \mathbb{R}^+$ , it holds

$$\text{Ent}_\pi(f) \leq \frac{C \log M}{\log \log M} \mathcal{E}_{\pi, P}(f, \log f). \quad (74)$$

This implies  $\rho_0(\pi, P) = \Omega\left(\frac{\log \log M}{\log M}\right)$ . By applying suitable scaling, we assume that  $f(1) = x$ ,  $f(2) = 1$ , and  $f(3) = y$  for some  $x, y > 0$ . Furthermore, we may assume without loss of generality that  $x \geq \frac{\beta}{M}$  where  $\beta > 0$  is some tiny absolute constant, due to [[TY24](#), Lemma 2.1]. More specifically, [[TY24](#), Lemma 2.1] shows that to establish a modified log-Sobolev inequality [Eq. \(74\)](#), it suffices to consider a restricted class of functions such that, among other restrictions,  $f(1) \geq \beta' \mathbb{E}_\pi f$  for an absolute constant  $\beta' > 0$ ; such a restriction immediately implies

$$x = f(1) \geq \beta' \left( \frac{Mx}{M+2} + \frac{1}{M+2} + \frac{y}{M+2} \right) \geq \frac{\beta'}{M+2} \geq \frac{\beta'}{3M}. \quad (75)$$

Hence, we can safely assume  $x \geq \frac{\beta}{M}$  where  $\beta = \beta'/3$ .

A direct calculation yields

$$(M+2) \cdot \text{Ent}_\pi(f) = Mx \log x + y \log y - (Mx + y + 1) \log \left( \frac{Mx + y + 1}{M+2} \right), \quad (76)$$

and

$$4(M+2) \cdot \mathcal{E}_{\pi, P}(f, \log f) = (x-1) \log x + (y-1) \log y. \quad (77)$$

The following claim is helpful.

**Claim 22.** *We have*

$$\frac{\text{Ent}_\pi(f)}{4 \mathcal{E}_{\pi, P}(f, \log f)} \leq \frac{(2x - y - 1) + y \log y + (y + 1) \log \left(\frac{1}{x}\right)}{(x - 1) \log x + (y - 1) \log y}. \quad (78)$$

---

<sup>3</sup>They in fact prove a more general statement that works for non-reversible  $(\pi, P)$ .

*Proof.* We rewrite the entropy as

$$\begin{aligned} (M+2) \cdot \text{Ent}_\pi(f) &= Mx \log x + y \log y - (Mx + y + 1) \log \left( \frac{Mx + y + 1}{M+2} \right) \\ &= (Mx + y + 1) \log \left( \frac{(M+2)x}{Mx + y + 1} \right) + y \log y + (y+1) \log \left( \frac{1}{x} \right). \end{aligned} \quad (79)$$

Notice that the first term above can be controlled by

$$\begin{aligned} (Mx + y + 1) \log \left( \frac{(M+2)x}{Mx + y + 1} \right) &= (Mx + y + 1) \log \left( 1 + \frac{2x - y - 1}{Mx + y + 1} \right) \\ &\leq (Mx + y + 1) \cdot \frac{2x - y - 1}{Mx + y + 1} = 2x - y - 1. \end{aligned} \quad (80)$$

The claim then follows.  $\square$

We consider two separate cases of  $(x, y)$  to establish [Eq. \(74\)](#).

**Case 1:**  $(x, y) \notin (\frac{1}{2}, \frac{3}{2}) \times (\frac{1}{2}, \frac{3}{2})$ . In this case, we have

$$(x-1) \log x + (y-1) \log y \geq \frac{1}{10}. \quad (81)$$

Since we have

$$2x - y - 1 \leq 2x - 1 \leq 2(x-1) \log x + 3 \leq 32((x-1) \log x + (y-1) \log y) \quad (82)$$

and also

$$y \log y \leq 2(y-1) \log y + 1 \leq 12((x-1) \log x + (y-1) \log y), \quad (83)$$

we deduce from [Claim 22](#) that

$$\frac{\text{Ent}_\pi(f)}{4 \mathcal{E}_{\pi, P}(f, \log f)} \leq 44 + \frac{(y+1) \log \left( \frac{1}{x} \right)}{(x-1) \log x + (y-1) \log y}. \quad (84)$$

Consider three subcases.

(i) If  $x \geq \frac{3}{2}$ , then  $\log(1/x) < 0$  and hence

$$\frac{\text{Ent}_\pi(f)}{4 \mathcal{E}_{\pi, P}(f, \log f)} \leq 44. \quad (85)$$

(ii) If  $x \leq \frac{1}{2}$ , then consider how large  $y$  is. If  $y \leq \frac{\log M}{\log \log M}$ , then by  $(x-1) \log x \geq \frac{1}{2} \log(1/x)$  and  $(y-1) \log y \geq 0$  we deduce that

$$\frac{(y+1) \log \left( \frac{1}{x} \right)}{(x-1) \log x + (y-1) \log y} \leq 2(y+1) \leq \frac{4 \log M}{\log \log M}. \quad (86)$$

If  $y > \frac{\log M}{\log \log M}$ , then by  $(x-1) \log x \geq 0$ ,  $\frac{y+1}{y-1} \leq 3$ , and  $x \geq \frac{\beta}{M}$  we deduce that

$$\frac{(y+1) \log \left( \frac{1}{x} \right)}{(x-1) \log x + (y-1) \log y} \leq \frac{(y+1) \log \left( \frac{1}{x} \right)}{(y-1) \log y} \leq \frac{3 \log \left( \frac{M}{\beta} \right)}{\log \left( \frac{\log M}{\log \log M} \right)} \leq \frac{C_0 \log M}{\log \log M}, \quad (87)$$

for some  $C_0 = C_0(\beta) > 0$  when  $M$  is sufficiently large.

(iii) If  $x \in (\frac{1}{2}, \frac{3}{2})$ , then by our assumption it must hold  $y \notin (\frac{1}{2}, \frac{3}{2})$ . Since  $\log(1/x) \leq 1$  when  $x > \frac{1}{2}$ , we have

$$\frac{(y+1) \log(\frac{1}{x})}{(x-1) \log x + (y-1) \log y} \leq \frac{y+1}{(y-1) \log y} \leq 10, \quad (88)$$

when  $y \notin (\frac{1}{2}, \frac{3}{2})$ .

Therefore, in all three subcases we have

$$\frac{\text{Ent}_\pi(f)}{4 \mathcal{E}_{\pi, P}(f, \log f)} \leq \frac{C \log M}{\log \log M} \quad (89)$$

where  $C = C(\beta) > 0$  is constant, whenever  $M$  is sufficiently large.

**Case 2:**  $(x, y) \in (\frac{1}{2}, \frac{3}{2}) \times (\frac{1}{2}, \frac{3}{2})$ . In this case, we have

$$(x-1) \log x + (y-1) \log y \geq \frac{1}{2}(x-1)^2 + \frac{1}{2}(y-1)^2, \quad (90)$$

and also

$$\begin{aligned} & (2x - y - 1) + y \log y + (y+1) \log\left(\frac{1}{x}\right) \\ & \leq 2(x-1) - (y-1) + y(y-1) + (y+1) \left(\frac{1}{x} - 1\right) \\ & = \frac{1}{x}(x-1)(2x-y-1) + (y-1)^2 \\ & = \frac{2}{x}(x-1)^2 - \frac{1}{x}(x-1)(y-1) + (y-1)^2 \\ & \leq 4(x-1)^2 + 2|(x-1)(y-1)| + (y-1)^2 \\ & \leq 5(x-1)^2 + 5(y-1)^2. \end{aligned} \quad (91)$$

By [Claim 22](#),

$$\frac{\text{Ent}_\pi(f)}{4 \mathcal{E}_{\pi, P}(f, \log f)} \leq 10. \quad (92)$$

Combining the two cases, we conclude that

$$\rho_0(\pi, P) = \Omega\left(\frac{\log \log M}{\log M}\right). \quad (93)$$

In fact, this lower bound on  $\rho_0(\pi, P)$  is asymptotically tight and can be achieved by, for example,  $f(1) = 1/M$ ,  $f(2) = 1$ , and  $f(3) = \log M$  as given in [\[Mün23\]](#). Therefore,

$$\rho_0(\pi, P) = \Theta\left(\frac{\log \log M}{\log M}\right) \quad (94)$$

as  $M \rightarrow \infty$ .

□

### 3.4 Modified log-Sobolev constant $\rho_0$ vs Poincaré constant $\lambda$

[BT06] shows that  $\rho_0(\pi, P) \leq 2\lambda(\pi, P)$ , and gave an example showing that the gap can be arbitrarily large. We record their example here.

**Example 23** (Expander graphs). Let  $G = (V, E)$  be an expander graph with bounded degree. Consider the lazy random walk on  $G$  (Definition 4). [BT06] shows that  $\rho_0(\pi, P) = \Theta\left(\frac{1}{\log|V|}\right) = \Theta\left(\frac{\lambda_1(\pi, P)}{\log|V|}\right)$ . Therefore as  $|V| \rightarrow \infty$ , we have  $\frac{\lambda_1(\pi, P)}{\rho_0(\pi, P)} \rightarrow \infty$ .

### 3.5 Other comparisons

[DSC96] shows that the spectral gap  $\lambda$  and the log-Sobolev constant  $\rho$  differ by at most a factor of  $O(\log(1/\pi_{\min}))$  where

$$\pi_{\min} = \min_{x \in \mathcal{X}: \pi(x) > 0} \pi(x). \quad (95)$$

More precisely, [DSC96, Corollary A.4] shows that, assuming  $\pi_{\min} \leq 1/2$ , it holds

$$\frac{\lambda}{2 + \log(1/\pi_{\min})} \leq \frac{(1 - 2\pi_{\min})\lambda}{\log(1/\pi_{\min} - 1)} \leq \rho \leq \frac{\lambda}{2}. \quad (96)$$

This in particular shows that all constants discussed in this paper, including also the half-step entropy contraction  $\alpha$ , the full-step entropy contraction  $\delta$ , and the modified log-Sobolev constant  $\rho_0$ , differ by at most a factor of  $O(\log(1/\pi_{\min}))$  from each other.

In a recent work [STY23], Salez, Tikhomirov, and Youssef establish a surprising and remarkable comparison between the modified log-Sobolev constant  $\rho_0$  and the log-Sobolev constant  $\rho$ . For a reversible Markov kernel  $P$  with respect to a probability measure  $\pi$ , define the sparsity parameter as

$$p_{\min} = \min_{(x,y) \in \mathcal{X}^2: P(x,y) > 0} P(x,y). \quad (97)$$

Then, [STY23, Theorem 1] shows that

$$\frac{\rho_0}{20 \log(1/p_{\min})} \leq \rho \leq \frac{\rho_0}{4}. \quad (98)$$

Hence, all entropy-related constants discussed in this paper, including also the half-step entropy contraction  $\alpha$  and the full-step entropy contraction  $\delta$ , differ by at most a factor of  $O(\log(1/p_{\min}))$  from each other.

### 3.6 Comments on several previous works

[DMLM03, Prop. 5.1] claims that

$$c\rho_0 \leq \alpha \quad (99)$$

for some universal constant  $c > 0$ . In fact, Eq. (99) fails for the random transposition model, showing that the claim must be incorrect.

**Example 24** (Random transposition model). Let  $n$  be a positive integer. Let  $G = (V, E)$  be the Cayley graph on the symmetric group  $S_n$  generated by transpositions. That is,  $V$  is the set of permutations of  $[n]$ , and there is an edge between  $\sigma, \tau \in V$  if and only if they differ by exactly two entries. The random transposition model is the lazy random walk on  $G$  (Definition 4). Considering the point measure at any  $x \in \mathcal{X}$  gives  $\alpha \leq \frac{\log 2}{\log(n!)} = \Theta\left(\frac{1}{n \log n}\right)$ . In fact, [BC24, Theorem 1.9] shows that this is tight, i.e.,  $\alpha = \frac{\log 2}{\log(n!)}$ . On the other hand, [GQ03] shows that  $\rho_0 = \Theta\left(\frac{1}{n}\right)$ . Therefore, as  $n \rightarrow \infty$ , we have  $\frac{\rho_0}{\alpha} \rightarrow \infty$ .

Our separations of  $\alpha$  vs  $\delta$  (Section 3.2) and  $\delta$  vs  $\rho_0$  (Section 3.3) also provide alternative counterexamples to the claim. We now explain briefly the issue in the proof of Eq. (99) in [DMLM03]. In their proof of Prop. 5.1, the authors apply a technical result, Lemma 5.2, which represents the (relative) entropy of a function  $f$  with expectation  $\mathbb{E}_n(f) = 1$  (where  $n$  denotes the underlying probability measure) as an integral of the covariance between  $f_t$  and  $\log(f_t)$  where  $f_t = e^{-t}f + (1 - e^{-t})$ ,  $t \in \mathbb{R}_{\geq 0}$  represents an interpolation between  $f$  and 1. However, in the actual application of Lemma 5.2, the measure  $n$  is a conditional probability measure under which the expectation of  $f$  is no longer 1, and hence the interpolation function  $f_t$  should be replaced by  $f_t = e^{-t}f + (1 - e^{-t})\mathbb{E}_n(f)$ ; this would require the function  $f_t$  to depend on the conditioning and the proofs following afterwards no longer work.

[Rag16, Prop. 4.3] claims that

$$\rho_0 \leq 1 - c(1 - \alpha) \tag{100}$$

for some universal constant  $c > 0$ . Our examples do not disprove the claim. However, the proof of [Rag16, Theorem 4.4] is a generalization of that of [DMLM03, Prop. 5.1], so it has the same error. In particular, the last display of the proof of [Rag16, Prop. 4.3] implies that  $c\rho_0 \leq \alpha$  for some constant  $c > 0$ , which we have shown to be incorrect. Therefore the proof of [Rag16, Prop. 4.3] is incorrect and does not establish Eq. (100). It is unclear whether Eq. (100) as stated is correct.

## 4 Extremal functions

In this section we discuss another difference between the continuous-time entropy contraction constants and the discrete-time entropy contraction constants. It is known ([BT06]) that for irreducible  $(\pi, P)$ , the log-Sobolev constant  $\rho(\pi, P)$  and the modified log-Sobolev constant  $\rho_0(\pi, P)$  satisfy a dichotomy: they are either equal to twice the Poincaré constant  $\lambda(\pi, P)$  or achieved at a full-support function. We show that this is no longer true for the discrete-time entropy contraction constants  $\alpha(\pi, K)$  and  $\delta(\pi, P)$  by providing explicit examples whose extremal functions have non-full support.

### 4.1 Log-Sobolev constant $\rho$ , modified log-Sobolev constant $\rho_0$ , Poincaré constant $\lambda$

[BT06] studies extremal functions for  $\rho$ ,  $\rho_0$ ,  $\lambda$ . The extremal functions for the Poincaré constant  $\lambda(\pi, P)$  are easy to describe. They are the (right) eigenfunctions of  $-L$  corresponding to the eigenvalue  $\lambda$ . In particular,  $\lambda$  is always achieved at some non-constant function  $f : \mathcal{X} \rightarrow \mathbb{R}$ .

For the log-Sobolev constant  $\rho$ , [BT06] shows that for any reversible  $(\pi, P)$ , either

- (i)  $\rho(\pi, P) = 2\lambda(\pi, P)$ , or
- (ii) there exists a non-constant function  $f : \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$  with  $\pi[f] = 1$  such that

$$\rho \text{Ent}_\pi(f) = \mathcal{E}_{\pi, P}(\sqrt{f}, \sqrt{f}). \tag{101}$$

Furthermore, such  $f$  satisfies the equation

$$-L\sqrt{f} = \rho\sqrt{f}\log f. \quad (102)$$

If  $(\pi, P)$  is irreducible, then  $f$  has full support.

For the modified log-Sobolev constant  $\rho_0(\pi, P)$ , [BT06] shows that for any reversible  $(\pi, P)$ , either

- (i)  $\rho_0(\pi, P) = 2\lambda(\pi, P)$ , or
- (ii) there exists a non-constant function  $f : \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$  with  $\pi[f] = 1$  such that

$$\rho_0 \text{Ent}_{\pi}(f) = \mathcal{E}_{\pi, P}(f, \log f). \quad (103)$$

Furthermore, such  $f$  satisfies the equation

$$-Lf - fL(\log f) = \rho_0 f \log f. \quad (104)$$

If  $(\pi, P)$  is irreducible, then  $f$  has full support.

## 4.2 Half-step entropy contraction $\alpha$ and full-step entropy contraction $\delta$

In this section we study the extremal distributions for  $\eta_{\text{KL}}$ , which includes both  $\alpha$  and  $\delta$ .

**Lemma 25** (Extremal distributions for  $\eta_{\text{KL}}$ ). *Let  $\pi$  be a distribution on  $\mathcal{X}$  and  $K : \mathcal{X} \rightarrow \mathcal{Y}$  be a Markov kernel. Then either*

- (i)  $\eta_{\text{KL}}(\pi, K) = \eta_{\chi^2}(\pi, K)$ , or
- (ii) there exists distribution  $\nu$  on  $\mathcal{X}$  such that

$$\eta_{\text{KL}}(\pi, K) = \frac{D(\nu K \| \pi K)}{D(\nu \| \pi)}. \quad (105)$$

*Proof.* WLOG assume that  $\pi$  has full support. If **Item (ii)** does not happen, then there exists a sequence  $\{f_n\}_n$  ( $f_n : \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$ ,  $\pi[f_n] = 1$ ) satisfying  $\|f_n - \mathbb{1}\|_{\infty} \rightarrow 0$  and  $\frac{\text{Ent}_{\pi K}(K_{\pi}^* f_n)}{\text{Ent}_{\pi}(f_n)} \rightarrow \eta_{\text{KL}}(\pi, K)$  as  $n \rightarrow \infty$ .

Write  $f_n = \mathbb{1} + \epsilon_n g_n$ , where  $\epsilon_n \geq 0$ ,  $\pi[g_n^2] = 1$ . Note that the space  $\mathcal{G} := \{g : \mathcal{X} \rightarrow \mathbb{R} : \pi[g] = 0, \pi[g^2] = 1\}$  is compact. By replacing the sequence  $\{f_n\}_n$  with a subsequence, we can WLOG assume that there exists  $g^* \in \mathcal{G}$  such that  $\|g_n - g^*\|_{\infty} \rightarrow 0$  as  $n \rightarrow \infty$ .

Now let us prove that

$$\lim_{n \rightarrow \infty} \frac{\text{Ent}_{\pi K}(K_{\pi}^* f_n)}{\text{Ent}_{\pi}(f_n)} = \frac{\text{Var}_{\pi K}(K_{\pi}^* g^*)}{\text{Var}_{\pi}(g^*)}. \quad (106)$$

If **Eq. (106)** holds, then

$$\eta_{\chi^2}(\pi, K) = \sup_{g \in \mathcal{G}} \frac{\text{Var}_{\pi K}(K_{\pi}^* g)}{\text{Var}_{\pi}(g)} \geq \eta_{\text{KL}}(\pi, K), \quad (107)$$

and **Item (i)** holds.

By Lemma 26,

$$\begin{aligned} & \|(1 + \epsilon_n g_n) \log(1 + \epsilon_n g_n) - (1 + \epsilon_n g^*) \log(1 + \epsilon_n g^*) - \epsilon_n (g_n - g^*)\|_\infty \\ &= O(\epsilon_n^2 (\epsilon_n + \|g_n - g^*\|_\infty)). \end{aligned} \quad (108)$$

Taking expectation and using triangle inequality, we get

$$|\text{Ent}_\pi(f_n) - \text{Ent}_\pi(1 + \epsilon_n g^*)| = O(\epsilon_n^2 (\epsilon_n + \|g_n - g^*\|_\infty)). \quad (109)$$

It is known that

$$\lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon^2} \text{Ent}_\pi(1 + \epsilon g^*) = \frac{1}{2} \text{Var}_\pi(g^*). \quad (110)$$

So

$$\lim_{n \rightarrow \infty} \frac{1}{\epsilon_n^2} \text{Ent}_\pi(f_n) = \frac{1}{2} \text{Var}_\pi(g^*). \quad (111)$$

Similarly

$$\lim_{n \rightarrow \infty} \frac{1}{\epsilon_n^2} \text{Ent}_{\pi K}(K_\pi^* f_n) = \frac{1}{2} \text{Var}_{\pi K}(K_\pi^* g^*). \quad (112)$$

This finishes the proof of Eq. (106).  $\square$

**Lemma 26.** *There exists  $C > 0$  such that for  $\epsilon, \epsilon' > 0$  small enough, for  $x, y, z, w \in \mathbb{R}$  satisfying  $|x|, |y|, |z|, |w| \leq \epsilon$ ,  $|x - z|, |y - w| \leq \epsilon\epsilon'$ , we have*

$$|(1 + x) \log(1 + y) - (1 + z) \log(1 + w) - (y - w)| \leq C\epsilon^2(\epsilon + \epsilon'). \quad (113)$$

*Proof.* First note that  $|\log(1 + y) - (y - y^2/2)| = O(\epsilon^3)$ . Then

$$\begin{aligned} & |(1 + x) \log(1 + y) - (1 + z) \log(1 + w) - (y - w)| \\ &= |(1 + x)(y - y^2/2) - (1 + z)(w - w^2/2) - (y - w)| + O(\epsilon^3) \\ &= |xy - zw| + |y^2 - w^2|/2 + |xy^2|/2 + |zw^2|/2 + O(\epsilon^3) \\ &= |(x - z)y + z(y - w)| + |(y - w)(y + w)|/2 + O(\epsilon^2(\epsilon + \epsilon')) \\ &= O(\epsilon^2(\epsilon + \epsilon')). \end{aligned} \quad (114)$$

$\square$

Unlike  $\rho$  and  $\rho_0$  where the extremal functions (if they exist) have full support, the extremal distributions for  $\eta_{\text{KL}}$  may have non-full support.

**Example 27** (Complete graph). Let  $n \geq 3$  be an integer. Let  $\mathcal{X} = [n]$  and  $\pi = \text{Unif}(\mathcal{X})$ . Let  $K : \mathcal{X} \rightarrow \mathcal{X}$  be the (non-lazy) random walk on the complete graph. That is,

$$K(x, y) = \frac{1}{n-1} \mathbb{1}\{x \neq y\} \quad (115)$$

for  $x, y \in \mathcal{X}$ . [GP23, Prop. 33] proves that  $\eta_{\text{KL}}(\pi, K) = \frac{\log n - \log(n-1)}{\log n}$ , and is achieved at and only at point distributions.



In the above example,  $K$  is not factorizable, so it corresponds to the half-step contraction coefficient  $\alpha$ . The next example shows the extremal distributions for the full-step contraction coefficient  $\delta$  can have non-full support.

**Example 28** (Complete bipartite graph). Let  $n \geq 3$  be an integer. Let  $G = K_{n,n}$  be the complete bipartite graph. That is,  $V = [2] \times [n]$ , and there is an edge between  $(1, i)$  and  $(2, j)$  for all  $i, j \in [n]$ . Let  $(\pi, P)$  be the random walk on  $G$  (Definition 4). Numerical computation suggests that  $\eta_{\text{KL}}(\pi, P) = \frac{\log n}{2 \log(2n)}$ , and equality is achieved at and only at point distributions. Proposition 29 proves this observation for  $n = 3$ .

**Proposition 29** (Complete bipartite graph). *Let  $\pi, P$  be as in Example 28. For  $n = 3$ ,  $\eta_{\text{KL}}(\pi, P) = \frac{\log n}{2 \log(2n)}$ , and equality is achieved at and only at point distributions.*

*Proof.* Let  $\nu$  be the point distribution on any  $x \in \mathcal{X}$ . Then

$$D(\nu \parallel \pi) = \log(2n), \quad (116)$$

$$D(\nu P \parallel \pi) = \frac{1}{2} \log n. \quad (117)$$

So

$$\eta_{\text{KL}}(\pi, P) \geq \frac{D(\nu P \parallel \pi)}{D(\nu \parallel \pi)} = \frac{\log n}{2 \log(2n)}. \quad (118)$$

Note that this holds for any  $n \geq 3$ .

We prove that for  $n = 3$ , the point distributions are the only maximizers. We represent a distribution  $\nu$  using a tuple  $(t, \nu_1, \nu_2)$ , where  $t \in [0, 1]$  and  $\nu_i$  ( $i = 1, 2$ ) is a distribution on  $\mathcal{X}_i = \{i\} \times [n] \subseteq \mathcal{X}$ . Given  $\nu$ , the corresponding tuple  $\phi(\nu) = (\nu(\mathcal{X}_1), \nu|_{\mathcal{X}_1}, \nu|_{\mathcal{X}_2})$ , where  $\nu|_{\mathcal{X}_i}$  denotes the conditional distribution (if  $\nu(\mathcal{X}_i) = 0$ , then choose an arbitrary distribution on  $\mathcal{X}_i$ ). Given a tuple  $(t, \nu_1, \nu_2)$ , the corresponding distribution  $\nu$  is  $\psi(t, \nu_1, \nu_2) = t\nu_1 + (1-t)\nu_2$ . Let  $\pi_i = \text{Unif}(\mathcal{X}_i)$  for  $i = 1, 2$ . By symmetry, we only need to consider the case  $\frac{1}{2} \leq t \leq 1$ .

Under the tuple parametrization, we have

$$D(\psi(t, \nu_1, \nu_2) \parallel \pi) = d_{\text{KL}}\left(t \parallel \frac{1}{2}\right) + tD(\nu_1 \parallel \pi_1) + (1-t)D(\nu_2 \parallel \pi_2), \quad (119)$$

$$\phi(\psi(t, \nu_1, \nu_2)P) = \left(\frac{1}{2}, t\nu_1 + (1-t)\pi_1, (1-t)\nu_2 + t\pi_2\right), \quad (120)$$

where  $d_{\text{KL}}(x \parallel y) = x \log \frac{x}{y} + (1-x) \log \frac{1-x}{1-y}$  is the binary KL divergence function. So for  $\nu = \psi(t, \nu_1, \nu_2)$ , we have

$$\frac{D(\nu P \parallel \pi)}{D(\nu \parallel \pi)} = \frac{1}{2} \cdot \frac{D(t\nu_1 + (1-t)\pi_1 \parallel \pi_1) + D((1-t)\nu_2 + t\pi_2 \parallel \pi_2)}{d_{\text{KL}}\left(t \parallel \frac{1}{2}\right) + tD(\nu_1 \parallel \pi_1) + (1-t)D(\nu_2 \parallel \pi_2)}. \quad (121)$$

When  $t = 1$ , Eq. (121) simplifies to

$$\frac{D(\nu P \parallel \pi)}{D(\nu \parallel \pi)} = \frac{1}{2} \cdot \frac{D(\nu_1 \parallel \pi_1)}{\log 2 + D(\nu_1 \parallel \pi_1)} \quad (122)$$

which is maximized when and only when  $\nu_1$  is a point measure, taking value  $\frac{\log n}{2 \log(2n)}$ . Note that when  $t = 1$  and  $\nu_1$  is a point measure,  $\nu = \psi(t, \nu_1, \nu_2)$  does not depend on  $\nu_2$  and is always a point measure.

Now assume that  $\frac{1}{2} \leq t < 1$ . By [GP23, Prop. 17], for fixed  $D(\nu_1 \|\pi_1)$  (resp.  $D(\nu_2 \|\pi_2)$ ),  $D(t\nu_1 + (1-t)\pi_1 \|\pi_1)$  (resp.  $D(\nu_2 + t\pi_2 \|\pi_2)$ ) is uniquely (up to permutation of the alphabet) at a distribution of form  $\left(x, \frac{1-x}{n-1}, \dots, \frac{1-x}{n-1}\right)$  where  $\frac{1}{n} \leq x \leq 1$ . Therefore, we can assume WLOG that  $\nu_1 = \left(x, \frac{1-x}{n-1}, \dots, \frac{1-x}{n-1}\right)$ ,  $\nu_2 = \left(y, \frac{1-y}{n-1}, \dots, \frac{1-y}{n-1}\right)$  for some  $x, y \in [\frac{1}{n}, 1]$ . In this case, we can simplify Eq. (121) as

$$\frac{D(\nu P \|\pi)}{D(\nu \|\pi)} = \frac{1}{2} \cdot \frac{f_n(tx + \frac{1-t}{n}) + f_n((1-t)y + \frac{t}{n})}{f_2(t) + tf_n(x) + (1-t)f_n(y)} =: F_n(t, x, y), \quad (123)$$

where

$$f_n(x) = D\left(\left(x, \frac{1-x}{n-1}, \dots, \frac{1-x}{n-1}\right) \parallel \text{Unif}([n])\right) = \log n + x \log x + (1-x) \log \frac{1-x}{n-1}. \quad (124)$$

Therefore, computing  $\eta_{\text{KL}}(\pi, P)$  is equivalent to computing

$$\sup_{\substack{\frac{1}{2} \leq t \leq 1 \\ \frac{1}{n} \leq x, y \leq 1 \\ (t, x, y) \neq (\frac{1}{2}, \frac{1}{n}, \frac{1}{n})}} F_n(t, x, y) \quad (125)$$

We have reduced the original problem of computing  $\eta_{\text{KL}}(\pi, K)$ , which is a priori an  $(2n - 1)$ -dimensional optimization problem, to a optimization problem with three real variables  $t \in [0, 1], x, y \in [\frac{1}{n}, 1]$ . The following lemma helps us reduce it further to two real variables.

**Lemma 30.** *For any  $n \geq 3$ , there exists  $t^* > \frac{1}{2}$  such that*

$$\sup_{\substack{0 < t < t^* \\ \frac{1}{n} < x \leq 1}} \frac{f_n(tx + \frac{1-t}{n})}{tf_n(x)} < \frac{\log n}{\log(2n)}. \quad (126)$$

*Proof.* By [GP23, Eqs. (23) and (27)], for fixed  $0 \leq t \leq 1$ , we have

$$\sup_{\frac{1}{n} < x \leq 1} \frac{f_n(tx + \frac{1-t}{n})}{tf_n(x)} \leq t^{\frac{1}{\log n}}. \quad (127)$$

For any  $t < \exp\left((\log n) \log \frac{\log n}{\log(2n)}\right)$ , we have

$$t^{\frac{1}{\log n}} < \frac{\log n}{\log(2n)}. \quad (128)$$

Furthermore, notice that  $\exp\left((\log n) \log \frac{\log n}{\log(2n)}\right) > \frac{1}{2}$  for all  $n \geq 3$ . So we can take  $t^*$  to be any number smaller than  $\exp\left((\log n) \log \frac{\log n}{\log(2n)}\right)$ .  $\square$

Note that all arguments until this point work for any  $n \geq 3$ . From now on we will use the assumption  $n = 3$ . For  $n = 3$ , we take  $t^* = 0.58 < \exp\left((\log 3) \log \frac{\log 3}{\log 6}\right)$  and apply Lemma 30 to Eq. (125). For  $\frac{1}{2} \leq t \leq t^*$ , we can apply Lemma 30 to  $(t, x)$  and  $(1-t, y)$  and get

$$\sup_{\substack{\frac{1}{2} \leq t \leq t^* \\ \frac{1}{3} \leq x, y \leq 1 \\ (t, x, y) \neq (\frac{1}{2}, \frac{1}{3}, \frac{1}{3})}} F_3(t, x, y) < \frac{\log 3}{2 \log 6} \quad (129)$$

For  $t^* < t < 1$ , we have

$$\sup_{\substack{t^* < t < 1 \\ \frac{1}{3} \leq x, y \leq 1}} F_3(t, x, y) \geq \lim_{t \rightarrow 1^-} F_3\left(t, 1, \frac{1}{3}\right) = \frac{\log 3}{2 \log 6}. \quad (130)$$

Applying [Lemma 30](#) to  $(1-t, y)$  we see that for any  $(t, x, y)$  with  $t^* < t < 1$  and  $F_3(t, x, y) \geq \frac{\log 3}{2 \log 6}$ , we have  $F_3(t, x, y) \leq F_3\left(t, x, \frac{1}{3}\right)$ . So

$$\sup_{\substack{t^* < t < 1 \\ \frac{1}{3} \leq x, y \leq 1}} F_3(t, x, y) = \sup_{\substack{t^* < t < 1 \\ \frac{1}{3} \leq x \leq 1}} F_3\left(t, x, \frac{1}{3}\right) = \sup_{\substack{t^* < t < 1 \\ \frac{1}{3} \leq x \leq 1}} \frac{1}{2} \cdot \frac{f_3\left(tx + \frac{1-t}{3}\right)}{f_2(t) + t f_3(x)}. \quad (131)$$

In the following, we prove that for  $t^* < t < 1$ ,  $\frac{1}{3} \leq x \leq 1$ , we have

$$G_3(t, x) = \frac{\log 6}{\log 3} \cdot f_3\left(tx + \frac{1-t}{3}\right) - (f_2(t) + t f_3(x)) < 0. \quad (132)$$

Note that [Eq. \(132\)](#) implies the desired result by [Eqs. \(129\)](#) and [\(131\)](#).

**Case 1:**  $t \geq 0.999$ ,  $x \geq 0.999$ . For simplicity of notation, write  $a = 1 - t$  and  $b = 1 - x$ . Then  $0 < a \leq 0.001$  and  $0 \leq b \leq 0.001$ . Let  $c = 1 - \left(tx + \frac{1-t}{3}\right) = \frac{2}{3}a + b - ab$ . Write  $g_n(u) = -(1-u) \log(1-u) - u \log \frac{u}{n} \geq 0$ . Then

$$\begin{aligned} & G_3(1-a, 1-b) \\ &= \frac{\log 6}{\log 3} \cdot (\log 3 - g_2(c)) - (\log 2 - g_1(a) + t(\log 3 - g_2(b))) \\ &\leq \frac{\log 6}{\log 3} \cdot (\log 3 - g_2(c)) - (\log 2 - g_1(a) - a \log 3 + \log 3 - g_2(b)) \\ &= g_3(a) + g_2(b) - \frac{\log 6}{\log 3} \cdot g_2(c). \end{aligned} \quad (133)$$

For  $0 \leq w \leq 0.001$ , we have

$$0.999w \leq -(1-w) \log(1-w) \leq w. \quad (134)$$

So

$$g_3(a) \leq a \log \frac{3e}{a}, \quad g_2(b) \leq b \log \frac{2e}{b}, \quad g_2(c) \geq c \left(0.999 + \log \frac{2}{c}\right) =: h(c). \quad (135)$$

Because  $h(c)$  is concave and increasing for  $0 \leq c \leq 0.002$ , we have

$$h(c) \geq \frac{1}{2}h\left(\frac{4}{3}a\right) + \frac{1}{2}h(2b(1-a)) \geq \frac{1}{2}h\left(\frac{4}{3}a\right) + \frac{1}{2}h(1.998b). \quad (136)$$

For  $0 < a \leq 0.001$  and  $0 \leq b \leq 0.001$ , we have

$$a \log \frac{3e}{a} < \frac{1}{2}h\left(\frac{4}{3}a\right), \quad b \log \frac{2e}{b} \leq \frac{1}{2}h(1.998b). \quad (137)$$

So  $G_3(t, x) < 0$  for  $0.999 \leq t < 1$ ,  $0.999 \leq x \leq 1$ . This finishes the proof for Case 1.

**Case 2:**  $t \leq 0.999$  or  $x \leq 0.999$ . Let  $A = ([t^*, 1] \times [\frac{1}{3}, 1]) \setminus ([0.999, 1] \times [0.999, 1])$ . Our goal is to prove that  $G_3(t, x) < 0$  for all  $(t, x) \in A$ . The proof strategy is as follows. We choose  $\epsilon, \delta > 0$  and a finite set  $A^* \subseteq A$  such that the following are true.

- (a)  $G_3(t, x) < -\epsilon$  for all  $(t, x) \in A^*$ ;
- (b) For any  $(t, x) \in A$ , there exists  $(t^*, x^*) \in A^*$  such that  $\max\{|t - t^*|, |x - x^*|\} \leq \delta$ ;
- (c) For any  $(t, x), (t', x') \in A$ , if  $\max\{|t - t'|, |x - x'|\} \leq \delta$ , then  $|G_3(t, x) - G_3(t', x')| \leq \epsilon$ .

If all three items hold, then they imply our goal.

We take  $\epsilon = 0.00078$ ,  $\delta = 10^{-5}$ ,  $A^* = A \cap (10^{-5}\mathbb{Z} \times 10^{-5}\mathbb{Z})$ . **Item (a)** is verified using a computer program by iterating over all points in  $A^*$ . **Item (b)** is immediate by our choice of  $A^*$ . It remains to prove **Item (c)**. Note that on  $A$ ,  $f_3\left(tx + \frac{1-t}{3}\right)$  and  $f_2(t) + tf_3(x)$  are non-decreasing in both  $t$  and  $x$ . By convexity and monotonicity,

$$\sup_{\substack{\frac{1}{3} \leq u, u' \leq 1 \\ |u - u'| \leq \delta}} |f_3(u) - f_3(u')| = |f_3(1) - f_3(1 - \delta)| \leq 0.00014, \quad (138)$$

$$\sup_{\substack{\frac{1}{2} \leq u, u' \leq 1 \\ |u - u'| \leq \delta}} |f_2(u) - f_2(u')| = |f_2(1) - f_2(1 - \delta)| \leq 0.00013. \quad (139)$$

So for  $(t, x), (t', x') \in A$  with  $\max\{|t - t^*|, |x - x^*|\} \leq \delta$ , we have

$$|G_3(t, x) - G_3(t, x')| \leq \max \left\{ \frac{\log 6}{\log 3} \cdot \left| f_3\left(tx + \frac{1-t}{3}\right) - f_3\left(tx' + \frac{1-t}{3}\right) \right|, \right. \quad (140)$$

$$\left. |(f_2(t) + tf_3(x)) - (f_2(t) + tf_3(x'))| \right\}$$

$$\leq \max \left\{ \frac{\log 6}{\log 3} \cdot 0.00014, 0.00014 \right\} \leq 0.00023,$$

$$|G_3(t, x') - G_3(t', x')| \leq \max \left\{ \frac{\log 6}{\log 3} \cdot \left| f_3\left(tx' + \frac{1-t}{3}\right) - f_3\left(t'x' + \frac{1-t'}{3}\right) \right|, \right. \quad (141)$$

$$\left. |(f_2(t) + tf_3(x')) - (f_2(t') + tf_3(x'))| \right\}$$

$$\leq \max \left\{ \frac{\log 6}{\log 3} \cdot 0.00014, 0.00013 + \delta \log 3 \right\} \leq 0.00023.$$

Therefore

$$|G_3(t, x) - G_3(t', x')| \leq 0.00023 + 0.00023 = 0.00046 < \epsilon. \quad (142)$$

This proves **Item (c)**, thus finishing the proof for Case 2.

**Eq. (132)** follows by combining the three cases. This finishes the proof that  $\eta_{\text{KL}}(\pi, P) = \frac{\log 3}{2 \log 6}$  and equality is achieved at and only at point distributions.  $\square$

Using the same proof strategy one could in principle prove the statement of **Proposition 29** for any given  $n \geq 3$  (assuming it is true). However it is unclear to us how to prove uniformly for all  $n \geq 3$ .

Finally, we provide sufficient conditions for the extremal distribution for  $\eta_{\text{KL}}$  to have full support. This result can be contrasted with **Examples 27** and **28**, and **[BC24]** where it is shown that the half-step entropy contraction for many natural chains (e.g., the random transposition model) has point measures as their only extremizers.

**Lemma 31.** *Let  $(\pi, P)$  be a reversible pair where  $P$  is irreducible. Suppose that either*

- (i)  $P(x, x) \geq \frac{1}{2}$  for all  $x \in \mathcal{X}$ , and  $\eta_{\text{KL}}(\pi, P) > \frac{1}{2}$ , or

(ii)  $P(x, y) > 0$  for all  $x, y \in \mathcal{X}$ .

Let  $\nu$  be a distribution on  $\mathcal{X}$  satisfying  $\eta_{\text{KL}}(\pi, P) = \frac{D(\nu P \| \pi P)}{D(\nu \| \pi)}$ . Then  $\nu$  has full support.

*Proof.* Let  $f = \frac{d\nu}{d\pi}$  be the Radon-Nikodym derivative. For the sake of contradiction assume that  $\text{supp } f \neq \mathcal{X}$ . Choose  $a \in \mathcal{X} - \text{supp } f$  and  $b \in \text{supp } f$  such that  $P(a, b) > 0$ . Because  $P$  is irreducible, such  $(a, b)$  always exists.

Let  $h = \mathbb{1}_a - \frac{\pi(a)}{\pi(b)} \mathbb{1}_b$  and  $f_\epsilon = f + \epsilon h$ . For small enough  $\epsilon > 0$ , we have  $f_\epsilon \geq 0$  and  $\pi[f_\epsilon] = 1$ . We prove that for  $\epsilon > 0$  small enough, we have

$$\frac{d \text{Ent}_\pi(Pf_\epsilon)}{d\epsilon \text{Ent}_\pi(f_\epsilon)} > 0. \quad (143)$$

Computation shows that

$$\frac{d \text{Ent}_\pi(Pf_\epsilon)}{d\epsilon \text{Ent}_\pi(f_\epsilon)} = \frac{1}{\text{Ent}_\pi(f_\epsilon)^2} (\pi[Ph \log(Pf_\epsilon)]\pi[f_\epsilon \log f_\epsilon] - \pi[h \log f_\epsilon]\pi[Pf_\epsilon \log(Pf_\epsilon)]). \quad (144)$$

Expanding near  $\epsilon = 0$  gives

$$\pi[h \log f_\epsilon] = \pi(a) \log \epsilon + O(1), \quad (145)$$

$$\pi[Ph \log(Pf_\epsilon)] = \sum_{j \in \mathcal{X}} \pi(j) \left( P(j, a) - P(j, b) \frac{\pi(a)}{\pi(b)} \right) \log(Pf_\epsilon(j)) \quad (146)$$

$$= \sum_{j \in \mathcal{X} - \text{supp}(Pf)} \pi(j) P(j, a) \log(P(j, a)\epsilon) + O(1)$$

$$= \pi(a) P(a, \mathcal{X} - \text{supp}(Pf)) \log \epsilon + O(1),$$

$$\pi[f_\epsilon \log f_\epsilon] = \text{Ent}_\pi(f) + o(1), \quad (147)$$

$$\pi[Pf_\epsilon \log(Pf_\epsilon)] = \text{Ent}_\pi(Pf) + o(1). \quad (148)$$

**Case Item (i).** Because  $P(a, b) > 0$ , we have  $a \in \text{supp}(Pf)$ . Because  $P$  is lazy, we have

$$P(a, \mathcal{X} - \text{supp}(Pf)) \leq 1 - P(a, a) \leq \frac{1}{2}. \quad (149)$$

By assumption,  $\frac{\text{Ent}_\pi(Pf)}{\text{Ent}_\pi(f)} = \eta_{\text{KL}}(\pi, P) > \frac{1}{2}$ . Therefore (143) holds for  $\epsilon > 0$  small enough.

**Case Item (ii).** In this case,  $P(a, \mathcal{X} - \text{supp}(Pf)) = 0$ . So (143) holds for  $\epsilon > 0$  small enough.  $\square$

## Acknowledgments

Y.G. is supported by the National Science Foundation under Grant No. DMS-1926686.

## References

- [AG76] Rudolf Ahlswede and Peter Gács. Spreading of sets in product spaces and hypercontraction of the Markov operator. *The Annals of Probability*, pages 925–939, 1976.

- [AJK<sup>+</sup>22] Nima Anari, Vishesh Jain, Frederic Koehler, Huy Tuan Pham, and Thuy-Duong Vuong. Entropic independence: optimal mixing of down-up random walks. In *Proceedings of the 54th Annual ACM SIGACT Symposium on Theory of Computing (STOC)*, pages 1418–1430, 2022.
- [BC24] Alexandre Bristiel and Pietro Caputo. Entropy inequalities for random walks and permutations. *Annales de l’Institut Henri Poincaré (B) Probabilités et statistiques*, 60(1):54–81, 2024.
- [BCC<sup>+</sup>22] Antonio Blanca, Pietro Caputo, Zongchen Chen, Daniel Parisi, Daniel Štefankovič, and Eric Vigoda. On mixing of Markov chains: coupling, spectral independence, and entropy factorization. *Electronic Journal of Probability*, 27:1–42, 2022.
- [BCP<sup>+</sup>21] Antonio Blanca, Pietro Caputo, Daniel Parisi, Alistair Sinclair, and Eric Vigoda. Entropy decay in the Swendsen-Wang dynamics on  $\mathbb{Z}^d$ . In *Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing*, pages 1551–1564, 2021.
- [BG99] Sergej G. Bobkov and Friedrich Götze. Exponential integrability and transportation cost related to logarithmic Sobolev inequalities. *Journal of Functional Analysis*, 163(1):1–28, 1999.
- [BSM03] Abraham Berman and Naomi Shaked-Monderer. *Completely positive matrices*. World Scientific, 2003.
- [BT06] Sergey G. Bobkov and Prasad Tetali. Modified logarithmic Sobolev inequalities in discrete settings. *Journal of Theoretical Probability*, 19(2):289–336, 2006.
- [Cap23] Pietro Caputo. Lecture notes on entropy and Markov chains. *Preprint*, 2023. Available at: <http://www.mat.uniroma3.it/users/caputo/entropy.pdf>.
- [CDPP09] Pietro Caputo, Paolo Dai Pra, and Gustavo Posta. Convex entropy decay via the Bochner-Bakry-Emery approach. *Annales de l’IHP Probabilités et statistiques*, 45(3):734–753, 2009.
- [Che03] Mu-Fa Chen. Variational formulas of Poincaré-type inequalities for birth-death processes. *Acta Mathematica Sinica*, 19(4):625–644, 2003.
- [Che05] Mu-Fa Chen. Poincaré-type inequalities in dimension one. *Eigenvalues, Inequalities, and Ergodic Theory*, pages 113–130, 2005.
- [CLV21] Zongchen Chen, Kuikui Liu, and Eric Vigoda. Optimal mixing of Glauber dynamics: Entropy factorization via high-dimensional expansion. In *Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing*, pages 1537–1550, 2021.
- [CMS24] Pietro Caputo, Florentin Münch, and Justin Salez. Entropy and curvature: beyond the Peres-Tetali conjecture. *arXiv preprint arXiv:2401.17148*, 2024.
- [DMLM03] Pierre Del Moral, Michel Ledoux, and Laurent Miclo. On contraction properties of Markov kernels. *Probability theory and related fields*, 126(3):395–420, 2003.
- [Dob56] Roland L. Dobrushin. Central limit theorem for nonstationary Markov chains. I. *Theory of Probability & Its Applications*, 1(1):65–80, 1956.

- [DSC96] Persi Diaconis and Laurent Saloff-Coste. Logarithmic Sobolev inequalities for finite Markov chains. *The Annals of Applied Probability*, 6(3):695–750, 1996.
- [GP23] Yuzhou Gu and Yury Polyanskiy. Non-linear log-Sobolev inequalities for the Potts semigroup and applications to reconstruction problems. *Communications in Mathematical Physics*, 404(2):769–831, 2023.
- [GQ03] Fuqing Gao and Jeremy Quastel. Exponential decay of entropy in the random transposition and Bernoulli-Laplace models. *The Annals of Applied Probability*, 13(4):1591–1600, 2003.
- [KM17] Tali Kaufman and David Mass. High dimensional random walks and colorful expansion. In *8th Innovations in Theoretical Computer Science Conference (ITCS 2017)*. Schloss-Dagstuhl-Leibniz Zentrum für Informatik, 2017.
- [LP17] David A. Levin and Yuval Peres. *Markov chains and mixing times*, volume 107. American Mathematical Soc., 2017.
- [Mic97] Laurent Miclo. Remarques sur l’hypercontractivité et l’évolution de l’entropie pour des chaînes de Markov finies. In *Séminaire de Probabilités XXXI*, pages 136–167. Springer, 1997.
- [MM62] John E. Maxfield and Henryk Minc. On the matrix equation  $X'X = A$ . *Proceedings of the Edinburgh Mathematical Society*, 13(2):125–129, 1962.
- [MN61] Marvin Marcus and Morris Newman. The permanent of a symmetric matrix. *Notices Amer. Math. Soc*, 8:595, 1961.
- [Mün23] Florentin Münch. Ollivier curvature, isoperimetry, concentration, and log-Sobolev inequality. *arXiv preprint arXiv:2309.06493*, 2023.
- [PW17] Yury Polyanskiy and Yihong Wu. Strong data-processing inequalities for channels and bayesian networks. In *Convexity and Concentration*, pages 211–249. Springer, 2017.
- [PW24] Yury Polyanskiy and Yihong Wu. *Information theory: From coding to learning*. Cambridge university press, 2024.
- [Rag16] Maxim Raginsky. Strong data processing inequalities and  $\Phi$ -Sobolev inequalities for discrete channels. *IEEE Transactions on Information Theory*, 62(6):3355–3389, 2016.
- [Rob01] Cyril Roberto. *Inégalités de Hardy et de Sobolev logarithmiques*. PhD thesis, PhD thesis, 2001.
- [Sin64] Richard Sinkhorn. A relationship between arbitrary positive matrices and doubly stochastic matrices. *The annals of mathematical statistics*, 35(2):876–879, 1964.
- [STY23] Justin Salez, Konstantin Tikhomirov, and Pierre Youssef. Upgrading MLSI to LSI for reversible Markov chains. *Journal of Functional Analysis*, 285(9):110076, 2023.
- [TY24] Konstantin Tikhomirov and Pierre Youssef. Regularized modified log-Sobolev inequalities and comparison of Markov chains. *The Annals of Probability*, 52(4):1201–1224, 2024.