

# CS395T: Continuous Algorithms, Part XVI

## Stochastic calculus

Kevin Tian

### 1 Drift-diffusion processes

In this lecture, we introduce the tools required to study a framework for sampling algorithm design based on discretizing *drift-diffusion processes*, which is able to exploit different structural aspects of the target density compared to the previous few lectures (for example, an appropriate notion of well-conditionedness for densities). We adopt a dual perspective of such processes, viewing them as both a stochastic evolution in particle space (i.e., a particle in  $\mathbb{R}^d$  following a stochastic trajectory), as well as a deterministic evolution in density space (i.e., an evolving measure in  $\mathcal{P}(\mathbb{R}^d)$ , the set of probability densities on  $\mathbb{R}^d$  which are absolutely continuous with respect to the Lebesgue measure). This perspective lets us adopt a rich set of analytical tools on both  $\mathbb{R}^d$  and  $\mathcal{P}(\mathbb{R}^d)$  to analyze the convergence and bound the discretization error of the resulting algorithms.

We mention that this section is largely structured to briefly set up the mathematical tools to rigorously study stochastic calculus, and most of the formalism can be ignored on a first read in the later sections. We start by introducing Itô calculus, beginning with *Brownian motion*.

**Definition 1** (Brownian motion). *We define Brownian motion in  $\mathbb{R}^d$ , denoted by  $\{\mathbf{B}_t\}_{t \geq 0}$ , to be a stochastic process, i.e., a random sequence of points in  $\mathbb{R}^d$  indexed by  $t \geq 0$  (thought of as time), satisfying the following properties.*

1.  $\mathbf{B}_0 = \mathbf{0}_d$ .
2.  $\{\mathbf{B}_t\}_{t \geq 0}$  is continuous with probability 1.
3. For all  $k \in \mathbb{N}$  and all  $\{t_i\}_{i=0}^k \subset \mathbb{R}_{\geq 0}$  with  $t_0 = 0 < t_1 < \dots < t_k$ , all of the random variables  $\mathbf{B}_{t_{i+1}} - \mathbf{B}_{t_i}$  for  $0 \leq i \leq k-1$ , are mutually independent.
4. For all  $0 \leq s \leq t$ ,  $\mathbf{B}_t - \mathbf{B}_s$  is distributed as  $\mathcal{N}(\mathbf{0}_d, (t-s)\mathbf{I}_d)$ .

For an existence proof of Brownian motion, see Chapter 7 of [Dur10]. The probability space that Brownian motion is defined on, as well as all stochastic processes we will study, is denoted  $\{\mathcal{F}_t\}_{t \geq 0}$ , a *filtration* satisfying  $\mathcal{F}_s \subseteq \mathcal{F}_t$  for all  $0 \leq s \leq t$ . We say that Brownian motion is a stochastic process *adapted to* the filtration  $\{\mathcal{F}_t\}_{t \geq 0}$ . Informally, we can think of  $\mathcal{F}_t$  as containing the information of the randomness used up to time  $t$ , i.e., the realizations of the random Brownian motion, so the inclusion condition means that more information is always available later in time.

To give some intuition for the properties of Brownian motion, we prove the reflection principle.

**Lemma 1** (Reflection principle). *Let  $\{B_t\}_{t \geq 0}$  be Brownian motion in  $\mathbb{R}$ . For all  $t \geq 0$  and  $a > 0$ ,*

$$\Pr \left[ \sup_{0 \leq s \leq t} B_s \geq a \right] = 2 \Pr [B_t \geq a].$$

*Proof.* Let  $\tau \geq 0$  be a random stopping time<sup>1</sup> corresponding to the first time the Brownian motion

<sup>1</sup>A stopping time  $\tau$  adapted to a filtration  $\{\mathcal{F}_t\}_{t \geq 0}$  is a random variable in  $\mathbb{R}_{\geq 0}$ , such that for all  $t \in \mathbb{R}_{\geq 0}$ , the event  $\tau \leq t$  is measurable with respect to  $\mathcal{F}_t$ . In other words, we can determine whether the event defining  $\tau$  has occurred using information in  $\mathcal{F}_t$ , e.g., the first time an  $\{\mathcal{F}_t\}_{t \geq 0}$ -measurable event occurs is a stopping time.

reaches  $a$ , i.e.,  $\tau$  is the random variable equal to  $\inf\{s \mid B_s = a\}$ . By independence of increments,

$$\begin{aligned} \Pr \left[ \sup_{0 \leq s \leq t} B_s \geq a \right] &= \Pr \left[ \sup_{0 \leq s \leq t} B_s \geq a \text{ and } B_t \geq a \right] + \Pr \left[ \sup_{0 \leq s \leq t} B_s \geq a \text{ and } B_t < a \right] \\ &= \Pr [B_t \geq a] + \Pr \left[ \sup_{0 \leq s \leq t} B_s \geq a \text{ and } B_t - B_\tau < 0 \right] \\ &= \Pr [B_t \geq a] + \frac{1}{2} \Pr \left[ \sup_{0 \leq s \leq t} B_s \geq a \right], \end{aligned}$$

as  $\sup_{0 \leq s \leq t} B_s \geq a$  implies that  $\tau \leq t$ , so  $B_t - B_\tau < 0$  with probability  $\frac{1}{2}$  regardless of the realization of  $\tau \leq t$ . Rearranging the above display gives the claim.  $\square$

Brownian motion is an example of a continuous *martingale* with respect to  $\{\mathcal{F}_t\}_{t \geq 0}$ . To define this, recall that a martingale  $\{\mathbf{x}_t\}_{t \geq 0}$  adapted to a filtration  $\{\mathcal{F}_t\}_{t \geq 0}$  satisfies the property that

$$\mathbb{E}[\mathbf{x}_t \mid \mathcal{F}_s] = \mathbf{x}_s \text{ for all } 0 \leq s \leq t.$$

A useful property of martingales is that by Jensen's inequality, for all convex  $\varphi$  and  $0 \leq s \leq t$ ,

$$\mathbb{E}\varphi(\mathbf{x}_t) = \mathbb{E}[\mathbb{E}[\varphi(\mathbf{x}_t) \mid \mathcal{F}_s]] \geq \mathbb{E}\varphi(\mathbf{x}_s).$$

One consequence of the martingale property (see Chapter 7.5, [Dur10]) is that for all bounded stopping times  $\tau$ ,  $\mathbb{E}\mathbf{x}_\tau = \mathbf{x}_0$ . We are now ready to introduce the Itô integral, which is defined with respect to a continuous process  $\{\mathbf{H}_t\}_{t \geq 0} \subset \mathbb{R}^{d \times d}$  adapted to  $\{\mathcal{F}_t\}_{t \geq 0}$ .<sup>2</sup> We can view  $\{\mathbf{H}_t\}_{t \geq 0}$  as a “reweighting” of Brownian motion, and we correspondingly define the Itô integral

$$\mathbf{x}_t := \int_0^t \mathbf{H}_s d\mathbf{B}_s,$$

which is a continuous martingale adapted to  $\{\mathcal{F}_t\}_{t \geq 0}$ . In particular, this means that

$$\mathbb{E} \left[ \int_s^t \mathbf{H}_u d\mathbf{B}_u \mid \mathcal{F}_s \right] = 0 \text{ for all } 0 \leq s \leq t.$$

It is useful to view the Itô integral as a stochastic process indexed by  $t$ . For example, choosing  $\mathbf{H}_t = \mathbf{I}_d$  for all  $t$  gives that  $\mathbf{x}_t = \mathbf{B}_t$  is simply Brownian motion.

An important characteristic of a stochastic process  $\{\mathbf{x}_t\}_{t \geq 0}$  is its *quadratic variation*, defined by

$$[\mathbf{x}]_t = \lim_{\|P\|_{\text{gap}} \rightarrow 0} \sum_{k=1}^{|P|} \|\mathbf{x}_{t_k} - \mathbf{x}_{t_{k-1}}\|_2^2, \quad (1)$$

the limit over meshes  $P = \{t_k\}_{k \in |P|} \subset [0, t]$  where  $\|P\|_{\text{gap}} = \max_{k \in |P|} |t_k - t_{k-1}|$ , and  $t_{|P|} = t$  and  $t_0 := 0$ . One can formally verify that when  $\mathbf{x}_t$  is an Itô integral driven by  $\{h_t\}_{t \geq 0}$ , we have that

$$[\mathbf{x}]_t = \int_0^t \|\mathbf{H}_s\|_F^2 ds, \quad (2)$$

thought of as the “total variance” accumulated thus far throughout the process. For example, if  $\{\mathbf{B}_t\}_{t \geq 0}$  is 1-dimensional Brownian motion then  $[\mathbf{B}]_t = t$ . The quadratic variation is useful in the context of a characterization known as the *Dambis-Dubins-Schwarz* theorem [Dam65, DS65], which says that in one dimension, all continuous martingales  $\{x_t\}_{t \geq 0}$  are distributionally identical to a time-changed Brownian motion  $\{B_{\tau(t)}\}_{t \geq 0}$ , where we let  $q(t) := [x]_t$  be the quadratic variation at time  $t$ . So, all continuous martingales in  $\mathbb{R}$  are characterized by their quadratic variation.

Intuitively, the expression (2) holds because of the heuristic “ $d\mathbf{B}_t^2 = dt$ ,” which makes sense in one dimension and can be extended to higher dimensions straightforwardly. This is because an infinitesimal advancement of Brownian motion (at timescale  $dt$ ) behaves like  $\mathcal{N}(0, dt)$ , so its square

<sup>2</sup>Formally, the continuity of  $\{\mathbf{H}_t\}_{t \geq 0}$  is not necessary, and it just needs to be *progressive*, a strengthening of being adapted that implies stopped processes are measurable. All of the Itô integrals that we encounter will be continuous and adapted processes, so we do not make this distinction for simplicity.

has expectation  $dt$ . Formalizing this intuition and using the martingale property of Itô integrals to deal with “cross-terms” establishes equivalence of the quadratic variation definitions (1) and (2).

Relatedly, we mention that one of the most important properties of the Itô integral is the *Itô isometry*, which holds if the right-hand side below is finite:

$$\mathbb{E} \left[ \left\| \int_0^t \mathbf{H}_s d\mathbf{B}_s \right\|_2^2 \right] = \mathbb{E} \left[ \int_0^t \|\mathbf{H}_s\|_{\mathbb{F}}^2 ds \right].$$

The handling of cross-terms when expanding the above expression is handled analogously to the intuition for (2), and the Itô isometry then follows by taking total expectations on both sides.

**Remark 1.** *Much of the previous discussion can in fact be substantially generalized via the formalism of local martingales and localizing sequences. Specifically, while martingales have  $\mathbb{E}\mathbf{x}_\tau = \mathbf{x}_0$  for all stopping times  $\tau$ , a local martingale  $\{\mathbf{x}_t\}_{t \geq 0}$  adapted to  $\{\mathcal{F}_t\}_{t \geq 0}$  is a stochastic process satisfying the weaker property that there is a sequence of stopping times  $\{\tau_n\}_{n \in \mathbb{N}}$  such that  $\lim_{n \rightarrow \infty} \tau_n \rightarrow \infty$  with probability 1, where  $\mathbf{x}_{\min\{t, \tau_n\}}$  is a martingale for all  $n \in \mathbb{N}$ .<sup>3</sup> These stopping times are called a localizing sequence, and are used to formalize Itô integrals, which as previously mentioned can be extended to handle progressive integrands, also through the machinery of localizing sequences. Localizing sequences are useful as a way to handle undesirable behavior as  $t \rightarrow \infty$ , for rigorously defining key concepts in stochastic calculus by using Itô integrals.<sup>4</sup>*

We also mention that Itô integrals are not the only method of formalizing a theory of stochastic calculus. In particular, the Stratonovich integral is a popular alternative in physics. Loosely speaking, the distinction is that the Stratonovich integral is defined with respect to a “midpoint” rule (familiar from Riemann integration), whereas the Itô integral is defined with respect to left endpoints. The key difference useful for our purposes is that the definition of the Itô integral makes it a martingale, whereas the Stratonovich integral is not in general. For the remainder of the lecture, we specialize Itô integrals to the setting of continuous adapted  $\{\mathbf{H}_t\}_{t \geq 0}$ .

Next, we define drift-diffusion processes on  $\mathbb{R}^d$ , adapted to the same filtration  $\{\mathcal{F}_t\}_{t \geq 0}$  that Brownian motion is adapted to. Informally, drift-diffusion processes can be thought of as modeling the deterministic and stochastic components in the time-evolution of a random particle. A drift-diffusion process  $\{\mathbf{x}_t\}_{t \geq 0}$  on  $\mathbb{R}^d$  is driven by a vector-valued function  $\boldsymbol{\mu} : \mathbb{R}^d \rightarrow \mathbb{R}^d$  and a matrix-valued function  $\boldsymbol{\sigma} : \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$ , and captured by the stochastic differential equation (SDE)

$$d\mathbf{x}_t = \boldsymbol{\mu}(x_t)dt + \boldsymbol{\sigma}(\mathbf{x}_t)d\mathbf{B}_t. \quad (3)$$

In principle, SDEs can be defined with respect to any progressive process (not just functions  $\boldsymbol{\mu}$  and  $\boldsymbol{\sigma}$  of  $x_t$ ), but (3) suffices for our purposes. We say such a stochastic process  $\{\mathbf{x}_t\}_{t \geq 0}$  is a drift-diffusion process, sometimes also called an Itô process, and we write it in integral form as

$$\mathbf{x}_t = \mathbf{x}_0 + \int_0^t \boldsymbol{\mu}(\mathbf{x}_s)ds + \int_0^t \boldsymbol{\sigma}(\mathbf{x}_s)d\mathbf{B}_s.$$

To establish existence and uniqueness of SDE solutions, the following characterization is helpful. We defer a proof to [Øk03]; the intuition for this result is similar to that for the Picard-Lindelöf theorem for establishing existence and uniqueness for the solution of ordinary differential equations (ODEs), which can be proven by defining an appropriate fixed-point iteration.

**Proposition 1** (Theorem 5.2.1, [Øk03]). *Suppose, for a SDE of the form (3), that  $\boldsymbol{\mu}$  and  $\boldsymbol{\sigma}$  both have a finite Lipschitz constant (with respect to  $\|\cdot\|_2$  and  $\|\cdot\|_{\mathbb{F}}$ , respectively). Then for all  $t \geq 0$ , there exists a unique solution  $\mathbf{x}_t$  to the SDE for every realization of  $\mathcal{F}_t$ . Moreover,  $\mathbf{x}_t$  is square-integrable in the sense that  $\mathbb{E} \int_0^t \|\mathbf{x}_s\|_2^2 ds < \infty$ .*

We conclude the section with the most useful property of drift-diffusion processes when performing computations: a stochastic calculus generalization of the chain rule.

<sup>3</sup>One example which may help see the difference is the “sticky Brownian motion” in  $\mathbb{R}$ , which is Brownian motion until the first time  $-1$  is reached, and is stuck at  $-1$  henceforth. This is not a martingale, because the stopping time  $\tau$  when  $-1$  is first reached is adapted to  $\{\mathcal{F}_t\}_{t \geq 0}$ , but  $\mathbb{E}[x_\tau] = -1$ . However, it is a local martingale.

<sup>4</sup>For instance, every bounded local martingale is a martingale, so the distinction is only due to limiting behavior.

**Proposition 2** (Itô's lemma). *Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be twice-differentiable, and suppose that  $\{\mathbf{x}_t\}_{t \geq 0}$  is a drift-diffusion process following the SDE (3). Then the stochastic process  $\{f(\mathbf{x}_t)\}_{t \geq 0}$  is also a drift-diffusion process, and follows the SDE*

$$df(\mathbf{x}_t) = \left( \langle \nabla f(\mathbf{x}_t), \boldsymbol{\mu}(\mathbf{x}_t) \rangle + \frac{1}{2} \langle \nabla^2 f(\mathbf{x}_t), \boldsymbol{\sigma}(\mathbf{x}_t) \boldsymbol{\sigma}(\mathbf{x}_t)^\top \rangle \right) dt + \langle \nabla f(\mathbf{x}_t), \boldsymbol{\sigma}(\mathbf{x}_t) d\mathbf{B}_t \rangle.$$

More generally, if  $\{\mathbf{x}_t\}_{t \geq 0}$  follows the SDE

$$d\mathbf{x}_t = \boldsymbol{\mu}_t dt + \boldsymbol{\sigma}_t d\mathbf{B}_t,$$

where  $\{\boldsymbol{\mu}_t, \boldsymbol{\sigma}_t\}_{t \geq 0} \subset \mathbb{R}^d \times \mathbb{R}^{d \times d}$  are continuous stochastic processes adapted to the same filtration  $\{\mathcal{F}_t\}_{t \geq 0}$  as  $\{\mathbf{B}_t\}_{t \geq 0} \subset \mathbb{R}^d$ , then  $\{f(\mathbf{x}_t)\}_{t \geq 0}$  is a stochastic process following the SDE

$$df(\mathbf{x}_t) = \left( \langle \nabla f(\mathbf{x}_t), \boldsymbol{\mu}_t \rangle + \frac{1}{2} \langle \nabla^2 f(\mathbf{x}_t), \boldsymbol{\sigma}_t \boldsymbol{\sigma}_t^\top \rangle \right) dt + \langle \nabla f(\mathbf{x}_t), \boldsymbol{\sigma}_t d\mathbf{B}_t \rangle.$$

We note that there is also a generalization of Proposition 2 to the more complicated setting where  $f$  is a time-dependent function, but we will not require it. The proof of Proposition 2 follows from a similar calculation as used to show the quadratic variation formula (2) holds, e.g., formalizing our aforementioned “ $dB_t^2 = dt$ ” argument. For example, a Taylor expansion gives that for vanishing  $\eta$ , and approximating  $\mathbf{x}_{t+\eta} \approx \mathbf{x}_t + \eta \boldsymbol{\mu}(\mathbf{x}_t) + \sqrt{\eta} \boldsymbol{\sigma}(\mathbf{x}_t) \boldsymbol{\xi}$  for  $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$ ,<sup>5</sup>

$$\begin{aligned} df(\mathbf{x}_t) &\approx f(\mathbf{x}_{t+\eta}) - f(\mathbf{x}_t) \\ &\approx \langle \nabla f(\mathbf{x}_t), \eta \boldsymbol{\mu}(\mathbf{x}_t) + \sqrt{\eta} \boldsymbol{\sigma}(\mathbf{x}_t) \boldsymbol{\xi} \rangle + \langle \nabla^2 f(\mathbf{x}_t), \eta \boldsymbol{\sigma}(\mathbf{x}_t) \boldsymbol{\xi} \boldsymbol{\xi}^\top \boldsymbol{\sigma}(\mathbf{x}_t)^\top \rangle \\ &\approx \eta \left( \langle \nabla f(\mathbf{x}_t), \boldsymbol{\mu}(\mathbf{x}_t) \rangle + \langle \nabla^2 f(\mathbf{x}_t), \boldsymbol{\sigma}(\mathbf{x}_t) \boldsymbol{\sigma}(\mathbf{x}_t)^\top \rangle \right) + \langle \nabla f(\mathbf{x}_t), \boldsymbol{\sigma}(\mathbf{x}_t) \cdot (\sqrt{\eta} \boldsymbol{\xi}) \rangle. \end{aligned}$$

We used the approximation  $\boldsymbol{\xi} \boldsymbol{\xi}^\top = \mathbf{I}_d$  in the last line, and dropped all terms of lower order than  $\eta$ . Unlike in the standard chain rule, there is an additional second-order component that persists, because  $\|\sqrt{\eta} \boldsymbol{\xi}\|_2^2$  scales with  $\eta$ , explaining the presence of the additional term in Proposition 2.

A particular drift-diffusion process of significant interest is the *Langevin diffusion*. Letting  $V : \mathbb{R}^d \rightarrow \mathbb{R}$  be a twice-differentiable function, the associated Langevin diffusion is the SDE

$$d\mathbf{x}_t = -\nabla V(\mathbf{x}_t) dt + \sqrt{2} d\mathbf{B}_t, \quad (4)$$

i.e., we take  $\boldsymbol{\mu} = -\nabla V$  and  $\boldsymbol{\sigma} = \sqrt{2} \mathbf{I}_d$  in (3). Intuitively, (4) can be thought of as a noisy gradient flow, with a deterministic drift  $-\nabla V(\mathbf{x}_t) dt$  (making function value progress) and a diffusion component  $\sqrt{2} d\mathbf{B}_t$ . For the SDE (4), Itô's lemma (Proposition 2) yields

$$df(\mathbf{x}_t) = (-\langle \nabla f(\mathbf{x}_t), \nabla V(\mathbf{x}_t) \rangle + \Delta f(\mathbf{x}_t)) dt + \sqrt{2} \langle \nabla f(\mathbf{x}_t), d\mathbf{B}_t \rangle, \quad (5)$$

where  $\Delta$  is the *Laplacian* operator given by the formula

$$\Delta f(\mathbf{x}) = \text{Tr}(\nabla^2 f(\mathbf{x})). \quad (6)$$

**Remark 2.** *There is a connection between the Laplacian operator defined in (6) and the Laplacian matrix associated with a graph, defined in Definition 5, Part II. At a high level, the Laplacian  $\mathbf{L}_G$  of an unweighted path graph  $G$  acts on each indicator vector  $\mathbf{e}_v$  for a vertex  $v$  of  $G$ , with neighbors  $u, w$ , by computing  $\mathbf{L}_G \mathbf{e}_v = 2\mathbf{e}_v - \mathbf{e}_u - \mathbf{e}_w$ , i.e., it computes the difference between a vertex value and “averages” around its neighbors. Similarly, the Laplacian operator  $\Delta$  computes a difference of a function value with an average in a small neighborhood, which is best-understood in one dimension:*

$$\begin{aligned} \eta f''(x) &\approx f'(x + \eta) - f'(x) \\ &\approx f(x + \eta) - f(x) - f(x) + f(x - \eta) = f(x + \eta) + f(x - \eta) - 2f(x). \end{aligned}$$

<sup>5</sup>Note that  $\sqrt{\eta} \boldsymbol{\xi}$  has covariance matrix  $\eta \mathbf{I}_d$ , as required by Definition 1 for advancing time by  $\eta$ .

## 2 Markov semigroups

We now take a dual view on stochastic differential equations, by assuming that  $\mathbf{x}_0$  is drawn from a density  $\pi_0 \in \mathcal{P}(\mathbb{R}^d)$ , and associating each time  $t \geq 0$  with a density  $\pi_t \in \mathcal{P}(\mathbb{R}^d)$  corresponding to the density of  $\mathbf{x}_t$  (called its law), which evolves from  $\mathbf{x}_0$  following a drift-diffusion process (3). That is, while the evolution of  $\mathbf{x}_t$  is stochastic, the evolution of its *distribution* can be viewed as a deterministic process following a partial differential equation (PDE) derived from (3).

To introduce this equation and its analysis, we first give some definitions from *Markov semigroup theory*, which studies “continuous-time Markov chains” such as drift-diffusion processes through the evolution of functions. For a drift-diffusion process  $\{\mathbf{x}_t\}_{t \geq 0}$ , we define an associated *Markov semigroup*  $\{P_t\}_{t \geq 0}$ , which acts on an appropriate set of functions  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  via the definition

$$P_t f(\mathbf{x}) := \mathbb{E}[f(\mathbf{x}_t) \mid \mathbf{x}_0 = \mathbf{x}]. \quad (7)$$

Formally, when an SDE (3) has a stationary density  $\pi^* : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$ , i.e., drawing  $\mathbf{x}_0 \sim \pi^*$  yields  $\mathbf{x}_t$  distributed according to  $\pi^*$  for all  $t \geq 0$ , an appropriate set of functions is

$$L^2(\pi^*) := \left\{ f : \mathbb{R}^d \rightarrow \mathbb{R} \mid \int f(\mathbf{x})^2 \pi^*(\mathbf{x}) d\mathbf{x} < \infty \right\}. \quad (8)$$

For brevity, we assume that all functions discussed in this section are sufficiently smooth<sup>6</sup> and in  $L^2(\pi^*)$ ; all arguments can be formalized for a broader family of functions in  $L^2(\pi^*)$  by taking appropriate limits with an approximating sequence of smooth functions.

The Markov (memoryless) property of stochastic processes, including drift-diffusion processes, says that for  $0 \leq s \leq t$ ,  $\mathbf{x}_t$  is independent of  $\mathcal{F}_s$  given  $\mathbf{x}_s$ . Because drift-diffusion processes are Markov, the law of iterated expectations shows that for all  $s, t \geq 0$ ,

$$P_{t+s} f = P_t P_s f = P_s P_t f \text{ for all } f : \mathbb{R}^d \rightarrow \mathbb{R}, \quad (9)$$

and  $P_0$  is clearly the identity operator, which justifies calling  $\{P_t\}_{t \geq 0}$  a semigroup. In fact, (9) shows that it is additionally a *commutative* semigroup, which will be used later.

Next, given a Markov semigroup  $\{P_t\}_{t \geq 0}$ , we define the associated *generator*  $\mathcal{L}$  which also acts on functions  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  via the following definition (which holds pointwise on  $\mathbb{R}^d$ ):

$$\mathcal{L}f := \lim_{\eta \rightarrow 0} \frac{P_\eta f - f}{\eta}. \quad (10)$$

The generator commutes with all elements of the semigroup  $\{P_t\}_{t \geq 0}$ , due to the property (9).

**Lemma 2.** For all  $t \geq 0$  and  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , we have pointwise on  $\mathbb{R}^d$ ,

$$\frac{\partial}{\partial t} P_t f = \mathcal{L} P_t f = P_t \mathcal{L} f. \quad (11)$$

*Proof.* By using the commutativity property in (9),

$$\mathcal{L} P_t f = \lim_{\eta \rightarrow 0} \frac{P_\eta - P_0}{\eta} P_t f = \lim_{\eta \rightarrow 0} \frac{P_{t+\eta} f - P_t f}{\eta} = \lim_{\eta \rightarrow 0} P_t \frac{P_\eta - P_0}{\eta} f = P_t \mathcal{L} f$$

holds pointwise in  $\mathbb{R}^d$ . The conclusion follows since  $\frac{d}{dt} P_t f = \lim_{\eta \rightarrow 0} \frac{P_{t+\eta} f - P_t f}{\eta}$  by definition.  $\square$

The expression (11) is sometimes called *Kolmogorov’s backward equation*, a name which is perhaps best justified by also introducing *Kolmogorov’s forward equation*. Let  $\pi_0 \in \mathcal{P}(\mathbb{R}^d)$  be a density. Because of the definition (7), we have that for all  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ ,

$$\mathbb{E}f(\mathbf{x}_t) = \int P_t f(\mathbf{x}) \pi_0(\mathbf{x}) d\mathbf{x} = \int f(\mathbf{x}) P_t^* \pi_0(\mathbf{x}) d\mathbf{x},$$

<sup>6</sup>When we say a function in this lecture is smooth, we mean that it has continuous derivatives of all orders.

where  $P_t^*$  is the *adjoint operator* to  $P_t$ .<sup>7</sup> This means that  $\pi_t$ , the law of  $\mathbf{x}_t$ , is given by  $P_t^* \pi_0$ , so  $P_t^*$  is formally the operator which advances  $\pi_0$  forward by time  $t$ . In other words, if we want to understand the density  $\pi_t$ , i.e., the law of  $x_t$ , we should aim to understand  $P_t^*$ . Indeed, appropriately taking adjoints of Lemma 2 yields Kolmogorov's forward equation,

$$\frac{\partial}{\partial t} \underbrace{P_t^* \pi_0}_{=\pi_t} = \mathcal{L}^* P_t^* \pi_0 = P_t^* \mathcal{L}^* \pi_0, \quad (12)$$

so  $P_t^*$  and  $\mathcal{L}^*$  also commute for all  $t \geq 0$ . In the case of drift-diffusion processes, there is in fact a convenient formula for  $P_t^*$ , enabling us to write a PDE for density evolution. We first remind the reader of the multivariate integration by parts formula, which holds for any differentiable function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  and vector field  $\mathbf{v} : \mathbb{R}^d \rightarrow \mathbb{R}^d$  with sufficiently fast decay at infinity,<sup>8</sup>

$$\int \langle \mathbf{v}(\mathbf{x}), \nabla f(\mathbf{x}) \rangle d\mathbf{x} = - \int f(\mathbf{x}) (\nabla \cdot \mathbf{v})(\mathbf{x}) d\mathbf{x}, \quad (13)$$

where  $\nabla \cdot$  is the divergence operator of a vector-valued function (for example,  $\nabla \cdot \nabla = \text{Tr} \nabla^2 = \Delta$ ).

**Proposition 3** (Fokker-Planck equation). *Let  $\{\mathbf{x}_t\}_{t \geq 0}$  follow the drift-diffusion process (3) from  $\mathbf{x}_0 \sim \pi_0 \in \mathcal{P}(\mathbb{R}^d)$ . Then for all  $t \geq 0$ , denoting the law of  $\mathbf{x}_t$  by  $\pi_t$ , we have*

$$\frac{\partial}{\partial t} \pi_t(\mathbf{x}) = -\nabla \cdot (\boldsymbol{\mu}(\mathbf{x}) \pi_t(\mathbf{x})) + \frac{1}{2} \sum_{i,j \in [d]} \frac{\partial^2}{\partial \mathbf{x}_i \partial \mathbf{x}_j} [\boldsymbol{\sigma}(\mathbf{x}) \boldsymbol{\sigma}(\mathbf{x})^\top \pi_t(\mathbf{x})]_{ij} \text{ for all } \mathbf{x} \in \mathbb{R}^d. \quad (14)$$

*Proof.* We first compute using Proposition 2 that for smooth  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , and  $\mathbf{x}_t$  following (3),

$$d\mathbb{E}f(\mathbf{x}_t) = \left( \langle \nabla f(\mathbf{x}_t), \boldsymbol{\mu}(\mathbf{x}_t) \rangle + \frac{1}{2} \langle \nabla^2 f(\mathbf{x}_t), \boldsymbol{\sigma}(\mathbf{x}_t) \boldsymbol{\sigma}(\mathbf{x}_t)^\top \rangle \right) dt,$$

where we drop the martingale term  $\langle \nabla f(\mathbf{x}_t), \boldsymbol{\sigma}(\mathbf{x}_t) d\mathbf{B}_t \rangle$  because it vanishes in expectation. Therefore, by the definition (10) of  $\mathcal{L}$ , we have

$$\mathcal{L}f(\mathbf{x}) = \langle \nabla f(\mathbf{x}), \boldsymbol{\mu}(\mathbf{x}) \rangle + \frac{1}{2} \langle \nabla^2 f(\mathbf{x}), \boldsymbol{\sigma}(\mathbf{x}) \boldsymbol{\sigma}(\mathbf{x})^\top \rangle, \quad (15)$$

so for all probability densities  $\pi \in \mathcal{P}(\mathbb{R}^d)$  (which must decay to zero as  $\mathbf{x} \rightarrow \infty$ ),

$$\begin{aligned} \int f(\mathbf{x}) \mathcal{L}^* \pi(\mathbf{x}) d\mathbf{x} &= \int \mathcal{L}f(\mathbf{x}) \pi(\mathbf{x}) d\mathbf{x} \\ &= \int \left( \langle \nabla f(\mathbf{x}), \boldsymbol{\mu}(\mathbf{x}) \rangle + \frac{1}{2} \langle \nabla^2 f(\mathbf{x}), \boldsymbol{\sigma}(\mathbf{x}) \boldsymbol{\sigma}(\mathbf{x})^\top \rangle \right) \pi(\mathbf{x}) d\mathbf{x} \\ &= \int \left( -f(\mathbf{x}) \nabla \cdot (\boldsymbol{\mu}(\mathbf{x}) \pi(\mathbf{x})) - \frac{1}{2} \langle \nabla f(\mathbf{x}), \nabla \cdot (\boldsymbol{\sigma}(\mathbf{x}) \boldsymbol{\sigma}(\mathbf{x})^\top \pi(\mathbf{x})) \rangle \right) d\mathbf{x} \\ &= \int f(\mathbf{x}) \left( -\nabla \cdot (\boldsymbol{\mu}(\mathbf{x}) \pi(\mathbf{x})) + \frac{1}{2} \sum_{i,j \in [d]} \frac{\partial^2}{\partial \mathbf{x}_i \partial \mathbf{x}_j} [\boldsymbol{\sigma}(\mathbf{x}) \boldsymbol{\sigma}(\mathbf{x})^\top \pi(\mathbf{x})]_{ij} \right) d\mathbf{x}, \end{aligned}$$

where we repeatedly applied the multivariate integration by parts formula (13). Because this holds for all test functions  $f$  and densities  $\pi \in \mathcal{P}(\mathbb{R}^d)$ , this means that pointwise,

$$\mathcal{L}^* \pi(\mathbf{x}) = -\nabla \cdot (\boldsymbol{\mu}(\mathbf{x}) \pi(\mathbf{x})) + \frac{1}{2} \sum_{i,j \in [d]} \frac{\partial^2}{\partial \mathbf{x}_i \partial \mathbf{x}_j} [\boldsymbol{\sigma}(\mathbf{x}) \boldsymbol{\sigma}(\mathbf{x})^\top \pi(\mathbf{x})]_{ij},$$

and the conclusion follows because  $\frac{\partial}{\partial t} \pi_t = \mathcal{L}^* \pi_t$  by Kolmogorov's forward equation (12).  $\square$

<sup>7</sup>The adjoint to an operator  $P$  on a Hilbert space is denoted  $P^*$ , and satisfies  $\langle Pf, g \rangle = \langle f, P^*g \rangle$ . Here, the relevant Hilbert space is sufficiently regular functions on  $\mathbb{R}^d$ , and the inner product is  $\langle f, g \rangle = \int f(\mathbf{x})g(\mathbf{x})d\mathbf{x}$ .

<sup>8</sup>The decay lets us discard the "boundary term" which typically arises in integration by parts formulas.

At this point, it is helpful to do an example. Consider the Langevin dynamics (4), where  $\boldsymbol{\sigma}(\mathbf{x}) = \sqrt{2}\mathbf{I}_d$  for all  $\mathbf{x} \in \mathbb{R}^d$ , so that  $\boldsymbol{\sigma}(\mathbf{x})\boldsymbol{\sigma}(\mathbf{x})^\top = 2\mathbf{I}_d$ , so the formula (14) reads

$$\frac{\partial}{\partial t}\pi_t(\mathbf{x}) = \nabla \cdot (\nabla V(\mathbf{x})\pi_t(\mathbf{x})) + \Delta\pi_t(\mathbf{x}). \quad (16)$$

The contribution of the component  $\Delta\pi_t(\mathbf{x})$  is the *heat equation*, and by going back through the calculations, we can check that it arose because of the presence of Brownian motion  $\sqrt{2}d\mathbf{B}_t$  in (4). Intuitively, this component is because of the dissipation of heat, where each point contributes symmetrically to its surroundings. On the other hand, the  $\nabla \cdot (\nabla V(\mathbf{x})\pi_t(\mathbf{x}))$  term is the contribution of the drift, and combining characterizes the stationary distribution of the Langevin dynamics.

**Theorem 1** (Langevin diffusion stationarity). *Suppose that  $V : \mathbb{R}^d \rightarrow \mathbb{R}$  is differentiable and  $\int \exp(-V(x))dx < \infty$ . Then, a stationary distribution for the Langevin dynamics (4) is given by the density  $\pi^* : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$  satisfying  $\pi^* \propto \exp(-V)$ .*

*Proof.* By the definition of the stationary distribution, we have that  $\frac{\partial}{\partial t}\pi_t$  vanishes pointwise when  $\pi_t = \pi^*$ . By reparameterizing  $U := -\log \pi^*$ , i.e.,  $\pi^* \propto \exp(-U)$ , (16) gives that pointwise

$$0 = -\nabla \cdot (\nabla V(\mathbf{x})\pi^*(\mathbf{x}) + \nabla\pi^*(\mathbf{x})) = -\nabla \cdot ((\nabla V(\mathbf{x}) - \nabla U(\mathbf{x}))\pi^*(\mathbf{x})),$$

which is solved by setting  $U = V$  up to addition by a universal constant.  $\square$

In general, it is not necessarily the case that the zero vector field  $\mathbf{v}$  is the only solution to  $\nabla \cdot (\mathbf{v}\pi^*) = 0$  pointwise, but we will soon give conditions (in Section 5) under which  $\pi^* \propto \exp(-V)$  is the unique stationary distribution for the Langevin dynamics. We pause here to draw some analogies between spectral graph theory, introduced in Part XIV, and the developments of this section. Intuitively,  $\mathcal{L}$  plays the role of an infinitesimal transition operator from functions to functions, and its adjoint acts on probability densities. By solving Kolmogorov's forward equation (12), we should heuristically expect the distribution  $\pi_t$  to follow the law

$$\pi_t = \exp(t\mathcal{L}^*)\pi_0.$$

Advancing time for  $t = 1$  units, this means that the state evolution of  $\pi_t$  roughly follows  $\pi_{t+1} = \exp(\mathcal{L}^*)\pi_t \approx (\text{id} + \mathcal{L}^*)\pi_t$ , so  $\exp(\mathcal{L}^*)$  acts like the discrete-time transition operator. We also expect that as  $t \rightarrow \infty$ ,  $\pi_t \rightarrow \pi^*$  (the stationary distribution of the semigroup  $\{P_t\}_{t \geq 0}$ ), so if we view  $\mathcal{L}^*$  as an infinite-dimensional operator, it should have a leading eigenfunction of  $\pi^*$  with eigenvalue 0 (so its exponential has eigenvalue 1), and if  $\pi^*$  is unique, then all other eigenvalues of  $\mathcal{L}^*$  should be negative (so that they rapidly decay with time after exponentiation). To formalize this spectral characterization of  $\mathcal{L}$  (which in the spectral graph theory case only made sense when the transition operator was reversible), it is helpful to introduce the analogous definition for Markov semigroups.

**Definition 2.** *Let  $\{P_t\}_{t \geq 0}$  be a Markov semigroup with generator  $\mathcal{L}$ , and stationary distribution  $\pi^* \in \mathcal{P}(\mathbb{R}^d)$ . We say  $\{P_t\}_{t \geq 0}$  is reversible if for all  $f, g \in L_2(\pi^*)$ , following notation (8),*

$$\int \mathcal{L}f(\mathbf{x})g(\mathbf{x})\pi^*(\mathbf{x})d\mathbf{x} = \int \mathcal{L}g(\mathbf{x})f(\mathbf{x})\pi^*(\mathbf{x})d\mathbf{x}.$$

In the spectral graph theory case, the transition operator  $\mathbf{P} = \mathbf{D}_G^{-1}\mathbf{A}_G$  corresponded to a degree-normalized adjacency matrix of a graph  $G$ , so reversibility corresponded to  $\mathbf{D}_G\mathbf{P} = \mathbf{P}^\top\mathbf{D}_G$ , or alternatively  $\mathbf{D}_G^{\frac{1}{2}}\mathbf{P}\mathbf{D}_G^{-\frac{1}{2}}$  being a symmetric matrix. In this case, the stationary distribution was  $\propto \mathbf{d}_G$ , the degrees of  $G$ . Analogously, we can interpret the condition in Definition 2 as saying

$$“\langle \mathcal{L}f, \mathbf{diag}(\pi^*)g \rangle = \langle \mathbf{diag}(\pi^*)f, \mathcal{L}g \rangle \iff \mathcal{L}^*\mathbf{diag}(\pi^*) = \mathbf{diag}(\pi^*)\mathcal{L},”$$

which also justifies the definition being an integration with respect to  $\pi^*$ , since we need to conjugate the generator by  $\pi^*$  to make it truly self-adjoint. If we let  $f = \mathbf{1}_E$  and  $g = \mathbf{1}_{E'}$  be indicators of subsets  $E$  and  $E'$  of  $\mathbb{R}^d$ , integrating Definition 2 for time  $t$  implies that

$$\begin{aligned} \int P_t\mathbf{1}_E(\mathbf{x})\mathbf{1}_{E'}(\mathbf{x})\pi^*(\mathbf{x})d\mathbf{x} &= \int P_t\mathbf{1}_{E'}(\mathbf{x})\mathbf{1}_E(\mathbf{x})\pi^*(\mathbf{x})d\mathbf{x} \\ \implies \Pr_{\mathbf{x}_0 \sim \pi^*}[\mathbf{x}_0 \in E', \mathbf{x}_t \in E] &= \Pr_{\mathbf{x}_0 \sim \pi^*}[\mathbf{x}_t \in E', \mathbf{x}_0 \in E]. \end{aligned}$$

This holds for all events  $E, E'$ , so we can verify that  $(\mathbf{x}_0, \mathbf{x}_t)$  has the same joint distribution as  $(\mathbf{x}_t, \mathbf{x}_0)$ . This can be interpreted as the fact that the reversed-time Markov semigroup, i.e., that which follows (12) but flips time, induces the same evolutions as the original Markov semigroup.

Finally, we wish to formalize our earlier intuition of a “spectral gap” driving all functions orthogonal to  $\pi^*$  rapidly to zero, so that the transition operator only preserves the stationary distribution. By using Jensen’s inequality with the definition (7), observe that for all functions  $f \in L^2(\pi^*)$ ,

$$(P_t f(\mathbf{x}))^2 = \mathbb{E}[f(\mathbf{x}_t) \mid \mathbf{x}_0 = \mathbf{x}]^2 \leq \mathbb{E}[f(\mathbf{x}_t)^2 \mid \mathbf{x}_0 = \mathbf{x}] = P_t(f^2)(\mathbf{x}), \text{ for all } \mathbf{x} \in \mathbb{R}^d. \quad (17)$$

Moreover, the above inequality is tight only when  $f$  is a constant function, corresponding to  $\mathcal{L}$  having a constant right eigenfunction (just as graph transition operators had a right eigenfunction of  $\mathbf{1}_V$ ). Our hope is to establish a quantitative variant of (17), where the inequality is strict for all other functions. When  $\{P_t\}_{t \geq 0}$  is reversible with stationary distribution  $\pi^*$ , we introduce the *carré du champ* operator  $\Gamma$ , and its integral  $\mathcal{E}$  (called the Dirichlet energy), defined as follows:

$$\Gamma(f, g)(\mathbf{x}) := -f(\mathbf{x})\mathcal{L}g(\mathbf{x}), \quad \mathcal{E}(f, g) := \int \Gamma(f, g)(\mathbf{x})\pi^*(\mathbf{x})d\mathbf{x}, \quad (18)$$

so that the reversibility assumption in Definition 2 gives  $\mathcal{E}(f, g) = \mathcal{E}(g, f)$ . As we explore in Section 5, the question of a spectral gap in (17) is really about understanding the spectrum of the symmetric operator  $\mathcal{E}$ . To be more concrete, we conclude the section by working out  $\mathcal{E}$  explicitly when  $\{P_t\}_{t \geq 0}$  is the Markov semigroup of the Langevin diffusion (4).

**Lemma 3.** *Following notation (18), if  $\{P_t\}_{t \geq 0}$  is the Markov semigroup corresponding to the Langevin diffusion (4) and  $\int \exp(-V(\mathbf{x}))d\mathbf{x} < \infty$ , then for  $\pi^* \in \mathcal{P}(\mathbb{R}^d)$  satisfying  $\pi^* \propto \exp(-V)$ ,*

$$\mathcal{E}(f, g) = \int \langle \nabla f(\mathbf{x}), \nabla g(\mathbf{x}) \rangle \pi^*(\mathbf{x})d\mathbf{x} \text{ for all } f, g \in L^2(\pi^*).$$

*Proof.* Recall from the calculation (15), specialized to the Langevin diffusion (4), that

$$\begin{aligned} \mathcal{E}(f, g) &= - \int \mathcal{L}f(\mathbf{x})g(\mathbf{x})\pi^*(\mathbf{x})d\mathbf{x} \\ &= \int \langle \nabla f(\mathbf{x}), \nabla V(\mathbf{x}) \rangle g(\mathbf{x})\pi^*(\mathbf{x})d\mathbf{x} + \int \Delta f(\mathbf{x})g(\mathbf{x})\pi^*(\mathbf{x})d\mathbf{x} \\ &= \int \langle \nabla f(\mathbf{x}), \nabla V(\mathbf{x}) \rangle g(\mathbf{x})\pi^*(\mathbf{x})d\mathbf{x} + \int \langle \nabla f(\mathbf{x}), \nabla(g \cdot \pi^*)(\mathbf{x}) \rangle d\mathbf{x} \\ &= \int \langle \nabla f(\mathbf{x}), \nabla V(\mathbf{x}) \rangle g(\mathbf{x})\pi^*(\mathbf{x})d\mathbf{x} - \int \langle \nabla f(\mathbf{x}), \nabla V(\mathbf{x}) \rangle g(\mathbf{x})\pi^*(\mathbf{x})d\mathbf{x} \\ &\quad + \int \langle \nabla f(\mathbf{x}), \nabla g(\mathbf{x}) \rangle \pi^*(\mathbf{x})d\mathbf{x} = \int \langle \nabla f(\mathbf{x}), \nabla g(\mathbf{x}) \rangle \pi^*(\mathbf{x})d\mathbf{x}, \end{aligned}$$

as claimed. The third line above used integration by parts (13) (with  $\Delta f = \nabla \cdot (\nabla f)$ ), and the fourth line used that  $\pi^*(\mathbf{x}) = \frac{\exp(-V(\mathbf{x}))}{Z}$  for a universal normalizing constant  $Z$ , so that

$$-\nabla V(\mathbf{x})\pi^*(\mathbf{x}) = -\nabla V(\mathbf{x}) \cdot \frac{\exp(-V(\mathbf{x}))}{Z} = \nabla \left( \frac{\exp(-V(\mathbf{x}))}{Z} \right) = \nabla \pi^*(\mathbf{x}).$$

□

Notice in particular that when  $f$  or  $g$  is a constant function, then  $\mathcal{E}(f, g) = 0$  by Lemma 3, so the constant function is an eigenfunction of  $\mathcal{E}$ . The question of a spectral gap then corresponds to characterizing the eigenvalues of all eigenfunctions of  $\mathcal{E}$  orthogonal to the constant function (which is exactly the space where probability densities are allowed to evolve, since they must stay densities). Lemma 3 also justifies the name “carré du champ,” which is French for “square of a field.” Here, we mean the vector field  $\nabla f$ , which is relevant because  $\mathcal{E}(f, f) = \int \|\nabla f(\mathbf{x})\|_2^2 \pi^*(\mathbf{x})d\mathbf{x}$ . Further, Lemma 3 proves that  $\mathcal{E}(f, g) = \mathcal{E}(g, f)$ , so the Langevin diffusion is indeed reversible.

**Remark 3.** *Markov semigroup theory extends to the discrete setting as well, where a continuous-time Markov chain is governed by a matrix specifying transition probabilities for a continuous-time jump process. The continuous-time analog of the “transition matrix” is a family of matrices  $\{\exp(t\mathbf{L})\}_{t \geq 0}$  for a generator matrix  $\mathbf{L}$ , specifying transition probabilities after  $t$  time has passed.*



### 3 Optimal transport

In this section, we take a brief digression to introduce a notion of distance between probability measures in  $\mathcal{P}(\mathbb{R}^d)$ , which is important for measuring convergence rates of stochastic processes to their stationary distribution. We begin by defining the relevant distance we will consider.

**Definition 3** (Wasserstein distance). *Let  $\mu, \pi \in \mathcal{P}(\mathbb{R}^d)$  have finite second moments, i.e.,  $\mathbb{E}_{\mathbf{x} \sim \mu}[\|\mathbf{x}\|_2^2] < \infty, \mathbb{E}_{\mathbf{x} \sim \pi}[\|\mathbf{x}\|_2^2] < \infty$ . We define the Wasserstein distance between  $\mu$  and  $\pi$  by*

$$W_2(\mu, \pi) := \inf_{\gamma \in \mathcal{C}(\mu, \pi)} \sqrt{\int \|\mathbf{x} - \mathbf{y}\|_2^2 \gamma(\mathbf{x}, \mathbf{y}) d\mathbf{x}d\mathbf{y}}, \quad (19)$$

where  $\mathcal{C}(\mu, \pi)$  is the set of all couplings of  $\mu$  and  $\pi$  (see Fact 1, Part XI).

In other words, the squared distance  $W_2^2(\mu, \pi)$  is induced by the coupling  $\gamma$  which minimizes the expected squared distance between  $(x, y) \sim \gamma$ . More generally, there is also a family of  $p$ -Wasserstein distances  $W_p(\mu, \pi)$ , which generalize Definition 3. In fact, historically it was the  $W_1$  distance which was first considered by Monge in the 1700s (who only considered a restricted set of deterministic mappings  $\mathbf{x} \rightarrow \mathbf{y}$ ), and later Kantorovich during World War II (who introduced the more general coupling definition), who were interested in the transportation of resources. The study of optimization problems of the form in Definition 3 is hence called *optimal transport*. For an extended introduction to optimal transport and its applications, we recommend the text [Vil08]; for a somewhat briefer summary of some of the key technical tools, see [MG10].

One useful observation for our purposes is the fact that  $W_2$  induces a metric on

$$\mathcal{P}_2(\mathbb{R}^d) := \left\{ \pi \in \mathcal{P}(\mathbb{R}^d) \mid \int \|\mathbf{x}\|_2^2 \pi(\mathbf{x}) d\mathbf{x} < \infty \right\},$$

the set of densities with finite second moment. Hence, it yields a distance on  $\mathcal{P}_2(\mathbb{R}^d)$ .

**Lemma 4.** *For all  $\mu, \nu, \pi \in \mathcal{P}_2(\mathbb{R}^d)$ , we have  $W_2(\mu, \pi) \leq W_2(\mu, \nu) + W_2(\nu, \pi)$ .*

*Proof.* Let  $\gamma_{\mu, \nu} \in \mathcal{C}(\mu, \nu)$  be the optimal coupling of  $\mu, \nu$  realizing the value  $W_2(\mu, \nu)$ , and similarly let  $\gamma_{\nu, \pi} \in \mathcal{C}(\nu, \pi)$  realize the value  $W_2(\nu, \pi)$ . We first claim that there exists a joint distribution  $(\mathbf{x}, \mathbf{y}, \mathbf{z}) \sim \gamma_{\mu, \nu, \pi}$  on  $\mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R}^d$ , such that the marginal distribution of  $(\mathbf{x}, \mathbf{y})$  is  $\gamma_{\mu, \nu}$ , and the marginal distribution of  $(\mathbf{y}, \mathbf{z})$  is  $\gamma_{\nu, \pi}$ . To see this, we can first draw  $\mathbf{y} \sim \nu$ , and then draw  $\mathbf{x}$  and  $\mathbf{z}$  from the conditional distributions of  $\gamma_{\mu, \nu} \mid \mathbf{y}$  and  $\gamma_{\nu, \pi} \mid \mathbf{y}$ , respectively.

Next, the marginal distribution of  $(\mathbf{x}, \mathbf{z})$  for  $(\mathbf{x}, \mathbf{y}, \mathbf{z}) \sim \gamma_{\mu, \nu, \pi}$  is a valid coupling of  $\mu, \pi$ , so

$$\begin{aligned} W_2(\mu, \pi) &\leq \sqrt{\int \|\mathbf{x} - \mathbf{z}\|_2^2 \gamma_{\mu, \nu, \pi}(\mathbf{x}, \mathbf{y}, \mathbf{z}) d\mathbf{x}d\mathbf{y}d\mathbf{z}} \\ &\leq \sqrt{\int (\|\mathbf{x} - \mathbf{y}\|_2 + \|\mathbf{y} - \mathbf{z}\|_2)^2 \gamma_{\mu, \nu, \pi}(\mathbf{x}, \mathbf{y}, \mathbf{z}) d\mathbf{x}d\mathbf{y}d\mathbf{z}} \\ &\leq \sqrt{\int \|\mathbf{x} - \mathbf{y}\|_2^2 \gamma_{\mu, \nu, \pi}(\mathbf{x}, \mathbf{y}, \mathbf{z}) d\mathbf{x}d\mathbf{y}d\mathbf{z}} + \sqrt{\int \|\mathbf{y} - \mathbf{z}\|_2^2 \gamma_{\mu, \nu, \pi}(\mathbf{x}, \mathbf{y}, \mathbf{z}) d\mathbf{x}d\mathbf{y}d\mathbf{z}} \\ &= \sqrt{\int \|\mathbf{x} - \mathbf{y}\|_2^2 \gamma_{\mu, \nu}(\mathbf{x}, \mathbf{y}) d\mathbf{x}d\mathbf{y}} + \sqrt{\int \|\mathbf{y} - \mathbf{z}\|_2^2 \gamma_{\nu, \pi}(\mathbf{y}, \mathbf{z}) d\mathbf{y}d\mathbf{z}} = W_2(\mu, \nu) + W_2(\nu, \pi). \end{aligned}$$

The second line used the triangle inequality, and the inequality in the third line follows by squaring both sides and bounding cross-terms. The last line used our construction of  $\gamma_{\mu, \nu, \pi}$ .  $\square$

We mention that while Lemma 4 shows  $W_2$  obeys the triangle inequality, it is not necessarily the case that  $W_2(\mu, \pi) = 0 \implies \mu = \pi$ . However, this conclusion is true except for on a set of measure zero, so we will treat  $W_2$  as a genuine metric under this set of equivalence classes.

We conclude the section by introducing the fundamental theorem of optimal transport, and providing a brief proof sketch. This result gives a remarkable characterization of the optimal transport

coupling  $\gamma$  inducing  $W_2$ , and establishes that strong duality holds for the corresponding problem. Specifically, note that  $\mathcal{C}(\mu, \pi)$  is a convex subset of the probability measures on  $\mathbb{R}^d \times \mathbb{R}^d$  (since taking convex combinations does not affect marginals), and further

$$\frac{1}{2}W_2^2(\mu, \pi) = \inf_{\gamma \in \mathcal{C}(\mu, \pi)} \int \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 \gamma(\mathbf{x}, \mathbf{y}) d\mathbf{x}d\mathbf{y} \quad (20)$$

is a linear optimization problem over  $\mathcal{C}(\mu, \pi)$  (i.e., the objective above is linear in  $\gamma$ ). Therefore, we may consider its Lagrangian dual, which is derived via

$$\begin{aligned} & \inf_{\gamma \in \mathcal{C}(\mu, \pi)} \int \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 \gamma(\mathbf{x}, \mathbf{y}) d\mathbf{x}d\mathbf{y} \\ = & \inf_{\gamma \in \mathcal{P}_2(\mathbb{R}^d \times \mathbb{R}^d)} \sup_{\substack{f \in L^1(\mu) \\ g \in L^1(\pi)}} \int \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 \gamma(\mathbf{x}, \mathbf{y}) d\mathbf{x}d\mathbf{y} + \int f(\mathbf{x})\mu(\mathbf{x})d\mathbf{x} - \int f(\mathbf{x})\gamma(\mathbf{x}, \mathbf{y})d\mathbf{x}d\mathbf{y} \\ & + \int g(\mathbf{y})\pi(\mathbf{y})d\mathbf{y} - \int g(\mathbf{y})\gamma(\mathbf{x}, \mathbf{y})d\mathbf{x}d\mathbf{y}, \end{aligned}$$

where  $L^1(\pi)$  is the set of functions with finite expectation with respect to  $\pi$  (analogously to (8)). The role of  $f$  and  $g$  above is to enforce the marginal constraints in the unconstrained problem over  $\gamma$ , forcing  $\gamma \in \mathcal{C}(\mu, \pi)$  at optimality. Now, if strong duality holds, the above expression equates to

$$\sup_{(f, g) \in \mathcal{D}(\mu, \pi)} \int f(\mathbf{x})\mu(\mathbf{x})d\mathbf{x} + \int g(\mathbf{y})\pi(\mathbf{y})d\mathbf{y}, \quad (21)$$

$$\text{where } \mathcal{D}(\mu, \pi) := \left\{ (f, g) \in L^1(\mu) \times L^1(\pi) \mid f(\mathbf{x}) + g(\mathbf{y}) \leq \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 \text{ for all } \mathbf{x}, \mathbf{y} \in \mathbb{R}^d \right\}.$$

We think of  $\mathcal{D}(\mu, \pi)$  as the set of dual feasible potentials  $(f, g)$  inducing the dual optimization problem. Incredibly, under mild conditions strong duality does hold, and we also have a complete characterization of the optimal coupling  $\gamma \in \mathcal{C}(\mu, \pi)$  and the optimal dual potentials  $(f, g)$ .

**Proposition 4** (Brenier's theorem). *If  $\mu, \pi \in \mathcal{P}_2(\mathbb{R}^d)$ , the values of (20) and (21) are equal (i.e., strong duality holds), and both values are realized, respectively by  $\gamma^* \in \mathcal{C}(\mu, \pi)$  and  $(f^*, g^*) \in \mathcal{D}(\mu, \pi)$ . Moreover, there exists a unique convex function  $\varphi : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  such that*

$$f^*(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|_2^2 - \varphi(\mathbf{x}), \quad g^*(\mathbf{y}) = \frac{1}{2} \|\mathbf{y}\|_2^2 - \varphi^*(\mathbf{y}),$$

and the optimal coupling  $\gamma^*$  is supported only on points  $(\mathbf{x}, \mathbf{y})$  where  $\mathbf{y} \in \partial\varphi(\mathbf{x})$ .

We remark that Alexandrov's theorem tells us that convex functions are differentiable almost everywhere, so the support is almost surely of the form  $(\mathbf{x}, \nabla\varphi(\mathbf{x}))$ . In other words, the optimal transport plan from  $\mu$  to  $\pi$  realizing  $W_2^2$  is actually a unique deterministic mapping  $\nabla\varphi$ , which sends  $\mathbf{x} \sim \mu$  to  $\nabla\varphi(\mathbf{x}) \sim \pi$ , except on a measure-zero subset. As a sanity check, we indeed have

$$\begin{aligned} f^*(\mathbf{x}) + g^*(\mathbf{y}) &= \frac{1}{2} \|\mathbf{x}\|_2^2 - \varphi(\mathbf{x}) + \frac{1}{2} \|\mathbf{y}\|_2^2 - \varphi^*(\mathbf{y}) \\ &= \frac{1}{2} \|\mathbf{x}\|_2^2 + \frac{1}{2} \|\mathbf{y}\|_2^2 - \varphi(\mathbf{x}) - \varphi^*(\nabla\varphi(\mathbf{x})) = \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2, \end{aligned}$$

for all  $(\mathbf{x}, \mathbf{y}) = (\mathbf{x}, \nabla\varphi(\mathbf{x}))$  in the support of  $\gamma^*$ , where we recall from Corollary 1, Part III, that

$$\varphi(\mathbf{x}) + \varphi^*(\nabla\varphi(\mathbf{x})) = \langle \nabla\varphi(\mathbf{x}), \mathbf{x} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle.$$

Therefore, all of the constraints in  $\mathcal{D}(\mu, \pi)$  are tight for the optimal dual potentials  $(f^*, g^*)$ . The induced convex function  $\varphi$  is often called a Brenier potential, and much of Proposition 4 was established by the works [Bre87, Bre91]. The proof of Proposition 4 is a bit tedious for this brief exposition; a summary of it can be found in Section 1.3 of [Che24]. However, we mention one crucial step of the proof which may be of broader interest to the reader.

**Fact 1** ([Roc70]). *If  $S \subseteq \mathbb{R}^d \times \mathbb{R}^d$  is cyclically monotone, i.e., for all  $n \in \mathbb{N}$  and all permutations  $\sigma : [n] \rightarrow [n]$ , we have for all  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i \in [n]} \subseteq S$ ,*

$$\sum_{i \in [n]} \langle \mathbf{x}_i, \mathbf{y}_i \rangle \geq \sum_{i \in [n]} \langle \mathbf{x}_i, \mathbf{y}_{\sigma(i)} \rangle, \quad (22)$$

then there exists a convex function  $\varphi : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  such that

$$S \subseteq \{(\mathbf{x}, \mathbf{y}) \mid \mathbf{x} \in \mathbb{R}^d, \mathbf{y} \in \partial\varphi(\mathbf{x})\}.$$

One can check that pairs  $(\mathbf{x}, \partial\varphi(\mathbf{x}))$  do indeed induce cyclically monotone sets, which follows from monotonicity of convex gradients (see Definition 2, Part IV). The characterization in Proposition 4 then follows because optimal transport plans induce cyclically monotone subsets. Indeed, if (22) did not hold for some pairs  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i \in [n]}$  in the support of the transport plan, then we could improve the  $W_2^2$  objective by pairing up the transported points differently:

$$\sum_{i \in [n]} \langle \mathbf{x}_i, \mathbf{y}_i \rangle \leq \sum_{i \in [n]} \langle \mathbf{x}_i, \mathbf{y}_{\sigma(i)} \rangle \implies \sum_{i \in [n]} \|\mathbf{x}_i - \mathbf{y}_i\|_2^2 \geq \sum_{i \in [n]} \|\mathbf{x}_i - \mathbf{y}_{\sigma(i)}\|_2^2.$$

**Remark 4** (Caffarelli's contraction theorem). *Under further structural assumptions on the measures  $\mu, \pi$  in Definition 3, one can say more about the optimal transport map  $\partial\varphi$  described in Proposition 4. For example, a famous result by [Caf00] states that if  $\mu \propto \exp(-V)$  and  $\pi \propto \exp(-W)$  where  $V$  is  $L$ -smooth<sup>9</sup> and  $W$  is  $m$ -strongly convex, then the optimal transport map is induced by a differentiable convex function  $\varphi$  which is  $\sqrt{L/\mu}$ -smooth. For a short proof of this fact and a generalization to entropic variants of optimal transport, see [CP23].*

We conclude this section with our first convergence analysis for the Langevin dynamics, albeit in continuous time (which does not give an implementable algorithm). We will later analyze a discretized version of the following result. Before stating our claim, we require one definition.

**Definition 4** (Strong logconcavity). *We say  $\pi \in \mathcal{P}(\mathbb{R}^d)$  is  $\mu$ -strongly logconcave if  $\pi \propto \exp(-V)$  for  $V : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ , such that  $V$  is  $\mu$ -strongly convex.*

For instance, the multivariate normal distribution  $\mathcal{N}(\mathbf{m}, \Sigma)$  for  $\mathbf{m} \in \mathbb{R}^d$ ,  $\Sigma \in \mathbb{S}_{>0}^{d \times d}$ , whose density is  $\propto \exp(-\frac{1}{2} \|\cdot - \mathbf{m}\|_{\Sigma^{-1}}^2)$ , is  $\mu$ -strongly logconcave for any value of  $\mu > 0$  with  $\Sigma^{-1} \succeq \mu \mathbf{I}_d$ .

**Theorem 2** (Wasserstein convergence of Langevin dynamics). *Let  $\{\mathbf{x}_t\}_{t \geq 0}$  follow the Langevin dynamics (4) with stationary distribution  $\pi^*$ , and assume that  $\pi^*$  is  $\mu$ -strongly logconcave. For all  $t \geq 0$ , the law of  $\mathbf{x}_t$ , denoted  $\pi_t \in \mathcal{P}(\mathbb{R}^d)$ , satisfies*

$$W_2^2(\pi_t, \pi^*) \leq \exp(-2\mu t) W_2^2(\pi_0, \pi^*).$$

*Proof.* Let  $\gamma_t$  be a coupling of  $(\pi_t, \pi^*)$  defined as follows. Let  $\gamma_0^*$  be the coupling of  $(\pi_0, \pi^*)$  which realizes  $W_2^2(\pi_0, \pi^*)$ , draw  $(\mathbf{x}_0, \mathbf{x}_0^*) \sim \gamma_0^*$ , and advance both points using the Langevin dynamics (4), sharing the same copy of Brownian motion  $\{\mathbf{B}_s\}_{0 \leq s \leq t}$ . It is clear that the marginals of  $\gamma_t$  are  $(\pi_t, \pi^*)$  respectively (the latter because  $\pi^*$  is stationary for (4)). Moreover, letting  $\{(\mathbf{x}_s, \mathbf{x}_s^*)\}_{s \in [0, t]}$  be advanced through (4) using the same copy of Brownian motion, so  $(\mathbf{x}_t, \mathbf{x}_t^*) \sim \gamma_t$ ,

$$\begin{aligned} \frac{d}{ds} \langle \mathbf{x}_s - \mathbf{x}_s^* \rangle &= \nabla V(\mathbf{x}_s^*) - \nabla V(\mathbf{x}_s) \\ \implies \frac{d}{ds} \|\mathbf{x}_s - \mathbf{x}_s^*\|_2^2 &= -2 \langle \nabla V(\mathbf{x}_s) - \nabla V(\mathbf{x}_s^*), \mathbf{x}_s - \mathbf{x}_s^* \rangle \leq -2\mu \|\mathbf{x}_s - \mathbf{x}_s^*\|_2^2 \\ \implies \|\mathbf{x}_t - \mathbf{x}_t^*\|_2^2 &\leq \exp(-2\mu t) \|\mathbf{x}_0 - \mathbf{x}_0^*\|_2^2. \end{aligned}$$

Above, the second inequality used Gronwall's inequality (Fact 1, Part II), and the first inequality used that strong convexity of  $V$  implies (e.g., by adding Eq. (9), Part II with  $x, x'$  interchanged)

$$\langle \nabla V(\mathbf{x}) - \nabla V(\mathbf{x}'), \mathbf{x} - \mathbf{x}' \rangle \geq \mu \|\mathbf{x} - \mathbf{x}'\|_2^2 \text{ for all } \mathbf{x}, \mathbf{x}' \in \text{dom}(V).$$

Finally, because  $\gamma_t$  is a coupling of  $(\pi_t, \pi^*)$  (so it has no better objective value for (19) than the optimal coupling), we conclude that

$$\begin{aligned} W_2^2(\pi_t, \pi^*) &\leq \mathbb{E}_{(\mathbf{x}_t, \mathbf{x}_t^*) \sim \gamma_t} \left[ \|\mathbf{x}_t - \mathbf{x}_t^*\|_2^2 \right] \\ &\leq \exp(-2\mu t) \mathbb{E}_{(\mathbf{x}_0, \mathbf{x}_0^*) \sim \gamma_0^*} \left[ \|\mathbf{x}_0 - \mathbf{x}_0^*\|_2^2 \right] = \exp(-2\mu t) W_2^2(\pi_0, \pi^*). \end{aligned}$$

□

<sup>9</sup>Here we mean smoothness in the sense of Definition 3, Part II.

We also give a basic initialization strategy under strong logconcavity. The following result shows that initializing at  $\pi_0$  set to a point mass at  $\mathbf{x}^*$  achieves bounded  $W_2^2(\pi_0, \pi^*)$ . The proof is tolerant to approximate minimizers; finding such a point typically does not dominate the cost of running a sampler (as state-of-the-art optimization rates are better than their sampling counterparts).

**Lemma 5.** *Let  $\pi^* \in \mathcal{P}(\mathbb{R}^d)$  satisfy  $\pi^* \propto \exp(-V)$ , and suppose  $V : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  is  $\mu$ -strongly convex and minimized at  $\mathbf{x}^*$ . Then,*

$$\mathbb{E}_{\mathbf{x} \sim \pi^*} \left[ \|\mathbf{x} - \mathbf{x}^*\|_2^2 \right] \leq \frac{2d}{\mu}.$$

*Proof.* Shifting  $V$  by a constant so  $\int \exp(-V(\mathbf{x}))d\mathbf{x} = 1$ , using  $\langle \nabla V(\mathbf{x}), \mathbf{x} - \mathbf{x}^* \rangle \geq \frac{\mu}{2} \|\mathbf{x} - \mathbf{x}^*\|_2^2$  which follows from strong logconcavity where  $V(\mathbf{x}) < \infty$ ,<sup>10</sup> and integrating by parts (using (13)),

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim \pi^*} \left[ \|\mathbf{x} - \mathbf{x}^*\|_2^2 \right] &= \int \|\mathbf{x} - \mathbf{x}^*\|_2^2 \exp(-V(\mathbf{x})) d\mathbf{x} \\ &\leq \frac{2}{\mu} \int \langle \nabla V(\mathbf{x}), \mathbf{x} - \mathbf{x}^* \rangle \exp(-V(\mathbf{x})) d\mathbf{x} \\ &= \frac{2}{\mu} \int (\nabla \cdot (\cdot - \mathbf{x}^*))(\mathbf{x}) \exp(-V(\mathbf{x})) d\mathbf{x} \\ &= \frac{2}{\mu} \int \text{Tr}(\mathbf{I}_d) \exp(-V(\mathbf{x})) d\mathbf{x} = \frac{2d}{\mu}. \end{aligned}$$

□

Theorem 2 has the appealing property of giving a linear rate of convergence, but the convergence guarantee is in the Wasserstein distance  $W_2$ , which may be less flexible in downstream applications. For instance,  $W_2$  is typically incomparable to the total variation distance  $D_{\text{TV}}$ , and when designing sampling algorithms to be used as subroutines (i.e., to be called multiple times),  $D_{\text{TV}}$  guarantees are preferable because they compose under the union bound. In Section 5, we introduce techniques for proving convergence of the Langevin dynamics and other stochastic differential equations in stronger error metrics such as the KL divergence (which implies bounds on  $D_{\text{TV}}$  by Pinsker’s inequality, as well as bounds on  $W_2$  under appropriate conditions discussed in Section 5).

## 4 Probability densities as a Riemannian manifold

### 4.1 Riemannian manifolds

In this section, we build upon our development of optimal transport to present a view of the space of square-integrable probability densities  $\mathcal{P}_2(\mathbb{R}^d)$  as a differentiable (in particular, Riemannian) manifold. It is first helpful to give some brief description of what differentiable manifolds are, and the sorts of calculations we can expect to arise from manipulating them. We informally do so here, focusing on presenting only the relevant material in a way that motivates the calculations specific to  $\mathcal{P}_2(\mathbb{R}^d)$ . For a more extended discussion on these topics, we suggest [Vis18] as a resource for learning about Riemannian manifolds in a way that is accessible to a computer science audience, as well as [dC92] for a more rigorous treatment from a mathematical perspective.

Roughly speaking, a manifold  $\mathcal{M}$  is a topological space such that every  $p \in \mathcal{M}$  has a open neighborhood  $U_p$  that “looks like” Euclidean space. Formally, in the finite-dimensional setting this means that there is a homeomorphism between  $U_p$  and an open set in  $\mathbb{R}^k$ . In many applications, this is interesting because  $\mathcal{M}$  is actually a subset of a higher-dimensional Euclidean space (e.g.,  $\mathbb{R}^d$  for  $d \geq k$ ). A canonical example is the surface of the unit sphere in  $\mathbb{R}^3$ , which is a 2-dimensional manifold (because neighborhoods of points are homeomorphic to balls in  $\mathbb{R}^2$ ). This definition extends straightforwardly to Hilbert spaces, which are inner product spaces (possibly of infinite dimension) that generalize Euclidean space. For instance,  $L^2(\mathbb{R}^d)$ , the set of all square-integrable functions in  $\mathbb{R}^d$  (i.e.,  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  such that  $\int f(\mathbf{x})^2 d\mathbf{x} < \infty$ ) is a Hilbert space equipped with the inner product

$$\langle f, g \rangle := \int f(\mathbf{x})g(\mathbf{x})d\mathbf{x}.$$

<sup>10</sup>When  $V : \mathbb{R}^d \rightarrow \mathbb{R}$  is unconstrained so that  $\nabla V(\mathbf{x}^*) = \mathbf{0}_d$  by first-order optimality, the constant factor can be sharpened by using  $\langle \nabla V(\mathbf{x}) - \nabla V(\mathbf{x}^*), \mathbf{x} - \mathbf{x}^* \rangle \geq \mu \|\mathbf{x} - \mathbf{x}^*\|_2^2$ .

In particular, note that  $\mathcal{P}_2(\mathbb{R}^d)$  is a subset of the space of square-integrable functions, and indeed we will view it as a manifold embedded in  $L^2(\mathbb{R}^d)$ . We now mention some definitions tailored to a specific family of manifolds which is convenient for calculation purposes, namely *Riemannian manifolds*. In a Riemannian manifold  $\mathcal{M}$ , every point  $p \in \mathcal{M}$  has an associated *tangent space*  $T_p\mathcal{M}$ . Informally, this should be viewed as the set of all possible “local velocities” of curves on  $\mathcal{M}$  passing through  $p$ . Notice that we treat points on  $\mathcal{M}$  (which should be thought of as representing position) differently than elements of a tangent space  $T_p\mathcal{M}$  (which should be thought of as representing velocities). Also, in general elements of  $T_p\mathcal{M}$  and  $T_q\mathcal{M}$  for  $p, q \in \mathcal{M}$ ,  $p \neq q$  are not directly comparable (e.g., we cannot simply add or subtract them directly, since they live in different spaces). To compare them, we should first transport  $T_p\mathcal{M}$  to  $T_q\mathcal{M}$  in an appropriate way along the manifold, which is done using a construction called the *Levi-Civita connection*. We will not explicitly describe this connection in the remainder of the lecture for brevity’s sake, but mention this caveat to caution the reader when performing calculations on manifolds.

There are two main properties of Riemannian manifolds that are convenient for our purposes. The first is that every point  $p \in \mathcal{M}$  comes with a local metric  $g_p$  which can be used to perform calculations specific to  $T_p\mathcal{M}$ . Specifically, there is a way to assign values  $g_p(u, v)$  such that  $g_p(u, v) = g_p(v, u)$  for all  $u, v \in T_p\mathcal{M}$ , and  $g_p(\cdot, w)$  is linear in its first argument for all  $w \in T_p\mathcal{M}$ . To help build intuition for this, when we view  $\mathbb{R}^d$  as a trivial manifold in  $\mathbb{R}^d$ ,  $g_p$  is simply the standard Euclidean inner product  $\langle \cdot, \cdot \rangle$  pointwise, i.e., it does not change depending on  $p \in \mathbb{R}^d$  (because the space is not “curved”). Another interesting example is the *Hessian manifold* of a self-concordant barrier function. For a subset  $\mathcal{K} \subseteq \mathbb{R}^d$  equipped with a self-concordant barrier function  $\varphi : \mathcal{K} \rightarrow \mathbb{R}$  (see Definition 2, Part X), we can define associated local metrics satisfying

$$g_{\mathbf{p}}(\mathbf{u}, \mathbf{v}) := \nabla^2\varphi(\mathbf{p})[u, v] \text{ for all } \mathbf{p} \in \mathcal{K}, \mathbf{u}, \mathbf{v} \in \mathbb{R}^d.$$

Note that in this example, each tangent space  $T_{\mathbf{p}}\mathcal{M}$  is  $d$ -dimensional (i.e., it is not actually lower-dimensional than its ambient space). The role of defining the local metrics  $g_{\mathbf{p}}$  above is to “curve space” in a way that can helpfully guide algorithm design, e.g., the Newton’s method we analyzed in Part X is just a discretization of a standard gradient flow from the perspective of the local metrics, which reweight the different directions in  $T_{\mathbf{p}}\mathcal{M}$  using  $\nabla^2\varphi(\mathbf{p})$ .<sup>11</sup>

We can correspondingly measure the “length” and “correlation” of elements in  $T_p\mathcal{M}$  by defining the norm and local inner product,

$$\|v\|_p := \sqrt{g_p(v, v)}, \quad \langle u, v \rangle_p := g_p(u, v), \text{ for all } u, v \in T_p\mathcal{M}. \quad (23)$$

The second main property we need is that there is a meaningful notion of differentiability associated with functions on the manifold,  $\mathcal{F} : \mathcal{M} \rightarrow \mathbb{R}$ . We begin by defining geodesics, which informally are length-minimizing curves along  $\mathcal{M}$  which do not leave the manifold. Concretely, given two points  $p, q \in \mathcal{M}$ , we denote their distance along the manifold by  $d_{\mathcal{M}}(p, q)$ , which is defined by

$$d_{\mathcal{M}}(p, q) := \inf_{\substack{\gamma : [0, 1] \rightarrow \mathcal{M} \\ \gamma(0) = p, \gamma(1) = q}} \left\{ \int_0^1 \|\dot{\gamma}(t)\|_{\gamma(t)} dt \right\}. \quad (24)$$

We pause to explain the above formula. The infimum is taken with respect to all curves  $\gamma : [0, 1] \rightarrow \mathcal{M}$ , which travel for one unit of time along  $\mathcal{M}$ , starting from  $\gamma(0) = p$  and ending at  $\gamma(1) = q$ . Moreover,  $\dot{\gamma}(t) \in T_{\gamma(t)}\mathcal{M}$  is the velocity vector of the curve at time  $t$ , whose length is measured in the local metric  $\|\cdot\|_{\gamma(t)}$ . So, the formula (24) measures the total magnitude of the distance accumulated over the curve, a meaningful generalization of the “length” of a curve. When the infimum is achieved, we call the length-minimizing curve  $\gamma$  which realizes the value of (24) the *geodesic* between  $p$  and  $q$ . Under mild conditions, geodesics actually have constant speed  $\|\dot{\gamma}(t)\|_{\gamma(t)}$  along their trajectories, so an equivalent definition is

$$d_{\mathcal{M}}(p, q) := \inf_{\substack{\gamma : [0, 1] \rightarrow \mathcal{M} \\ \gamma(0) = p, \gamma(1) = q}} \left\{ \sqrt{\int_0^1 \|\dot{\gamma}(t)\|_{\gamma(t)}^2 dt} \right\}. \quad (25)$$

<sup>11</sup>A helpful example to consider is  $\mathcal{M} = \mathbb{R}_{>0}^d$ , equipped with the self-concordant barrier function  $\varphi(\mathbf{p}) := -\sum_{i \in [d]} \log \mathbf{p}_i$ . In this example,  $\nabla^2\varphi(\mathbf{p}) = \mathbf{diag}(\frac{1}{\mathbf{p}_i})$ , so as  $\mathbf{p}_i$  approaches any boundary of the orthant for any coordinate  $i \in [d]$ , space in  $\mathcal{M}$  is curved in a way such that the weight assigned to that coordinate blows up.

For instance, geodesics in  $\mathbb{R}^d$  are just straight-line paths (which clearly have a constant-speed parameterization, i.e., move along the line at a fixed speed), and geodesics along the unit sphere example mentioned earlier follow arcs in the way that one would expect.

We can now finally define our notion of differentiability. Given a function  $\mathcal{F} : \mathcal{M} \rightarrow \mathbb{R}$ , the gradient of  $\mathcal{F}$  at  $p \in \mathcal{M}$  is denoted  $\nabla_{\mathcal{M}}\mathcal{F}(p) \in T_p\mathcal{M}$ , and is the element of  $T_p\mathcal{M}$  such that for all curves  $\{\gamma(t)\}_{t \in \mathbb{R}}$  passing through  $p$  at time  $t = 0$ , i.e. with  $\gamma(0) = p$ , it holds that

$$\frac{\partial}{\partial t}\mathcal{F}(\gamma(t)) \Big|_{t=0} = \langle \nabla_{\mathcal{M}}\mathcal{F}(p), \dot{\gamma}(0) \rangle_p. \quad (26)$$

Namely, this is just the generalization of the chain rule in  $\mathbb{R}^d$ : if  $x(t)$  is a smooth curve parameterized by time  $t \in \mathbb{R}$ , then we have  $\frac{d}{dt}f(x(t)) = \langle \nabla f(x(t)), \frac{d}{dt}x(t) \rangle$ . The formula (26) extends this definition to hold for all curves passing through  $p$ , which may not be fully comparable (due to their different trajectories) except locally at  $p$ . In the remainder of the section, we focus on describing how to do these calculations for an appropriate Riemannian manifold associated with  $\mathcal{P}_2(\mathbb{R}^d)$ .

## 4.2 Wasserstein space

The punchline of the above digression, and the developments of Section 3, is that  $\mathcal{P}_2(\mathbb{R}^d)$  can naturally be viewed as a manifold  $\mathcal{M}$ , equipped with the distance function (see (24))

$$d_{\mathcal{M}}(\mu, \pi) := W_2(\mu, \pi), \text{ for all } \mu, \pi \in \mathcal{M} := \mathcal{P}_2(\mathbb{R}^d). \quad (27)$$

We call the manifold  $\mathcal{P}_2(\mathbb{R}^d)$  (referred to as  $\mathcal{M}$  for the remainder of this section for brevity), equipped with the distance  $W_2$ , *Wasserstein space* for short. Importantly, there are succinct characterizations of tangent spaces, local metrics, geodesics, and gradients of functionals on  $\mathcal{M}$ . The goal of this section is to state these characterizations, as well as some rough sketches for their proofs. For a more extended proof overview, we defer to Section 1.3 of [Che24], and we defer formal derivations of these results to the excellent resource [AGS08].

To introduce the Riemannian structure of  $\mathcal{M}$ , it is helpful to first adopt a dual view on curves along Wasserstein space, just as we did in Section 2. In particular, in Section 2 we showed that every *particle flow* governed by an SDE on particles in  $\mathbb{R}^d$  induces a deterministic *density flow* along  $\mathcal{M}$ . In the case of Section 2, we considered stochastic particle flows, but this connection works fine for deterministic particle flows as well,<sup>12</sup> summarized in the following fundamental result.

**Lemma 6** (Continuity equation). *Let  $\{\mathbf{v}_t\}_{t \in \mathbb{R}}$  be a family of vector fields on  $\mathbb{R}^d$ , i.e., for all  $t \in \mathbb{R}$ ,  $\mathbf{v}_t : \mathbb{R}^d \rightarrow \mathbb{R}^d$  is a corresponding vector field. Let  $\mathbf{x}_0 \sim \pi_0$  and suppose that  $\mathbf{x}_t$  follows the ODE  $\frac{d}{dt}\mathbf{x}_t = \mathbf{v}_t(\mathbf{x}_t)$  for all  $t \geq 0$ . Then  $\pi_t$ , the law of  $\mathbf{x}_t$ , follows the following PDE pointwise on  $\mathbb{R}^d$ :*

$$\frac{\partial}{\partial t}\pi_t + \nabla \cdot (\pi_t \mathbf{v}_t) = 0. \quad (28)$$

*Proof.* For all test functions  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ ,

$$\begin{aligned} \int f(\mathbf{x}) \frac{\partial}{\partial t} \pi_t(\mathbf{x}) d\mathbf{x} &= \frac{\partial}{\partial t} \left( \int f(\mathbf{x}) \pi_t(\mathbf{x}) d\mathbf{x} \right) = \frac{\partial}{\partial t} \mathbb{E}_{\mathbf{x} \sim \pi_t} [f(\mathbf{x})] \\ &= \mathbb{E} \left\langle \nabla f(\mathbf{x}_t), \frac{d}{dt} \mathbf{x}_t \right\rangle = \mathbb{E} \langle \nabla f(\mathbf{x}_t), \mathbf{v}_t(\mathbf{x}_t) \rangle \\ &= \int \langle \nabla f(\mathbf{x}), \mathbf{v}_t(\mathbf{x}) \rangle \pi_t(\mathbf{x}) d\mathbf{x} = - \int f(\mathbf{x}) \nabla \cdot (\mathbf{v}_t(\mathbf{x}) \pi_t(\mathbf{x})) d\mathbf{x}, \end{aligned} \quad (29)$$

where we used integration by parts (13) in the last line. As this holds for all  $f$ , the claim follows.  $\square$

**Tangent spaces.** Lemma 6 suggests that if the connection between vector fields in  $\mathbb{R}^d$  (corresponding to particle flows) and curves on  $\mathcal{M}$  (corresponding to density flows) goes both ways, then every geodesic on  $\mathcal{M}$  is interpretable by studying a corresponding family of vector fields which

<sup>12</sup>It is unsurprising in light of Proposition 4, which shows that Wasserstein distances are induced by deterministic maps in particle space, that the most interesting particle flows from a geodesic perspective are deterministic.

induces it. This viewpoint turns out to be formalizable, and in fact the relevant vector fields turn out to be those induced by gradients of functions. More precisely, for  $\pi \in \mathcal{M}$ , the tangent space is

$$T_\pi \mathcal{M} = \left\{ \nabla \cdot (\pi \nabla \psi) \mid \nabla \psi \in \overline{\left\{ \nabla \psi \mid \psi : \mathbb{R}^d \rightarrow \mathbb{R} \text{ is smooth, } \int \|\nabla \psi(\mathbf{x})\|_2^2 \pi(\mathbf{x}) d\mathbf{x} < \infty \right\}} \right\}, \quad (30)$$

where  $\overline{S}$  denotes the closure of a set  $S$ .<sup>13</sup> In other words, when moving along the manifold  $\mathcal{M}$  (following a trajectory in density space), the resulting movement in particle space is induced by a vector field given by the gradient of a smooth function, or a limit of such gradients. Moreover, each curve  $\{\pi_t\}_{t \in \mathbb{R}} \in \mathcal{M}$  has a time derivative  $\frac{\partial}{\partial t} \pi_t$  given by a continuity equation (28), parameterized by a family of vector fields  $\{\mathbf{v}_t\}_{t \geq 0}$ , such that  $\mathbf{v}_t = \nabla \psi_t$  for a function  $\psi_t : \mathbb{R}^d \rightarrow \mathbb{R}$  almost surely.

We have already seen two pieces of evidence suggesting this should be the case. First, vector fields in Wasserstein space yield  $W_2$  distances between probability distributions (27), and we know that “shortest paths” (i.e., optimal transport plans) are given by gradient vector fields from Brenier’s theorem, Proposition 4. Second, when performing the computation in (29) to derive the continuity equation, the vector field  $\mathbf{v}_t$  only interacts with gradients of functions  $\nabla f$ . So, projecting  $\mathbf{v}_t$  into the subspace of function gradients<sup>14</sup> decreases its norm measured in the local metric, without affecting the continuity equation (28). To give an example, we showed that the Langevin dynamics (4) induce a curve in Wasserstein space given by the Fokker-Planck equation (16), i.e.,  $\frac{\partial}{\partial t} \pi_t = \nabla \cdot (\pi_t \nabla V) + \Delta \pi_t$  pointwise. This is consistent with (30), as it gives the continuity equation

$$\begin{aligned} \frac{\partial}{\partial t} \pi_t &= \nabla \cdot (\pi_t \nabla V + \nabla \pi_t) = \nabla \cdot (\pi_t (\nabla V + \nabla \log \pi_t)) \\ &= \nabla \cdot \left( \pi_t \nabla \log \frac{\pi_t}{\pi^*} \right) \in T_{\pi_t} \mathcal{M}. \end{aligned} \quad (31)$$

We now mention some additional calculations relevant to Wasserstein space.

**Local metrics.** First, the way to define the local metric (following notation (23)) between  $\nabla \cdot (\pi \nabla \psi), \nabla \cdot (\pi \nabla \varphi) \in T_\pi \mathcal{M}$  for smooth functions  $\psi, \varphi : \mathbb{R}^d \rightarrow \mathbb{R}$  is the formula

$$\langle \nabla \cdot (\pi \nabla \psi), \nabla \cdot (\pi \nabla \varphi) \rangle_\pi := \int \langle \nabla \psi(\mathbf{x}), \nabla \varphi(\mathbf{x}) \rangle \pi(\mathbf{x}) d\mathbf{x}. \quad (32)$$

For elements of  $T_\pi \mathcal{M}$  which are obtained by taking the limit of smooth gradients, we similarly define their inner product in the local metric by taking limits of the above formula.

**Geodesics.** For two probability densities in Wasserstein space,  $\pi, \mu \in \mathcal{M}$ , we can characterize the geodesic  $\gamma : [0, 1] \rightarrow \mathcal{M}$  joining  $\pi = \gamma(0)$  and  $\mu = \gamma(1)$  as follows. Let  $(\mathbf{x}_0, \mathbf{x}_1)$  be drawn from the optimal coupling  $c$  inducing  $W_2^2(\pi, \mu)$ , which exists and is unique by Proposition 4. Then, the density  $\gamma(t) \in \mathcal{M}$  is given by the law of  $(1-t)\mathbf{x}_0 + t\mathbf{x}_1$ . This result states that the geodesic  $\gamma$  is (in particle space) traced out by the straight-line interpolation within the optimal coupling, and a simple calculation using the formulas (25), (32) shows that indeed,

$$d_{\mathcal{M}}(\pi, \mu) = \sqrt{\int \|\mathbf{x}_0 - \mathbf{x}_1\|_2^2 c(\mathbf{x}_0, \mathbf{x}_1) d\mathbf{x}_0 d\mathbf{x}_1} = W_2(\pi, \mu),$$

as claimed. This geodesic is also called McCann’s displacement interpolation.

**Gradients.** Finally, for a functional  $\mathcal{F} : \mathcal{M} \rightarrow \mathbb{R}$ ,<sup>15</sup> we compute its Wasserstein gradient  $\nabla_{\mathcal{M}} \mathcal{F}(\pi)$  (i.e., its gradient over the differentiable manifold  $\mathcal{M}$  in the sense of (26)), at a point  $\pi \in \mathcal{M}$ , as follows. Let  $\{\pi_t\}_{t \in \mathbb{R}} \subset \mathcal{M}$  be an arbitrary smooth curve satisfying  $\pi_0 = \pi$ . By the discussion following Lemma 6, there is a family of functions  $\{\psi_t : \mathbb{R}^d \rightarrow \mathbb{R}\}_{t \in \mathbb{R}}$ , such that

$$\frac{\partial}{\partial t} \pi_t = \nabla \cdot (\pi_t \nabla \psi_t), \text{ for all } t \in \mathbb{R}. \quad (33)$$

<sup>13</sup>Formally, the closure here is with respect to the vector topology induced by the Hilbert space  $L^2(\pi)$ .

<sup>14</sup>Square-integrable gradients of smooth functions form a subspace within the Hilbert space  $L^2(\pi)$ , because this space is closed under linear combinations (i.e.,  $a\nabla f + b\nabla g$  is the gradient of  $af + bg$ ).

<sup>15</sup>Following the literature, we call functions defined over  $\mathcal{M}$  “functionals,” as they are functions of functions. Indeed, points in  $\mathcal{M}$  are probability densities, which are functions from  $\mathbb{R}^d \rightarrow \mathbb{R}$ .

Moreover, the element  $\nabla_{\mathcal{M}}\mathcal{F}(\pi)$  is an element of  $T_{\pi}\mathcal{M}$ , so it is identifiable with a gradient  $\nabla\psi$  for  $\psi : \mathbb{R}^d \rightarrow \mathbb{R}$  following (30). Let us further suppose that the functional  $\mathcal{F}$  satisfies the following assumption: for any smooth curve  $\{\pi_t\}_{t \in \mathbb{R}} \in \mathcal{M}$ , we have

$$\frac{\partial}{\partial t}\mathcal{F}(\pi_t)|_{t=0} = \int (\delta\mathcal{F}(\pi)(\mathbf{x})) \left( \frac{\partial}{\partial t}\pi_t(\mathbf{x})|_{t=0} \right) d\mathbf{x}, \quad (34)$$

for a scalar-valued function  $\delta\mathcal{F}(\pi) : \mathbb{R}^d \rightarrow \mathbb{R}$ . We will see an example of a functional  $\mathcal{F}$  satisfying (34) shortly; the function  $\delta\mathcal{F}(\pi)$  is called the *first variation* of  $\mathcal{F}$  at  $\pi$ . Continuing, we see that

$$\begin{aligned} \frac{\partial}{\partial t}\mathcal{F}(\pi_t)|_{t=0} &= \int (\delta\mathcal{F}(\pi)(\mathbf{x})) (\nabla \cdot (\pi \nabla \psi_0)(\mathbf{x})) d\mathbf{x} \\ &= - \int \langle \nabla \delta\mathcal{F}(\pi)(\mathbf{x}), \nabla \psi_0(\mathbf{x}) \rangle \pi(\mathbf{x}) d\mathbf{x} = \langle \nabla \cdot (-\pi \nabla \delta\mathcal{F}(\pi)), \nabla \cdot (\pi \nabla \psi_0) \rangle_{\pi}. \end{aligned}$$

The first equation used (33), the second was integration by parts (13), and the last used our local metric formula (32). Comparing the above display with the definition of the manifold gradient in (26), we see that whenever (34) holds for a functional  $\mathcal{F}$ , the manifold gradient is simply given by

$$\nabla_{\mathcal{M}}\mathcal{F}(\pi) = \nabla \cdot (-\pi \nabla \delta\mathcal{F}(\pi)) \equiv -\nabla \delta\mathcal{F}(\pi). \quad (35)$$

We use the above notation to mean  $-\nabla \delta\mathcal{F}(\pi)$  is the vector field followed in particle space pointwise, inducing a density flow of the form  $\frac{\partial}{\partial t}\pi_t|_{t=0} = \nabla \cdot (-\pi \nabla \delta\mathcal{F}(\pi))$  for curves  $\{\pi_t\}_{t \in \mathbb{R}}$  passing through  $\pi_0 = \pi$  (see Lemma 6).

To give an example, suppose that  $V : \mathbb{R}^d \rightarrow \mathbb{R}$  satisfies  $\int \exp(-V(\mathbf{x})) d\mathbf{x} < \infty$ , and let  $\pi^* \in \mathcal{P}(\mathbb{R}^d)$  be the density with  $\pi^* \propto \exp(-V)$ . Finally, consider the functional

$$\mathcal{F}(\pi) := D_{\text{KL}}(\pi \| \pi^*), \text{ for all } \pi \in \mathcal{M}. \quad (36)$$

We begin by computing the first variation of  $\mathcal{F}$ , in the sense of (34). Observe that

$$\begin{aligned} \frac{\partial}{\partial t}\mathcal{F}(\pi_t) &= \frac{\partial}{\partial t} \left( \int \pi_t(\mathbf{x}) \log \left( \frac{\pi_t(\mathbf{x})}{\pi^*(\mathbf{x})} \right) d\mathbf{x} \right) = \int \left( \log \left( \frac{\pi_t(\mathbf{x})}{\pi^*(\mathbf{x})} \right) + 1 \right) \left( \frac{\partial}{\partial t}\pi_t(\mathbf{x}) \right) d\mathbf{x} \\ &= \int \log \left( \frac{\pi_t(\mathbf{x})}{\pi^*(\mathbf{x})} \right) \left( \frac{\partial}{\partial t}\pi_t(\mathbf{x}) \right) d\mathbf{x} + \frac{\partial}{\partial t} \int \pi_t(\mathbf{x}) d\mathbf{x} \\ &= \int \log \left( \frac{\pi_t(\mathbf{x})}{\pi^*(\mathbf{x})} \right) \left( \frac{\partial}{\partial t}\pi_t(\mathbf{x}) \right) d\mathbf{x}, \end{aligned} \quad (37)$$

since  $\int \pi_t(\mathbf{x}) d\mathbf{x} = 1$  for all  $\pi_t \in \mathcal{M}$ . Therefore, comparing (37) to (34) shows that

$$\nabla \delta\mathcal{F}(\pi)(\mathbf{x}) = \nabla \log \left( \frac{\pi(\mathbf{x})}{\pi^*(\mathbf{x})} \right) = \nabla V(\mathbf{x}) + \nabla \log \pi(\mathbf{x}), \quad (38)$$

up to a universal additive constant. In other words, following (35), the negated manifold gradient of the functional  $\mathcal{F} = D_{\text{KL}}(\cdot \| \pi^*)$  at  $\pi \in \mathcal{M}$  is the vector field given pointwise by  $\nabla V + \nabla \log \pi$ .

Perhaps surprisingly (as first observed by the landmark result of [JKO98]), this is the same gradient field induced by the continuity equation of the Langevin dynamics, as computed in (31)! Therefore, following the gradient flow along  $\mathcal{M}$ , i.e., defining  $\frac{\partial}{\partial t}\pi_t = -\nabla_{\mathcal{M}}\mathcal{F}(\pi_t)$ , is equivalent to a particle flow governed by the Langevin dynamics. We summarize this observation as follows.

**Theorem 3** ([JKO98]). *Let  $\pi^* \in \mathcal{P}_2(\mathbb{R}^d)$  satisfy  $\pi^* \propto \exp(-V)$  for  $V : \mathbb{R}^d \rightarrow \mathbb{R}$ , and let  $\mathcal{F}(\pi) := D_{\text{KL}}(\pi \| \pi^*)$  for  $\pi \in \mathcal{P}_2(\mathbb{R}^d)$ . Drawing  $\mathbf{x}_0 \sim \pi_0$  for a density  $\pi_0 \in \mathcal{P}(\mathbb{R}^d)$ , and letting  $\pi_t$  be the law of  $\mathbf{x}_t$  for  $\{\mathbf{x}_t\}_{t \geq 0}$  following the Langevin dynamics (4), the curve  $\{\pi_t\}_{t \geq 0}$  follows the curve*

$$\frac{\partial}{\partial t}\pi_t = -\nabla_{\mathcal{M}}\mathcal{F}(\pi_t),$$

where  $\nabla_{\mathcal{M}}$  is the manifold gradient defined in (26), and  $\mathcal{M}$  is the Riemannian manifold given by  $\mathcal{P}_2(\mathbb{R}^d)$  equipped with the distance function  $W_2$ . In other words, the Langevin dynamics in particle space induce the gradient flow of the KL divergence to  $\pi^*$  in the space of measures.



We conclude this section with one additional manifold gradient formula, applied to the functional

$$\mathcal{F}(\pi) := \frac{1}{2} W_2^2(\pi, \pi^*),$$

for a target density  $\pi^* \in \mathcal{M}$ . The manifold gradient in this case is induced by the vector field

$$\begin{aligned} \nabla_{\mathcal{M}} \mathcal{F}(\pi)(\mathbf{x}) &\equiv \mathbf{x} - \nabla \varphi(\mathbf{x}), \\ \text{where } (\mathbf{x}, \nabla \varphi(\mathbf{x})) &\in \mathcal{C}(\pi, \pi^*) \text{ for } \mathbf{x} \sim \pi \text{ is the coupling inducing } W_2^2(\pi, \pi^*). \end{aligned} \quad (39)$$

This is intuitive along geodesics, given our earlier discussion of McCann’s displacement interpolation and following (34); more generally we defer proving (39) to [Vil08], Theorem 23.9.

## 5 Functional inequalities

In this section, we introduce *functional inequalities*, which characterize certain spectral properties of a Markov semigroup. As a first example, we define and interpret the Poincaré inequality.

**Definition 5** (Poincaré inequality). *We say that  $\pi^* \in \mathcal{P}(\mathbb{R}^d)$  satisfies a Poincaré inequality with constant  $C_{\text{PI}}$  if for all differentiable  $f \in L^2(\pi^*)$ ,*

$$\text{Var}_{\pi^*}[f] \leq C_{\text{PI}} \int \|\nabla f(\mathbf{x})\|_2^2 \pi^*(\mathbf{x}) d\mathbf{x}, \quad (40)$$

where

$$\text{Var}_{\pi^*}[f] := \int \left( f(\mathbf{x}) - \int f(\mathbf{y}) \pi^*(\mathbf{y}) d\mathbf{y} \right)^2 \pi^*(\mathbf{x}) d\mathbf{x}.$$

As a first example, suppose that  $\pi^*$  satisfies a Poincaré constant with constant  $C_{\text{PI}}$ . Then plugging in each  $f(\mathbf{x}) = \langle \mathbf{v}, \mathbf{x} \rangle$ , and defining  $\text{Cov}_{\pi^*} := \mathbb{E}_{\mathbf{x} \sim \pi^*}[\mathbf{x}\mathbf{x}^\top] - \mathbb{E}_{\mathbf{x} \sim \pi^*}[\mathbf{x}]\mathbb{E}_{\mathbf{x} \sim \pi^*}[\mathbf{x}]^\top$ , (40) reads

$$\mathbf{v}^\top \text{Cov}_{\pi^*} \mathbf{v} \leq C_{\text{PI}} \|\mathbf{v}\|_2^2 \text{ for all } \mathbf{v} \in \mathbb{R}^d \implies \text{Cov}_{\pi^*} \preceq C_{\text{PI}} \mathbf{I}_d.$$

Thus, Poincaré inequalities imply covariance bounds, suggesting that they are related to concentration of  $\pi^*$ . In fact, as we will see, much stronger concentration follows from Poincaré inequalities.

We now explain how Definition 5 has a natural interpretation as providing a spectral gap on the Dirichlet energy operator  $\mathcal{E}$  defined in (18). Recall from Lemma 3 that when  $\mathcal{E}$  corresponds to the Langevin dynamics (4) with stationary density  $\pi^*$ , we have  $\mathcal{E}(f, g) = \int \langle \nabla f(\mathbf{x}), \nabla g(\mathbf{x}) \rangle \pi^*(\mathbf{x}) d\mathbf{x}$  for  $f, g \in L^2(\pi^*)$ , so that the right-hand side of (40) is simply  $C_{\text{PI}} \cdot \mathcal{E}(f, f)$ . Moreover, we know that constant functions are vanishing eigenfunctions of  $\mathcal{E}$ , since they have pointwise zero gradients. Hence, viewing  $L^2(\pi^*)$  as a Hilbert space with the inner product  $\langle f, g \rangle = \int f(\mathbf{x})g(\mathbf{x})\pi^*(\mathbf{x})d\mathbf{x}$  (and defining  $\|f\|^2 := \langle f, f \rangle$ ), the condition (40) is equivalent to the statement that

$$\frac{\mathcal{E}(f, f)}{\|f\|^2} \geq \frac{1}{C_{\text{PI}}} \text{ for all } f \in L^2(\pi^*) \text{ with } \langle f, c \rangle = 0,$$

where  $c$  is any constant-valued function. This is precisely imposing a spectral gap on the operator  $\mathcal{E}$ , since it bounds its second-smallest eigenvalue. As we will see, we can formally use Poincaré inequalities to bound the decay of any eigenfunctions orthogonal to constant functions to derive convergence to  $\pi^*$  in the  $\chi^2$  distance. We also mention that our definition (40) is specialized to the Dirichlet energy operator  $\mathcal{E}$  arising from the Langevin dynamics. Indeed, other reversible Markov semigroups induce alternative compatible formulations of the Poincaré inequality, by setting the right-hand side of (40) to the generalization of  $\mathcal{E}$  derived in Lemma 3.

In general, the  $\chi^2$  divergence between probability distributions (whose convergence proofs are facilitated by Poincaré inequalities, as statements about “variance decay”) dominates the KL divergence  $D_{\text{KL}}$ . However, in important cases of interest the initial  $\chi^2$  divergence can be exponentially larger than  $D_{\text{KL}}$  (see discussion at the end of Section 2.2, Part XII), so we would like a means to directly argue about the convergence of a semigroup in  $D_{\text{KL}}$ . This “entropy decay” strengthening of the Poincaré inequality is the log-Sobolev inequality, which we next define.

**Definition 6** (Log-Sobolev inequality). *We say that  $\pi^* \in \mathcal{P}(\mathbb{R}^d)$  satisfies a log-Sobolev inequality with constant  $C_{\text{LSI}}$  if for all differentiable  $f \in L^2(\pi^*)$ ,*

$$\text{Ent}_{\pi^*}[f^2] \leq 2C_{\text{LSI}} \int \|\nabla f(\mathbf{x})\|_2^2 \pi^*(\mathbf{x}) d\mathbf{x}, \quad (41)$$

where

$$\text{Ent}_{\pi^*}[f] := \int f(\mathbf{x}) \log(f(\mathbf{x})) \pi^*(\mathbf{x}) d\mathbf{x} - \left( \int f(\mathbf{x}) \pi^*(\mathbf{x}) d\mathbf{x} \right) \log \left( \int f(\mathbf{x}) \pi^*(\mathbf{x}) d\mathbf{x} \right).$$

We show that Definition 6 is indeed a strengthening of Definition 5.

**Lemma 7.** *If  $\pi^* \in \mathcal{P}(\mathbb{R}^d)$  satisfies a log-Sobolev inequality with constant  $C$ , it also satisfies a Poincaré inequality with constant  $C$ .*

*Proof.* Let  $f \in L^2(\pi^*)$ , and without loss of generality suppose that  $\mathbb{E}_{\pi^*}[f] = \int f(\mathbf{x}) \pi^*(\mathbf{x}) d\mathbf{x} = 0$ , else we may subtract an appropriate constant from  $f$  pointwise which does not affect the calculation (40). Next, consider taking  $f \leftarrow 1 + \epsilon f$  in (41) for a small constant  $\epsilon$ . We compute that

$$\lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon^2} \text{Ent}_{\pi^*}[(1 + \epsilon f)^2] = 2 \int f(\mathbf{x})^2 \pi^*(\mathbf{x}) d\mathbf{x},$$

which follows by Lebesgue's dominated convergence theorem, the assumption  $\mathbb{E}_{\pi^*}[f] = 0$ , and

$$\begin{aligned} (1 + \epsilon c)^2 \log((1 + \epsilon c)^2) &= 2(1 + 2\epsilon c) \left( \epsilon c - \frac{\epsilon^2 c^2}{2} \right) + O(\epsilon^3) = 2\epsilon c + 3\epsilon^2 c^2, \\ (1 + \epsilon^2 c) \log(1 + \epsilon^2 c) &= \epsilon^2 c + O(\epsilon^3). \end{aligned}$$

Meanwhile, it is clear that

$$\lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon^2} \int \|\epsilon \nabla f(\mathbf{x})\|_2^2 \pi^*(\mathbf{x}) d\mathbf{x} = \int \|\nabla f(\mathbf{x})\|_2^2 \pi^*(\mathbf{x}) d\mathbf{x}.$$

Combining in (41) and taking limits yields the conclusion.  $\square$

We also mention one alternative formulation of (41), which is easier to apply in certain settings.

**Lemma 8.** *If  $\pi^* \in \mathcal{P}(\mathbb{R}^d)$  satisfies a log-Sobolev inequality with constant  $C_{\text{LSI}}$ , then for any  $\pi \in \mathcal{P}(\mathbb{R}^d)$  such that  $\frac{\pi}{\pi^*} < \infty$  almost surely,*

$$D_{\text{KL}}(\pi \| \pi^*) \leq \frac{C_{\text{LSI}}}{2} \int \left\| \nabla \log \left( \frac{\pi(\mathbf{x})}{\pi^*(\mathbf{x})} \right) \right\|_2^2 \pi(\mathbf{x}) d\mathbf{x}, \quad (42)$$

where

$$D_{\text{KL}}(\pi \| \pi^*) := \int \pi(\mathbf{x}) \log \left( \frac{\pi(\mathbf{x})}{\pi^*(\mathbf{x})} \right) d\mathbf{x}.$$

*Proof.* By plugging in  $f = \sqrt{\frac{\pi}{\pi^*}}$  into (41), the left-hand side reads

$$\text{Ent}_{\pi^*}[f^2] = \int \frac{\pi(\mathbf{x})}{\pi^*(\mathbf{x})} \log \left( \frac{\pi(\mathbf{x})}{\pi^*(\mathbf{x})} \right) \pi^*(\mathbf{x}) d\mathbf{x} = D_{\text{KL}}(\pi \| \pi^*),$$

since  $\log(\int \pi(\mathbf{x}) d\mathbf{x}) = \log 1 = 0$ . Moreover, the right-hand side of (41) is

$$\begin{aligned} 2C_{\text{LSI}} \int \|\nabla f(\mathbf{x})\|_2^2 \pi^*(\mathbf{x}) d\mathbf{x} &= \frac{C_{\text{LSI}}}{2} \int \left\| \nabla \frac{\pi(\mathbf{x})}{\pi^*(\mathbf{x})} \right\|_2^2 \frac{\pi^*(\mathbf{x})^2}{\pi(\mathbf{x})} d\mathbf{x} \\ &= \frac{C_{\text{LSI}}}{2} \int \left\| \nabla \log \left( \frac{\pi(\mathbf{x})}{\pi^*(\mathbf{x})} \right) \right\|_2^2 \pi(\mathbf{x}) d\mathbf{x}, \end{aligned}$$

and combining yields the desired (42).  $\square$

The term  $\int \|\nabla \log(\frac{\pi(\mathbf{x})}{\pi^*(\mathbf{x})})\|_2^2 \pi(\mathbf{x}) d\mathbf{x}$  in the right-hand side of (42) is sometimes called the *Fisher information*, an important quantity in statistics. For now, we mention that under our earlier calculations (35) and (38), letting  $\mathcal{F}(\pi) := D_{\text{KL}}(\pi \| \pi^*)$  so that  $\mathcal{F}(\pi^*) = 0$ , (42) alternatively reads

$$\mathcal{F}(\pi) - \mathcal{F}(\pi^*) \leq \frac{C_{\text{LSI}}}{2} \langle \nabla_{\mathcal{M}} \mathcal{F}(\pi), \nabla_{\mathcal{M}} \mathcal{F}(\pi) \rangle_{\pi}. \quad (43)$$

This is the functional analog of the gradient domination condition in Corollary 3, Part II which allowed us to prove linear convergence rates for gradient flow. We will show shortly that log-Sobolev inequalities establish a similar linear rate of decay for the KL divergence.

## 5.1 Convergence from functional inequalities

We now formalize our earlier statements that the Poincaré inequality implies a linear rate of decay on an appropriate notion of variance under the Langevin dynamics, and that the log-Sobolev inequality implies a similar decay on the KL divergence. We begin by proving variance decay, reminding the reader of the definition of the  $\chi^2$  divergence between distributions:

$$\chi^2(\pi \| \pi^*) := \text{Var}_{\pi^*} \left[ \frac{\pi}{\pi^*} \right] = \int \left( \frac{\pi(x)}{\pi^*(x)} \right)^2 \pi^*(x) dx - 1.$$

**Lemma 9.** *Let  $\{\mathbf{x}_t\}_{t \geq 0}$  follow the Langevin dynamics (4) with stationary distribution  $\pi^*$ , and assume that  $\pi^*$  satisfies a Poincaré inequality with constant  $C_{\text{PI}}$ . For all  $t \geq 0$ , the law of  $\mathbf{x}_t$ , denoted  $\pi_t \in \mathcal{P}(\mathbb{R}^d)$ , satisfies*

$$\chi^2(\pi_t \| \pi^*) \leq \exp\left(-\frac{2t}{C_{\text{PI}}}\right) \chi^2(\pi_0 \| \pi^*).$$

*Proof.* It suffices to show that  $\frac{d}{dt} \chi^2(\pi_t \| \pi^*) \leq -\frac{2}{C_{\text{PI}}} \chi^2(\pi_t \| \pi^*)$ , at which point the conclusion follows from Grönwall's inequality (Fact 1, Part II). We compute

$$\begin{aligned} \frac{d}{dt} \chi^2(\pi_t \| \pi^*) &= \frac{\partial}{\partial t} \int \left( \frac{\pi_t(\mathbf{x})}{\pi^*(\mathbf{x})} - 1 \right)^2 \pi^*(\mathbf{x}) d\mathbf{x} \\ &= 2 \int \left( \frac{\pi_t(\mathbf{x})}{\pi^*(\mathbf{x})} \right) \left( \frac{\partial}{\partial t} \frac{\pi_t(\mathbf{x})}{\pi^*(\mathbf{x})} \right) \pi^*(\mathbf{x}) d\mathbf{x} \\ &= 2 \int \left( \frac{\pi_t(\mathbf{x})}{\pi^*(\mathbf{x})} \right) \left( \frac{\mathcal{L}^* \pi_t(\mathbf{x})}{\pi^*(\mathbf{x})} \right) \pi^*(\mathbf{x}) d\mathbf{x} = -2\mathcal{E} \left( \frac{\pi_t}{\pi^*}, \frac{\pi_t}{\pi^*} \right), \end{aligned}$$

where we used Kolmogorov's backward equation (11) in the third equality, which implies  $\frac{\partial}{\partial t} \pi_t(x) = \mathcal{L}^* \pi_t(\mathbf{x})$ . Finally, using the definition (40), we have the desired

$$-2\mathcal{E} \left( \frac{\pi_t}{\pi^*}, \frac{\pi_t}{\pi^*} \right) \leq -\frac{2}{C_{\text{PI}}} \text{Var}_{\pi^*} \left[ \frac{\pi_t}{\pi^*} \right] = -\frac{2}{C_{\text{PI}}} \chi^2(\pi_t \| \pi^*),$$

where we recalled the formula for  $\mathcal{E}$  specialized to the Langevin diffusion, stated in Lemma 3.  $\square$

We prove a similar convergence for the KL divergence under a log-Sobolev inequality.

**Lemma 10.** *Let  $\{\mathbf{x}_t\}_{t \geq 0}$  follow the Langevin dynamics (4) with stationary distribution  $\pi^*$ , and assume that  $\pi^*$  satisfies a log-Sobolev inequality with constant  $C_{\text{LSI}}$ . For all  $t \geq 0$ , the law of  $\mathbf{x}_t$ , denoted  $\pi_t \in \mathcal{P}(\mathbb{R}^d)$ , satisfies*

$$D_{\text{KL}}(\pi_t \| \pi^*) \leq \exp\left(-\frac{2t}{C_{\text{LSI}}}\right) D_{\text{KL}}(\pi_0 \| \pi^*).$$

*Proof.* As in the proof of Lemma 9, it suffices to show that  $\frac{d}{dt} D_{\text{KL}}(\pi_t \| \pi^*) \leq -\frac{2}{C_{\text{LSI}}} D_{\text{KL}}(\pi_t \| \pi^*)$ . We previously computed in (38) that the KL divergence satisfies

$$\frac{\partial}{\partial t} D_{\text{KL}}(\pi_t \| \pi^*) = \int \log \left( \frac{\pi_t(\mathbf{x})}{\pi^*(\mathbf{x})} \right) \left( \frac{\mathcal{L}^* \pi_t(\mathbf{x})}{\pi^*(\mathbf{x})} \right) \pi^*(\mathbf{x}) d\mathbf{x} = -\mathcal{E} \left( \frac{\pi_t}{\pi^*}, \log \left( \frac{\pi_t}{\pi^*} \right) \right).$$

The conclusion follows analogously to Lemma 9, using (41), upon realizing

$$\begin{aligned} \mathcal{E}\left(\frac{\pi_t}{\pi^*}, \log\left(\frac{\pi_t}{\pi^*}\right)\right) &= \int \left\langle \nabla \log\left(\frac{\pi_t(\mathbf{x})}{\pi^*(\mathbf{x})}\right), \nabla\left(\frac{\pi_t(\mathbf{x})}{\pi^*(\mathbf{x})}\right) \right\rangle \pi^*(\mathbf{x}) d\mathbf{x} \\ &= \int \left\| \nabla \log\left(\frac{\pi_t(\mathbf{x})}{\pi^*(\mathbf{x})}\right) \right\|_2^2 \pi(\mathbf{x}) d\mathbf{x} \geq \frac{2}{C_{\text{LSI}}} D_{\text{KL}}(\pi_t \| \pi^*). \end{aligned}$$

□

## 5.2 Functional inequalities from strong logconcavity

We have established in Lemmas 9 and 10 that functional inequalities give us a way of proving convergence of Markov semigroups (in relative variance, i.e.,  $\chi^2$ , or relative entropy, i.e.,  $D_{\text{KL}}$ ). These results can be thought of as continuous-time generalizations of the spectral gap arguments we used in the discrete-time setting, e.g., Parts XIV and XV. However, we have not yet given examples of densities which actually satisfy Definitions 5 or 6, which is the purpose of this section.

We first establish a ‘‘Hessian-reweighted’’ variant of (40) for logconcave densities, which streamlines our presentation. Specifically, recall the Prékopa-Leindler inequality (Theorem 4, Part I), which says (changing function names from  $(f, g, h) \rightarrow (u, v, w)$  for convenience) that for  $\lambda \in [0, 1]$  and  $u, v, w : \mathbb{R}^d \rightarrow \mathbb{R}$  with  $w((1 - \lambda)\mathbf{x} + \lambda\mathbf{x}') \geq u(\mathbf{x})^{1-\lambda}v(\mathbf{x}')^\lambda$  for all  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ ,

$$\int w(\mathbf{x}) d\mathbf{x} \geq \left( \int u(\mathbf{x}) d\mathbf{x} \right)^{1-\lambda} \left( \int v(\mathbf{x}) d\mathbf{x} \right)^\lambda. \quad (44)$$

The following result is known as the *Brascamp-Lieb inequality* [BL76], and our proof is based on an elegant strategy from [BL00] which applies the Prékopa-Leindler inequality (44).

**Proposition 5** (Brascamp-Lieb inequality). *Let  $\pi^* \in \mathcal{P}(\mathbb{R}^d)$  be logconcave, and let  $\pi^* \propto \exp(-V)$  where  $V : \mathbb{R}^d \rightarrow \mathbb{R}$  is strictly convex and twice-continuously differentiable. Then for all twice-differentiable  $f \in L^2(\pi^*)$ ,*

$$\text{Var}_{\pi^*}[f] \leq \int \|\nabla f(\mathbf{x})\|_{(\nabla^2 V(\mathbf{x}))^{-1}}^2 \pi^*(\mathbf{x}) d\mathbf{x}.$$

*Proof.* Throughout the proof, assume without loss that  $\int \exp(-V(\mathbf{x})) d\mathbf{x} = 1$ , and that there exists  $\alpha > 0$  such that  $\nabla^2 V \succeq \alpha \mathbf{I}_d$  pointwise; this latter assumption can be removed using a limiting argument. Our strategy is to apply the Prékopa-Leindler inequality (44) with  $\lambda = \frac{1}{2}$ , and

$$\begin{aligned} u(\mathbf{x}) &:= \exp(2\delta f(\mathbf{x}) - V(\mathbf{x})), \quad v(\mathbf{x}) := \exp(-V(\mathbf{x})), \quad w(\mathbf{x}) := \exp(g^{\delta f}(\mathbf{x}) - V(\mathbf{x})), \\ \text{where } \delta \rightarrow 0 \text{ and } g^{\delta f}(\mathbf{y}) &:= \sup_{\mathbf{y}=\frac{1}{2}(\mathbf{x}+\mathbf{x}')} \left\{ \delta f(\mathbf{x}) - \left( \frac{1}{2}V(\mathbf{x}) + \frac{1}{2}V(\mathbf{x}') - V(\mathbf{y}) \right) \right\}. \end{aligned}$$

By the definition of  $g^{\delta f}$ ,  $w(\frac{1}{2}(\mathbf{x} + \mathbf{x}')) \geq u(\mathbf{x})^{\frac{1}{2}}v(\mathbf{x}')^{\frac{1}{2}}$  for all  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ , so (44) gives

$$\int \exp(2\delta f(\mathbf{x})) \pi^*(\mathbf{x}) d\mathbf{x} \leq \left( \int \exp(g^{\delta f}(\mathbf{x})) \pi^*(\mathbf{x}) d\mathbf{x} \right)^2. \quad (45)$$

Next, we claim that as  $\delta \rightarrow 0$ , we have that for all  $\mathbf{y} \in \mathbb{R}^d$ ,

$$g^{\delta f}(\mathbf{y}) = \delta f(\mathbf{y}) + \frac{\delta^2}{2} (\nabla^2 V(\mathbf{y}))^{-1} [\nabla f(\mathbf{y}), \nabla f(\mathbf{y})] + o(\delta^2). \quad (46)$$

To see this, essentially the idea is to treat  $\delta f$  as a small perturbation of the strongly convex function  $V$ , and see how the minimizer varies. For fixed  $\delta, \mathbf{y}$ , let us parameterize the optimization problem governing  $g^{\delta f}(\mathbf{y})$  by  $\mathbf{h} \in \mathbb{R}^d$  such that  $\mathbf{x} = \mathbf{y} + \mathbf{h}$  and  $\mathbf{x}' = \mathbf{y} - \mathbf{h}$ . Then a necessary condition for the optimal  $\mathbf{h}$  is

$$\delta \nabla f(\mathbf{y} + \mathbf{h}) - \frac{1}{2} \nabla V(\mathbf{y} + \mathbf{h}) + \frac{1}{2} \nabla V(\mathbf{y} - \mathbf{h}) = 0.$$

Moreover, under our assumption that  $\nabla^2 V$  is pointwise lower bounded, this gives

$$\delta \nabla f(\mathbf{y}) + \delta \nabla^2 f(\mathbf{y}) \mathbf{h} = \nabla^2 V(\mathbf{y}) \mathbf{h} + \mathbf{r} \text{ for } \|\mathbf{r}\|_2 = o(\delta^2) \implies \|\mathbf{h}\|_2 = O(\delta),$$

where the asymptotic notation above hides local smoothness parameters of  $f, V$  at  $\mathbf{y}$ , that vanish as  $\delta \rightarrow 0$ . Plugging this estimate in above, we have shown that at optimality,

$$\mathbf{h} = \delta \nabla^2 V(\mathbf{y})^{-1} \nabla f(\mathbf{y}) + \mathbf{r} \text{ for } \|\mathbf{r}\|_2 = o(\delta).$$

Finally, we have from a second-order Taylor expansion of  $V$  that

$$\delta f(\mathbf{y} + \mathbf{h}) - \left( \frac{1}{2} V(\mathbf{y} + \mathbf{h}) + \frac{1}{2} V(\mathbf{y} - \mathbf{h}) - V(\mathbf{y}) \right) = \delta f(\mathbf{y}) + \frac{1}{2} \nabla^2 V(\mathbf{y})[\mathbf{h}, \mathbf{h}] + o(\delta^2),$$

so plugging in our optimal  $\mathbf{h}$ , we have established (46). Now, expanding (45),

$$\begin{aligned} & 1 + 2\delta \int f(\mathbf{x}) \pi^*(\mathbf{x}) d\mathbf{x} + 2\delta^2 \int f(\mathbf{x})^2 \pi^*(\mathbf{x}) d\mathbf{x} - o(\delta^2) \\ &= \int \exp(2\delta f(\mathbf{x})) \pi^*(\mathbf{x}) d\mathbf{x} \\ &\leq \left( \int \exp \left( \delta f(\mathbf{x}) + \frac{\delta^2}{2} (\nabla^2 V(\mathbf{x}))^{-1} [\nabla f(\mathbf{x}), \nabla f(\mathbf{x})] \right) \pi^*(\mathbf{x}) d\mathbf{x} \right)^2 + o(\delta^2) \\ &\leq 1 + 2\delta \int f(\mathbf{x}) \pi^*(\mathbf{x}) + \delta^2 \int f(\mathbf{x})^2 \pi^*(\mathbf{x}) d\mathbf{x} + \delta^2 \left( \int f(\mathbf{x}) \pi^*(\mathbf{x}) d\mathbf{x} \right)^2 \\ &\quad + \delta^2 \int (\nabla^2 V(\mathbf{x}))^{-1} [\nabla f(\mathbf{x}), \nabla f(\mathbf{x})] \pi^*(\mathbf{x}) d\mathbf{x} + o(\delta^2). \end{aligned}$$

where the second line used a second-order Taylor expansion of  $\exp$ , the third line combined (45) and (46), and the last line also used a second-order Taylor expansion of  $\exp$  and  $(1 + \cdot)^2$ . At this point, cancelling like terms and normalizing by  $\delta^2$  for  $\delta \rightarrow 0$  gives the claimed result:

$$\left( \int f(\mathbf{x})^2 \pi^*(\mathbf{x}) d\mathbf{x} - \left( \int f(\mathbf{x}) \pi^*(\mathbf{x}) d\mathbf{x} \right)^2 \right) \leq \int (\nabla^2 V(\mathbf{x}))^{-1} [\nabla f(\mathbf{x}), \nabla f(\mathbf{x})] \pi^*(\mathbf{x}) d\mathbf{x}.$$

□

Proposition 5 implies a Poincaré inequality for *strongly logconcave* densities (Definition 4).

**Corollary 1.** *Any  $\mu$ -strongly logconcave density satisfies a Poincaré inequality with constant  $\frac{1}{\mu}$ .*

*Proof.* Apply Proposition 5 with  $(\nabla^2 V(\cdot))^{-1} \preceq \frac{1}{\mu} \mathbf{I}_d$ , where  $V$  is the negative log-density. □

Although one may naturally next conjecture that the ‘‘Hessian-reweighted’’ variant of the log-Sobolev inequality (41) holds, i.e.,

$$\text{Ent}_{\pi^*} [f^2] \leq 2 \int \|\nabla f(\mathbf{x})\|_{(\nabla^2 V(\mathbf{x}))^{-1}}^2 \pi^*(\mathbf{x}) d\mathbf{x},$$

it turns out that there are one-dimensional counterexamples (Section 2, [BL00]). Nonetheless, we have quite a substantial generalization of Corollary 1.

**Theorem 4** (Proposition 3.1, [BL00]). *Let  $\pi^* \in \mathcal{P}(\mathbb{R}^d)$  be  $\mu$ -strongly logconcave with respect to  $\|\cdot\|$ , i.e.,  $\pi^* \propto \exp(-V)$  where  $V$  is  $\mu$ -strongly convex with respect to  $\|\cdot\|$  (Definition 4, Part II). Then for all differentiable  $f \in L^2(\pi^*)$ ,*

$$\text{Ent}_{\pi^*} [f^2] \leq \frac{2}{\mu} \int \|\nabla f(\mathbf{x})\|_*^2 \pi^*(\mathbf{x}) d\mathbf{x}.$$

We do not prove Theorem 4 here, but mention that it follows from a variant of our proof of Proposition 5, instead applied with  $\lambda \rightarrow 0$ . The key fact used is that

$$\frac{d}{dx} x^{1+\lambda} \Big|_{\lambda=0} = \lim_{\lambda \rightarrow 0} \frac{x^{1+\lambda} - x}{\lambda} = x \log(x).$$

and much of the rest of the proof is the same. Thus, Theorem 4 implies a log-Sobolev generalization of Corollary 1 holds, for all norms (not just  $\ell_2$ ). This also shows that the Poincaré variant in Corollary 1 holds in all norms as well, using the reduction in Lemma 7.

### 5.3 Further consequences

In this section, we mention a few further consequences of functional inequalities which are often useful in applications. First, we show that they imply concentration of Lipschitz functions. The following proof strategy is attributed by [Led99] to Ira Herbst in an unpublished letter to Leonard Gross, and hence has come to be known as ‘‘Herbst’s argument.’’

**Lemma 11** (Proposition 2.3, [Led99]). *If  $\pi^* \in \mathcal{P}(\mathbb{R}^d)$  satisfies a log-Sobolev inequality with constant  $C_{\text{LSI}}$ , then for any 1-Lipschitz function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , we have that*

$$\Pr_{\mathbf{x} \sim \pi^*} [f(\mathbf{x}) \geq \mathbb{E}_{\pi^*}[f] + c] \leq \exp\left(-\frac{c^2}{2C_{\text{LSI}}}\right) \text{ for all } c > 0.$$

*Proof.* The basic idea is to apply the log-Sobolev inequality with  $f^2 \leftarrow \exp(\lambda f)$ , the *moment-generating function* of  $f$ , treated as a random variable. Let

$$H(\lambda) := \mathbb{E}_{\pi^*} [\exp(\lambda f)] \implies \lambda H'(\lambda) = \mathbb{E}_{\pi^*} [\lambda f \exp(\lambda f)].$$

Then, by applying the log-Sobolev inequality,

$$\begin{aligned} \lambda H'(\lambda) - H(\lambda) \log H(\lambda) &\leq 2C_{\text{LSI}} \int \left\| \nabla \exp\left(\frac{\lambda}{2} f(\mathbf{x})\right) \right\|_2^2 \pi^*(\mathbf{x}) d\mathbf{x} \\ &= \frac{\lambda^2 C_{\text{LSI}}}{2} \int \|\nabla f(\mathbf{x})\|_2^2 \exp(\lambda f(\mathbf{x})) \pi^*(\mathbf{x}) d\mathbf{x} \leq \frac{\lambda^2 C_{\text{LSI}}}{2} H(\lambda). \end{aligned} \tag{47}$$

Now consider the function  $K(\lambda) := \frac{1}{\lambda} \log H(\lambda)$ , where we define

$$K(0) := \lim_{\lambda \rightarrow 0} \frac{1}{\lambda} \log (\mathbb{E}_{\pi^*} [\exp(\lambda f)]) = \lim_{\lambda \rightarrow 0} \frac{H(\lambda) - H(0)}{\lambda} = H'(0) = \mathbb{E}_{\pi^*} [f],$$

using  $\log(1+x) = x + o(x^2)$  for  $x \rightarrow 0$ . Moreover,

$$K'(\lambda) = -\frac{1}{\lambda^2} \log H(\lambda) + \frac{H'(\lambda)}{\lambda H(\lambda)} = \frac{\lambda H'(\lambda) - H(\lambda) \log H(\lambda)}{\lambda^2 H(\lambda)} \leq \frac{C_{\text{LSI}}}{2},$$

where the last inequality used (47). Therefore, for all  $\lambda \geq 0$ , by integrating  $K$  we have

$$K(\lambda) \leq K(0) + \frac{\lambda C_{\text{LSI}}}{2},$$

which yields the moment-generating function bound

$$\mathbb{E}_{\pi^*} [\exp(\lambda f)] = H(\lambda) = \exp(\lambda K(\lambda)) \leq \exp\left(\lambda \mathbb{E}_{\pi^*} [f] + \frac{\lambda^2 C_{\text{LSI}}}{2}\right).$$

The remainder of the proof follows analogously to Theorem 1, Part VI:

$$\begin{aligned} \Pr_{\mathbf{x} \sim \pi^*} [f(\mathbf{x}) \geq \mathbb{E}_{\pi^*}[f] + c] &= \Pr_{\mathbf{x} \sim \pi^*} [\exp(\lambda f(\mathbf{x})) \geq \exp(\lambda \mathbb{E}_{\pi^*}[f] + \lambda c)] \\ &\leq \exp\left(-\lambda c + \frac{\lambda^2 C_{\text{LSI}}}{2}\right) = \exp\left(-\frac{c^2}{2C_{\text{LSI}}}\right), \text{ for } \lambda = \frac{c}{C_{\text{LSI}}}. \end{aligned}$$

We remark that the same arguments extend straightforwardly to any norm  $\|\cdot\| \neq \|\cdot\|_2$ .  $\square$

It turns out that with some additional effort, one can also derive tail bounds from a Poincaré inequality. The proof is similar to Lemma 11 and we omit it.

**Lemma 12** (Proposition 2.13, [Led99]). *If  $\pi^* \in \mathcal{P}(\mathbb{R}^d)$  satisfies a Poincaré inequality with constant  $C_{\text{PI}}$ , then for any 1-Lipschitz function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , we have that*

$$\Pr_{\mathbf{x} \sim \pi^*} [f(\mathbf{x}) \geq \mathbb{E}_{\pi^*}[f] + c] \leq 3 \exp\left(-\frac{c}{\sqrt{C_{\text{PI}}}}\right) \text{ for all } c > 0.$$

There is an interesting qualitative difference between Lemmas 11 and 12: random variables satisfying log-Sobolev inequalities exhibit *sub-Gaussian* concentration of Lipschitz functions, whereas if they satisfy Poincaré inequalities, then Lipschitz functions are merely *sub-exponential*.

To provide some justification, recall from the localization lemma (Proposition 4, Part XV) that extreme examples of logconcave densities are exponential, which induce piecewise-linear log-densities (whose Hessians do not have curvature). Intuitively, this is consistent with our understanding of Poincaré and log-Sobolev inequalities: all logconcave densities exhibit weak forms of Poincaré inequalities, via Proposition 5, but we have only claimed that log-Sobolev holds under some strict amount of positive curvature, e.g., strong logconcavity (Theorem 4).

Our final consequence is that a log-Sobolev inequality for  $\pi^*$  implies that  $D_{\text{KL}}(\cdot\|\pi^*)$  grows quadratically with respect to  $W_2(\cdot, \pi^*)$ . This is known as the Otto-Villani theorem [FO00], and is essentially a *quadratic growth* bound, a weakening of strong convexity (Remark 3, Part II).

**Lemma 13.** *If  $\pi^* \in \mathcal{P}_2(\mathbb{R}^d)$  satisfies a log-Sobolev inequality with constant  $C_{\text{LSI}}$ , then*

$$\frac{1}{2C_{\text{LSI}}} W_2^2(\pi, \pi^*) \leq D_{\text{KL}}(\pi\|\pi^*) \text{ for all } \pi \in \mathcal{P}_2(\mathbb{R}^d).$$

*Proof.* We instead prove that if  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  satisfies the gradient domination condition

$$f(\mathbf{x}) - f(\mathbf{x}^*) \leq \frac{1}{2\mu} \|\nabla f(\mathbf{x})\|_2^2 \text{ for all } \mathbf{x} \in \mathbb{R}^d, \quad (48)$$

where  $\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$  is unique, then the quadratic growth bound

$$\frac{\mu}{2} \|\mathbf{x} - \mathbf{x}^*\|_2^2 \leq f(\mathbf{x}) - f(\mathbf{x}^*) \text{ for all } \mathbf{x} \in \mathbb{R}^d \quad (49)$$

also holds. The lemma statement follows by appropriately generalizing this proof strategy to the functional  $f \leftarrow D_{\text{KL}}(\cdot\|\pi^*)$  on Wasserstein space, and applying the observation (43). This is done using the formalism of *Otto calculus* [FO00, Ott01], which extends Section 4; we defer further proof details to Proposition 1 of [FO00], as well as an alternative perspective in [GLRT20].

To show that gradient domination implies quadratic growth on  $\mathbb{R}^d$ , let  $\mathbf{x}_0 \leftarrow \mathbf{x}$ , and define a trajectory  $\{\mathbf{x}_t\}_{t \geq 0}$  that follows the gradient flow ODE  $d\mathbf{x}_t = -\nabla f(\mathbf{x}_t)dt$ . We argued in Section 4, Part II that under the gradient domination condition (48),  $\mathbf{x}_t$  converges exponentially fast in function value, so  $\mathbf{x}_t \rightarrow \mathbf{x}^*$  as  $t \rightarrow \infty$ . Next, consider the potential

$$\Phi_t := \|\mathbf{x}_t - \mathbf{x}_0\|_2 + \sqrt{\frac{2(f(\mathbf{x}_t) - f(\mathbf{x}^*))}{\mu}},$$

applied to this trajectory. Note that

$$\begin{aligned} \frac{d\Phi_t}{dt} &\leq \|\nabla f(\mathbf{x}_t)\|_2 + \sqrt{\frac{2}{\mu}} \cdot \frac{d}{dt} \sqrt{f(\mathbf{x}_t) - f(\mathbf{x}^*)} \\ &\leq \|\nabla f(\mathbf{x}_t)\|_2 - \frac{\|\nabla f(\mathbf{x}_t)\|_2^2}{\sqrt{2\mu}(f(\mathbf{x}_t) - f(\mathbf{x}^*))} \leq 0, \end{aligned}$$

where the first line used the triangle inequality, and the second applied the assumption (48). Thus we have that  $\Phi_0 \geq \Phi_t$  for all  $t \geq 0$ , and rearranging after taking  $t \rightarrow \infty$  gives the claim (49).  $\square$

The conclusion of Lemma 13 is sometimes called (Talagrand's) transportation inequality, so log-Sobolev inequalities imply transportation inequalities. Interestingly, transportation inequalities in turn imply Poincaré inequalities via a linearization argument similar to Lemma 7 (see Section 3, [CE17]), showing that these three properties of a density  $\pi^*$  form a natural hierarchy.

## Source material

Portions of this lecture are based on reference material in [Øk03, AGS08, Vil08, Dur10, MG10, vH16, Vis18, Che24], as well as the author’s own experience working in the field. We would like to mention that much of the exposition is patterned off of the excellent resource [Che24], which we highly recommend as a deeper dive into the topics covered in this lecture.

## References

- [AGS08] Luigi Ambrosio, Nicola Gigli, and Giuseppe Savaré. *Gradient Flows in Metric Spaces and in the Space of Probability Measures*. 2008.
- [BL76] Herm Jan Brascamp and Elliott H Lieb. On extensions of the brunn-minkowski and prékopa-leindler theorems, including inequalities for log concave functions, and with an application to the diffusion equation. *Journal of Functional Analysis*, 22(4):366–389, 1976.
- [BL00] Sergey G Bobkov and Michel Ledoux. From brunn-minkowski to brascamp-lieb and to logarithmic sobolev inequalities. *GAFSA, Geometric and Functional Analysis*, 10:1028–1052, 2000.
- [Bre87] Yann Brenier. Décomposition polaire et réarrangement monotone des champs de vecteurs. *C.R. Acad. Sci. Paris Sér I Math.*, 305:805–808, 1987.
- [Bre91] Yann Brenier. Polar factorization and monotone rearrangement of vector-valued functions. *Comm. Pure Appl. Math.*, 44:375–417, 1991.
- [Caf00] Luis A. Caffarelli. Monotonicity properties of optimal transportation and the fkg and related inequalities. *Communications in Mathematical Physics*, 214(3):547–563, 2000.
- [CE17] Dario Cordero-Erausquin. Transport inequalities for log-concave measures, quantitative forms, and applications. *Canadian Journal of Mathematics*, 69(3):481–501, 2017.
- [Che24] Sinho Chewi. *Log-Concave Sampling*. 2024.
- [CP23] Sinho Chewi and Aram-Alexandre Pooladian. An entropic generalization of Caffarelli’s contraction theorem via covariance inequalities. *Reports. Mathematical*, 361:1471–1482, 2023.
- [Dam65] K. E. Dambis. On the decomposition of continuous submartingales. *Theor. Probab. Appl.*, 10:401–410, 1965.
- [dC92] Manfredo P. do Carmo. *Riemannian geometry*. 1992.
- [DS65] Lester E. Dubins and Gideon Schwarz. On continuous martingales. *Proceedings of the National Academy of Sciences*, 53(5):913–916, 1965.
- [Dur10] Rick Durrett. *Probability: Theory and Examples*. 2010.
- [FO00] Cédric Villani Felix Otto. Generalization of an inequality by talagrand and links with the logarithmic sobolev inequality. *Journal of Functional Analysis*, 173(2):361–400, 2000.
- [GLRT20] Ivan Gentil, Christian Léonard, Luigia Ripani, and Luca Tamanini. An entropic interpolation proof of the hwi inequality. *Stochastic Processes and their Applications*, 130(2):907–923, 2020.
- [JKO98] Richard Jordan, David Kinderlehrer, and Felix Otto. The variational formulation of the fokker-planck equation. *SIAM J. Math. Anal.*, 29(1):1–17, 1998.
- [Led99] Michel Ledoux. *Concentration of measure and logarithmic Sobolev inequalities*. Séminaire de probabilités XXXIII, 1999.
- [MG10] Robert J. McCann and Nestor Guillen. *Five Lectures on Optimal Transportation: Geometry, Regularity and Applications*. 2010.



- [Ott01] Felix Otto. The geometry of dissipative evolution equations: the porous medium equation. *Comm. in Partial Differential Equations*, 26(1–2):101–174, 2001.
- [Roc70] R. T. Rockafellar. On the maximal monotonicity of subdifferential mappings. *Pacific Journal of Mathematics*, 33:209–216, 1970.
- [vH16] Ramon van Handel. *Probability in High Dimension*. 2016.
- [Vil08] Cédric Villani. *Optimal transport, old and new*. 2008.
- [Vis18] Nisheeth K. Vishnoi. Geodesic convex optimization: Differentiation on manifolds, geodesics, and convexity. *CoRR*, abs/1806.06373, 2018.
- [Øk03] Bernt Øksendal. *Stochastic Differential Equations: An Introduction with Applications*. 2003.