

# STOCHASTIC DYNAMICS AND MACHINE LEARNING

— NOTES FOR A SEMINAR COURSE

YU CAO

*Shanghai Jiao Tong University*

Spring, 2025

## Abstract:

This material is part of a seminar course/mini-course designed for undergraduate students in the School of Mathematical Sciences at Shanghai Jiao Tong University in Spring 2025. A basic understanding of ordinary differential equations (ODEs), partial differential equations (PDEs), and probability theory is assumed. Familiarity with stochastic processes is preferable, but key concepts will be briefly reviewed in Chapter 1.

Over the course of these four weeks, the lectures will cover the following topics:

- Chapter 1: A review of basic stochastic processes and key preliminary results.
- Chapters 2: An exploration of neural dynamics in control problems. As a remark, this chapter does not aim to discuss the mathematics of control theory, but rather a discussion of backpropagation, together with its continuous-time and stochastic generalizations.
- Chapter 3: An exploration of neural dynamics in generative problems, including discussions on normalizing flows and probability flows, along with simple applications.
- Chapters 4: Fundamental results in Langevin dynamics, its origin from open classical systems and applications in sampling.
- Chapter 5: Discuss how one can use Langevin dynamics for generative models, specifically, the diffusion model.

The source code of some examples will be available at [https://gitee.com/yucaoyc/Spring25\\_StocDynML](https://gitee.com/yucaoyc/Spring25_StocDynML).

If you spot any errors or have suggestions for improvement, please don't hesitate to email me at [yucao@sjtu.edu.cn](mailto:yucao@sjtu.edu.cn).

# Contents

<b>1</b>	<b>Introduction and Preliminaries</b>	<b>4</b>
1.1	Introduction	4
1.2	Deterministic Dynamics and Evolution Operators	4
1.3	Stochastic Differential Equations and Fokker-Planck Equations	5
1.3.1	Brownian motion	5
1.3.2	Itô and Stratonovich integrals	6
1.3.3	Stochastic differential equations (SDEs)	7
1.3.4	Fokker-Planck equation	10
1.4	A Heuristic Introduction to Machine Learning	11
1.4.1	Discussion on the origins of data and loss functions	12
1.4.2	Examples of architecture: residual neural network and ODEs	13
1.4.3	Optimization algorithms: SGD	13
1.5	Notations	14
1.6	Further Readings	14
<b>2</b>	<b>Neural Dynamics for Control Problems</b>	<b>15</b>
2.1	Introduction	15
2.2	Backpropagation	16
2.3	Neural ODEs	17
2.3.1	Proof of Theorem 2.3	17
2.4	Neural SDEs	19
2.4.1	Proof of Theorem 2.7	20
2.4.2	Proof of lemmas	21
2.5	Further Readings	22
<b>3</b>	<b>Neural Dynamics for Generative Problems: (I) Normalizing Flows</b>	<b>23</b>
3.1	Background	23
3.2	Connection to Optimal Transport	23
3.3	Normalizing Flows	25
3.3.1	General architecture	25
3.3.2	Coupling flow and RealNVP	26
3.3.3	Loss functions for learning without data	26
3.3.4	Loss functions for learning with data	27
3.4	ODE Flows	27
3.4.1	Evolution of log-density function	27
3.4.2	Connection to (discrete) normalizing flows	28
3.4.3	Training of ODE flow for learning tasks without data	28
3.4.4	Training of ODE flow for learning tasks with data	29
3.5	Further Readings	29

<i>CONTENTS</i>	3
<b>4 Langevin Sampling</b>	<b>30</b>
4.1 Generalized Langevin Equation and Underdamped Langevin Dynamics . . . . .	30
4.2 Over-damped Langevin Sampling . . . . .	31
4.3 Further Readings . . . . .	32
<b>5 Neural Dynamics for Generative Problems: (II) Diffusion Models</b>	<b>33</b>
5.1 SDE-based Diffusion Models . . . . .	34
5.2 ODE-based Models . . . . .	35
5.3 Further Readings . . . . .	37
<b>A Erratum (which has been corrected)</b>	<b>40</b>

# Chapter 1

## Introduction and Preliminaries

### 1.1 Introduction

The study of dynamical systems and machine learning has been very closely connected [8]. This is a seminar course for stochastic dynamics (including deterministic dynamics) and recent machine learning techniques. Before we delve into topics, we shall review results about dynamical systems, e.g, notions like evolution operator, and review basic concepts from Stochastic Differential Equations (SDEs). Then we will provide a very brief introduction to machine learning. As a remark, this mini-course does not aim to be an introductory course for machine learning.

### 1.2 Deterministic Dynamics and Evolution Operators

Suppose we consider a deterministic dynamics,

$$\partial_t \rho(t) = \mathcal{L}(t)\rho(t),$$

where the state space of  $\rho(t)$ , denoted as  $\mathcal{S}$ , can be a finite-dimensional vector space, or a proper functional space. The operator  $\mathcal{L}(t) : \mathcal{S} \rightarrow \mathcal{S}$  can be time-dependent. Then the solution can be written as

$$\rho(t) = U(t, s) (\rho(s)), \quad U(t, s) = \mathcal{T} e^{\int_s^t \mathcal{L}(r) dr},$$

where  $\mathcal{T}$  is the time chronological ordering operator, namely

$$\mathcal{T}(A(t_1)A(t_2)) = \begin{cases} A(t_1)A(t_2) & \text{if } t_1 \geq t_2 \\ A(t_2)A(t_1) & \text{if } t_1 < t_2 \end{cases}$$

for appropriate time-dependent operators  $t \mapsto A(t)$ .  $U(t, s)$  can be approximated by

$$U(t, s) \approx e^{\Delta t \mathcal{L}(t_{N-1})} e^{\Delta t \mathcal{L}(t_{N-2})} \dots e^{\Delta t \mathcal{L}(t_0)},$$

where  $t_k = s + k\Delta t$  and  $\Delta t = \frac{t-s}{N}$ . This is well-known as Lie's first-order algorithm. Generalizations include Strang splitting scheme (2 nd order), Lie-Trotter-Suzuki splitting schemes (arbitrary order), and et al.

These abstract results are valid for both general ODEs and PDEs. We remark that if  $\mathcal{L}(t)$  is anti-Hermitian (as in the case of Schrödinger's equation in quantum mechanics), then  $U(t, s)$  represents unitary dynamics.

**Lemma 1.1.**

$$\frac{d}{dt} U(t, s) = \mathcal{L}(t)U(t, s), \quad \frac{d}{ds} U(t, s) = U(t, s)(-\mathcal{L}(s)). \quad (1.1)$$

A common situation is that  $\mathcal{L}$  can be split into two parts:

$$\mathcal{L}(t) = \mathcal{L}_A(t) + \mathcal{L}_B(t).$$

A useful tool is called the Duhamel's principle:

**Theorem 1.2** (Duhamel's principle).

$$\rho(t) = U_A(t, s)\rho(s) + \int_s^t U_A(t, r) \mathcal{L}_B(r) \rho(r) dr, \quad (1.2)$$

where  $U_A(t, s) := \mathcal{T}\left(e^{\int_s^t \mathcal{L}_A(r) dr}\right)$  is the evolution operator for  $\mathcal{L}_A$  only.

*Proof.* Hint: Duhamel's principle can be proved by introducing integrating factors; consider estimating  $\frac{d}{dt}(U_A(s, t)\rho(t))$  and write it in the integral form.  $\square$

## 1.3 Stochastic Differential Equations and Fokker-Planck Equations

This section is prepared by partially referring to the textbook [12].

### 1.3.1 Brownian motion

Back into 1820s, Robert Brown (a Scottish botanist) discovered the motion of pollen particle in fluid<sup>1</sup>. Later Albert Einstein discovered the equation of Brownian motion in 1905 (the other two most famous results in the same year by Einstein is special relativity and light quanta) [32]. The mathematical theory of Brownian motion was later formulated by Wiener [12]. Brownian motion is arguably the most famous continuous-time stochastic process. But what is a stochastic process then?

**Definition 1.3** (An informal definition of stochastic process). Suppose there is a probability space  $(\Omega, \mathcal{F}, \mathcal{P})$ , where  $\Omega$  is the event space,  $\mathcal{F}$  is the filtration, and  $\mathcal{P}$  is the probability measure<sup>2</sup>. A stochastic process  $\{X(t)\}_{t \in \mathfrak{T}}$  is a time-parameterized family of random variables, where  $\mathfrak{T}$  is the time index, which can be either discrete or continuous.

For instance, the location of a pollen inside the fluid at time  $t$  is a random variable, which can be modeled via  $X(t)$ ; for the event  $\omega \in \Omega$ ,  $X(t, \omega)$  is the location of pollen particle if  $\omega$  “happens”. The stock price of amazon in year 2030 is random, which can be modeled via a random variable  $X(t = 2030)$ ; the stock price of Tencent in year 2030 is also random, which can be modeled via a random variable  $\tilde{X}(t = 2030)$ ; clearly, these stock prices can be imagined to share the same probability space  $\Omega$  (which can be thought of as all possibilities of stocks on earth, even though we cannot literally list all outcomes on a paper). This example explains that for the same space  $\Omega$ , we can have different stochastic processes.

**Definition 1.4** (Brown motion). Brownian motion is a stochastic process with the following properties:

1.  $B(t) - B(s)$  ( $t > s$ ) is independent of all histories up to time  $s$ ;
2.  $B(t) - B(s)$  is a normal random variable  $\mathcal{N}(0, t - s)$ ;
3.  $B(t)$  is continuous in time  $t$ .

Therefore, the transition probability

$$P(y, t|x, s) = P(B(t) = y|B(s) = x) = \frac{1}{\sqrt{2\pi(t-s)}} \exp\left(-\frac{(y-x)^2}{2(t-s)}\right).$$

As you may observe in Figure 1.1, the Brownian motion path is a bit “rough”.

<sup>1</sup>See e.g., <https://www.britannica.com/science/Brownian-motion>

<sup>2</sup> $\mathcal{F}$  is called the filtration, which can be regarded as a technical part to adapt this notion into the rigorous measure theory; we will disregard the details for simplicity.

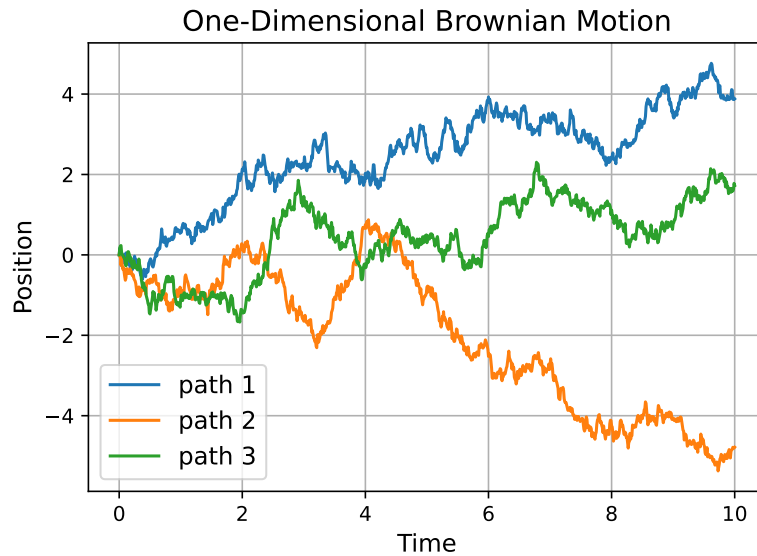


Figure 1.1: Example of Brownian motion trajectories

**Theorem 1.5** (Properties of Brownian motion; see [12, Chp. 3, page 64]).

- (i) Brownian motion trajectories are continuous but not differentiable.
- (ii) Brownian motion trajectories have infinite total variation.

### 1.3.2 Itô and Stratonovich integrals

Now we need to make sense of stochastic calculus in a heuristic way. For any given continuous function  $f$ , the integral of Brownian motion can be clearly approximated as follows:

$$\int_0^T f(B(s)) \, ds = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=0}^{N-1} f(B(s_k)), \quad (s_k = k\Delta s, \Delta s = \frac{T}{N}).$$

As Brownian motion  $s \mapsto B(s, \omega)$  is a continuous function (also uniformly continuous), the above definition (in the pointwise limit sense) is mathematically valid.

*Remark.* To understand why we need this, imagine the following artificial situation. Suppose a certain stock price exactly follows the Brownian motion, and someone bet with you that for every minute in the next hour, for each outcome  $B(s)$  you will earn  $f(B(s))\Delta s$  dollars (of course if the value is negative, it means you lose dollars). The time interval is  $\frac{1}{60} = \Delta s$ , and your total asset is exactly  $\frac{1}{N} \sum_{k=0}^{N-1} f(B(s_k))$ . Of course, it is reasonable to model this problem using continuous time as time interval  $\Delta s$  is small. This will motivate the definition of the integral  $\int_a^b f(B(s)) \, ds$  on the interval  $[a, b]$ .

**Exercise 1.6.** Compute the probability of events  $\left\{ \omega : \int_0^1 B(s, \omega) \, ds > 1 \right\}$ .

*Proof.* Hint: first prove that  $\int_0^1 B(s, \omega) \, ds$  is a normal random variable with mean 0 and variance 1/3.  $\square$

However, things become tricky when we want to take the integral of “ $\int_0^T X(t) \times dB(t)$ ”. In this case,

$dB(t)$  kind of plays the role as time-rescaling. Suppose that  $G(t)$  is a given differentiable function, then

$$\begin{aligned} \int_0^T X(t) dG(t) &= \int_0^T X(t) G'(t) dt \approx \sum_{k=0}^{N-1} X(s_k) \int_{s_k}^{s_{k+1}} G'(s) ds \\ &= \sum_{k=0}^{N-1} X(s_k) (G(s_{k+1}) - G(s_k)). \end{aligned} \quad (1.3)$$

Similarly, we can also approximate the integral via

$$\begin{aligned} \int_0^T X(t) dG(t) &= \int_0^T X(t) G'(t) dt \approx \sum_{k=0}^{N-1} \frac{X(s_k) + X(s_{k+1})}{2} \int_{s_k}^{s_{k+1}} G'(s) ds \\ &= \sum_{k=0}^{N-1} \frac{X(s_k) + X(s_{k+1})}{2} (G(s_{k+1}) - G(s_k)). \end{aligned} \quad (1.4)$$

**Exercise 1.7.** If  $G$  has finite total variation (e.g., when  $G$  is continuously differentiable), then the approximations in (1.3) and (1.4) are the same in the large  $N$  limit.

However, the above two approximations are distinct for Brownian motions (as Brownian motion has infinite total variation).

$$\begin{aligned} \text{It\^o integral : } \int_0^T X(t) dB(t) &= \lim_{N \rightarrow \infty} \sum_{k=0}^{N-1} X(s_k) (B(s_{k+1}) - B(s_k)) \\ \text{Stratonovich integral : } \int_0^T X(t) \circ dB(t) &= \lim_{N \rightarrow \infty} \sum_{k=0}^{N-1} \frac{X(s_k) + X(s_{k+1})}{2} (B(s_{k+1}) - B(s_k)) \end{aligned} \quad (1.5)$$

Using the above definitions, we can verify the following two results.

**Example 1.8.** For It\^o and Stratonovich integrals, we have the following results respectively:

$$\mathbb{E} \left( \int_0^T B(t) dB(t) \right) = 0, \quad \mathbb{E} \left( \int_0^T B(t) \circ dB(t) \right) = \frac{T}{2}.$$

Therefore, these two definitions are clearly different.

**Lemma 1.9** (Properties of It\^o integral). Suppose that  $X(t)$  is  $\mathcal{F}(t)$ -adapted, then

- **(Zero mean)**  $\mathbb{E} \left( \int_0^T X(t) dB(t) \right) = 0.$
- **(It\^o isometry)**  $\mathbb{E} \left( \int_0^T X(t) dB(t) \right)^2 = \int_0^T \mathbb{E}(X(t)^2) dt.$

These can be proved (at least formally) directly using definitions (1.5). The proofs are left as exercises.

### 1.3.3 Stochastic differential equations (SDEs)

**Definition 1.10** (Stochastic Differential Equations, SDEs). Suppose a stochastic process  $X(t)$  satisfies the following:

$$X(t) = X(0) + \int_0^t b(X(s), s) ds + \int_0^t \sigma(X(s), s) dB(s) \quad \text{with probability one,} \quad (1.6)$$

where  $b : \mathbb{R}^d \times \mathbb{R} \rightarrow \mathbb{R}^d$ ,  $\sigma : \mathbb{R}^d \times \mathbb{R} \rightarrow \mathbb{R}^{d \times d}$ ,<sup>3</sup> then  $X(t)$  is called the solution to the following SDE

$$dX(t) = b(X(t), t) dt + \sigma(X(t), t) dB(t). \quad (1.7)$$

If the integral inside (1.6) is replaced by Stratonovich integrals, then the corresponding SDE is denoted as

$$dX(t) = b(X(t), t) dt + \sigma(X(t), t) \circ dB(t). \quad (1.8)$$

<sup>3</sup>Diffusion coefficient  $\sigma$  does not necessarily needs to be a square matrix, but we use simpler cases without loss of generality.



**Theorem 1.11** (Relationship between Itô and Stratonovich SDEs).

(i) Suppose  $X(t)$  follows the Stratonovich SDEs

$$dX(t) = b(X(t), t) dt + \sigma(X(t), t) \circ dB(t),$$

then  $X(t)$  also follows the following Itô SDE (with a different drift function):

$$dX(t) = (b(X(t), t) + \theta(X(t), t)) dt + \sigma(X(t), t) dB(t), \quad (1.9)$$

where  $\theta_i(X(t), t) := \frac{1}{2} \sum_{j,k} \partial_{x_k} \sigma_{i,j}(X(t), t) \sigma_{k,j}(X(t), t)$ .

(ii) For 1D case ( $d = 1$ ), one has the following simplified expression:

$$dX(t) = \left( b(X(t), t) + \frac{1}{2} \partial_x \sigma(X(t), t) \sigma(X(t), t) \right) dt + \sigma(X(t), t) dB(t).$$

*Proof.* For Stratonovich SDEs,

$$\begin{aligned} \Delta X(t) &= b(X(t), t) \Delta t + \frac{\sigma(X(t + \Delta t), t + \Delta t) + \sigma(X(t), t)}{2} \Delta B(t) + o(\Delta t) \\ &= b(X(t), t) \Delta t + \sigma(X(t), t) \Delta B(t) + \frac{\sigma(X(t + \Delta t), t + \Delta t) - \sigma(X(t), t)}{2} \Delta B(t) + o(\Delta t) \\ &= b(X(t), t) \Delta t + \sigma(X(t), t) \Delta B(t) + \theta(X(t), t) \Delta t + o(\Delta t), \end{aligned}$$

where

$$\begin{aligned} \theta_i(X(t), t) &= \sum_j \frac{\sigma_{i,j}(X(t + \Delta t), t + \Delta t) - \sigma_{i,j}(X(t), t)}{2} (\Delta B(t))_j \\ &= \frac{1}{2} \sum_{j,k,\ell} \partial_{x_k} \sigma_{i,j}(X(t), t) \sigma_{k,\ell}(X(t), t) (\Delta B(t))_\ell (\Delta B(t))_j + o(\Delta t) \\ &= \frac{1}{2} \sum_{j,k} \partial_{x_k} \sigma_{i,j}(X(t), t) \sigma_{k,j}(X(t), t) \Delta t + o(\Delta t) \end{aligned}$$

where in the last line, we used the fact that  $(dB_j(t) dB_k(t)) = \delta_{j,k} dt$  (which can be more rigorously established using the law of large numbers). The proof is complete after simplifying notations.  $\square$

**Theorem 1.12** (Itô's formula and Stratonovich's formula). Suppose  $f \in C^2(\mathbb{R}^d \times \mathbb{R})$ .

(i) Suppose  $X(t)$  follows the Itô SDE (1.7), then

$$\begin{aligned} df(X(t), t) &= \partial_t f(X(t), t) dt + \nabla_x f(X(t), t) \cdot dX(t) + \frac{1}{2} (\sigma \sigma^\top)(X(t), t) : \nabla_x^2 f(X(t), t) dt \\ &\equiv \partial_t f(X(t), t) dt + \nabla_x f(X(t), t) \cdot b(X(t), t) dt + (\nabla_x f(X(t), t))^\top \sigma(X(t), t) dB(t) \\ &\quad + \frac{1}{2} (\sigma \sigma^\top)(X(t), t) : \nabla_x^2 f(X(t), t) dt. \end{aligned} \quad (1.10)$$

We will denote  $D(x, t) := \frac{1}{2} (\sigma \sigma^\top)(x, t)$  and the notation  $D : \nabla^2 f$  means  $\sum_{i,j=1}^d D_{i,j} \nabla_{x_i} \nabla_{x_j} f$ .

(ii) Suppose  $X(t)$  follows the Stratonovich SDE (1.8), then

$$\begin{aligned} df(X(t), t) &= \partial_t f(X(t), t) dt + \nabla_x f(X(t), t) \cdot dX(t) \\ &\equiv \partial_t f(X(t), t) dt + \nabla_x f(X(t), t) \cdot b(X(t), t) dt + \left( \nabla_x f(X(t), t) \right)^\top \sigma(X(t), t) \circ dB(t). \end{aligned} \quad (1.11)$$

*Proof.* Below are heuristic proofs for illustrating main ideas.

Proof of Part (i). By Taylor's expansion,

$$\begin{aligned}
& f(X(t + \Delta t), t + \Delta t) - f(X(t), t) \\
&= \partial_t f(X(t), t) \Delta t + \nabla_x f(X(t), t) \cdot (X(t + \Delta t) - X(t)) \\
&\quad + \frac{1}{2} (X(t + \Delta t) - X(t))^\top \nabla_x^2 f(X(t), t) (X(t + \Delta t) - X(t)) + o(\Delta t) \\
&= \partial_t f(X(t), t) \Delta t + \nabla_x f(X(t), t) \cdot (X(t + \Delta t) - X(t)) \\
&\quad + \frac{1}{2} \Delta B(t)^\top \sigma(X(t), t)^\top \nabla_x^2 f(X(t), t) (\sigma(X(t), t) \Delta B(t)) + o(\Delta t).
\end{aligned}$$

The cross product terms  $(\Delta B(t))_i (\Delta B(t))_j$  ( $i \neq j$ ) has negligible contributions. Hence,

$$\begin{aligned}
& f(X(t + \Delta t), t + \Delta t) - f(X(t), t) \\
&= \partial_t f(X(t), t) \Delta t + \nabla_x f(X(t), t) \cdot (X(t + \Delta t) - X(t)) \\
&\quad + \frac{1}{2} \sum_{j, \ell} \left( \nabla_x^2 f(X(t), t) \right)_{\ell, j} \left( \sigma \sigma^\top(X(t), t) \right)_{\ell, j} \Delta t + o(\Delta t).
\end{aligned}$$

By writing the l.h.s as  $df(X(t), t) \Delta t$ , and by matching terms, one can obtain (1.7).

Proof of Part (ii). Suppose  $X(t)$  follows the Stratonovich SDE (1.8), then it follows the SDE (1.9). By part (i), we have

$$\begin{aligned}
df(X(t), t) &= \partial_t f(X(t), t) dt + \nabla_x f(X(t), t) \cdot \left( b(X(t), t) + \theta(X(t), t) \right) dt \\
&\quad + \left( \nabla_x f(X(t), t) \right)^\top \sigma(X(t), t) dB(t) \\
&\quad + \frac{1}{2} (\sigma \sigma^\top)(X(t), t) : \nabla_x^2 f(X(t), t) dt.
\end{aligned} \tag{1.12}$$

Let us transform this back into Stratonovich SDE. Imagine that the corresponding Stratonovich SDE has the form

$$\begin{aligned}
df(X(t), t) &= \partial_t f(X(t), t) dt + \nabla_x f(X(t), t) \cdot \left( b(X(t), t) + \theta(X(t), t) \right) dt \\
&\quad + \frac{1}{2} (\sigma \sigma^\top)(X(t), t) : \nabla_x^2 f(X(t), t) dt + \kappa(X(t), t) dt \\
&\quad + \left( \nabla_x f(X(t), t) \right)^\top \sigma(X(t), t) \circ dB(t),
\end{aligned} \tag{1.13}$$

where  $\kappa$  is a term to be determined. By following the same calculation as in Theorem 1.11, we know that the relevant Ito's SDE for the last equation is

$$\begin{aligned}
df(X(t), t) &= \partial_t f(X(t), t) dt + \nabla_x f(X(t), t) \cdot \left( b(X(t), t) + \theta(X(t), t) \right) dt \\
&\quad + \frac{1}{2} (\sigma \sigma^\top)(X(t), t) : \nabla_x^2 f(X(t), t) dt + \kappa(X(t), t) dt \\
&\quad + \left( \nabla_x f(X(t), t) \right)^\top \sigma(X(t), t) dB(t) \\
&\quad + \frac{1}{2} \sum_{i, j, k} \partial_{x_k} (\partial_{x_i} f \sigma_{i, j}) \sigma_{k, j} \Big|_{(X(t), t)} dt.
\end{aligned} \tag{1.14}$$

By matching terms in (1.12) and (1.14), we have

$$\kappa(X(t), t) + \frac{1}{2} \sum_{i, j, k} \partial_{x_k} (\partial_{x_i} f \sigma_{i, j}) \sigma_{k, j} \Big|_{(X(t), t)} = 0,$$

which means,

$$\kappa = \frac{1}{2} \sum_{i,j,k} \left( -\partial_{i,k} f \sigma_{i,j} \sigma_{k,j} - \partial_i f \partial_k \sigma_{i,j} \sigma_{k,j} \right) = -D : \nabla^2 f - (\nabla_x f) \cdot \theta.$$

Therefore, (1.13) can be simplified to (1.11).  $\square$

**Example 1.13.** A widely used example is the Ornstein-Uhlenbeck process:

$$dX(t) = aX(t) dt + \sigma dB(t), \tag{1.15}$$

whose solution can be explicitly written as

$$X(t) = e^{at} X_0 + \sigma \int_0^t e^{a(t-s)} dB(s).$$

The mean is  $\mathbb{E}[X(t)] = e^{at} \mathbb{E}[X_0]$ , the variance is

$$\text{Var}(X(t)) = e^{2at} \text{Var}(X_0) + \sigma^2 \int_0^t e^{2a(t-s)} ds. \tag{1.16}$$

*Proof.* By Lemma 1.9 and Theorem 1.12.  $\square$

**Exercise 1.14.** Write code to simulate the above Ornstein-Uhlenbeck (1.15) and verify (1.16) numerically.

**Exercise 1.15** (1D Geometric Brownian Motion). The geometric Brownian motion is the SDE behind Black–Scholes equation which is widely used as a typical model in finance:

$$dX(t) = \alpha X(t) dt + \sigma X(t) \circ dB(t). \tag{1.17}$$

(i) Prove that the SDE (1.17) is equivalent to the following Itô's SDE form

$$dX(t) = \left( \alpha + \frac{1}{2} \sigma^2 \right) X(t) dt + \sigma X(t) dB(t).$$

(ii) Prove that the explicit solution of the above SDE (1.17) is

$$X(t) = X(0) \exp(\alpha t + \sigma B(t)) \equiv X(0) \exp\left(\int_0^t \alpha + \sigma \circ dB(s)\right).$$

### 1.3.4 Fokker-Planck equation

The probability density function of  $X(t)$ , denoted as  $p(x, t)$ , satisfies a certain type of PDEs, well-known as Fokker-Planck equations. Let us consider the Itô SDE for example. Suppose  $g$  is any compactly supported smooth function, then

$$\int g(x) p(x, t) dx = \mathbb{E}(g(X(t))).$$

Let us take derivatives on both sides,

$$\begin{aligned} & \left( \int g(x) \partial_t p(x, t) dx \right) \Delta t + o(\Delta t) = \mathbb{E} \left( \Delta g(X(t)) \right) \\ & = \mathbb{E} \left( \nabla_x g(X(t)) \cdot b(X(t), t) + \frac{1}{2} (\sigma \sigma^\top)(X(t), t) : \nabla_x^2 g(X(t)) \right) \Delta t + o(\Delta t) \\ & = \left( \int \nabla_x g(x) \cdot b(x, t) p(x, t) + D_{i,j}(x, t) \nabla_{x_i, x_j} g(x) p(x, t) dx \right) \Delta t + o(\Delta t) \\ & = \left( \int -g(x) \nabla_x \cdot (b(x, t) p(x, t)) + g(x) \nabla_{x_i, x_j} (D_{i,j}(x, t) p(x, t)) dx \right) \Delta t + o(\Delta t). \end{aligned}$$

In the limit of small  $\Delta t$  and since  $g$  is arbitrary, one has:

**Theorem 1.16.** The evolution of  $p(t, \cdot)$  for Itô SDE (1.7) satisfies the following PDE

$$\begin{aligned} \partial_t p(x, t) &= -\nabla_x \cdot (b(x, t)p(x, t)) + \sum_{i,j} \nabla_{x_i, x_j} (D_{i,j}(x, t)p(x, t)) \\ &\equiv -\nabla_x \cdot (b(x, t)p(x, t)) + \nabla_x^2 : (D(x, t)p(x, t)). \end{aligned} \quad (1.18)$$

**Example 1.17.** When  $\sigma(x, t) = \sqrt{2} \mathbf{I}_d$  is independent of time and space, then the Fokker-Planck equation can be simplified as

$$\partial_t p(x, t) = -\nabla_x \cdot (b(x, t)p(x, t)) + \Delta p(x, t).$$

**Exercise 1.18.** Verify that when  $b(x, t) = -\nabla U(x)$  and  $\sigma(x, t) = \sqrt{2} \mathbf{I}_d$ , where  $U$  is an appropriate potential function, one stationary solution of the Fokker-Planck equation is  $p_\infty(x) = e^{-U(x)}/\mathcal{Z}$  where  $\mathcal{Z} = \int_{\mathbb{R}^d} e^{-U(x)} dx$  is a normalizing constant. This Itô SDE is well-known as the overdamped Langevin dynamics, which has been widely used for sampling in statistics (see Chapter 4).

## 1.4 A Heuristic Introduction to Machine Learning

Machine learning problems are typically formulated in terms of an optimization problem

$$\min_{\theta} L(\theta) = \sum_{i=1}^N \omega_i f_i(\theta). \quad (1.19)$$

where  $\omega_i \in \mathbb{R}$  is the weight,  $f_i$  is a family of functions, and  $\theta$  is the parameter to find.

The success of machine learning relies on at least the following:

- high-quality data  $x_i$ ; the data can come from either realistic problems (for many AI products), or comes from generative process via some existing processes (perhaps more common in AI for Science);
- a good way to turn the original problems into an optimization problem with effective parameterization (e.g., a neural network like below); this means a good choice of  $f$  overall;
- an efficient optimization algorithm.

We will go through these ingredients for different tasks in the following chapters, as there is no universal practical way to design  $f$  and find  $x_i$ .

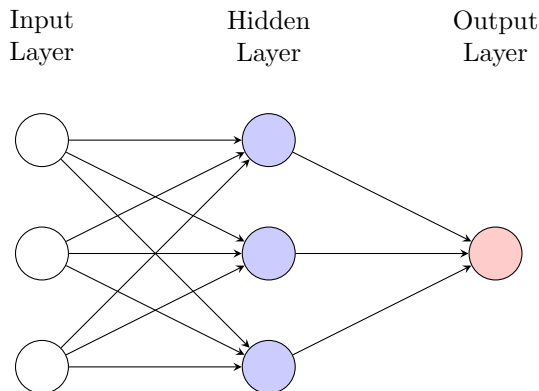


Figure 1.2: A basic structure of a two-layer neural network.

### 1.4.1 Discussion on the origins of data and loss functions

The data can come from either the realistic problems or from generated data.

**Example 1.19** (Classification Problems). Suppose  $x_i$  represents image data, and suppose  $x \mapsto g(x, \theta)$  is a parameterized mapping that decides the probability of certain events to occur, e.g., the probability that the medical image contains a cancer. Suppose we have some labeled data  $y_i \in \{0, 1\}$  that represent whether an image contains a cancer (this is where the data comes from). Therefore, to minimize the gap between prediction and the actual data, one may want to minimize the mean-square error

$$\frac{1}{N} \sum_{i=1}^N \underbrace{|g(x_i, \theta) - y_i|^2}_{\text{denoted as } f_i(\theta)}.$$

Of course, one may use any monotone function  $\varphi$ , and design a loss function as

$$\frac{1}{N} \sum_{i=1}^N \underbrace{\varphi(g(x_i, \theta) - y_i)}_{\text{denoted as } f_i(\theta)}.$$

**Example 1.20** (Deep Ritz [9]). Suppose one wants to solve the Poisson's equation

$$\begin{cases} -\Delta u = \phi & \text{in } \Omega; \\ u = 0, & \text{on } \partial\Omega. \end{cases}$$

Then one can show that it is equivalent to solve the following optimization problem

$$\min_u \frac{1}{2} \int_{\Omega} |\nabla u(x)|^2 dx - \int_{\Omega} u(x)\phi(x) dx.$$

This PDE model finds wide applications, e.g., in electrostatics. <sup>4</sup> To incorporate the boundary conditions, we may minimize the following in practice

$$\min_u \frac{1}{2} \int_{\Omega} |\nabla u(x)|^2 dx - \int_{\Omega} u(x)\phi(x) dx + \lambda \int_{\partial\Omega} |u(x)|^2 dx.$$

where  $\lambda$  is a hyper-parameter for penalty term arising from boundary conditions. Then we can parameterize the solution via a class of functions  $u_{\theta}$ , and also deterministically or randomly pick sample points  $x_i \in \Omega$  and  $z_i \in \partial\Omega$  to approximate the integral via

$$\min_{\theta} \sum_i \omega_i \left( \frac{1}{2} |\nabla u_{\theta}(x_i)|^2 - u_{\theta}(x_i)\phi(x_i) \right) + \lambda \sum_i \bar{\omega}_i (u_{\theta}(z_i))^2,$$

where  $\omega_i, \bar{\omega}_i$  is the weight in the quadrature. This again falls into the above form (1.19). <sup>5</sup>

**Example 1.21** (Physics-Informed Neural Network [25]). Suppose one wants to solve the differential equation

$$\begin{cases} \mathcal{L}(u) = 0, & \text{in } \Omega; \\ u = 0, & \text{on } \partial\Omega. \end{cases}$$

Imagine that we pick sample points  $\{x_i\}_{i=1}^N$  from the domain  $\Omega$  and  $\{z_i\}_{i=1}^M$  from the boundary  $\partial\Omega$ , then one may minimize the following

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N (\mathcal{L}(u_{\theta})(x_i))^2 + \lambda \frac{1}{M} \sum_{i=1}^M (u_{\theta}(z_i) - 0)^2,$$

and  $\lambda$  is a hyper-parameter. This again falls into the above form (1.19). <sup>6</sup>

<sup>4</sup> [https://en.wikipedia.org/wiki/Poisson%27s\\_equation](https://en.wikipedia.org/wiki/Poisson%27s_equation).

<sup>5</sup> See [https://github.com/yucaoyc/math6008\\_num\\_pde/blob/main/chp08-1D-Poisson.ipynb](https://github.com/yucaoyc/math6008_num_pde/blob/main/chp08-1D-Poisson.ipynb) for illustrative codes.

<sup>6</sup> see the above footnote.

### 1.4.2 Examples of architecture: residual neural network and ODEs

There are many architectures, from the simplest fully connected neural network, to convolutional neural network (CNN), to Recurrent Neural Network (RNN), to Residual Network (ResNet), and to more recent ones like Transformers. The study of their performances and their expressibility is an active research field.

Though perhaps not a one-to-one correspondence, there is a tight connection between dynamical systems and neural network architectures. We follow the discussion in [8] below. The deep Residual network takes the following form

$$y_\ell = h(z_\ell) + \mathcal{F}(z_\ell, W_\ell), \quad z_{\ell+1} = g(y_\ell), \quad \text{or in the form } z_{\ell+1} = g(h(z_\ell) + \mathcal{F}(z_\ell, W_\ell)),$$

where  $z_\ell$  are intermediate layers, and  $y_\ell$  are auxiliary variables. Suppose  $\mathcal{F}$  has a small magnitude, and  $g = h$  are identity mapping, then

$$z_{\ell+1} = z_\ell + \mathcal{F}(z_\ell, W_\ell),$$

can be regarded as Euler-discretization of a certain ODE dynamics. There are other connections — normalizing flow can be regarded as discretization of ODE/probability flows; see later chapters.

### 1.4.3 Optimization algorithms: SGD

The Gradient Descent (GD) is perhaps one of the easiest algorithms to solve an optimization problem. However, GD algorithm is well-known to fail to find a good estimates of global minimum when the loss function is not convex. For machine learning problems, typically, we use Stochastic Gradient Descent (SGD), or other more complicated algorithms [26]. We shall simply explain SGD for illustration.

Suppose the loss function is (1.19) and suppose the weight  $\omega_i = 1/N$  for simplicity. The SGD refers to the following iterative algorithms

$$\theta_{k+1} = \theta_k - h_k \nabla_\theta f_{i_k}(\theta_k),$$

where the index  $i_k$  is randomly drawn with probability  $\omega_i = 1/N$ , and  $h_k > 0$  is called the learning rate. Note that  $\nabla f_{i_k}(\theta)$  is random variable with unbiased mean  $\nabla_\theta L(\theta)$ :

$$\mathbb{E}[\nabla_\theta f_{i_k}(\theta)] = \nabla_\theta L(\theta).$$

Of course, there is an intermediate region between GD and SGD, which is called Mini-batch gradient descent — which is less fluctuation than SGD, and has more exploring power than GD [26].

**Theorem 1.22** (See e.g., [27]). Assume that

- The loss function  $L$  is bounded from below by  $L^*$  (the global minimum value);
- The gradient of loss function  $\nabla L$  is  $L$ -Lipschitz;
- $\mathbb{E}[\|\nabla f_i(\theta)\|^2] \leq \sigma^2$  for any  $\theta$ .

Then

$$\min_{k=0,1,\dots,K-1} \mathbb{E}\|\nabla L(\theta_k)\|^2 \leq \frac{L(\theta_0) - L^*}{\sum_{k=0}^{K-1} h_k} + \frac{L\sigma^2}{2} \frac{\sum_{k=0}^{K-1} h_k^2}{\sum_{k=0}^{K-1} h_k}.$$

If  $h_k = h$  is independent of  $k$ , then the error scales like  $\mathcal{O}(\frac{1}{K} + h)$ .

*Proof.* See slides 3-6 from [27]. □

**Exercise 1.23.** When  $h_k = h/k$  and  $h_k = h/\sqrt{k}$ , please find the scaling of error with respect to  $K$ . Answers can be found in [27].

**Theorem 1.24** (See e.g., [27]). Under the same assumption as Theorem 1.22, if we further assume that  $\nabla^2 L(\theta) \geq \mu \mathbf{I}$  for any  $\theta$  (namely, the strongly convex assumption holds) and assume that  $h_k = h < \frac{2}{\mu}$  (small enough), then

$$\mathbb{E}[\|\theta_k - \theta^*\|^2] \leq (1 - 2h\mu)^k \|\theta_0 - \theta^*\|^2 + \frac{h\sigma^2}{2\mu}, \quad \forall k.$$

This suggests that for a large enough  $k$ ,  $\theta_k$  is close to  $\theta^*$  with distance  $\sqrt{h\sigma^2/(2\mu)}$ . By Markov inequality,

$$\mathbb{P}(\|\theta_k - \theta^*\| \geq \epsilon) \leq \frac{\mathbb{E}[\|\theta_k - \theta^*\|^2]}{\epsilon^2} \approx \frac{h\sigma^2}{2\mu\epsilon^2}.$$

For instance, when  $\epsilon = 10\sqrt{h\sigma^2/(2\mu)}$ , one has  $\mathbb{P}(\|\theta_k - \theta^*\| \geq \epsilon) \leq 1\%$  for large enough  $k$ .

*Proof.*

$$\begin{aligned} \mathbb{E}[\|\theta_{k+1} - \theta^*\|^2] &= \mathbb{E}[\|\theta_k - \theta^* - h\nabla_{\theta} f_{i_k}(\theta_k)\|^2] \\ &= \mathbb{E}\left[\|\theta_k - \theta^*\|^2 - 2h(\theta_k - \theta^*) \cdot \nabla_{\theta} f_{i_k}(\theta_k) + h^2 \|\nabla_{\theta} f_{i_k}(\theta_k)\|^2\right] \\ &\leq \mathbb{E}[\|\theta_k - \theta^*\|^2 - 2h\mu\|\theta_k - \theta^*\|^2] + h^2\sigma^2 \\ &\leq (1 - 2h\mu)\mathbb{E}[\|\theta_k - \theta^*\|^2] + h^2\sigma^2. \end{aligned}$$

By the iterative relation,

$$\begin{aligned} \mathbb{E}[\|\theta_k - \theta^*\|^2] &\leq (1 - 2h\mu)^k \|\theta_0 - \theta^*\|^2 + h^2\sigma^2((1 - 2h\mu)^0 + (1 - 2h\mu)^1 + \dots + (1 - 2h\mu)^{k-1}) \\ &\leq (1 - 2h\mu)^k \|\theta_0 - \theta^*\|^2 + \frac{h\sigma^2}{2\mu}. \end{aligned}$$

□

## 1.5 Notations

By default, we take the column-based vector form rather than row-based vector form. Suppose  $z \in \mathbb{R}^d \mapsto f(z) \in \mathbb{R}^d$ , then  $\nabla_z f(z)$  is a matrix whose  $(i, j)$ <sup>th</sup> element is  $\partial_{z_j} f_i(z)$ . Suppose  $z \in \mathbb{R}^d \mapsto F(z) \in \mathbb{R}^{d \times d}$ , then  $\nabla_z F(z)$  is a 3-fold tensor whose  $(i, j, k)$ <sup>th</sup> element is  $\partial_{z_k} F_{i,j}(z)$ . For any arbitrary vector  $v$ , we use notation

$$\nabla F v := \left[ \sum_k \partial_{z_k} F_{i,j} v_k \right]_{i,j}$$

which is a matrix.

## 1.6 Further Readings

- **Theory of Stochastic Calculus:** For further details on fundamental results of stochastic calculus, one may refer to e.g., [12, 24].
- **Numerical Methods of Stochastic Calculus:** For further details on numerical methods for stochastic dynamics, one may refer to [13].
- **Deep Learning:** For background on machine learning (in particular, deep learning), one may refer to e.g., [4]. For more optimization algorithms and practical challenges, one may refer to e.g., [26]. One may also refer to course materials e.g., by Mark Schmidt (UBC) at <https://www.cs.ubc.ca/~schmidtm/Courses/540-W19/>, and by Martin Jaggi and Nicolas Flammarion (EPFL) at [https://github.com/epfml/OptML\\_course](https://github.com/epfml/OptML_course).
- **Learning to Code a ML Project:** 深入浅出 PyTorch. <https://datawhalechina.github.io/thorough-pytorch/> This is a Chinese book that introduce you to code in PyTorch for machine learning problems.

# Chapter 2

# Neural Dynamics for Control Problems

## 2.1 Introduction

**Problem:** For many control problems, our goal is to find a certain dynamics, such that the resulting dynamics at a certain time  $T$  satisfies a certain condition, either matching a certain value, or minimizes a certain target value (namely, loss function).

The dynamics itself may have additional constraint itself.

- **(Deterministic Network).** Suppose we consider a neural network with layer depth  $\ell$  and width  $N_d$  and let us consider a simple setup:

$$f = f_\ell \circ f_{\ell-1} \circ \cdots \circ f_1$$

where each  $f_i : \mathbb{R}^{N_d} \rightarrow \mathbb{R}^{N_d}$  are parameterized functions. For simplicity, let us denote

$$y_1 = f_1(x), \quad y_2 = f_2(y_1), \quad \cdots, \quad y_\ell = f_\ell(y_{\ell-1}).$$

For the control problem, one wants the final output state to satisfy certain criterions, e.g., one may want to minimize

$$\min_{\theta} L(y_\ell) \tag{2.1}$$

where  $L$  is a certain loss function, and  $\theta$  are parameters inside functions  $f_i$ .

- **(Neural ODE).** If each layer takes the form of Residual Network and each  $f_i \approx \mathbf{I}$ , then one may regard this network as an approximation of an ordinary differential equation (ODE). Suppose we consider  $z(t, \theta) : \mathbb{R} \times \mathbb{R}^{N_p} \rightarrow \mathbb{R}^{N_d}$ :

$$\frac{dz(t, \theta)}{dt} = f(z(t, \theta), t, \theta), \tag{2.2}$$

where  $\theta$  is the control parameter, and  $(z, t, \theta) \mapsto f(z, t, \theta)$  is the ODE flow for consideration. The target is to find such a function  $f$  such that the output state  $z(T, \cdot)$  satisfies certain conditions.

- **(Neural SDE).** We can further extend the setup to the probability space  $\Omega$ . Suppose we consider  $X(t, \theta, \omega) : \mathbb{R} \times \mathbb{R}^{N_p} \times \Omega \rightarrow \mathbb{R}^{N_d}$ :

$$dX(t, \theta) = b(X(t, \theta), t, \theta) dt + \sigma(X(t, \theta), t, \theta) \circ dB(t), \tag{2.3}$$

where  $\theta$  is the control parameter, and  $(z, t, \theta) \mapsto b(z, t, \theta)$  is the drift and  $(z, t, \theta) \mapsto \sigma(z, t, \theta)$  is the diffusion coefficient. The target is again to find such a function  $b$  and  $\sigma$ , such that the output state  $X(T, \cdot)$  satisfies certain conditions.



## 2.2 Backpropagation

Our goal is to find a parameter to solve (2.1). In order to use many gradient-based optimization algorithm, we have to be able to compute the gradient of  $L(y_\ell)$  with respect to the parameters inside the layer. One widely used efficient method is called the backpropagation [4, Chapter 8], which we will review below. Suppose  $f_i$  has parameters  $\theta_j^{(i)}$  where the index  $j$  characterize the index of parameters for this layer.

By the chain rule, when we compute, e.g.,

$$\frac{dL(y_\ell)}{d\theta_\lambda^{(1)}} = \sum_{j_1, j_2, \dots, j_\ell=1, \dots, N_d} \frac{\partial L}{\partial (y_\ell)_{j_\ell}}(y_\ell) \cdot \frac{\partial (y_\ell)_{j_\ell}}{\partial (y_{\ell-1})_{j_{\ell-1}}}(y_{\ell-1}) \cdots \frac{\partial (y_1)_{j_1}}{\partial \theta_\lambda^{(1)}}(x). \quad (2.4)$$

For the forward process, one compute  $y_j$  in the forward order; to compute the gradient, one may start with  $\frac{\partial L}{\partial (y_\ell)_{j_\ell}}(y_\ell)$  and backwardly update the gradient value in a backward manner.

**Example 2.1.** Let us consider the fully-connected neural network without the bias term for simplicity:

$$y_k = f_k(y_{k-1}) \equiv \varsigma(W^{(k)}y_{k-1}),$$

where  $W^{(k)}$  is the weight matrix with size  $N_d \times N_d$  and  $\varsigma$  is an activation function. Then

$$\frac{\partial (y_k)_{j_k}}{\partial (y_{k-1})_{j_{k-1}}}(y_{k-1}) = \varsigma'(W^{(k)}y_{k-1})_{j_k} W_{j_k, j_{k-1}}^{(k)}.$$

- As one needs to save all intermediate steps  $\{y_k\}_{k=1,2,\dots,\ell}$ , one requires  $\mathcal{O}(N_d\ell)$  total memory cost. As we also need to store  $\left[\frac{dL(y_\ell)}{d\theta_\lambda^{(1)}}\right]_{\lambda=1,2,\dots,N_d^2}$ , the memory cost is  $\mathcal{O}(N_d^2\ell)$ , which is essentially the number of parameters. In any case, the memory clearly scales linearly with respect to the depth  $\ell$ .
- The total computational cost for the backpropagation of (2.4) is

$$\mathcal{O}(N_d^2\ell) = \mathcal{O}(\text{total number of parameters}),$$

where the cost of  $N_d^2$  comes from matrix-vector multiplication for each layer and inside (2.4).

**Exercise 2.2** (Comparison with FDM). Please compare the cost of backpropagation with the central finite difference method, e.g.,

$$\frac{dL(y_\ell)}{d\theta_\lambda^{(1)}} = \frac{L(y_\ell(\theta_\lambda^{(1)} + \epsilon)) - L(y_\ell(\theta_\lambda^{(1)} - \epsilon))}{2\epsilon} + \mathcal{O}(\epsilon^2).$$

In particular, prove that to compute the gradient of the loss with respect to all parameters in weight, one requires  $\mathcal{O}(N_d^4\ell^2) = \mathcal{O}(\text{total number of parameters}^2)$ ; see the discussion e.g., in [4, Chapter 8].

### Automatic Differentiation

Automatic differentiation is a well-implemented computer technique to compute gradient automatically based on the above chain rule. The forward-mode automatic differentiation computes (2.4) from right to left, and the reverse-mode automatic differentiation computes (2.4) from left to right. Whether one should use forward-mode or reverse-mode automatic differentiation depends on the dimension of input and output variables. For machine learning problems, one typically has a large number of parameters where the output dimension is small (or perhaps one). Then generally, reverse-mode automatic differentiation is preferred. But for a specific problem, it is better to analyze it case-by-case.

## 2.3 Neural ODEs

Our goal is to find the parameter  $\theta$  such that

$$\min_{\theta} L(z(T, \theta)),$$

where  $L$  is a scalar-valued function. Suppose that we use a gradient-based optimization method, e.g., gradient descent method, then we need to estimate  $\frac{dL(z(T, \theta))}{d\theta}$ . A core ingredient in Neural ODE [5] is to find an efficient way to compute this quantity.

**Theorem 2.3** (Results in [5]). We consider the ODE in (2.2) with the initial time  $T_0$  and the final time  $T$ . Suppose that  $z(T_0, \theta) = z_0$  is fixed, and  $L(\cdot)$ ,  $f(\cdot, \cdot, \cdot)$  are continuously differentiable. Then

$$\begin{cases} (\nabla_{\theta} L(z(T, \theta)))^{\top} = \int_{T_0}^T a(s)^{\top} \partial_{\theta} f(z(s, \theta), s, \theta) ds \\ \frac{d}{ds} a(s) = -(\partial_z f(z(s, \theta), s, \theta))^{\top} a(s), & s \in [T_0, T]; \\ a(T) = \nabla L(z(T, \theta)). \end{cases} \quad (2.5)$$

where  $\partial_{\theta} f(z(s, \theta), s, \theta) = [\partial_{\theta_1} f(z(s, \theta), s, \theta) \ \cdots \ \partial_{\theta_{N_p}} f(z(s, \theta), s, \theta)] \in \mathbb{R}^{N_d \times N_p}$ , and  $a(s) \in \mathbb{R}^{N_d \times 1}$ . Alternatively, one could compute  $\nabla_{\theta} L(z(T, \theta))$  via running the following system of ODEs for time  $s \in [T_0, T]$ ,

$$\begin{cases} \frac{dy(s)}{ds} = -(\partial_{\theta} f(z(s, \theta), s, \theta))^{\top} a(s), & y(T) = 0_{N_p \times 1}; \\ \frac{d}{ds} a(s) = -(\partial_z f(z(s, \theta), s, \theta))^{\top} a(s), & a(T) = \nabla L(z(T, \theta)). \end{cases} \quad (2.6)$$

The output  $y(T_0) = \nabla_{\theta} L(z(T, \theta))$ .

**Exercise 2.4.** How does the storage resource and computational time scale with respect to  $N_d$  and  $N_p$ ?

The following proof is presented in a way different from the original paper [5]. It fully uses the language of differential equations, echoing the proposal in [8]. It avoids to discuss any derivatives along the propagation of (discretized) deep neural networks. In [8, Section 1.1], this problem was also briefly discussed in the language of functional derivatives.

### 2.3.1 Proof of Theorem 2.3

By direct computation,

$$\partial_{\theta_j} L(z(T, \theta)) = \sum_i \nabla_i L(z(T, \theta)) \frac{\partial z_i(T, \theta)}{\partial \theta_j}. \quad (2.7)$$

Therefore, we need to compute  $\frac{\partial z_i(T, \theta)}{\partial \theta_j}$ . Let us take derivative in (2.2) with respect to  $\theta_j$ ,

$$\frac{d^2 z(t, \theta)}{dt d\theta_j} = \partial_{\theta_j} f(z(t, \theta), t, \theta) + \nabla_z f(z(t, \theta), t, \theta) \frac{dz(t, \theta)}{d\theta_j}.$$

Let us denote  $h(t, \theta) = \frac{dz(t, \theta)}{d\theta_j}$ . Then one has

$$\frac{d}{dt} h(t, \theta) = \partial_{\theta_j} f(z(t, \theta), t, \theta) + \nabla_z f(z(t, \theta), t, \theta) h(t, \theta).$$

Let us denote the evolution/solution operator  $U(t_f, t_i) = e^{\int_{\mathcal{T}}^{t_f} \partial_z f(z(s, \theta), s, \theta) ds}$  where  $\mathcal{T}$  is the chronological time-ordering operator. Then one can show that

$$\frac{d}{dt} (U(T_0, t) h(t, \theta)) = U(T_0, t) \left( -\nabla_z f(z(t, \theta), t, \theta) \right) h(t, \theta)$$

$$\begin{aligned}
& + U(T_0, t) \left( \partial_{\theta_j} f(z(t, \theta), t, \theta) + \nabla_z f(z(t, \theta), t, \theta) h(t, \theta) \right) \\
& = U(T_0, t) \partial_{\theta_j} f(z(t, \theta), t, \theta).
\end{aligned}$$

By rewriting this in the integral form, one has

$$U(T_0, t)h(t, \theta) - U(T_0, T_0)h(T_0, \theta) = \int_{T_0}^t U(T_0, s) \partial_{\theta_j} f(z(s, \theta), s, \theta) ds.$$

Since we considered  $z(T_0, \theta) = z_0$  is fixed, thus  $h(T_0, \theta) = 0$ , and therefore

$$h(t, \theta) = \int_{T_0}^t U(t, s) \partial_{\theta_j} f(z(s, \theta), s, \theta) ds.$$

One can also obtain this result by applying the Duhamel's principle stated in Theorem 1.2.

By plugging this equation into (2.7), one has

$$\begin{aligned}
\partial_{\theta_j} L(z(T, \theta)) &= \sum_i \nabla_i L(z(T, \theta)) \frac{\partial z_i(T, \theta)}{\partial \theta_j} \\
&= \nabla_i L(z(T, \theta)) \cdot h(T, \theta) \\
&= \int_{T_0}^T \nabla_i L(z(T, \theta)) \cdot \left( U(T, s) \partial_{\theta_j} f(z(s, \theta), s, \theta) \right) ds.
\end{aligned}$$

Let us denote  $a(s)^\top = \nabla L(z(T, \theta))^\top U(T, s)$  (an  $1 \times N_d$  time-dependent row vector), then

$$\partial_{\theta_j} L(z(T, \theta)) = \int_{T_0}^T a(s)^\top \partial_{\theta_j} f(z(s, \theta), s, \theta) ds.$$

If we adopt the matrix form,

$$(\nabla_{\theta} L(z(T, \theta)))^\top = \int_{T_0}^T a(s)^\top \partial_{\theta} f(z(s, \theta), s, \theta) ds,$$

where  $\partial_{\theta} f(z(s, \theta), s, \theta) = [\partial_{\theta_1} f(z(s, \theta), s, \theta) \quad \cdots \quad \partial_{\theta_{N_p}} f(z(s, \theta), s, \theta)] \in \mathbb{R}^{N_d \times N_p}$ .

The dynamics of  $a(s)$  can be easily verified as follows:

$$\frac{d}{ds} a(s)^\top = \nabla L(z(T, \theta))^\top U(T, s) (-\partial_z f(z(s, \theta), s, \theta)) = -a(s)^\top \partial_z f(z(s, \theta), s, \theta).$$

Or in a column form as in (2.5). The terminal condition of  $a$  can be easily checked, as well as (2.6).

**Example 2.5.** For Neural ODE in (2.2), if we consider a special family  $f(z, t, \theta) = -iH(t, \theta)z$ , where  $H(t, \theta)$  is a  $\theta$ -parameterized time-dependent Hermitian operator, then for the forward process, one has

$$\frac{dz(s, \theta)}{dt} = -iH(s, \theta)z(s, \theta).$$

For illustration, let us consider the following pedagogical example:

$$H(t, \theta) = \frac{\omega}{2} \sigma_Z + u_{\theta}(t) \sigma_X \equiv \begin{bmatrix} \frac{\omega}{2} & u_{\theta}(t) \\ u_{\theta}(t) & \frac{\omega}{2} \end{bmatrix}$$

with the constraint  $\|u_{\theta}(t)\| \leq u_{\max}$  and with initial condition  $z(0, \theta) = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ . We want to find a control  $u_{\theta}$

such that  $z(T, \theta) \approx \begin{bmatrix} 0 \\ 1 \end{bmatrix}$  as much as possible. We choose  $\omega = 1$ ,  $u_{\max} = 2$ , and parameterize the  $u_{\theta}$  via fully-connected neural network. The training results for  $T = 5$  and  $T = 8$  are visualized in Figure below. For the case  $T = 8$ , the final target state is almost identical to the target state (up to a constant rescaling) with Fidelity as high as 0.999.

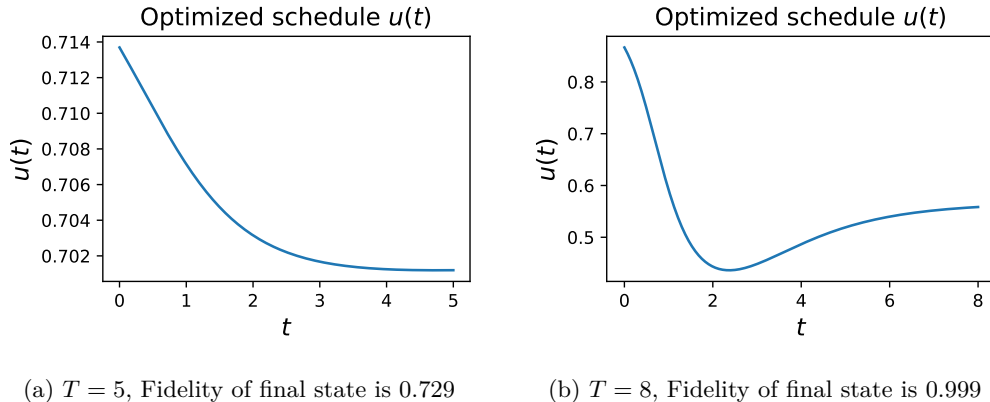


Figure 2.1: Training result for Example 2.5. Fidelity measures the similarity between the final state  $z(T, \theta)$  and the target state. If the number of Fidelity is close to 1, it means the similarity is very high and the control is very efficient.

## 2.4 Neural SDEs

Lastly, we will discuss how to work on backpropagation for an SDE. We will mainly focus on the most fundamental mathematical facts without getting deep into various applications. In this section, we focus on the strong solution of SDE. The following notes is prepared partially based on [17]. In order to extend the backpropagation (adjoint equation) from ODEs to SDEs, it is easier to work with the Stratonovich SDE. For the Itô SDE, one may convert it into Stratonovich SDE by Theorem 1.11, and then apply the following procedures.

In the adjoint equation of ODEs, one needs to run dynamics backwardly. The first result below is to show how one can reverse the dynamics of Stratonovich SDE, provided that one is given access to the Brownian motion paths.

**Lemma 2.6** (Reverse of Stratonovich SDE). Suppose  $\{X(t)\}_{t \in [0, T]}$  satisfies the Stratonovich SDE in (1.8) (also copied below herein):

$$dX(t) = b(X(t), t) dt + \sigma(X(t), t) \circ dB(t).$$

Let  $\tilde{X}(t) := X(T - t)$ , then  $\tilde{X}$  satisfies the following Stratonovich integral

$$d\tilde{X}(t) = -b(\tilde{X}(t), T - t) dt + \sigma(\tilde{X}(t), T - t) \circ dB(T - t). \quad (2.8)$$

A direct proof will be provided in § 2.4.2.

Similar to Theorem 2.3 for Neural ODEs, one can show the following:

**Theorem 2.7** (Gradient estimates for neural SDE; see also [17]). Suppose that the initial time is  $T_0$  and the final time is  $T$ . Suppose that  $X(t, \theta)$  follows the SDE (2.3) (also copied below herein)

$$dX(t, \theta) = b(X(t, \theta), t, \theta) dt + \sigma(X(t, \theta), t, \theta) \circ dB(t), \quad X(T_0, \theta) = x_0 \text{ is fixed.} \quad (2.9)$$

Suppose that  $L(\cdot), b(\cdot, \cdot, \cdot), \sigma(\cdot, \cdot, \cdot)$  are continuously differentiable. For simplicity, let us consider the case dimension  $N_d = 1$ . Then

$$\begin{cases} \nabla_{\theta} L(X(T, \theta)) = \int_{T_0}^T a(s) \nabla_{\theta} b(X(s, \theta), s, \theta) ds + a(s) \nabla_{\theta} \sigma(X(s, \theta), s, \theta) \circ dB(s), \\ da(s) = -a(s) \partial_x b(X(s, \theta), s, \theta) ds - a(s) \partial_x \sigma(X(s, \theta), s, \theta) \circ dB(s). \end{cases} \quad (2.10)$$

Alternatively, one could compute  $\nabla_{\theta} L(X(T, \theta))$  via running the following system of SDEs for time  $s \in [T_0, T]$ ,

$$\begin{cases} dy(s) = -a(s) \nabla_{\theta} b(X(s, \theta), s, \theta) ds - a(s) \nabla_{\theta} \sigma(X(s, \theta), s, \theta) \circ dB(s), & y(T) = 0_{N_p} \\ da(s) = -a(s) \partial_x b(X(s, \theta), s, \theta) ds - a(s) \partial_x \sigma(X(s, \theta), s, \theta) \circ dB(s), & a(T) = L'(X(T, \theta)) \end{cases} \quad (2.11)$$

The output  $y(T_0) = \nabla_{\theta} L(X(T, \theta))$ .

**Exercise 2.8.** Generalize Theorem 2.7 to a general dimension  $N_d$ .

### 2.4.1 Proof of Theorem 2.7

The proof is essentially the same as that for Theorem 2.3. As mentioned above, for simplicity, we shall only consider the 1D case. By direct computation,

$$\partial_{\theta_j} L(X(T, \theta)) = L'(X(T, \theta)) \frac{\partial X(T, \theta)}{\partial \theta_j}. \quad (2.12)$$

Therefore, we need to compute  $\frac{\partial X(T, \theta)}{\partial \theta_j}$ . For the time being, let us fix  $j$ , denote  $h(t, \theta) = \frac{dX(t, \theta)}{d\theta_j}$ , and let us take derivative in (2.9) with respect to  $\theta_j$ . Then

$$\begin{aligned} dh(t, \theta) = & \left( \partial_{\theta_j} b(X(t, \theta), t, \theta) + \partial_x b(X(t, \theta), t, \theta) h(t, \theta) \right) dt \\ & + \left( \partial_{\theta_j} \sigma(X(t, \theta), t, \theta) + \partial_x \sigma(X(t, \theta), t, \theta) h(t, \theta) \right) \circ dB(t). \end{aligned} \quad (2.13)$$

Since the path of Stratonovich SDE is reversible for any given Brownian motion path, let us denote the evolution/solution operator  $U(s_f, s_i)$  by the following: suppose we fix  $\theta$  (and thus fix  $X(t, \theta)$ ), and then for any  $h_0$ , we define  $U(s_f, s_i)h_0$  as the solution at time  $s_f$  of the following SDE

$$d\tilde{h}(t) = \partial_x b(X(t, \theta), t, \theta) \tilde{h}(t) dt + \partial_x \sigma(X(t, \theta), t, \theta) \tilde{h}(t) \circ dB(t), \quad \tilde{h}(s_i) = h_0.$$

**Lemma 2.9.** For arbitrary  $s_i, s_f \in [T_0, T]$ , the evolution operator  $U(s_f, s_i)$  has the following form:

$$U(s_f, s_i) = \exp \left( \int_{s_i}^{s_f} \partial_x b(X(t, \theta), t, \theta) dt + \partial_x \sigma(X(t, \theta), t, \theta) \circ dB(t) \right). \quad (2.14)$$

A direct proof will be provided in § 2.4.2.

*Remark.* This result is essentially a time-dependent generalization of Exercise 1.15.

**Lemma 2.10.** For a fixed  $s_i$  and consider any  $t \in [T_0, T]$ , then

$$dU(t, s_i) = \partial_x b(X(t, \theta), t, \theta) U(t, s_i) dt + \partial_x \sigma(X(t, \theta), t, \theta) U(t, s_i) \circ dB(t).$$

For a fixed  $s_f$ , and consider any  $t \in [T_0, T]$ , then

$$dU(s_f, t) = -U(s_f, t) \partial_x b(X(t, \theta), t, \theta) dt - U(s_f, t) \partial_x \sigma(X(t, \theta), t, \theta) \circ dB(t).$$

*Proof.* The proof of Lemma 2.9 can be easily modified to prove the conclusion. A detailed proof is skipped.  $\square$

*Remark.* This lemma resembles Lemma 1.1.

Then one can directly show that

$$d(U(T_0, t)h(t, \theta)) = U(T_0, t) \partial_{\theta_j} b(X(t, \theta), t, \theta) dt + U(T_0, t) \partial_{\theta_j} \sigma(X(t, \theta), t, \theta) \circ dB(t).$$

By rewriting this in the integral form, one has

$$U(T_0, t)h(t, \theta) - h(T_0, \theta) = \int_{T_0}^t U(T_0, s) \partial_{\theta_j} b(X(s, \theta), s, \theta) ds + U(T_0, s) \partial_{\theta_j} \sigma(X(s, \theta), s, \theta) \circ dB(s).$$

Since we considered  $X(T_0, \theta) = X_0$  is fixed, thus  $h(T_0, \theta) = 0$ , and therefore

$$h(t, \theta) = \int_{T_0}^t U(t, s) \partial_{\theta_j} b(X(s, \theta), s, \theta) ds + U(t, s) \partial_{\theta_j} \sigma(X(s, \theta), s, \theta) \circ dB(s).$$

By plugging this equation into (2.12), one has

$$\begin{aligned}
& \partial_{\theta_j} L(z(T, \theta)) \\
&= L'(X(T, \theta)) \frac{\partial X(T, \theta)}{\partial \theta_j} \\
&= L'(X(T, \theta)) h(T, \theta) \\
&= \int_{T_0}^T L'(X(T, \theta)) \left( U(T, s) \partial_{\theta_j} b(X(s, \theta), s, \theta) ds + U(T, s) \partial_{\theta_j} \sigma(X(s, \theta), s, \theta) \circ dB(s) \right).
\end{aligned}$$

Let us denote  $a(s) = L'(X(T, \theta)) U(T, s)$ , then

$$\partial_{\theta_j} L(z(T, \theta)) = \int_{T_0}^T a(s) \partial_{\theta_j} b(X(s, \theta), s, \theta) ds + a(s) \partial_{\theta_j} \sigma(X(s, \theta), s, \theta) \circ dB(s)$$

If we adopt the matrix form, we can obtain the formula of  $\nabla_{\theta} L(z(T, \theta))$  as in (2.10). The dynamics of  $a(s)$  in (2.10) can be easily verified by Lemma 2.10.

## 2.4.2 Proof of lemmas

*Proof of Lemma 2.6.* Let us divide the time interval into  $N$  equally spaced grid points and let  $t_k = k\Delta t$  where  $\Delta t = T/N$ .

$$\begin{aligned}
& \tilde{X}(t_{k+1}) - \tilde{X}(t_k) \\
&= X(T - t_{k+1}) - X(T - t_k) \\
&= -b(X(T - t_{k+1}), T - t_{k+1}) \Delta t - \frac{\sigma(X(T - t_{k+1}), T - t_{k+1}) + \sigma(X(T - t_k), T - t_k)}{2} (B(T - t_k) - B(T - t_{k+1})) \\
&\quad + o(\Delta t) \\
&= -b(\tilde{X}(t_k), T - t_k) \Delta t + \frac{\sigma(\tilde{X}(t_{k+1}), T - t_{k+1}) + \sigma(\tilde{X}(t_k), T - t_k)}{2} (B(T - t_{k+1}) - B(T - t_k)) + o(\Delta t),
\end{aligned}$$

which leads into (2.8).  $\square$

*Proof of Lemma 2.9.* For simplicity of notations, we shall suppress the dependence on  $\theta$ . Let us consider an arbitrary  $h_0$  and by definition,

$$\begin{aligned}
& U(s_f + \Delta t, s_i) h_0 - U(s_f, s_i) h_0 \\
&= U(s_f + \Delta t, s_f) U(s_f, s_i) h_0 - U(s_f, s_i) h_0 \quad (\text{Denote } \tilde{h}(s_f) = U(s_f, s_i) h_0) \\
&= \partial_x b(X(s_f), s_f) \tilde{h}(s_f) \Delta t + \frac{1}{2} \left( \partial_x \sigma(X(s_f), s_f) \tilde{h}(s_f) + \partial_x \sigma(X(s_f + \Delta t), s_f + \Delta t) \tilde{h}(s_f + \Delta t) \right) \Delta B(s_f) + o(\Delta t) \\
&= \partial_x b(X(s_f), s_f) \tilde{h}(s_f) \Delta t + \frac{1}{2} \left( \partial_x \sigma(X(s_f), s_f) \tilde{h}(s_f) + \partial_x \sigma(X(s_f + \Delta t), s_f + \Delta t) \tilde{h}(s_f) \right) \Delta B(s_f) \\
&\quad + \frac{1}{2} \partial_x \sigma(X(s_f + \Delta t), s_f + \Delta t) \Delta B(s_f) (\tilde{h}(s_f + \Delta t) - \tilde{h}(s_f)) + o(\Delta t) \\
&= \partial_x b(X(s_f), s_f) \tilde{h}(s_f) \Delta t + \frac{1}{2} \left( \partial_x \sigma(X(s_f), s_f) + \partial_x \sigma(X(s_f + \Delta t), s_f + \Delta t) \right) \Delta B(s_f) \tilde{h}(s_f) \\
&\quad + \frac{1}{2} \left( \partial_x \sigma(X(s_f), s_f) \Delta B(s_f) \right)^2 \tilde{h}(s_f) + o(\Delta t).
\end{aligned}$$

Let us denote the expression on the right hand side of (2.14) as  $\bar{U}$  to distinguish expressions:

$$\begin{aligned}
& \bar{U}(s_f + \Delta t, s_i) h_0 - \bar{U}(s_f, s_i) h_0 \\
&= \exp \left( \partial_x b(X(s_f), s_f) \Delta t + \frac{1}{2} \partial_x \sigma(X(s_f), s_f) \Delta B(s_f) + \frac{1}{2} \partial_x \sigma(X(s_f + \Delta t), s_f + \Delta t) \Delta B(s_f) + o(\Delta t) \right) \bar{U}(s_f, s_i) h_0
\end{aligned}$$

$$\begin{aligned}
& -\bar{U}(s_f, s_i)h_0 \\
= & \left( \partial_x b(X(s_f), s_f)\Delta t + \frac{1}{2}\partial_x \sigma(X(s_f), s_f)\Delta B(s_f) + \frac{1}{2}\partial_x \sigma(X(s_f + \Delta t), s_f + \Delta t)\Delta B(s_f) \right) \bar{U}(s_f, s_i)h_0 \\
& + \frac{1}{2} \left( \partial_x \sigma(X(s_f), s_f)\Delta B(s_f) \right)^2 \bar{U}(s_f, s_i)h_0 + o(\Delta t)
\end{aligned}$$

As the above two dynamics have the same discretization and match with each other, we can obtain (2.14) as a limit.  $\square$

## 2.5 Further Readings

Below we list a few related topics for possible exploration:

- [15] gave a more detailed discussion on the relation between control theory and deep learning in lecture notes for a summer school at Peking University (2020).
- When the problem involves taking higher-order derivatives, e.g.,

$$\min_{\theta} \mathbb{E}_x g(\nabla_x p(\theta, x), \nabla_x^{(2)} p(\theta, x), \dots, \nabla_x^{(k)} p(\theta, x)),$$

and  $g$  is some non-linear non-negative function, then an efficient way to take derivatives becomes important. An example is to solve PDEs using PINNs (cf. Example 1.21). For Laplacian operator, the forward method can provide speed-up over backpropagation [16].

- The above results can be further generalized to Neural Jump Process, and Levy process. This is an interesting topic for exploration on your own.

## Chapter 3

# Neural Dynamics for Generative Problems: (I) Normalizing Flows

### 3.1 Background

Generative problem, in a mathematical way, means to generate samples according to a certain probability measure  $p_X$  (which could be extremely complex). Here we list two examples among many:

- Each black/white image (with  $n$  pixels) can be represented by a data point in  $\mathbb{R}^n$ . A task in machine learning is to find a way to generate samples according to  $p_X$ , provided that one is given a certain amount of data samples  $\{x_i\}_{i=1}^N$ . This is a reminiscent of the first type of data source, namely, data directly coming from realistic problems.
- For  $N$  interacting particles in statistical physics with potential  $U$ , the equilibrium distribution (also known as the Boltzmann distribution) takes the following form  $p_X(x) = e^{-U(x)}/\mathcal{Z}$  where  $\mathcal{Z} = \int e^{-U}$  is known as the partition function or normalization constant (depending on which research field that you come from). One goal is to find a way to generate samples according to such a distribution, which could enable the estimation of many physical quantities for equilibrium dynamics. This is a reminiscent of the second type of data source: the data is not provided in general and needs to be generated artificially in order to formulate the generative problem as a machine learning problem.

A widely studied tool for sample generation is normalizing flows. When normalizing flows have infinitely many layers and for each layer, the update is infinitesimally small, the corresponding mathematical object is simply ODEs, and also refers to as probability flow in some contexts. The following notes are prepared mainly by referring to review papers [14, 23] and [4, Chapter 18]. This topic is also tightly connected to a field called optimal transport theory; one may refer to e.g., [2, 33].

For this mini-course, we will only cover some basic concepts about normalizing flows and probability flows, as well as a little from optimal transport theory. In this notes, we will not specifically distinguish the notation of probability measures and probability densities.

### 3.2 Connection to Optimal Transport

As mentioned above, the goal is to find a way to generate samples according to a target measure  $p_X$ . One widely used approach is to find a mapping  $f$  such that  $f(Z)$  follows the target distribution  $p_X$ . Designing a function  $f$  to generate samples using  $f(Z)$  is a very classical idea, with enormous applications. For instance, *Box-Muller transform* (a backend algorithm for generating normal random variables in our computers) uses the following algorithm: generate two independent  $U_1$  and  $U_2$  following uniform distributions on  $(0, 1)$  and then let

$$Z_0 = \sqrt{-2 \ln(U_1)} \cos(2\pi U_2);$$



$$Z_1 = \sqrt{-2 \ln(U_1)} \sin(2\pi U_2).$$

Then  $(Z_0, Z_1)$  are independent standard normal random variables.<sup>1</sup>

The prior distribution  $p_Z$  contains certain degree of freedom, but it is typically chosen as simple distributions like uniform distributions or Gaussians. The mapping  $f$  is what we need to find. Such a mapping is also well-known as the *Monge map* in some literatures, and this name dated back to the problem of finding a shortest path to move a pile of sand into another pile of sand in 18th century.

In general, if we consider general probability measures (rather than probability densities), there are examples such that one cannot find such a mapping. For instance, let  $p_Z(z) = \delta_0(z)$  (a Dirac delta distribution at 0), and  $p_X(x) = \frac{1}{2}\delta_1(x) + \frac{1}{2}\delta_{-1}(x)$  (half probability at 1 and half probability at  $-1$ ). It is easy to see that one cannot find such a mapping  $f$ . Of course, in real problems, life is not that tough and we have or we can assume the existence of distributions, which has the following result.

**Theorem 3.1** (Existence of Monge map [2, Theorem 2.3]). Suppose that the initial measure  $\mu_Z \in \mathcal{P}_2(\mathbb{R}^n)$  is absolutely continuous with respect to Lebesgue measure, and  $\mu_X \in \mathcal{P}_2(\mathbb{R}^n)$ , then there exists a Monge map  $f$  such that  $\mu_X = f\#\mu_Z$ .<sup>2</sup>

Due to the above existential theorem, and also for many realistic problems, we can assume the existence of probability densities, we shall stick with this setup from now on. Let us consider the 1D case.

**Example 3.2** (1D case).

- Let us assume that the support of  $p_X, p_Z$  is simply  $\mathbb{R}$ , which means their corresponding cumulative probability density functions  $P_X$  and  $P_Z$  are strictly monotone functions and is invertible. Suppose that  $f$  is monotone increasing, then

$$f(x) = P_X^{-1} \circ P_Z(x) \tag{3.1}$$

is a mapping that satisfies the condition. This can be easily checked: with such a choice

$$\mathbb{P}(f(Z) \leq a) = \int_{\mathbb{R}} \chi_{\{f(z) \leq a\}} p_Z(z) dz = \int_{\mathbb{R}} \chi_{\{P_Z(z) \leq P_X(a)\}} p_Z(z) dz = P_X(a) = P(X \leq a),$$

which verifies that  $f(Z)$  and  $X$  have the same distribution.

- Conversely, if one is given a monotone increasing and differentiable mapping  $f$  and the prior distribution  $p_Z$ , the probability density function of  $X := f(Z)$  can be computed

$$p_X(x) = p_Z(f^{-1}(x)) (f^{-1})'(x).$$

The proof is straightforward by (3.1).

Such a result can be generalized to arbitrary dimensions. If  $f$  is invertible, we can similarly compute the density of  $p_X$ :

**Theorem 3.3.** Suppose the mapping  $f$  is invertible, and let  $X = f(Z)$ , then

$$p_X(x) = p_Z(f^{-1}(x)) |J_{f^{-1}}(x)|, \quad J_{f^{-1}}(x) := \det(\nabla f^{-1}(x)), \tag{3.2}$$

and we also have

$$p_X(f(z)) |J_f(z)| = p_Z(z). \tag{3.3}$$

<sup>1</sup>See [https://en.wikipedia.org/wiki/Box%E2%80%93Muller\\_transform](https://en.wikipedia.org/wiki/Box%E2%80%93Muller_transform)

<sup>2</sup>The notation  $\mathcal{P}_2(\mathbb{R}^n)$  means the space of probability measures with finite second moment. The original statement of [2, Theorem 2.3] contains a lot more interesting results, and we only selectively pick what we need.

*Proof.* For an arbitrary test function  $\phi \in C_c^\infty(\mathbb{R}^n)$ , we have

$$\begin{aligned} \int \phi(x) p_X(x) \, dx &= \mathbb{E}[\phi(X)] = \mathbb{E}[\phi(f(Z))] = \int \phi(f(z)) p_Z(z) \, dz \\ &= \int \phi(x) p_Z(f^{-1}(x)) |J_{f^{-1}}(x)| \, dx. \end{aligned}$$

Since  $\phi$  is arbitrary, one has the above result. The second equation can be obtained similarly.  $\square$

*Remark.* (Monge-Ampere equation). If  $p_X = f\#p_Z$ , namely, (3.3) holds, and suppose that the mapping  $f = \nabla\varphi$  with  $\varphi$  being convex, then

$$\det(\nabla^2\varphi) p_X(\nabla\varphi) = p_Z. \quad (3.4)$$

This is well-known as Monge-Ampere equation.

### 3.3 Normalizing Flows

The next task is how to design some architectures for us to learn such a mapping  $f$ .

#### 3.3.1 General architecture

For normalizing flows, a core ingredient is to represent  $f$  as the composition of a sequence of mapping  $f_i$ :

$$f(z) = f_\ell \circ f_{\ell-1} \circ \cdots \circ f_1(z),$$

in same way as deep neural networks. The second new ingredient is that if we also want to compute the density of  $p_X$ , we'd better require  $f$  to be invertible (cf Theorem 3.3), which becomes a requirement for each  $f_i$  to be invertible. If so,

$$f^{-1}(x) = f_1^{-1} \circ f_2^{-1} \circ \cdots \circ f_\ell^{-1}(x).$$

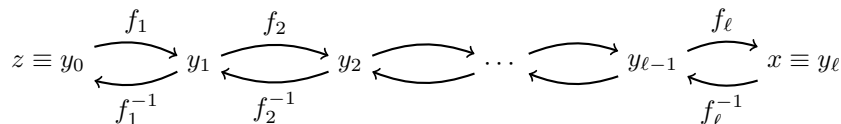


Figure 3.1: A visualization of the forward and backward directions in normalizing flows

Then we can compute the density, which involves computing the gradient of a composition of two invertible mapping.

**Lemma 3.4.** Suppose  $h_1$  and  $h_2$  are mappings, and  $h = h_1 \circ h_2$ , then

$$\det(\nabla h(x)) = \det(\nabla h_1(h_2(x))) \det(\nabla h_2(x)) \implies J_{h_1 \circ h_2}(x) = J_{h_1}(h_2(x)) J_{h_2}(x). \quad (3.5)$$

*Proof.* By chain rules,  $\nabla h(x) = \nabla h_1(h_2(x)) \nabla h_2(x)$ , which immediately leads into the above result.  $\square$

With such a result, one can show that

$$\begin{aligned} p_X(x) &\stackrel{(3.2)}{=} p_Z(f^{-1}(x)) |J_{f^{-1}}(x)| \stackrel{(3.5)}{=} p_Z(f^{-1}(x)) \left| J_{f_1^{-1}}(y_1) J_{f_2^{-1}}(y_2) \cdots J_{f_\ell^{-1}}(x) \right| \\ &= p_Z(f^{-1}(x)) \left| \prod_{i=1}^{\ell} J_{f_i^{-1}}(y_i) \right|. \end{aligned} \quad (3.6)$$

We have just shown that with such a mapping  $f$  available, we can surely generate samples according to  $X = f(Z)$  and also compute the density of  $p_X$ .

The remaining questions become:

(Q1) how to find a parameterized invertible mapping?

(Q2) how to design such a mapping such that the determinant of Jacobian is easy to compute?

(Q3) the universality (expressiveness) and efficiency of such a design?

### 3.3.2 Coupling flow and RealNVP

A widely used idea to design invertible mapping is through the coupling flows. For a  $n$  dimensional problem,  $x = (x_1, x_2) \in \mathbb{R}^d \times \mathbb{R}^{n-d}$ , one defines such a mapping:

$$\begin{cases} y^{(1)} = h(x^{(1)}, \theta(x^{(2)})); \\ y^{(2)} = x^{(2)}. \end{cases}$$

where  $h(\cdot, \theta)$  is an invertible function. The inverse is simply

$$\begin{cases} y^{(1)} = h^{-1}(y^{(1)}, \theta(y^{(2)})); \\ x^{(2)} = y^{(2)}. \end{cases}$$

Namely, we fix certain coordinates as conditioner, and we only update other coordinates, a bit like coordinate-wise gradient descent idea.

A widely used example is RealNVP [7] (standing for real-valued non-volume preserving):

$$f(x) = \begin{cases} y^{(1)} = x^{(1)} \odot \exp(s(x^{(2)})) + t(x^{(2)}); \\ y^{(2)} = x^{(2)}. \end{cases}$$

$s, t : \mathbb{R}^{n-d} \times \mathbb{R}^d$  are parameterized functions, and  $\odot$  is the element-wise product. The inverse function is

$$f(y)^{-1} = \begin{cases} x^{(1)} = (y^{(1)} - t(y^{(2)})) \odot \exp(-s(y^{(2)})); \\ x^{(2)} = y^{(2)}. \end{cases}$$

In such a formalism, one does not require  $s, t$  to be invertible. The Jacobian

$$J_f(x) = \det(\nabla f(x)) = \exp\left(\sum_j s(x^{(2)})_j\right)$$

In order for the component  $x^{(2)}$  to also evolve, one can change the partition (a collection of coordinates that remain unchanged) above so that all coordinates are overall “coupled”.

### 3.3.3 Loss functions for learning without data

Given a parameterized density  $X_\theta \sim p_{X,\theta} = f_\theta \# p_Z$ , we want to match this with  $X \sim p_X$ . Therefore, our optimization problem can be setup as

$$\min_{\theta} \text{Dist}(f_\theta \# p_Z, p_X).$$

A commonly used “metric” is KL divergence, which leads into

$$\begin{aligned} \min_{\theta} D_{\text{KL}}(f_\theta \# p_Z, p_X) &= \int p_{X,\theta}(x) \log p_{X,\theta}(x) - p_{X,\theta}(x) \log p_X(x) \, dx \\ &= \mathbb{E}_{X_\theta \sim f_\theta(Z)} \left[ \log p_{X,\theta}(X_\theta) - \log p_X(X_\theta) \right]. \end{aligned}$$

Suppose we already know that our target distribution has a potential  $U$  with  $p_X(x) = e^{-U(x)}/\mathcal{Z}$ , then the optimization problem becomes

$$\min_{\theta} \mathbb{E}_{X_{\theta} \sim f_{\theta}(Z)} \left[ \log p_{X,\theta}(X_{\theta}) + U(X_{\theta}) \right]. \quad (3.7)$$

The term involving  $\mathcal{Z}$  can be dropped out as it is independent of  $\theta$ . This formalism can help to learn a target distribution without given data. To discretize the above expectation, we can take samples

$$\min_{\theta} \frac{1}{M} \sum_{i=1}^M \log p_{X,\theta}(f_{\theta}(z_i)) + U(f_{\theta}(z_i)),$$

where  $z_i$  are *i.i.d.* samples drawn from the base distribution  $p_Z$ .

### 3.3.4 Loss functions for learning with data

Another possibility is to use

$$\min_{\theta} D_{\text{KL}}(p_X, f_{\theta} \# p_Z) = \int p_X \log(p_X) - p_X \log(p_{X,\theta}).$$

As  $p_X$  is the true distribution, independent of  $\theta$ , this is equivalent to estimate

$$\max_{\theta} \mathbb{E}_{X \sim p_X} \left[ \log(p_{X,\theta}(X)) \right].$$

Suppose we are given data samples  $\{x_i\}_{i=1}^N$  rather than the potential function, we can use

$$\max_{\theta} \frac{1}{N} \sum_{i=1}^N \log p_{X,\theta}(x_i) \stackrel{(3.2)}{=} \frac{1}{N} \sum_{i=1}^N \log p_Z(f_{\theta}^{-1}(x_i)) + \log |J_{f_{\theta}^{-1}}(x_i)|.$$

Then everything is computable. This formalism can also be interpreted as maximizing the log-likelihood for data conditioned on the model parameter  $\theta$ .

## 3.4 ODE Flows

As mentioned above, the NFs structure can help to more easily compute the density functions. When the layers become infinitely deep, the normalizing flow structure can turn into an ODE system or say a “continuous normalizing flow”. Moreover, ODE is naturally time-reversible (meaning invertible flow map).

### 3.4.1 Evolution of log-density function

More specifically, suppose we consider the ODE:

$$\frac{d}{dt} Y(t) = F(Y(t), t), \quad Y(0) \text{ is random.}$$

Though the initial condition is random, each trajectory is still deterministic. The time- $t$  flow map  $y \mapsto \phi(y, t)$  also solves the ODE

$$\frac{d}{dt} \phi(y, t) = F(\phi(y, t), t), \quad \phi(y, 0) = y.$$

The probability density function of  $Y(t)$ , denoted as  $p(y, t)$ , solves the following PDE

$$\partial_t p(y, t) \stackrel{(1.18)}{=} -\nabla_y \cdot (F(y, t)p(y, t)) = -(\nabla_y \cdot F(y, t))p(y, t) - F(y, t) \cdot \nabla_y p(y, t). \quad (3.8)$$

This can be obtained by the formula of Fokker-Planck equation in (1.18) as ODE can be regarded as a special SDE without diffusion coefficient. Then by the chain rules, the log-density function can be computed as follows:

$$\begin{aligned} \frac{d}{dt} \log p(Y(t), t) &= \frac{\partial_t p(Y(t), t) + \nabla_y p(Y(t), t) \cdot \frac{dY(t)}{dt}}{p(Y(t), t)} \\ &\stackrel{(3.8)}{=} \frac{-(\nabla_y \cdot F)(Y(t), t)p(Y(t), t) - F(Y(t), t) \cdot \nabla_y p(Y(t), t) + \nabla_y p(Y(t), t) \cdot F(Y(t), t)}{p(Y(t), t)} \\ &= -(\nabla_y \cdot F)(Y(t), t). \end{aligned}$$

Hence,

$$\log p(Y(T), T) - \log p(Y(0), 0) = \int_0^T -(\nabla_y \cdot F)(Y(t), t) dt. \quad (3.9)$$

### 3.4.2 Connection to (discrete) normalizing flows

Suppose the initial distribution is  $p_Z$ , we have

$$p(Y(T), T) = p_Z(Y(0)) \exp\left(-\int_0^T (\nabla_y \cdot F)(Y(t), t) dt\right).$$

If we choose grid points  $0 = t_0 < t_1 < \dots < t_\ell = T$ , grid size  $\Delta t = \frac{T}{\ell}$  and suppose that  $\ell \gg 1$ . If we approximately choose  $f_i^{-1}(x) = 1 - F(x, t_i)\Delta t + \mathcal{O}(\Delta t^2)$  and let mapping  $f(x) = f_\ell \circ f_{\ell-1} \circ \dots \circ f_1(x)$  and denote  $\hat{y}_k = f_k \circ f_{k-1} \circ \dots \circ f_1(x)$  for  $k = 1, 2, \dots, \ell$ . Then

$$\begin{aligned} \det(\nabla f_i^{-1}(x)) &= \det\left(I - \nabla F(x, t_i)\Delta t\right) + \mathcal{O}(\Delta t^2) \\ &= 1 - \text{tr}(\nabla F(x, t_i)\Delta t) + \mathcal{O}(\Delta t^2) \\ &= \exp\left(-\text{tr}(\nabla F(x, t_i)\Delta t)\right) + \mathcal{O}(\Delta t^2), \end{aligned}$$

and the last equation becomes

$$\begin{aligned} p(\hat{y}_\ell, T) &= p_Z(\hat{y}_0) \prod_{i=1}^{\ell} \exp\left(-\Delta t (\nabla_y \cdot F)(\hat{y}_i, t_i)\right) + \mathcal{O}(\Delta t) \\ &= p_Z(\hat{y}_0) \prod_{i=1}^{\ell} \det(\nabla f_i^{-1}(\hat{y}_i)) + \mathcal{O}(\Delta t). \end{aligned}$$

This reduces to (3.6) for the discrete-time case. In the above, we used a fact that  $\det(I + \epsilon A) = 1 + \text{tr}(A)\epsilon + \mathcal{O}(\epsilon^2)$  for any square matrix  $A$  and small enough  $\epsilon$ .

### 3.4.3 Training of ODE flow for learning tasks without data

Suppose we consider the problem of learning probability distribution without data. Then similar to (3.7), we need to optimize:

$$\min_{\theta} \mathbb{E} \left[ \log p(Y(T, \theta), T, \theta) + U(Y(T, \theta)) \right].$$

By (3.9), if we denote

$$\begin{aligned} \frac{d}{dt} Y(t, \theta) &= F(Y(t, \theta), t, \theta), & Y(0, \theta) &= Y_0; \\ \frac{d}{dt} J(t, \theta) &= -\nabla_y \cdot F(Y(t, \theta), t, \theta), & J(0, \theta) &= \log p_Z(Y(0)), \end{aligned} \quad (3.10)$$

and denote  $\bar{Y} = (Y, J) \in \mathbb{R}^{n+1}$ , the optimization problem above can be written as

$$\min_{\theta} \mathbb{E} \left[ V(\bar{Y}(T, \theta)) \right], \quad V(\bar{Y}) := \bar{Y}_{n+1} + U(\bar{Y}_{1:n})$$

Finding an optimal parameter  $\theta$  for the augmented dynamics (3.10) with terminal loss  $V$  is essentially the task for neural ODEs which had been revisited in Chapter 2.3.

### 3.4.4 Training of ODE flow for learning tasks with data

Similar to the discrete-time case, one can train the ODE flow with log-likelihood as loss (or equivalently the KL divergence as discussed above):

$$\max_{\theta} \mathbb{E}_{X \sim p_X} [\log p(X, T, \theta)] = \frac{1}{N} \sum_{i=1}^N \log p(x_i, T, \theta).$$

Using (3.9), we could simply solve the following system of ODEs:

$$\begin{aligned} \frac{d}{dt} Y(t, \theta) &= F(Y(t, \theta), t, \theta), & Y(T, \theta) &= x_i; \\ \frac{d}{dt} J(t, \theta) &= -(\nabla_y \cdot F)(Y(t, \theta), t), & J(0, \theta) &= \log p_Z(Y(0, \theta)). \end{aligned}$$

This can also be formulated as a standard neural ODE problem with a little effort.

## 3.5 Further Readings

Below we list a few related topics for reading and possible exploration:

- Boltzmann Generator [22]: use normalizing flows to sample from the Boltzmann distributions of many-body particles.
- The Benamou-Brenier formalism [3] (or say a continuous-time dynamical perspective) in optimal transport, which connects to the ODE flow formalism.
- We only cover a few narrow aspects of normalizing flows. More can be found in review papers [14, 23] and references therein, as well as [4, Chapter 18], e.g., we haven't covered auto-regressive flow above.
- We haven't discussed the estimation of the complexity of computing Jacobian, which is left as exercises; see e.g., [4, Chapter 18].

# Chapter 4

## Langevin Sampling

The open system studies the interaction of the system with its environment. As the environment could have a huge amount of degrees of freedom (e.g., a large amount of water molecules), we won't want to, and we are unable to directly study all degrees of freedom of the environment. Therefore, this challenge leads into a topic called *open system*. For instance, in the example of the discovery of Brownian motion, the pollen particle is the system that we want to study, and all water molecules form an environment. The study of pollen particle without modeling all water molecules is a typical example for the study of open classical systems.

Among many models, *Langevin dynamics* is an important dynamical equation to characterize the interaction of system and environment, with applications in various research fields, in particular with recent applications in machine learning. Later we will show why compared with normalizing flows and ODE flows above (cf. Chp. 3.3.3), Langevin dynamics can be used as a sampling algorithm without any training. The study of Langevin dynamics and Langevin sampling is a very large topic and we will only very concisely touch the very basic concepts only.

### 4.1 Generalized Langevin Equation and Underdamped Langevin Dynamics

For those who are interested in its origin in physics, one may refer to classical papers [20, 34] in last 60-70s or [24, Chapter 8] for the introduction of deriving Langevin dynamics in open classical systems. The *Generalized Langevin equation (GLE)* takes the following form [34, Eq. (22)]

$$\begin{aligned}\dot{Q}(t) &= P(t) \\ \dot{P}(t) &= -\nabla U(Q(t)) - \int_0^t dt' \zeta(t')P(t-t') + \mathcal{F}(t)\end{aligned}$$

where  $U$  is the potential function,  $\zeta$  is the correlation function, the noise  $\mathcal{F}(t)$  has mean zero and correlation  $\langle \mathcal{F}(t)\mathcal{F}(t') \rangle = k_B \mathbb{T} \zeta(t-t')$ , where  $k_B$  is the Boltzmann constant,  $\mathbb{T}$  is the temperature.

Suppose that the correlation function decays sufficiently fast, e.g.,  $\zeta(t) \approx 2\gamma \frac{e^{-t^2/(2\delta^2)}}{\sqrt{2\pi\delta^2}}$  ( $\delta \ll 1$ ) is approximately a delta measure, then we may model the above dynamics “equivalently” using the following SDE, also known as the *underdamped Langevin dynamics*:

$$\begin{aligned}dQ(t) &= P(t) dt \\ dP(t) &= -\nabla U(Q(t))dt - \gamma P(t) dt + \sqrt{2\gamma\beta^{-1}} dW(t),\end{aligned}\tag{4.1}$$

where  $\gamma > 0$  represents the damping coefficient, and  $\beta = \frac{1}{k_B \mathbb{T}}$  means the inverse temperature. When the temperature  $\mathbb{T}$  is large,  $\beta^{-1} \gg 1$ , and the fluctuation due to Brownian motion is more significant; when the temperature  $\mathbb{T}$  is very low,  $\beta^{-1} \ll 1$ , and the influence of random noise due to environment is small. Hence, at least, the role of  $\beta$  is indeed compatible with thermodynamics.

For underdamped Langevin dynamics (4.1), the convergence analysis of the relevant Fokker-Planck equation, the limiting behavior, as well as the performances of various numerical discretization schemes have been active research topics in the last decade and there are still open questions along this line.

## 4.2 Over-damped Langevin Sampling

For simplicity, let us consider *over-damped Langevin dynamics* instead:<sup>1</sup>

$$dX(t) = -\nabla U(X(t)) dt + \sqrt{2\beta^{-1}} dW(t). \quad (4.2)$$

The derivation of this SDE from underdamped Langevin dynamics (4.1) in the limit of  $\gamma \rightarrow \infty$  could be found in e.g., [24, Chapter 6.5]. This is perhaps also where the name “overdamped” comes from.

The Fokker-Planck equation of this SDE is the following:

$$\partial_t p(x, t) \stackrel{(1.18)}{=} \nabla \cdot (\nabla U(x)p(x, t)) + \beta^{-1} \Delta p(x, t). \quad (4.3)$$

It can be easily checked that

$$p_\infty(x) \propto e^{-\beta U(x)},$$

is a stationary solution to the above PDE.

*Remark.* Of course, to reach such an equilibrium, one may use time-dependent potential  $U(\cdot, t)$ . Designing a dynamics (namely, designing drift and diffusion terms) so that it reaches a certain equilibrium faster is a topic of interest for research.

Due to such a convergence behavior, one may design a numerical scheme to simulate the above SDE to time  $T$  where  $T \gg 1$ , so that  $p(\cdot, T) \approx p_\infty$ . If so, then  $X(T)$  are samples approximately from  $p_\infty$ . To simulate the SDE, one can use a discrete time-step  $\Delta t = \frac{T}{N}$  where  $\Delta t \ll 1$ , and discretize the SDE using Euler-Maruyama method:

$$\hat{X}_{k+1} = \hat{X}_k - \nabla U(\hat{X}_k) \Delta t + \sqrt{2\beta^{-1}\Delta t} \mathcal{N}(0, 1), \quad (4.4)$$

where  $\mathcal{N}(0, 1)$  are standard random variables, e.g., generated by Box-Muller transform discussed in Chapter 3.2. Fortunately, most numerical packages have already provided reliable interface to a generator of normal random variables, like numpy in Python.

In practice, we don't have to re-draw samples every time. If one wants to estimate the expected value in the following form

$$\int_{\mathbb{R}^n} f(x) p_\infty(x) dx = \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T f(X(t)) dt$$

where the equality holds due to ergodicity. One may further approximate the integral via

$$\lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T f(X(t)) dt \approx \frac{1}{T} \int_0^T f(X(t)) dt \approx \frac{1}{N} \sum_{i=1}^N f(\hat{X}_k). \quad (4.5)$$

where the time step is  $\Delta t = \frac{T}{N}$  is supposed to be small.

Langevin dynamics, though originated from statistical physics, can be used for other statistical problems, and machine learning problems. A common task is that one wants to sample from a certain distribution  $q$ , one may define  $U(x) := -\beta^{-1} \log q(x)$  and thus one can utilize such a dynamics above (4.2) to sample from  $q$  (provided that one is given an expression of  $q$  up to some multiplicative prefactor).

<sup>1</sup>the notation of the position  $Q$  has been changed to  $X$  below



### 4.3 Further Readings

- [24]: Chapter 6 and Chapter 8 provide a detailed mathematical description of Langevin dynamics and its Fokker-Planck equations.
- A more comprehensive introduction of this sampling topic can be found in [6].
- Algorithmic problem: Due to the hardness for Langevin dynamics to converge when a tall energy barrier exists under low temperature  $\beta \gg 1$ , it is still an active research topic to explore more efficient algorithms.



Diffusion-based models and other variants mostly take the second approach (if not all). The dynamical system can of course be (a) deterministic mapping; (b) Markov Chains; (c) ODE flows; (d) SDEs. To give some background, in the earlier stage of developing diffusion models, the process of generating samples from  $p_X$  utilized discrete-time Markov process [28, 30, 10]. Later, an SDE-based formulation was available in [31] which unifies previous Markov process based models. Not too much time later, diffusion-like models based on ODEs become a popular alternatives, and frameworks like flow-matching [18] and rectified flow [19] are introduced; these frameworks can be further unified within stochastic interpolants [1]. More recently, for the consideration of faster generation, the advantage of discrete mapping becomes more appreciated and e.g., consistency model [29] was introduced partially under this background. So far, the diffusion-based models and many variants have become a powerful tool in data synthesis, e.g., text-to-image generation and et al.

In the following notes, for simplicity, we will not follow historical development, and one may check the review papers for an overarching understanding of this field and many other important works. Due to the time limitation in this mini-course, we will only cover the most basic mathematical objects behind these models and only focus on continuous-time dynamics.

## 5.1 SDE-based Diffusion Models

### Step 1: specify the path

In the following, we will adopt the continuous-time formalism based on [31] for easier illustration. In this case, we consider the following OU process

$$dX(t) = b(X(t)) dt + \sqrt{2} dW_t, \quad b(x) = -x. \quad (5.2)$$

This is a very special Langevin dynamics with potential  $U(x) = |x|^2/2$ . Therefore, the limiting distribution is  $p_\infty = \mathcal{N}(0, \mathbf{I}_n)$  (see Exercise 1.18). With a time  $T$  large enough, we know that  $p(\cdot, T) \approx \mathcal{N}(0, \mathbf{I}_n) := p_Z$ . For the time being, let us worry not the gap between  $p_Z$  and  $p(\cdot, T)$ . What this process really did is to define a path in the space of probability distributions. The associated Fokker-Planck equation is

$$\partial_t p(x, t) = -\nabla \cdot (b(x)p(x, t)) + \Delta p(x, t),$$

and if we use  $\tilde{p}(x, t) := p(x, T - t)$  for the time-reversed process, then clearly

$$\partial_t \tilde{p}(x, t) = \nabla \cdot (b(x)\tilde{p}(x, t)) - \Delta \tilde{p}(x, t).$$

### Step 2: specify the reverse-time SDE

In order to reverse the process for data generation (driving a random variable following  $p_Z$  to a random variable following  $p_X$ ), it is natural to image that one also uses an SDE

$$dY(t) = A(Y(t), t) dt + \sqrt{2h} dW_t, \quad Y(0) \sim p(\cdot, T) \approx p_Z. \quad (5.3)$$

Since this SDE also need to follow the same path (in the sense of associated probability distributions), we have

$$\nabla \cdot (b(x)\tilde{p}(x, t)) - \Delta \tilde{p}(x, t) = \partial_t \tilde{p}(x, t) = -\nabla \cdot (A(x, t)\tilde{p}(x, t)) + h \Delta \tilde{p}(x, t).$$

By matching terms, we need

$$b(x)\tilde{p}(x, t) - \nabla \tilde{p}(x, t) = -A(x, t)\tilde{p}(x, t) + h\nabla \tilde{p}(x, t)$$

which leads into

$$A(x, t) = -b(x) + (h + 1) \frac{\nabla \tilde{p}(x, t)}{\tilde{p}(x, t)} = -b(x) + (h + 1) \nabla \log \tilde{p}(x, t).$$

### Step 3: Training of the score function

Therefore, the unknown part is  $\nabla \log p(x, t)$ , which is also known as the *score function*. A natural idea is to parameterize  $s_\theta(x, t) \approx \nabla \log p(x, t)$  and use the mean-square error

$$\begin{aligned}
 \min_{\theta} L_{\text{SDE}}(\theta) &:= \int_0^T \int_{\mathbb{R}^n} \|\nabla \log p(x, t) - s_\theta(x, t)\|^2 p(x, t) \, dt \\
 &= \int_0^T \int_{\mathbb{R}^n} \|s_\theta(x, t)\|^2 p(x, t) - 2s_\theta(x, t) \cdot \nabla_x \log p(x, t) p(x, t) \, dx \, dt + C \\
 &= \int_0^T \int_{\mathbb{R}^n} \|s_\theta(x, t)\|^2 p(x, t) - 2s_\theta(x, t) \cdot \nabla_x p(x, t) \, dx \, dt + C \\
 &= \int_0^T dt \int_{\mathbb{R}^n} dx \int_{\mathbb{R}^n} dy \|s_\theta(x, t)\|^2 P(x, t|y, 0)p(y, 0) - 2s_\theta(x, t) \cdot \nabla_x P(x, t|y, 0)p(y, 0) + C \\
 &= \int_0^T dt \mathbb{E}_{X_0 \sim p(\cdot, 0)} \mathbb{E}_{X_t \sim P(\cdot, t|X_0, 0)} \left[ \|s_\theta(X_t, t)\|^2 - 2s_\theta(X_t, t) \cdot \nabla_x \log P(X_t, t | X_0, 0) \right] + C \\
 &= \int_0^T dt \mathbb{E}_{X_0 \sim p(\cdot, 0)} \mathbb{E}_{X_t \sim P(\cdot, t|X_0, 0)} \left[ \|s_\theta(X_t, t) - \nabla_x \log P(X_t, t | X_0, 0)\|^2 \right] + C
 \end{aligned}$$

where “+C” means some terms independent of the parameter  $\theta$ . Therefore, we only need to solve

$$\min_{\theta} \int_0^T dt \mathbb{E}_{X_0 \sim p(\cdot, 0)} \mathbb{E}_{X_t \sim P(\cdot, t|X_0, 0)} \left[ \|s_\theta(X_t, t) - \nabla_x \log P(X_t, t | X_0, 0)\|^2 \right], \quad (5.4)$$

where the transition probability for OU process is easy to obtain: for a fixed  $y$ , the conditional probability  $P(x, t|y, 0)$  follows the normal random variable with mean  $e^{-t}y$  and covariance matrix  $c(t)\mathbf{I}_n$  and the scalar-valued function

$$c(t) = 2 \int_0^t e^{-2(t-s)} \, ds.$$

see Example 1.13; hence,

$$P(x, t|y, 0) = \frac{1}{\sqrt{2\pi c(t)^n}} \exp\left(-\frac{\|x - e^{-t}y\|^2}{2c(t)}\right) \implies \nabla_x \log P(x, t|y, 0) = -\frac{x - e^{-t}y}{c(t)},$$

which is easily computable.

### Summary

One can train a parameterized function  $(x, t) \in \mathbb{R}^n \times \mathbb{R} \rightarrow s_\theta(x, t) \in \mathbb{R}^n$  using the loss function in (5.4) and then use the SDE (5.3) to generate samples, also copied below

$$\begin{aligned}
 dY(t) &= A_\theta(Y(t), t) \, dt + \sqrt{2h} \, dW_t, & Y(0) &\sim p(\cdot, T) \approx p_Z \\
 A_\theta(x, t) &= -b(x) + (h+1)s_\theta(x, T-t).
 \end{aligned} \quad (5.5)$$

When we pick  $h = 0$  in (5.5), then the backward dynamics is simply an ODE.

## 5.2 ODE-based Models

### Specify the loss and the ODE

Suppose that the probability distribution at time  $t = T$  is (or very close to) the prior distribution  $p_Z$  (e.g., chosen as  $p_Z = \mathcal{N}(0, \mathbf{I}_n)$  for simplicity), and at time  $t = 0$ , it is the target distribution  $p_X$ . As discussed earlier, there are infinitely many paths in the probability space to connect them, let us say a smooth path

$\{p(\cdot, t)\}_{t=0}^T$  in the functional space  $\mathcal{P}_2(\mathbb{R}^n)$ , then there exists a velocity field  $v$  satisfying (5.1). One only needs to parameterize the velocity field using  $(x, t) \mapsto v_\theta(x, t)$  so that

$$\begin{aligned} \min_{\theta} L_{\text{ODE}}(\theta) &:= \int_0^T \int_{\mathbb{R}^n} \|v(x, t) - v_\theta(x, t)\|^2 p(x, t) \, dt \\ &= \int_0^T \mathbb{E} \left[ \|v(X(t), t) - v_\theta(X(t), t)\|^2 \right] \end{aligned}$$

where  $X(t) \sim p(\cdot, t)$ . Therefore, the next task becomes to find some examples of such  $X(t)$  and also  $v(X(t), t)$  is easily computable.

### Specify the path

A widely used example is the following, studied in e.g., [18, 19, 1]:

$$X(t) = \alpha(t)X + \beta(t)Z \tag{5.6}$$

where  $X \sim p_X$  and  $Z \sim p_Z$ , and  $\alpha(0) = 1, \alpha(T) = 0, \beta(0) = 0, \beta(T) = 1$  are used to satisfy terminal conditions. It has been shown in [1] that:

**Lemma 5.1** ([1, Theorem 2.6]). For (5.6), let  $p(\cdot, t)$  be the associated probability distributions and thus specifies the velocity field  $v(\cdot, t)$  to satisfy (5.1). Then

$$v(x, t) = \mathbb{E}[\dot{X}(t)|X(t) = x] \equiv \mathbb{E}[\dot{\alpha}(t)X + \dot{\beta}(t)Z | X(t) = x]. \tag{5.7}$$

*Proof.* For any test function  $\phi \in C_c^\infty(\mathbb{R}^n)$ ,

$$\frac{d}{dt} \mathbb{E}[\phi(X(t))] = \mathbb{E}[\nabla \phi(X(t)) \cdot (\dot{\alpha}(t)X + \dot{\beta}(t)Z)]$$

and in the weak form

$$\begin{aligned} \frac{d}{dt} \mathbb{E}[\phi(X(t))] &= \frac{d}{dt} \int \phi(x)p(x, t) \, dx = \int \phi(x) \partial_t p(x, t) \, dx \\ &= - \int \phi(x) \nabla \cdot (v(x, t)p(x, t)) \, dx = \int \nabla \phi(x) v(x, t) p(x, t) \, dx \\ &= \mathbb{E}[\nabla \phi(X(t)) v(X(t), t)] \end{aligned}$$

Since the test function  $\phi$  can be arbitrary, this proves (5.7). □

Therefore, one can write the loss function as

$$\begin{aligned} L_{\text{ODE}}(\theta) &= \int_0^T \mathbb{E} \left[ \|v_\theta(X(t), t)\|^2 - 2 \langle v(X(t), t), v_\theta(X(t), t) \rangle \right] + C \\ &= \int_0^T \mathbb{E} \left[ \|v_\theta(X(t), t)\|^2 \right] - 2 \mathbb{E} \left[ \langle v(X(t), t), v_\theta(X(t), t) \rangle \right] + C \\ &= \int_0^T \mathbb{E} \left[ \|v_\theta(X(t), t)\|^2 \right] - 2 \mathbb{E} \left[ \mathbb{E} \left[ \langle v(X(t), t), v_\theta(X(t), t) \rangle | X(t) \right] \right] + C \\ &= \int_0^T \mathbb{E} \left[ \|v_\theta(X(t), t)\|^2 \right] - 2 \mathbb{E} \left[ \mathbb{E} \left[ \langle \dot{\alpha}(t)X + \dot{\beta}(t)Z, v_\theta(X(t), t) \rangle | X(t) \right] \right] + C \\ &= \int_0^T \mathbb{E} \left[ \|v_\theta(X(t), t)\|^2 \right] - 2 \mathbb{E} \left[ \langle \dot{\alpha}(t)X + \dot{\beta}(t)Z, v_\theta(X(t), t) \rangle \right] + C \\ &= \int_0^T \mathbb{E} \left[ \|v_\theta(X(t), t) - (\dot{\alpha}(t)X + \dot{\beta}(t)Z)\|^2 \right] + C. \end{aligned}$$

**Summary**

Therefore, one can solve the following optimization problem

$$\min_{\theta} \int_0^T \mathbb{E}_{X \sim p_X, Z \sim p_Z} \left[ \left\| v_{\theta}(X(t), t) - (\dot{\alpha}(t)X + \dot{\beta}(t)Z) \right\|^2 \right], \quad (5.8)$$

and then simulate the following ODE systems (with time-already reversed):

$$\frac{d}{dt}Y(t) = -v_{\theta}(Y(t), T-t) \quad Y(0) \sim p_Z, \quad (5.9)$$

for the data generation using the output  $Y(T)$ .

**5.3 Further Readings**

- [4, Chapter 20]: a textbook style introduction to this topic.
- Classifier-free Diffusion Guidance [11]: generating samples conditioned on certain label (or say class) information, which leads into the application like text-to-image generation.
- Stochastic Interpolants [1]: a very unified framework of many mathematical models.
- Consistency Models [29]: provide a one-shot generation (aka fast) for a data generation trajectory.

# Bibliography

- [1] Michael S. Albergo, Nicholas M. Boffi, and Eric Vanden-Eijnden. Stochastic Interpolants: A Unifying Framework for Flows and Diffusions, 2023. arXiv:2303.08797 [cond-mat].
- [2] Luigi Ambrosio and Giuseppe Savaré. Chapter 1 - Gradient Flows of Probability Measures. In C. M. Dafermos and E. Feireisl, editors, *Handbook of Differential Equations: Evolutionary Equations*, volume 3, pages 1–136. 2007.
- [3] Jean-David Benamou and Yann Brenier. A computational fluid mechanics solution to the Monge-Kantorovich mass transfer problem. *Numerische Mathematik*, 84(3):375–393, 2000.
- [4] Christopher M. Bishop and Hugh Bishop. *Deep Learning: Foundations and Concepts*. Springer International Publishing, 2024.
- [5] Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural ordinary differential equations. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, NIPS’18, pages 6572–6583. Curran Associates Inc., 2018.
- [6] Sinho Chewi. *Log-Concave Sampling*. 2024. <https://chewisinho.github.io/main.pdf>.
- [7] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real NVP. In *International Conference on Learning Representations*, 2017.
- [8] Weinan E. A Proposal on Machine Learning via Dynamical Systems. *Communications in Mathematics and Statistics*, 5(1):1–11, 2017.
- [9] Weinan E and Bing Yu. The Deep Ritz Method: A Deep Learning-Based Numerical Algorithm for Solving Variational Problems. *Communications in Mathematics and Statistics*, 6(1):1–12, 2018.
- [10] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising Diffusion Probabilistic Models. In *Advances in Neural Information Processing Systems*, volume 33, pages 6840–6851, 2020.
- [11] Jonathan Ho and Tim Salimans. Classifier-Free Diffusion Guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021.
- [12] Fima C Klebaner. *Introduction to stochastic calculus with applications*. World Scientific Publishing Company, 2012.
- [13] Peter E. Kloeden and Eckhard Platen. *Numerical Solution of Stochastic Differential Equations*. Springer, Berlin, Heidelberg, 1992.
- [14] Ivan Kobyzev, Simon J.D. Prince, and Marcus A. Brubaker. Normalizing Flows: An Introduction and Review of Current Methods. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(11):3964–3979, 2021.
- [15] Qianxiao Li. Dynamical Systems and Machine Learning. 2020. URL: <https://www.math.pku.edu.cn/amel/docs/20200719122925684287.pdf>.

- [16] Ruichen Li, Haotian Ye, Du Jiang, Xuelan Wen, Chuwei Wang, Zhe Li, Xiang Li, Di He, Ji Chen, Weiluo Ren, and Liwei Wang. A computational framework for neural network-based variational Monte Carlo with Forward Laplacian. *Nature Machine Intelligence*, 6(2):209–219, 2024.
- [17] Xuechen Li, Ting-Kam Leonard Wong, Ricky T. Q. Chen, and David Duvenaud. Scalable Gradients for Stochastic Differential Equations. In *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, pages 3870–3882. PMLR, 2020. arXiv:2001.01328.
- [18] Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. In *The Eleventh International Conference on Learning Representations*, 2023.
- [19] Xingchao Liu, Chengyue Gong, and qiang liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. In *The Eleventh International Conference on Learning Representations*, 2023.
- [20] Hazime Mori. Transport, Collective Motion, and Brownian Motion. *Progress of Theoretical Physics*, 33(3):423–455, March 1965.
- [21] Radford M. Neal. Annealed importance sampling. *Statistics and Computing*, 11(2):125–139, 2001.
- [22] Frank Noé, Simon Olsson, Jonas Köhler, and Hao Wu. Boltzmann generators: Sampling equilibrium states of many-body systems with deep learning. *Science*, 365(6457):eaaw1147, 2019.
- [23] George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji Lakshminarayanan. Normalizing Flows for Probabilistic Modeling and Inference. *Journal of Machine Learning Research*, 22(57):1–64, 2021.
- [24] Grigorios A Pavliotis. Stochastic processes and applications. *Texts in applied mathematics*, 60, 2014.
- [25] M. Raissi, P. Perdikaris, and G. E. Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics*, 378:686–707, 2019.
- [26] Sebastian Ruder. An overview of gradient descent optimization algorithms, 2017. arXiv:1609.04747.
- [27] Mark Schmidt. CPSC 540: Machine Learning: SGD Convergence Rate, 2019. <https://www.cs.ubc.ca/~schmidtm/Courses/540-W19/L11.pdf>.
- [28] Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37*, pages 2256–2265, 2015.
- [29] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency Models. In *Proceedings of the 40th International Conference on Machine Learning*, pages 32211–32252, 2023.
- [30] Yang Song and Stefano Ermon. Generative Modeling by Estimating Gradients of the Data Distribution. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- [31] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-Based Generative Modeling through Stochastic Differential Equations. In *International Conference on Learning Representations*, 2021.
- [32] Unknown. Einstein’s random walk. *Physics World*, 2005. URL: <https://physicsworld.com/a/einsteins-random-walk/>.
- [33] Cédric Villani. *Optimal Transport: Old and New*. Springer-Verlag, Berlin Heidelberg, 2009.
- [34] Robert Zwanzig. Nonlinear generalized langevin equations. *Journal of Statistical Physics*, 9(3):215–220, 1973.



## Appendix A

# Erratum (which has been corrected)

- A missing argument ( $r$ ) inside the Theorem (1.2).
- The function  $f$  in the original version really means the loss function  $L$  inside Section 1.4.3 SGD.
- Two typos inside Section 2.1 for the output state.
- The normalizing constant has been changed to  $\mathcal{Z}$  to avoid notation conflicts.