

---

# Deep Learning

---

Mingming Li

First Created: September 29, 2020  
Last Modified: March 6, 2023

---

# Contents

<b>Contents</b>	<b>iii</b>
<b>List of Figures</b>	<b>vii</b>
<b>List of Tables</b>	<b>ix</b>
<b>I Theory</b>	<b>1</b>
<b>1 Probability</b>	<b>3</b>
1.1 Markov chain . . . . .	3
1.1.1 Transitions . . . . .	3
1.1.2 Discrete-time Markov chain . . . . .	3
1.1.3 Continous-time Markov chain . . . . .	3
1.1.4 Markov kernel . . . . .	4
1.1.5 Measurable space . . . . .	4
1.1.6 $\sigma$ -algebra . . . . .	4
1.2 Conditional probability . . . . .	4
1.3 Posterior probability . . . . .	5
1.4 Guassian distribution . . . . .	5
1.5 Laplace distribution . . . . .	6
1.6 Monte Carlo method . . . . .	6
1.7 Variational Bayesian methods . . . . .	6

1.8	Bayesian probability . . . . .	7
1.9	Binomial distribution . . . . .	7
1.10	Kullback–Leibler divergence . . . . .	8
1.11	Jensen’s inequality . . . . .	8
<b>2</b>	<b>Learning</b>	<b>9</b>
<b>3</b>	<b>Model</b>	<b>11</b>
3.1	Capacity . . . . .	11
3.1.1	The no free lunch theorem . . . . .	11
3.1.2	Universal approximation theorem . . . . .	12
3.2	Overfitting and underfitting . . . . .	12
<b>4</b>	<b>Training (optimization)</b>	<b>13</b>
4.1	SGD . . . . .	13
4.2	Momentum . . . . .	14
4.3	Nesterov momentum . . . . .	14
4.4	AdaGrad . . . . .	14
4.5	Adam . . . . .	15
<b>5</b>	<b>Cost function</b>	<b>17</b>
5.1	Regression cost function . . . . .	17
5.1.1	Mean squared error . . . . .	17
5.1.2	Mean absolute error . . . . .	17
5.2	Classification cost function . . . . .	17
<b>6</b>	<b>Full Connected Networks</b>	<b>19</b>
<b>7</b>	<b>CNN</b>	<b>21</b>
7.1	Convolution . . . . .	21
7.2	Properties . . . . .	22
7.2.1	Sparse interaction . . . . .	22

7.2.2	Parameter sharing . . . . .	22
7.2.3	Equivariant representations . . . . .	22
7.3	Pooling . . . . .	23
<b>8</b>	<b>Metric</b>	<b>25</b>
8.1	Precision . . . . .	25
8.2	Recall . . . . .	25
8.3	Accuracy . . . . .	25
8.4	F-score . . . . .	26
<b>9</b>	<b>Regularization</b>	<b>27</b>
9.1	Parameter norm penalties . . . . .	27
9.1.1	$L^2$ Parameter Regularization . . . . .	27
9.1.2	$L^1$ Regularization . . . . .	27
9.2	Dataset Augmentation . . . . .	28
9.3	Early stopping . . . . .	28
9.4	Droupout . . . . .	28
<b>10</b>	<b>Activation functions</b>	<b>29</b>
10.1	Desirable features . . . . .	29
10.2	Sigmoid . . . . .	29
10.3	Tanh . . . . .	29
10.4	ReLU . . . . .	29
10.5	Leaky ReLU . . . . .	30
10.6	Swish . . . . .	30
<b>11</b>	<b>Machine learning process</b>	<b>31</b>
<b>II</b>	<b>Models</b>	<b>33</b>
<b>12</b>	<b>Diffusion</b>	<b>35</b>

<b>III Tools</b>	<b>37</b>
<b>13 PyTorch</b>	<b>39</b>
<b>14 NumPy</b>	<b>41</b>
14.1 Load data . . . . .	41
<b>IV Practice</b>	<b>43</b>
<b>15 Preprocessing</b>	<b>45</b>
15.1 Not valued based data . . . . .	45
<b>V Projects</b>	<b>47</b>
<b>VI Papers</b>	<b>49</b>
<b>16 Deep Unsupervised Learning using Nonequilibrium Thermodynamics[1]</b>	<b>51</b>
16.1 Algorithm . . . . .	51
16.1.1 Forward trajectory . . . . .	51
16.1.2 Reverse Trajectory . . . . .	51
16.2 Model probability . . . . .	52
16.3 Training . . . . .	52
<b>Bibliography</b>	<b>55</b>

# List of Figures

1.1	Gaussian distribution . . . . .	5
1.2	Laplace distribution . . . . .	6
1.3	Binomial distribution . . . . .	7
1.4	Jensen's inequality . . . . .	8
6.1	Full connected neural network . . . . .	19
7.1	Convolution operation . . . . .	22





# List of Tables

8.1 Confusion matrix . . . . .	25
--------------------------------	----



# **Part I**

# **Theory**



# Probability

## 1.1 Markov chain

A **Markov chain** or **Markov process** is a stochastic model describing a sequence of possible events in which the probability of each event depends only on the state attained in the previous event.

### 1.1.1 Transitions

The changes of state of the system are called **transitions**. The probabilities associated with various state changes are called **transition probabilities**. The process is characterized by a state space, a transition matrix describing the probabilities of particular transitions, and an initial state (or initial distribution) across the state space.

### 1.1.2 Discrete-time Markov chain

A **discrete-time Markov chain** is a sequence of random variables  $X_1, X_2, X_3, \dots$  with the Markov property, namely that the probability of moving to the next state depends only on the present state and not on the previous states:

$$P(X_{n+1} = x \mid X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) = P(X_{n+1} = x \mid X_n = x_n) \quad (1.1)$$

if both conditional probabilities are well defined, that is, if  $P(X_1 = x_1, \dots, X_n = x_n) > 0$ .

The possible values of  $X_i$  form a countable set  $S$  called the state space of the chain.

### 1.1.3 Continuous-time Markov chain

A **continuous-time Markov chain**  $X(t)$  is defined by two components: a jump chain, and a set of holding time parameters  $\lambda_i$ . The jump chain consists of a countable set of states  $S \subset \{0, 1, 2, \dots\}$  along with transition probabilities  $p_{ij}$ . We assume  $p_{ii} = 0$ , for all non-absorbing states  $i \in S$ . We assume

- 1 if  $X(t) = i$ , the time until the state changes has Exponential ( $\lambda_i$ ) distribution;
- 2 if  $X(t) = i$ , the next state will be  $j$  with probability  $p_{ij}$ .

The process satisfies the Markov property. That is, for all  $0 \leq t_1 < t_2 < \dots < t_n < t_{n+1}$ , we have

$$\begin{aligned} P(X(t_{n+1}) = j \mid X(t_n) = i, \mathbf{X}(t_{n-1}) = i_{n-1}, \dots, X(t_1) = i_1) \\ = P(X(t_{n+1}) = j \mid \mathbf{X}(t_n) = i) \end{aligned}$$

### 1.1.4 Markov kernel

Let  $(X, \mathcal{A})$  and  $(Y, \mathcal{B})$  be **measurable spaces**. A **Markov kernel** with source  $(X, \mathcal{A})$  and target  $(Y, \mathcal{B})$  is a map  $\kappa : \mathcal{B} \times X \rightarrow [0, 1]$  with the following properties:

- 1 For every (fixed)  $B \in \mathcal{B}$ , the map  $x \mapsto \kappa(B, x)$  is  $\mathcal{A}$ -measurable
- 2 For every (fixed)  $x \in X$ , the map  $B \mapsto \kappa(B, x)$  is a **probability measure** on  $(Y, \mathcal{B})$

In other words it associates to each point  $x \in X$  a probability measure  $\kappa(dy \mid x) : B \mapsto \kappa(B, x)$  on  $(Y, \mathcal{B})$  such that, for every measurable set  $B \in \mathcal{B}$ , the map  $x \mapsto \kappa(B, x)$  is measurable with respect to the  $\sigma$ -algebra  $\mathcal{A}$ .

### 1.1.5 Measurable space

A **measurable space** consists of a set and a  **$\sigma$ -algebra**, which defines the subsets that will be measured.

Consider a set  $X$  and a  $\sigma$ -algebra  $\mathcal{A}$  on  $X$ . Then the tuple  $(X, \mathcal{A})$  is called a measurable space.

### 1.1.6 $\sigma$ -algebra

A  **$\sigma$ -algebra** (also  $\sigma$ -field) on a set  $X$  is a nonempty collection  $\Sigma$  of subsets of  $X$  closed under complement, countable unions, and countable intersections.

Let  $X$  be some set, and let  $P(X)$  represent its power set. Then a subset  $\Sigma \subseteq P(X)$  is called a  $\sigma$ -algebra if it satisfies the following three properties:

- 1  $X$  is in  $\Sigma$ , and  $X$  is considered to be the universal set in the following context.
- 2  $\Sigma$  is closed under complementation: If  $A$  is in  $\Sigma$ , then so is its complement,  $X \setminus A$ .
- 3  $\Sigma$  is closed under countable unions: If  $A_1, A_2, A_3, \dots$  are in  $\Sigma$ , then so is  $A = A_1 \cup A_2 \cup A_3 \cup \dots$ .

From these properties, it follows that the  $\sigma$ -algebra is also closed under countable intersections (by applying De Morgan's laws).

## 1.2 Conditional probability

Given two events  $A$  and  $B$ .

$$P(A \mid B) = \frac{P(A \cap B)}{P(B)} \tag{1.2}$$

$P(A|B)$  stands for the probability of  $A$  happening given that  $B$  happened.  $P(A \cap B)$  stands for the probability of  $A$  and  $B$  happening at the same time.  $P(B)$  stands for the probability of  $B$  happening.

### 1.3 Posterior probability

The **posterior probability** is a type of conditional probability that results from updating the prior probability with information summarized by the likelihood via an application of Bayes' rule.

In variational Bayesian methods, the posterior probability is the probability of the parameters  $\theta$  given the evidence  $X$ , and is denoted  $p(\theta | X)$ . It contrasts with the likelihood function, which is the probability of the evidence given the parameters:  $p(X | \theta)$ .

The two are related as follows: Given a prior belief that a probability distribution function is  $p(\theta)$  and that the observations  $x$  have a likelihood  $p(x | \theta)$ , then the posterior probability is defined as

$$p(\theta | x) = \frac{p(x | \theta)}{p(x)} p(\theta) \quad (1.3)$$

where  $p(x)$  is the normalizing constant and is calculated as

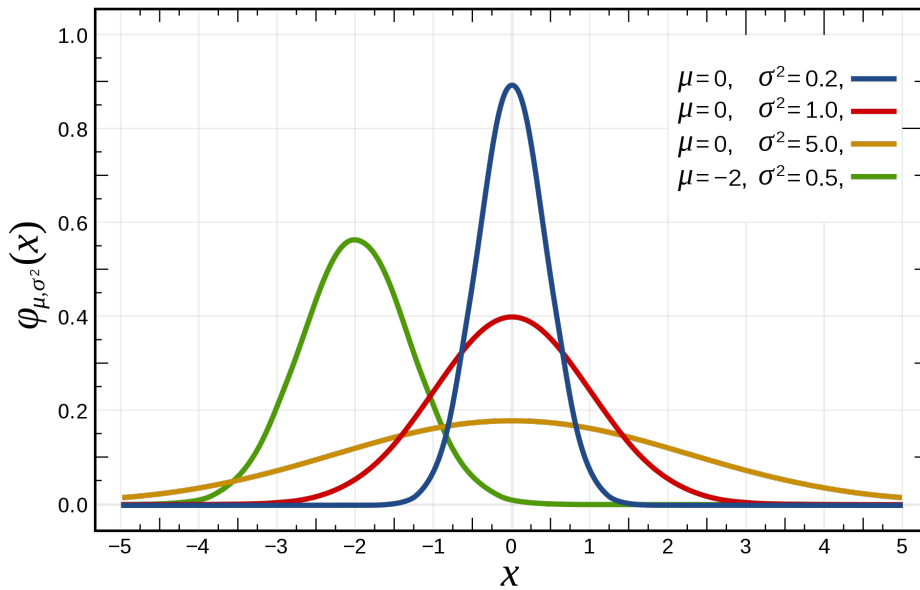
$$p(x) = \int p(x | \theta) p(\theta) d\theta \quad (1.4)$$

### 1.4 Gaussian distribution

**Gaussian distribution** (normal distribution) is a type of continuous probability distribution for a real-valued random variable. The general form of its probability density function is

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left( \frac{x-\mu}{\sigma} \right)^2} \quad (1.5)$$

The parameter  $\mu$  is the mean or expectation of the distribution (and also its median and mode), while the parameter  $\sigma$  is its standard deviation. The variance of the distribution is  $\sigma^2$ . The Figure 1.1 shows the Gaussian distributions.



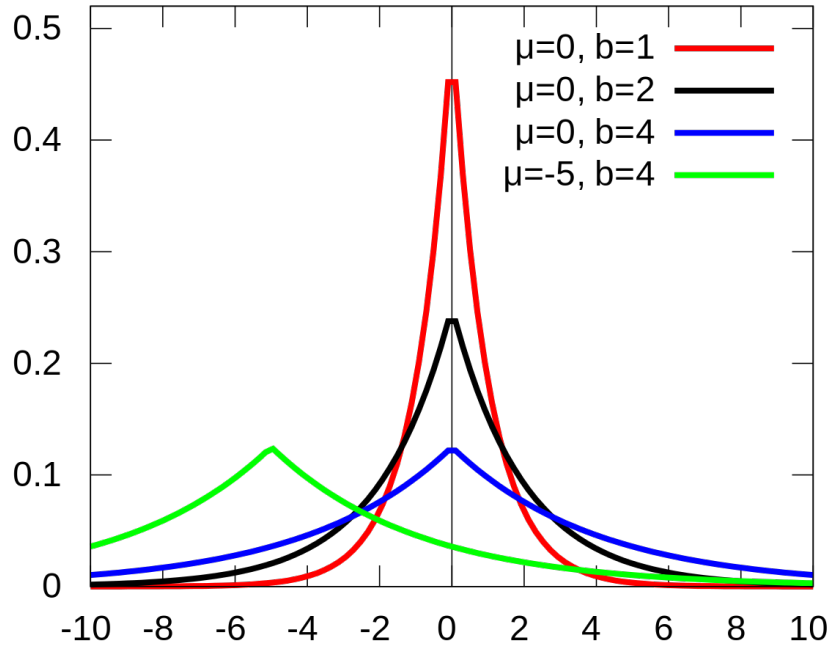
**Figures 1.1:** Gaussian distribution

## 1.5 Laplace distribution

A random variable has a Laplace  $(\mu, b)$  distribution if its probability density function is

$$f(x) = \frac{1}{2b} \exp\left(-\frac{|x - \mu|}{b}\right) \quad (1.6)$$

Here  $\mu$  is a location parameter and  $b$  is a scale parameter. Figure 1.2 shows the Laplace distributions.



**Figures 1.2:** Laplace distribution

## 1.6 Monte Carlo method

**Monte Carlo methods** are a broad class of computational algorithms that rely on repeated random sampling to obtain numerical results. The underlying concept is to use randomness to solve problems that might be deterministic in principle.

Monte Carlo methods vary, but tend to follow a particular pattern:

- 1 Define a domain of possible inputs
- 2 Generate inputs randomly from a probability distribution over the domain
- 3 Perform a deterministic computation on the inputs
- 4 Aggregate the results

## 1.7 Variational Bayesian methods

It is a simplifying that makes the intractable events into tractable events.

For example, in variational inference, the posterior distribution over a set of unobserved variables  $\mathbf{Z} =$



$\{Z_1 \dots Z_n\}$  given some data  $\mathbf{X}$  is approximated by a so-called variational distribution,  $Q(\mathbf{Z})$ :

$$P(\mathbf{Z} | \mathbf{X}) \approx Q(\mathbf{Z}) \quad (1.7)$$

The distribution  $Q(\mathbf{Z})$  is restricted to belong to a family of distributions of simpler form than  $P(\mathbf{Z} | \mathbf{X})$  (e.g. a family of Gaussian distributions), selected with the intention of making  $Q(\mathbf{Z})$  similar to the true posterior,  $P(\mathbf{Z} | \mathbf{X})$ .

The similarity (or dissimilarity) is measured in terms of a dissimilarity function  $d(Q; P)$  and hence inference is performed by selecting the distribution  $Q(\mathbf{Z})$  that minimizes  $d(Q; P)$ . The most common type of variational Bayes uses the Kullback-Leibler divergence (KL-divergence) of  $Q$  from  $P$  as the choice of dissimilarity function.

## 1.8 Bayesian probability

Bayesian probability is an interpretation of the concept of probability, in which, instead of frequency or propensity of some phenomenon, probability is interpreted as reasonable expectation representing a state of knowledge or as quantification of a personal belief.

## 1.9 Binomial distribution

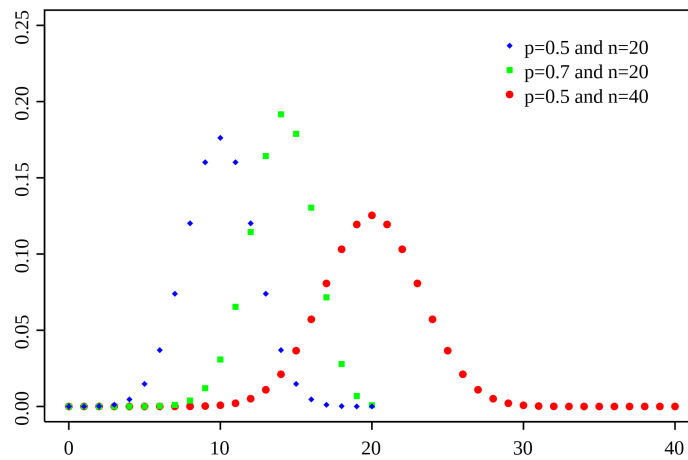
In general, if the random variable  $X$  follows the binomial distribution with parameters  $n \in \mathbb{N}$  and  $p \in [0, 1]$ , we write  $X \sim B(n, p)$ . The probability of getting exactly  $k$  successes in  $n$  independent Bernoulli trials is given by the probability mass function:

$$f(k) = \binom{n}{k} p^k (1-p)^{n-k} \quad (1.8)$$

for  $k = 0, 1, 2, \dots, n$ , where

$$\binom{n}{k} = \frac{n!}{k!(n-k)!} \quad (1.9)$$

The figure 1.3 shows the distribution.



**Figures 1.3:** Binomial distribution

### 1.10 Kullback–Leibler divergence

The Kullback–Leibler divergence, denoted  $D_{\text{KL}}(P\|Q)$ , is a type of statistical distance: a measure of how one probability distribution  $P$  is different from a second, reference probability distribution  $Q$ .

For discrete probability distributions  $P$  and  $Q$  defined on the same sample space,  $\mathcal{X}$ , the relative entropy from  $Q$  to  $P$  is defined <sup>[11]</sup> to be

$$D_{\text{KL}}(P\|Q) = \sum_{x \in \mathcal{X}} P(x) \log \left( \frac{P(x)}{Q(x)} \right). \quad (1.10)$$

which is equivalent to

$$D_{\text{KL}}(P\|Q) = - \sum_{x \in \mathcal{X}} P(x) \log \left( \frac{Q(x)}{P(x)} \right) \quad (1.11)$$

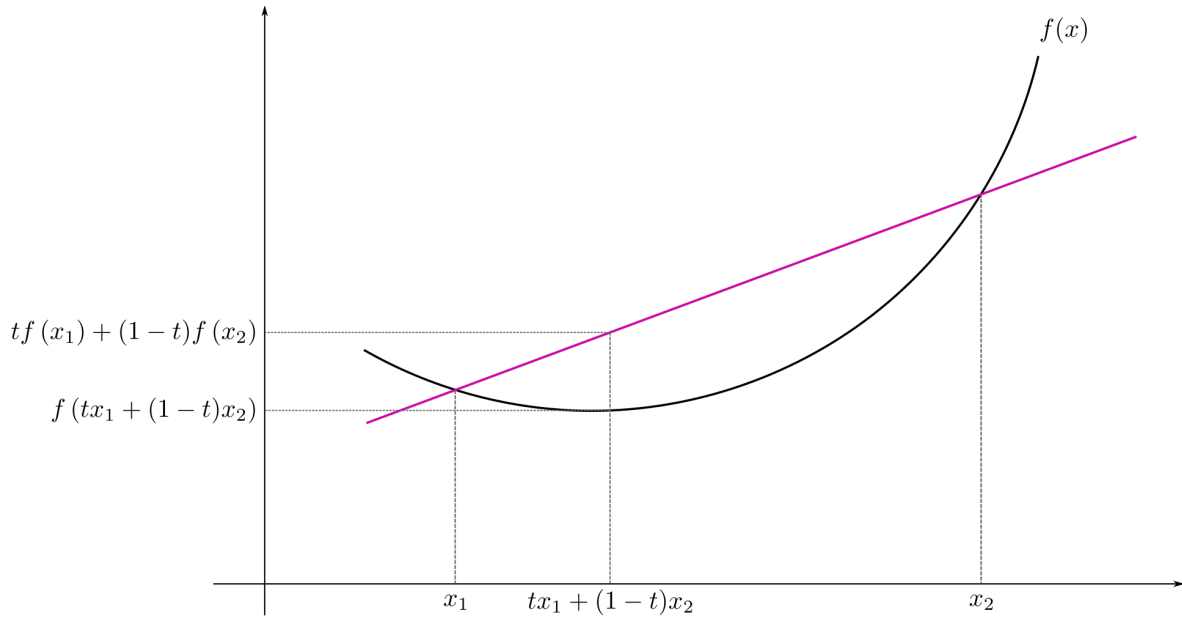
### 1.11 Jensen's inequality

Jensen's inequality generalizes the statement that a secant line of a convex function lies above its graph.

Thus, Jensen's inequality is

$$f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2). \quad (1.12)$$

Shown in Figure 1.4.



**Figures 1.4:** Jensen's inequality

In the context of probability theory, it is generally stated in the following form: if  $X$  is a random variable and  $\varphi$  is a convex function, then

$$\varphi(\mathbb{E}[X]) \leq \mathbb{E}[\varphi(X)] \quad (1.13)$$

The difference between the two sides of the inequality,  $\mathbb{E}[\varphi(X)] - \varphi(\mathbb{E}[X])$ , is called the Jensen gap.

# Chapter 2

## Learning

A machine learning algorithm is an algorithm that is able to learn from data. But what do we mean by learning? Mitchell provides a succinct definition: “A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ ”

Now, how does the learn happen? The model or algorithm learns by adjusting the parameters contained in it.

How to adjust the parameters? Refer to Chapter 4.



# Chapter 3

## Model

What is a model?

You can think a model as a function in math. For example,

$$ax + by = c$$

In this function,  $a, b$  and  $d$  are the parameters,  $x$  and  $y$  are input from data. Through adjusting the parameters  $a, b$  and  $d$ , we can implement the learning process. This is also called training in machine learning.

### 3.1 Capacity

The capacity is the function space. It defines the limitation that can learn from data.

For example,

$$ax + by = c$$

This function can only learn linear function from data.

$$ax_1^2 + bx_2 = c$$

This function can learn curve function from data.

#### 3.1.1 The no free lunch theorem

For any algorithms (functions)  $a_1$  and  $a_2$ , at iteration step  $m$

$$\sum P(d_m^y|f, m, a_1) = \sum P(d_m^y|f, m, a_2) \quad (3.1)$$

where  $d_m^y$  denotes the ordered set of size  $m$  of the cost values  $y$  associated to input values  $x \in X$ ,  $f : X \rightarrow Y$  is the function being optimized and  $P(d_m^y|f, m, a)$  is the conditional probability of obtaining a given sequence of cost values from algorithm  $a$  run  $m$  times on function  $f$ .

The no free lunch theorem implies that we must design our machine learning algorithms to perform well on a specific task but not a universal task.

### 3.1.2 Universal approximation theorem

A feed-forward network with a single hidden layer containing a finite number of neurons can approximate continuous functions on compact subsets of  $R^n$ , under mild assumptions on the activation function.

Let  $\varphi : R \rightarrow R$  be a nonconstant, bounded, and continuous functions. Let  $I_m$  denote the  $m$ -dimensional unit hypercube  $[0, 1]^m$ . The space of real-value continuous function on  $I_m$  is denoted by  $C(I_m)$ . Then, given any  $\varepsilon > 0$  and function  $f \in C(I_m)$ , there exist an integer  $N$ , real constants  $v_i, b_i \in R$  and real vectors  $\omega_i \in R^m$  for  $i = 1, \dots, N$ , such that we may define:

$$F(x) = \sum_{i=1}^N v_i \varphi(\omega_i^T x + b_i) \quad (3.2)$$

as an approximate realization of the function  $f$ ; that is

$$|F(x) - f(x)| < \varepsilon \quad (3.3)$$

for all  $x \in I_m$ . In other words, functions of the form  $F(x)$  are dense in  $C(I_m)$ .

This still holds when replacing  $I_m$  with any compact subset of  $R^m$ .

Kurt Hornik showed in 1991 that it is not the specific choice of the activation function, but rather the multi-layer feedforward architecture itself which gives neural networks the potential of being universal approximators. The output units are always assumed to be linear. For notational convenience, only the single output case will be shown. The general case can easily be deduced from the single output case.

In 2017 Lu et al. proved universal approximation theorem for width-bounded deep neural networks.

This is the base of deep learning.

## 3.2 Overfitting and underfitting

We train model on training data but use test data (not used to train the model) to test out model. The ability to perform on test data is called **generalization**. We can use model on test data because we assume that the train data and the test has the same probability distribution (i.e. they have relationship).

The error on training data is called **training error**. The error on test data is called **test error**. Underfitting occurs when the model is not able to obtain a sufficient low error value on the training set. Overfitting occurs when the gap between the training error and test error is too large.

We can control whether a model is more likely to overfit or underfit by altering its **capacity**. The overfitting and underfitting happen because of the mismatch of capacity and the data.

# Training (optimization)

Training is the process of learning. By inputting the data, we adjust the parameters in the model to achieve better performance. In machine learning this is also called the optimization.

## 4.1 SGD

Stochastic gradient descent (SGD) is a very basic and important algorithm.

The cost function used by a machine learning algorithm often decomposes as a sum over training examples of some per-example loss function.

$$J(\theta) = \mathbb{E}_{\mathbf{x}, y \sim \hat{p}_{data}} L(\mathbf{x}, y, \theta) = \frac{1}{m} \sum_{i=1}^m L(\mathbf{x}^{(i)}, y^{(i)}, \theta), \quad (4.1)$$

Where  $L$  is the loss function,  $\mathbf{x}$  and  $y$  are the input data and labels,  $\theta$  is the parameters in the model,  $m$  is the number of samples.

The gradient is

$$\nabla_{\theta} J(\theta) = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} L(\mathbf{x}^{(i)}, y^{(i)}, \theta). \quad (4.2)$$

The computational cost of this operation is  $O(m)$ . As the training set size grows to billions of examples, the time to take a single gradient step becomes prohibitively long.

The insight of SGD is that the gradient is an expectation. The expectation may be approximately estimated using a small set of samples. Specifically, on each step of the algorithm, we can sample a **minibatch** of examples  $\mathbb{B} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m')}\}$  drawn uniformly from the training set. The minibatch size  $m'$  is typically chosen to be a relatively small number of examples, ranging from one to a few hundred.

The estimator of the gradient is formed as

$$\mathbf{g} = \frac{1}{m'} \nabla_{\theta} \sum_{i=1}^{m'} L(\mathbf{x}^{(i)}, y^{(i)}, \theta) \quad (4.3)$$

using examples from the minibatch  $\mathbb{B}$ . The stochastic gradient descent algorithm then follows the estimated gradient downhill:

$$\theta \leftarrow \theta - \epsilon g, \quad (4.4)$$

where  $\epsilon$  is the learning rate.

## 4.2 Momentum

The method of momentum is designed to accelerate learning in SGD. The momentum algorithm accumulates an exponentially decaying moving average of past gradients and continues to move in their direction.

$$v \leftarrow \alpha v - \epsilon g \quad (4.5)$$

$$\theta \leftarrow \theta + v. \quad (4.6)$$

Where  $g$  is the gradient,  $\alpha \in [0, 1]$  determines how quickly the contributions of previous gradients exponentially decay.

## 4.3 Nesterov momentum

Nesterov Momentum is a variant of the Momentum algorithm. The difference between Nesterov momentum and standard momentum is where the gradient is evaluated. With Nesterov momentum the gradient is evaluated after the current velocity is applied.

$$g = \frac{1}{m'} \nabla_{\theta} \sum_{i=1}^{m'} L(x, y, \theta + \alpha v) \quad (4.7)$$

## 4.4 AdaGrad

AdaGrad is designed to converge rapidly when applied to a convex function. Comparing to SGD, AdaGrad algorithm individually adapts the learning rates of all model parameters by scaling them inversely proportional to the square root of the sum of all of their historical squared values.

$$\theta = \theta - \frac{\epsilon}{\sqrt{\delta \mathbf{I} + \text{diag}(G)}} \odot g \quad (4.8)$$

where  $\theta$  is the parameter to be updated,  $\epsilon$  is the initial learning rate,  $\delta$  is some small quantity that used to avoid the division of zero,  $\mathbf{I}$  is the identity matrix,  $g$  is the gradient estimate.

$$G = \sum_{\tau=1}^t g_{\tau} g_{\tau}^T \quad (4.9)$$

AdaGrad shrinks the learning rate according to the entire history of the squared gradient and may have made the learning rate too small before arriving at such a convex structure.



## 4.5 Adam

The RMSProp algorithm modifies AdaGrad to perform better in the non-convex setting by changing the gradient accumulation into an exponentially weighted moving average. RMSProp uses an exponentially decaying average to discard history from the extreme past so that it can converge rapidly after finding a convex bowl.

$$\mathbf{r} \leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g} \tag{4.10}$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{g} \tag{4.11}$$

Where  $\mathbf{g}$  is the gradient,  $\boldsymbol{\theta}$  is the parameters in a model,  $\mathbf{r}$  is initialized to 0..



## Cost function

**Loss** is the difference between the predicted value and the actual value. The Function used to quantify this loss during the training phase in the form of a single real number is known as **Loss Function**. **Cost function** refers to an average of the loss functions over an entire training dataset. Cost function helps us reach the optimal solution. The cost function is the technique of evaluating “the performance of our algorithm/model”.

There are many cost functions in machine learning and each has its use cases depending on whether it is a regression problem or classification problem.

### 5.1 Regression cost function

Regression models deal with predicting a continuous value. They are calculated on the distance-based error.

$$\text{Error} = y - y' \quad (5.1)$$

#### 5.1.1 Mean squared error

$$\text{MSE} = \frac{\sum_{i=0}^n (y - y')^2}{n} \quad (5.2)$$

#### 5.1.2 Mean absolute error

$$\text{MSE} = \frac{\sum_{i=0}^n |y - y'|}{n} \quad (5.3)$$

### 5.2 Classification cost function

In classification, we usually use one hot to encode the labels. This can eliminate the affect of the distance.

The cross entropy of the distribution  $q$  relative to distribution  $p$  over a given set is defined as

$$H(p, q) = -E_p[\log q] \quad (5.4)$$

Where the  $E$  is the expected value operator respected to the probability  $q$ .

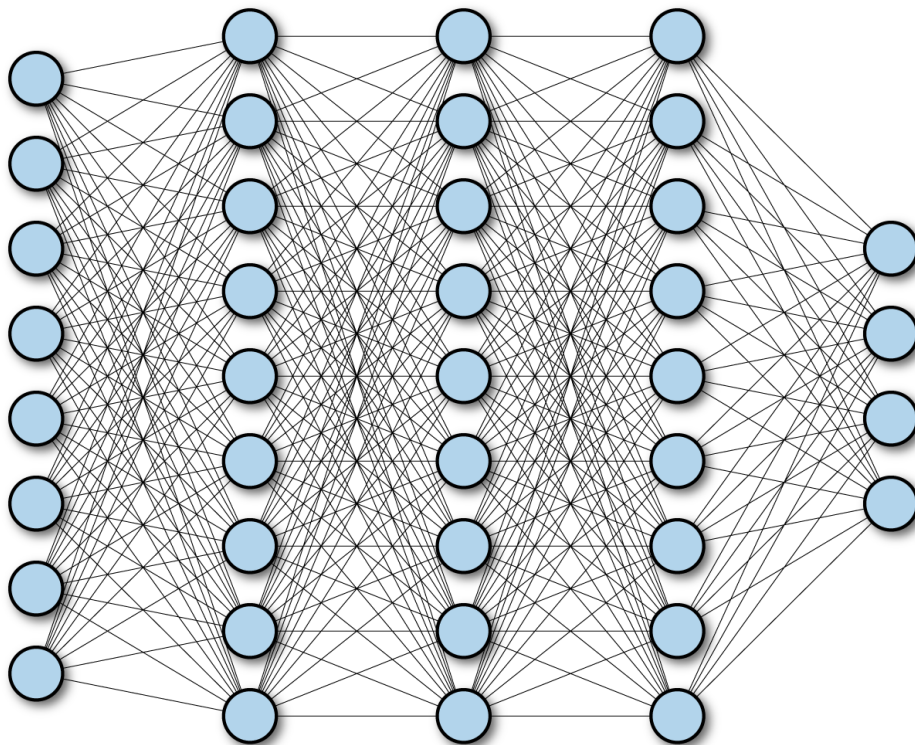
For discrete probability distribution  $p$  and  $q$  with the same support  $\mathcal{X}$  this means

$$H(P, Q) = - \sum_{x \in \mathcal{X}} p(x) \log q(x) \quad (5.5)$$

## Full Connected Networks

A fully connected neural network consists of a series of fully connected layers that connect every neuron in one layer to every neuron in the other layer. Full connected neural network layers use matrix multiplication by a matrix of parameters with a separate parameter describing the interaction between each input unit and each output unit. This means that every output unit interacts with every input unit.

Figure 6.1 show the full connects neural network.



**Figures 6.1:** Full connected neural network



# CNN

CNN stands for convolutional neural network. Convolutional networks are neural networks that have convolutional layers. A typical convolutional layer consists of three stages:

- 1 convolution stage: affine transform
- 2 detector stage: nonlinearty
- 3 pooling stage

## 7.1 Convolution

$$s(t) = \int x(a)w(t-a)da. \quad (7.1)$$

This operation is called **convolution**. The convolution operation is typically denoted with an asterisk:

$$s(t) = (x * w)(t). \quad (7.2)$$

In convolutional network terminology, the first argument (in this example, the function  $x$ ) to the convolution is often referred to as the **input**, and the second argument (int this example, the function  $w$ ) as the **kernel**. The output is sometimes referred to as the **feature map**.

If we assume that  $x$  and  $w$  are defined only on integer  $t$ , we can define the discrete convolution:

$$s(t) = (x * w)(t) = \sum_{a=-\infty}^{\infty} x(a)w(t-a). \quad (7.3)$$

We often use convolutions over more than one axis at a time. For example, if we use a two-dimensional image  $I$  as our input, we probably also want to use a two-dimensional kernel  $K$ :

$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(m, n)K(i-m, j-n). \quad (7.4)$$

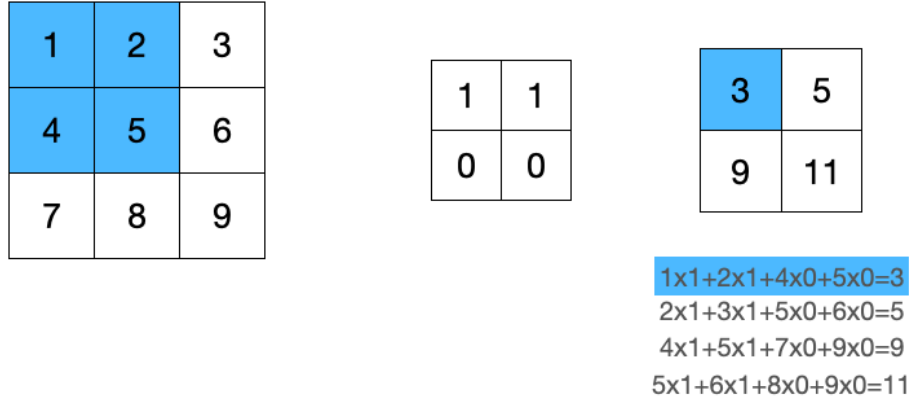
The following formula can be used to calculate the output dimension.

$$h_o = \frac{h_i - h_k}{h_s} + 1 \quad (7.5)$$

$$w_o = \frac{w_i - w_k}{w_s} + 1 \quad (7.6)$$

where  $h_o$  is the output height,  $h_i$  is the input height,  $h_k$  is the kernel height,  $h_s$  is the stride height,  $w_o$  is the output width,  $w_i$  is the input width,  $w_k$  is the kernel width,  $w_s$  is the stride width.

The convolution operation is shown in Figure 7.1.



**Figures 7.1:** Convolution operation

## 7.2 Properties

CNN leverages three important ideas:

- ♥ sparse interaction.
- ♥ parameter sharing.
- ♥ equivariant representations.

### 7.2.1 Sparse interaction

This is accomplished by making the kernel smaller than the input.

### 7.2.2 Parameter sharing

In convolutional layers, the same parameter defined in one kernel are used at every position of the input.

### 7.2.3 Equivariant representations

In the case of convolution, the particular form of a parameter sharing causes the layer to have a property called **equivariance** to translation. To say a function is equivariant means that if the input changes, the output changes in the same way.



### 7.3 Pooling

A pooling function replaces the output of the net at a certain location with a summary statistic of the nearby outputs. For example, the max pooling operation reports the maximum output within a rectangular neighborhood. Pooling helps to make the representation approximately invariant to small translations of the input. Invariant to translation means that if we translate the input by a small amount, the values of most of the pooled outputs do not change.

The following formula can be used to calculate the output dimension.

$$h_o = \frac{h_i - h_k}{h_s} + 1 \quad (7.7)$$

$$w_o = \frac{w_i - w_k}{w_s} + 1 \quad (7.8)$$

where  $h_o$  is the output height,  $h_i$  is the input height,  $h_k$  is the pooling height,  $h_s$  is the stride height,  $w_o$  is the output width,  $w_i$  is the input width,  $w_k$  is the pooling width,  $w_s$  is the stride width.



## Metric

		Real value	
		positive	negative
Predict value	positive	true positive	false positive
	negative	false negative	true negative

**Tables 8.1:** Confusion matrix

### 8.1 Precision

$$\text{precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (8.1)$$

Of all the predicted positive values, the ratio of the true positive values (the real value is positive and the predicted value is positive).

### 8.2 Recall

$$\text{recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (8.2)$$

Of all the real positive values, the ratio of the true positive values.

### 8.3 Accuracy

$$\text{accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} \quad (8.3)$$

Of all the values, the ratio of correctly predicted values. The disadvantage is that it is not suitable for unbalanced data.

## 8.4 F-score

$$F = \frac{(\alpha^2 + 1) \times \text{precision} \times \text{recall}}{\alpha^2 \times \text{precision} + \text{recall}} \quad (8.4)$$

When determining the value of the parameter  $\alpha$ , if we pay more attention to recall (compared to precision), we should choose a larger  $\alpha$ . The F-1 score is the expression when  $\alpha = 1$

# Regularization

We train model on training data but use test data (not used to train the model) to test out model. The ability to perform on test data is called **generalization**. We can use model on test data because we assume that the train data and the test has the same probability distribution (i.e. they have relationship).

**Regularization** is any modification we make to a learning algorithm that is intended to reduce its generalization error.

In practice, it is very difficult to find a model with the right number of parameters. Indeed, we often use a larger model that has been regularized appropriately.

## 9.1 Parameter norm penalties

Parameter norm penalties ( $\Omega(\theta)$ ) can be added to the object function  $J$  to limit the capacity of the model.

$$\tilde{J}(\theta; X, y) = J(\theta, X, y) + \alpha\Omega(\theta) \quad (9.1)$$

### 9.1.1 $L^2$ Parameter Regularization

The  $L^2$  parameter norm penalty commonly known as **weight decay**.

$$\Omega(\theta) = \frac{1}{2}\|w\|_2^2 \quad (9.2)$$

Where  $w$  is the model parameter matrix.

### 9.1.2 $L^1$ Regularization

$L^1$  regularization on the model parameter  $w$  is defined as

$$\Omega(\theta) = \|w\|_1 \quad (9.3)$$

## 9.2 Dataset Augmentation

The best way to make a machine learning model generalize better is to train it on more data. Of course, in practice, the amount of data we have is limited. One way to get around this problem is to create fake data and add it to the training set. Dataset augmentation has been a particular effective technique for a specific classification problem: object recognition.

## 9.3 Early stopping

Early stopping is used to avoid overfit.

The algorithm terminates when no parameters have improved over the best recorded validation error for some pre-specified number of iterations. This strategy is known as **early stopping**. It is probably the most commonly used form of regularization in deep learning.

## 9.4 Droupout

# Chapter 10

## Activation functions

Activation functions are usually added into neural network in order to help the network learn complex patterns. Because many of the patterns we want to learn are non-linear, we usually add non-linear activation functions into neural network to add the ability to learn non-linear patterns. Its function is to determine what information can be passed to the next neuron.

### 10.1 Desirable features

- ♥ Differentiable. This is needed for optimization.
- ♥ Low computational expense.
- ♥ Zero-centered.

### 10.2 Sigmoid

$$\text{sigmoid}(x) = \frac{1}{1 + e^{-x}} \quad (10.1)$$

It is no longer used in recent deep learning models because it is computationally expensive, causes vanishing gradient problem and not zero-centered.

### 10.3 Tanh

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{e^{2x} - 1}{e^{2x} + 1} \quad (10.2)$$

Comparing to sigmoid, it solve the not zero-centered problem.

### 10.4 ReLU

$$\text{relu}(x) = \max(0, x) \quad (10.3)$$

This is a widely used activation function.

### 10.5 Leaky ReLU

$$\text{lrelu}(x) = \max(\alpha x, x) \tag{10.4}$$

### 10.6 Swish

$$\text{swish}(x) = x * \text{sigmoid}(x) = x * (1 + e^{-x})^{-1} \tag{10.5}$$



# Chapter 1

## Machine learning process

- 1 confirm your target
- 2 observe the data
- 3 preprocess the data if necessary
- 4 select the model
- 5 select the cost function
- 6 select the optimizer
- 7 train the model
- 8 fine-tune the model
- 9 predict the data using the model
- 10 deploy the model



## **Part II**

# **Models**



# Chapter 12

## Diffusion

Diffusion models are generative models, meaning that they are used to generate data similar to the data on which they are trained. Fundamentally, diffusion models work by destroying training data through the successive addition of Gaussian noise, and then learning to recover the data by reversing this noising process. After training, we can use the diffusion model to generate data by simply passing randomly sampled noise through the learned denoising process.



## **Part III**

### **Tools**





# Chapter 13

## PyTorch



# Chapter 14

## NumPy

### 14.1 Load data

```
1 import numpy as np
2
3 data = np.load('foo.npz')
```



## **Part IV**

# **Practice**



# Chapter 15

## Preprocessing

Preprocessing is not necessary. We preprocess the data because we can not get what we want from the data directly.

### 15.1 Not valued based data

The computer can only process number-based data. For not number-based data, we need to convert them into number-based data. For example, convert a string into a number list. For null values, we can drop them or compute them with other not null values.

```
1 from sklearn.preprocessing import LabelBinarizer
2
3 data = ['dog', 'cat', 'dog', 'horse']
4 print('data: ', data)
5 binarizer = LabelBinarizer()
6 binarizer.fit(data)
7 print(binarizer.transform(data))
8 print(binarizer.transform(['dog']))
9
10 """
11 data:  ['dog', 'cat', 'dog', 'horse']
12 [[0 1 0]
13  [1 0 0]
14  [0 1 0]
15  [0 0 1]]
16 [[0 1 0]]
17 """
```





## **Part V**

# **Projects**



## **Part VI**

## **Papers**



# Chapter 16

## Deep Unsupervised Learning using Nonequilibrium Thermodynamics[1]

The algorithm consists of two trajectories: forward (inference) diffusion process and inverse (generative) diffusion process. The forward diffusion process converts complex data distribution into a simple, tractable distribution and the inverse diffusion process learn a finite-time reversal of the forward diffusion process which define a generative model distribution.

### 16.1 Algorithm

#### 16.1.1 Forward trajectory

$$\pi(\mathbf{y}) = \int d\mathbf{y}' T_{\pi}(\mathbf{y} | \mathbf{y}'; \beta) \pi(\mathbf{y}') \quad (16.1)$$

Where  $\pi(\mathbf{y})$  is a well behaved (analytically tractable) distribution.  $T_{\pi}(\mathbf{y} | \mathbf{y}'; \beta)$  is a Markov diffusion kernel.  $\beta$  is the diffusion rate.

$$q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)}) = T_{\pi}(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)}; \beta_t) \quad (16.2)$$

The forward trajectory, corresponding to starting at the data distribution and performing  $T$  steps of diffusion is

$$q(\mathbf{x}^{(0 \dots T)}) = q(\mathbf{x}^{(0)}) \prod_{t=1}^T q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)}) \quad (16.3)$$

Where  $q(\mathbf{x}^{(0)})$  is the initial data distribution.

#### 16.1.2 Reverse Trajectory

The generative distribution will be trained to describe the same trajectory, but in reverse,

$$p(\mathbf{x}^{(T)}) = \pi(\mathbf{x}^{(T)}) \quad (16.4)$$

$$p(\mathbf{x}^{(0:T)}) = p(\mathbf{x}^{(T)}) \prod_{t=1}^T p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}) \quad (16.5)$$

For both Gaussian and binomial diffusion, for continuous diffusion (limit of small step size  $\beta$ ) the reversal of the diffusion process has the identical functional form as the forward process. Since  $q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})$  is a Gaussian (binomial) distribution, and if  $\beta_t$  is small, then  $q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})$  will also be a Gaussian (binomial) distribution. The longer the trajectory the smaller the diffusion rate  $\beta$  can be made.

During learning only the mean and covariance for a Gaussian diffusion kernel, or the bit flip probability for a binomial kernel, need be estimated.

## 16.2 Model probability

The probability the generative model assigns to the data is

$$p(\mathbf{x}^{(0)}) = \int d\mathbf{x}^{(1:T)} p(\mathbf{x}^{(0:T)}) \quad (16.6)$$

Naively this integral is intractable – but taking a cue from annealed importance sampling and the Jarzynski equality, we instead evaluate the relative probability of the forward and reverse trajectories, averaged over forward trajectories

$$\begin{aligned} p(\mathbf{x}^{(0)}) &= \int d\mathbf{x}^{(1:T)} p(\mathbf{x}^{(0:T)}) \frac{q(\mathbf{x}^{(1:T)} | \mathbf{x}^{(0)})}{q(\mathbf{x}^{(1:T)} | \mathbf{x}^{(0)})} \\ &= \int d\mathbf{x}^{(1:T)} q(\mathbf{x}^{(1:T)} | \mathbf{x}^{(0)}) \frac{p(\mathbf{x}^{(0:T)})}{q(\mathbf{x}^{(1:T)} | \mathbf{x}^{(0)})} \\ &= \int d\mathbf{x}^{(1:T)} q(\mathbf{x}^{(1:T)} | \mathbf{x}^{(0)}) \cdot p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})} \end{aligned} \quad (16.7)$$

## 16.3 Training

Training amounts to maximizing the model log likelihood,

$$\begin{aligned} L &= \int d\mathbf{x}^{(0)} q(\mathbf{x}^{(0)}) \log p(\mathbf{x}^{(0)}) \\ &= \int d\mathbf{x}^{(0)} q(\mathbf{x}^{(0)}) \cdot \log \left[ \frac{\int d\mathbf{x}^{(1:T)} q(\mathbf{x}^{(1:T)} | \mathbf{x}^{(0)}) \cdot p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})}}{p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})}} \right], \end{aligned} \quad (16.8)$$

which has a lower bound provided by Jensen's inequality,

$$L \geq \int d\mathbf{x}^{(0:T)} q(\mathbf{x}^{(0:T)}) \cdot \log \left[ p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})} \right]. \quad (16.9)$$

for our diffusion trajectories this reduces to,

$$L \geq K \quad (16.10)$$

$$\begin{aligned}
K = & - \sum_{t=2}^T \int d\mathbf{x}^{(0)} d\mathbf{x}^{(t)} q(\mathbf{x}^{(0)}, \mathbf{x}^{(t)}) \\
& D_{KL}(q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}) \| p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})) \\
& + H_q(\mathbf{X}^{(T)} | \mathbf{X}^{(0)}) - H_q(\mathbf{X}^{(1)} | \mathbf{X}^{(0)}) - H_p(\mathbf{X}^{(T)})
\end{aligned} \tag{16.11}$$

where the entropies and KL divergences can be analytically computed. The derivation of this bound parallels the derivation of the log likelihood bound in variational Bayesian methods.

If the forward and reverse trajectories are identical, corresponding to a quasi-static process, then the inequality in Equation 16.10 becomes an equality.

Training consists of finding the reverse Markov transitions which maximize this lower bound on the log likelihood,

$$\hat{p}(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}) = \operatorname{argmax}_{p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})} K. \tag{16.12}$$





# Bibliography

- [1] Jascha Sohl-Dickstein. Deep unsupervised learning using nonequilibrium thermodynamics. *arXiv preprint arxiv:arXiv:1503.03585*, 2015.

