

Mathematics, Machine Learning and Deep Learning Notes

Jinming Su

Last update: April 24, 2019

Contents

1	Mathematical Foundation	1
1.1	Probability theory and mathematical statistics	1
1.1.1	How to get expected value and variance?	1
1.1.2	Discrete probability distribution	1
1.1.3	Continuous probability distribution	2
1.1.4	Sample mean and sample variance	3
1.2	Prior and posterior	4
2	Machine Learning	5
2.1	Why can't perceptron solve XOR problem?	5
2.1.1	Definition of perceptron	5
2.1.2	Learning Algorithm	5
2.1.3	Dual form of perceptron learning algorithm	6
2.2	How to get the update rule of parameters of backpropagation in gradient descent algorithm?	7
2.3	Generative model and discriminative model	7
2.4	Support vector machine	7
3	Deep Network	9
3.1	How does backpropagation work?	9
3.2	Why does batch normalization work?	11
3.2.1	Dataset shift and covariate shift	11
3.2.2	Internal covariate shift	11
3.2.3	Batch normalization	12
3.3	Why does residual learning work?	14
4	Contents specifically referenced	15

Todo list

□	Note in Eq. 1.5, maybe $P(X = 1)$ is equivalent $P(X^2 = 1^2)$, which ensures the establishment of this equation. . .	2
□	There exists another solution directly through derivation of the probability mass function of Binomial distribution. .	2
□	Note there exists Taylor Expansion $e^x = 1 + \frac{x}{1!} + \frac{x^2}{2!} + \frac{x^3}{3!} + \cdots + \cdots = \sum_{i=0}^{\infty} \frac{x^i}{i!}$	2
□	Todo: Taylor's formula.	2
□	The derivation of the expectation and invariance of normal distribution requires multiple integration operations. . .	3
□	The derivation of the expectation and invariance of exponential distribution requires multiple integration operations.	3
□	Todo: Beta distribution.	3
□	The convergence of this algorithm about perceptron learning can be proved, but it is not within the scope of this note at present.	6
□	Todo: SVM	7
□	Todo: Random forest	7
□	A more precise definition of covariate can't be found.	11
□	PCA whitening and ZCA Whitening.	11
□	Why is multi-layer linear network equivalent to the representation ability of a single layer linear network.	12
□	There exists \sum_i for y_i because gradient from each y_i will backpropagate for γ and the effect of summation is achieved.	13

Chapter 1

Mathematical Foundation

1.1 Probability theory and mathematical statistics

Probability theory mainly focuses on the probability of occurrence of a single event, while **mathematical statistics** is more inclined to statistics. It focuses on the sampling probability of a group and the possible interval of occurrence of this probability.

In the following introduction, these two concepts are introduced without distinction.

1.1.1 How to get expected value and variance?

X is a random variable whose values are X_1, X_2, \dots, X_n . $P(X_1), P(X_2), \dots, P(X_n)$ are the probability corresponding to these values. The expected value of X can be denoted as $E(X)$, and the variance is denoted as $Var(X)$. Then,

$$\begin{aligned} E(X) &= \sum_{i=1}^n X_i P(X_i) \text{ for discrete variable} \\ &= \int_X x f(x) dx \text{ for continuous variable} \end{aligned} \quad (1.1)$$

and

$$\begin{aligned} Var(X) &= \sum_{i=1}^n (X_i - E(X))^2 P(X_i) \text{ for discrete variable} \\ &= \int_X (x - E(X))^2 f(x) dx \text{ for continuous variable} \end{aligned} \quad (1.2)$$

For discrete variable,

$$\begin{aligned} Var(X) &= \sum_{i=1}^n (X_i - E(X))^2 P(X_i) \\ &= E[(X - E(X))^2] \\ &= E[X^2 + E(X)^2 - 2XE(X)] \\ &= E(X^2) + E(X)^2 - 2E(X)E(X) \\ &= E(X^2) - E(X)^2 \end{aligned} \quad (1.3)$$

1.1.2 Discrete probability distribution

Bernoulli distribution is the discrete probability distribution of a random variable which takes the value 1 with probability p and the value 0 with probability $q = 1 - p$. We denote Bernoulli distribution as $B(1, p)$. Mathematically, if X is a random variable with $B(1, p)$, then $P(X = 1) = p, P(X = 0) = q = 1 - p$. The probability mass function f of this distribution over possible outcomes k , is

$$f(k; p) = \begin{cases} p & \text{if } k = 1, \\ q = 1 - p & \text{if } k = 0. \end{cases} \quad (1.4)$$

The expected value and invariance of a Bernoulli variable X are

$$\begin{cases} E(X) = P(X = 1) \cdot 1 + P(X = 0) \cdot 0 = p, \\ E(X^2) = P(X = 1) \cdot 1^2 + P(X = 0) \cdot 0^2 = P(X = 1) = p, \\ Var(X) = E(X^2) - E(X)^2 = p - p^2 = p(1 - p) = pq. \end{cases} \quad (1.5)$$

Note in Eq. 1.5, maybe $P(X = 1)$ is equivalent $P(X^2 = 1^2)$, which ensures the establishment of this equation.

Binomial distribution with parameters n and p is the discrete probability distribution of the number of successes in a sequence of n independent experiments. Each experiment is a Bernoulli trial. 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 trials is given by the probability mass function:

$$f(k; n, p) = P(k; n, p) = P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}. \quad (1.6)$$

The expected value and invariance of a Binomial variable X are

$$\begin{cases} E(X) = E(X_1 + X_2 + \cdots + X_n) \\ \quad = E(X_1) + E(X_2) + \cdots + E(X_n) \\ \quad = p + p + \cdots + p \\ \quad = np, \\ Var(X) = Var(X_1) + Var(X_1) + \cdots + Var(X_1) \\ \quad = nVar(X_1) \\ \quad = np(1 - p). \end{cases} \quad (1.7)$$

There exists another solution directly through derivation of the probability mass function of Binomial distribution.

Poisson distribution is a discrete probability distribution that expresses the probability of a given number k of events occurring in a fixed interval of time or space if these events occur with a known constant rate λ and independently of the time since the last events. If X is a Poisson variable with the average number of events λ , we write $X \sim Poisson(\lambda)$. The probability mass function is

$$f(k; n, \lambda) = P(X = k) = \frac{e^{-\lambda} \lambda^k}{k!}. \quad (1.8)$$

The expected value and invariance of a Poisson variable X are

$$\begin{cases} E(X) = \sum_{i=0}^{\infty} i P(X = i) \\ \quad = \sum_{i=1}^{\infty} i \frac{e^{-\lambda} \lambda^i}{i!} = \lambda e^{-\lambda} \sum_{i=1}^{\infty} \frac{\lambda^{i-1}}{(i-1)!} = \lambda e^{-\lambda} \sum_{i=0}^{\infty} \frac{\lambda^i}{i!} \\ \quad = \lambda e^{-\lambda} e^{\lambda} \\ \quad = \lambda \\ E(X^2) = \lambda + \lambda^2 \\ Var(X) = \lambda + \lambda^2 - \lambda^2 \\ \quad = \lambda \end{cases} \quad (1.9)$$

Note there exists Taylor Expansion $e^x = 1 + \frac{x}{1!} + \frac{x^2}{2!} + \frac{x^3}{3!} + \cdots + \cdots = \sum_{i=0}^{\infty} \frac{x^i}{i!}$.

Todo: Taylor's formula.

1.1.3 Continuous probability distribution

For **continuous uniform distribution**, all intervals of the same length on the distribution's support are equally probable. The support is defined by the two parameters, a and b , which are its minimum and maximum values. The distribution is often

abbreviated $U(a, b)$. The probability density function of the continuous uniform distribution is

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{for } a \leq x \leq b, \\ 0 & \text{for } x < a \text{ or } x > b. \end{cases} \quad (1.10)$$

The expected value and invariance of a Uniform variable X are

$$\left\{ \begin{array}{l} E(X) = \int_a^b x \frac{1}{b-a} \\ \quad = \frac{x^2}{2(b-a)} \Big|_a^b = \frac{b^2 - a^2}{2(b-a)} = \frac{(b+a)(b-a)}{2(b-a)} \\ \quad = \frac{b+a}{2} \\ E(X^2) = \int_a^b x^2 \frac{1}{b-a} \\ \quad = \frac{x^3}{3(b-a)} \Big|_a^b = \frac{b^3 - a^3}{3(b-a)} = \frac{(b-a)(b^2 + a^2 + ab)}{3(b-a)} \\ \quad = \frac{a^2 + b^2 + ab}{3} \\ Var(X) = E(X^2) - E(X)^2 = \frac{a^2 + b^2 + ab}{3} - \left(\frac{b+a}{2}\right)^2 \\ \quad = \frac{4a^2 + 4b^2 + 4ab}{12} - \frac{3a^2 + 3b^2 + 6ab}{12} \\ \quad = \frac{(a-b)^2}{12} \end{array} \right. \quad (1.11)$$

The probability density function of the **normal distribution** is

$$f(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \quad (1.12)$$

where μ is the expectation and σ is the standard deviation. If a random variable X is distributed normally with mean μ and variance σ^2 , one may write $X \sim N(\mu, \sigma^2)$.

The derivation of the expectation and invariance of normal distribution requires multiple integration operations.

Exponential distribution is the probability distribution that describes the time between events in a Poisson point process, *i.e.* a process in which event occur continuously and independently at a constant average rate. The distribution is often abbreviated $Exponential(\lambda)$. The probability density function of an exponential distribution is

$$f(x; \lambda) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0, \\ 0 & x < 0. \end{cases} \quad (1.13)$$

The expected value and invariance of a Exponential variable X are

$$\left\{ \begin{array}{l} E(X) = \frac{1}{\lambda}, \\ Var(X) = \frac{1}{\lambda^2}. \end{array} \right. \quad (1.14)$$

The derivation of the expectation and invariance of exponential distribution requires multiple integration operations.

Todo: Beta distribution.

1.1.4 Sample mean and sample variance

Sample mean is defined as

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i. \quad (1.15)$$

Sample Variance is defined as

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2. \quad (1.16)$$

The above definitions follow the theorem: sample mean is unbiased estimation of population mean and sample variance is unbiased estimation of population variance.

proof:

$$\begin{aligned}
 E[\bar{X}] &= E\left[\frac{1}{n} \sum_{i=1}^n X_i\right] \\
 &= \frac{1}{n} \sum_{i=1}^n E[X_i] \\
 &= \frac{1}{n} nE(X) \\
 &= E(X) \\
 E[S^2] &= E\left[\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2\right] \\
 &= \frac{1}{n-1} E\left[\sum_{i=1}^n (X_i^2 + \bar{X}^2 - 2X_i\bar{X})\right] \\
 &= \frac{1}{n-1} (E[\sum_{i=1}^n X_i^2] + n\bar{X}^2 - 2n\bar{X}^2) \\
 &= \frac{1}{n-1} (nE[X^2] - nE[\bar{X}^2]) \\
 &= \frac{n}{n-1} (Var[X] + E[X]^2 - Var[\bar{X}] - E[\bar{X}]^2) \\
 &= \frac{n}{n-1} (Var[X] - Var[\bar{X}]) \\
 &= \frac{n}{n-1} (Var[X] - \frac{1}{n} Var[X]) \\
 &= Var[X]
 \end{aligned} \tag{1.17}$$

1.2 Prior and posterior

The prior probability is the probability of a cause inferred from experience, denoted as $P(\theta)$. The posterior probability is the probability of the cause estimated from the result, denoted as $P(\theta|x)$. The posterior probability is defined as

$$P(\theta|x) = \frac{P(x|\theta)P(\theta)}{P(x)} \tag{1.18}$$

where $P(x|\theta)$ represents likelihood of x . In fact, likelihood is the function of parameters. Likelihood is equal to the probability of the result occurring caused by a cause (parameters) based on the cause. Likelihood is the viewpoint of the frequency. In the frequency, the parameter is a true value, not a random variable, so the parameter has no distribution and no probability.

Chapter 2

Machine Learning

2.1 Why can't perceptron solve XOR problem?

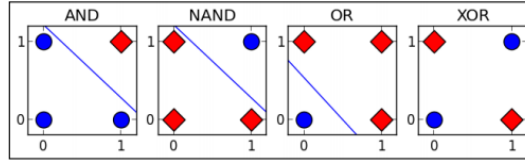


Figure 2.1: XOR problem.

Linear classification models can't classify linear non-separable problems. Perceptron is a linear classification model, and XOR problem is a linear non-separable problem, so perceptron can't solve the XOR problem.

2.1.1 Definition of perceptron

Suppose the input space is $\mathcal{X} \subseteq \mathbb{R}^n$, and the output space is $\mathcal{Y} = \{+1, -1\}$. For an example $x \in \mathcal{X}$ where x is an n -dimensional vector $(x_1, x_2, \dots, x_n)^T$, $y \in \mathcal{Y}$ represents the category of x . Then, we get the perceptron model from the input space to output space:

$$f(\vec{x}) = \text{sign}(\vec{w}^T \vec{x} + b), \quad (2.1)$$

where \vec{w} is weight and b is bias. And sign is a sign function, *i.e.*

$$\text{sign}(x) = \begin{cases} +1, & x \geq 0 \\ -1, & x < 0 \end{cases} \quad (2.2)$$

There exists a geometric interpretation that a hyperplane whose normal vector is \vec{w} and intercept is b .

2.1.2 Learning Algorithm

Given a linear separable dataset $T = (x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$, where $x_i \in \mathcal{X} \subseteq \mathbb{R}^n$, and $y_i \in \mathcal{Y} = \{+1, -1\}$, $i = 1, 2, \dots, N$, we can construct a perceptron model to classify this dataset.

(1) construct the loss function. we define the loss function as the total distance from misclassified points to hyperplane. Toward this end, the distance from any point x_0 to hyperplane:

$$\frac{1}{\|\vec{w}\|} |\vec{w}^T \cdot \vec{x}_0 + b| \quad (2.3)$$

where $\|\vec{w}\|$ means the L_2 norm of \vec{w} .

As for misclassified points, $\vec{w}^T \cdot \vec{x}_0 + b > 0$ when $y_i = -1$, and $\vec{w}^T \cdot \vec{x}_0 + b < 0$ when $y_i = +1$. So the total distance from all the misclassified points to hyperplane is

$$- \frac{1}{\|\vec{w}\|} \sum_{\vec{x}_i \in M} y_i (\vec{w}^T \cdot \vec{x}_i + b), \quad (2.4)$$

where M is the set of misclassified points.

In Eq. 2.6, $\frac{1}{\|\vec{w}\|}$ can be ignored. The reason is (a) $\frac{1}{\|\vec{w}\|}$ doesn't affect positive or negative judgment of $y_i(\vec{w}^T \cdot \vec{x}_i + b)$ and (b) $\frac{1}{\|\vec{w}\|}$ doesn't affect the final optimization result of Eq. 2.6. The final optimization result is there are no points with wrong classification, which leads to the loss of 0. Toward this end, the loss function of preceptron is

$$L(\vec{w}, b) = - \sum_{\vec{x}_i \in M} y_i(\vec{w}^T \cdot \vec{x}_i + b), \quad (2.5)$$

and the optimization object is

$$\min_{\vec{w}, b} L(\vec{w}, b) = \min_{\vec{w}, b} - \sum_{\vec{x}_i \in M} y_i(\vec{w}^T \cdot \vec{x}_i + b). \quad (2.6)$$

(2) We adopt stochastic gradient descent algorithm to train the perceptron model. Suppose the set of misclassified points is fixed, thus the gradient of loss function $L(\vec{w}, b)$:

$$\begin{aligned} \nabla_{\vec{w}} L(\vec{w}, b) &= - \sum_{\vec{x}_i \in M} y_i \vec{x}_i \\ \nabla_b L(\vec{w}, b) &= - \sum_{\vec{x}_i \in M} y_i \end{aligned} \quad (2.7)$$

So, the strategies of \vec{w} and b based on one misclassified point are:

$$\begin{aligned} \vec{w} &\leftarrow \vec{w} + \eta y_i \vec{x}_i \\ b &\leftarrow b + \eta y_i \end{aligned} \quad (2.8)$$

where η is the learning rate.

(3) The standard form of perceptron learning algorithm:

input: linear separable training set: $T = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, where $\vec{x}_i \in \mathcal{X} \subseteq \mathbb{R}^n$ and $y_i \in \mathcal{Y} = \{-1, +1\}$, $i = 1, 2, \dots, N$. In addition, the learning rate is denoted as η ($0 < \eta \leq 1$).

output: \vec{w}, b . And the perceptron model $f(\vec{x}) = \text{sign}(\vec{w}^T \vec{x} + b)$.

(a) choose initial values, \vec{w}_0, b_0 , and usually $\vec{w}_0 = 0, b_0 = 0$;

(b) choose one example (\vec{x}_i, y_i) from the training set;

(c) if $y_i(\vec{w}^T \cdot \vec{x}_i + b) \leq 0$:

$$\begin{aligned} \vec{w} &\leftarrow \vec{w} + \eta y_i \vec{x}_i \\ b &\leftarrow b + \eta y_i \end{aligned} \quad (2.9)$$

(d) go to (b) until there are no misclassified points in training set.

(4)

The convergence of this algorithm about perceptron learning can be proved, but it is not within the scope of this note at present.

2.1.3 Dual form of perceptron learning algorithm

From Eq. 2.8, \vec{w} and b are updated for many times. Suppose the number of updates about each example is n_i ($i = 1, 2, \dots, N$), and we define $\alpha_i = n_i \eta$. Then, the final learned \vec{w} and b can be expressed as

$$\begin{aligned} \vec{w} &= \sum_{i=1}^N \alpha_i y_i \vec{x}_i \\ b &= \sum_{i=1}^N \alpha_i y_i \end{aligned} \quad (2.10)$$

Then, we can obtain the dual form of perceptron learning algorithm.

input: linear separable training set: $T = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, where $\vec{x}_i \in \mathcal{X} \subseteq \mathbb{R}^n$ and $y_i \in \mathcal{Y} = \{-1, +1\}$, $i = 1, 2, \dots, N$. In addition, the learning rate is denoted as η ($0 < \eta \leq 1$).

output: $\vec{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_N)^T, b$. And the perceptron model $f(\vec{x}) = \text{sign}(\sum_{j=1}^N \alpha_j y_j \vec{x}_j^T \vec{x} + b)$.

- (a) choose initial values, $\vec{\alpha}, b_0$, and usually $\vec{\alpha} = 0, b_0 = 0$;
- (b) choose one example (\vec{x}_i, y_i) from the training set;
- (c) if $y_i(\sum_{j=1}^N \alpha_j y_j \vec{x}_j^T \vec{x} + b) \leq 0$:

$$\begin{aligned}\vec{\alpha}_i &\leftarrow \alpha_i + \eta \\ b &\leftarrow b + \eta y_i\end{aligned}\tag{2.11}$$

- (d) go to (b) until there are no misclassified points in training set.

2.2 How to get the update rule of parameters of backpropagation in gradient descent algorithm?

An intuitive and imprecise proof is as follows. Assuming that there is a differentiable loss function \mathcal{L} on training set with respect to parameter $\vec{w} = (w_1, w_2, \dots, w_n)$ of the training model. Then, the total increment of \mathcal{L} is

$$\begin{aligned}\Delta \mathcal{L} &= \frac{\partial \mathcal{L}}{\partial w_1} \Delta w_1 + \frac{\partial \mathcal{L}}{\partial w_2} \Delta w_2 + \dots + \frac{\partial \mathcal{L}}{\partial w_n} \Delta w_n + o(\rho) \\ &\approx \frac{\partial \mathcal{L}}{\partial w_1} \Delta w_1 + \frac{\partial \mathcal{L}}{\partial w_2} \Delta w_2 + \dots + \frac{\partial \mathcal{L}}{\partial w_n} \Delta w_n\end{aligned}\tag{2.12}$$

Define $\nabla \mathcal{L}_{\vec{w}} = (\frac{\partial \mathcal{L}}{\partial w_1}, \frac{\partial \mathcal{L}}{\partial w_2}, \dots, \frac{\partial \mathcal{L}}{\partial w_n})^T$, which represents the vector of gradients. Therefore,

$$\Delta \mathcal{L} \approx \nabla \mathcal{L}_{\vec{w}}^T \cdot \Delta \vec{w}.\tag{2.13}$$

Our goal is to reduce the loss function \mathcal{L} , which is equivalent to making the total increment $\Delta \mathcal{L}$ negative. Toward this end, construct

$$\Delta \vec{w} = -\eta \nabla \mathcal{L}_{\vec{w}},\tag{2.14}$$

which ensure that $\Delta \mathcal{L} \approx \nabla \mathcal{L}_{\vec{w}}^T \cdot \Delta \vec{w} = -\eta \|\nabla \mathcal{L}_{\vec{w}}\|_2^2 \leq 0$. So, the update rule of parameter \vec{w} is

$$\vec{w} \leftarrow \vec{w} - \eta \nabla \mathcal{L}_{\vec{w}}.\tag{2.15}$$

In other word, in gradient descent algorithm, the direction of parameter update is opposite to its gradient direction.

2.3 Generative model and discriminative model

- (1) Generative model: perceptron

2.4 Support vector machine

Todo: SVM

Todo: Random forest

Chapter 3

Deep Network

3.1 How does backpropagation work?

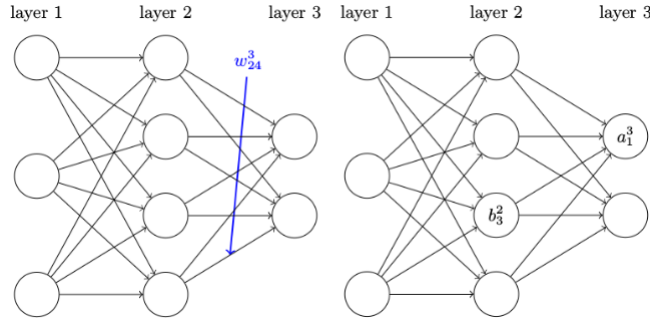


Figure 3.1: Neural network.

First, symbol definition is carried out. We use w^l_{jk} to denote the weight for the connection from the k th neuron in the $(l-1)$ th layer to the j th neuron in the l th layer. Moreover, we use b^l_j for the bias of j th neuron in l th layer and we use a^l_j for the activation of the j th neuron in l th layer. We also use z^l_j as the weighted input of j th neuron in l th layer. Then

$$\begin{aligned} z^l_j &= \sum_k w^l_{jk} a^{l-1}_k + b^l_j, \\ a^l_j &= \sigma\left(\sum_k w^l_{jk} a^{l-1}_k + b^l_j\right) \\ &= \sigma(z^l_j). \end{aligned} \tag{3.1}$$

where σ is the activation function, such as Sigmoid function. Then, Eq. 3.1 can be expressed in a vector form. In detail, we use \vec{a}^l for the activation vector in l th layer. \vec{z}^l is the same. Moreover, we use \vec{w}^l for the weight matrix and \vec{b}^l for the bias vector in l th layer. Then,

$$\begin{aligned} \vec{z}^l &= \vec{w}^l \vec{a}^{l-1} + \vec{b}^l, \\ \vec{a}^l &= \sigma(\vec{w}^l \vec{a}^{l-1} + \vec{b}^l) \\ &= \sigma(\vec{z}^l). \end{aligned} \tag{3.2}$$

The understanding of backpropagation is about how to change the weights and biases in a network to change the cost function. Ultimately, this means computing the partial derivatives $\frac{\partial \mathcal{L}}{\partial w^l_{jk}}$ and $\frac{\partial \mathcal{L}}{\partial b^l_j}$. To compute those, we first introduce an intermediate quantity δ^l_j , which we call the error of the j th neuron in the l th layer. In the actual calculation process, backpropagation will give us a prelude to compute δ^l_j , and then will relate δ^l_j to $\frac{\partial \mathcal{L}}{\partial w^l_{jk}}$ and $\frac{\partial \mathcal{L}}{\partial b^l_j}$.

We define the error of the j th neuroon in the l th layer by

$$\delta_j^l = \frac{\partial \mathcal{L}}{\partial z_j^l}. \quad (3.3)$$

Then, we can get the error of the output layer δ^L , and the components of δ^L are given by

$$\begin{aligned} \delta_j^L &= \frac{\partial \mathcal{L}}{\partial z_j^L} = \frac{\partial \mathcal{L}}{\partial a_j^L} \frac{\partial a_j^L}{\partial z_j^L} \\ &= \frac{\partial \mathcal{L}}{\partial a_j^L} \sigma'(z_j^L). \end{aligned} \quad (BP1)$$

Moreover, we can give the matrix-based form of Eq. BP1 by

$$\vec{\delta^L} = \nabla_{\vec{a^L}} \otimes \sigma'(\vec{z^L}). \quad (3.4)$$

Then, we can obtain the error δ^l in terms of the error in the next layer δ^{l+1} . In particular,

$$\begin{aligned} \delta_j^l &= \frac{\partial \mathcal{L}}{\partial z_j^l} = \sum_k \frac{\partial \mathcal{L}}{\partial z_k^{l+1}} \frac{\partial z_k^{l+1}}{\partial z_j^l} = \sum_k \delta_k^{l+1} \frac{\partial z_k^{l+1}}{\partial z_j^l} \\ &= \sum_k \delta_k^{l+1} \frac{\partial \sum_j w_{kj}^{l+1} \sigma(z_j^l) + b_k^l}{\partial z_j^l} \\ &= \sum_k \delta_k^{l+1} w_{kj}^{l+1} \sigma'(z_j^l) \end{aligned} \quad (BP2)$$

Moreover, we can give the matrix-based form of Eq. BP2 by

$$\vec{\delta^l} = ((\vec{w^{l+1}})^T \delta^{l+1}) \otimes \sigma'(\vec{z^l}). \quad (3.5)$$

Next, we can obtain the rate of change of the cost with respect to any weight in the network. In particular,

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial w_{jk}^l} &= \frac{\partial \mathcal{L}}{\partial z_j^l} \frac{\partial z_j^l}{\partial w_{jk}^l} = \frac{\partial \mathcal{L}}{\partial z_j^l} \frac{\partial \sum_k w_{jk}^l a_k^{l-1} + b_j^l}{\partial w_{jk}^l} \\ &= \frac{\partial \mathcal{L}}{\partial z_j^l} a_k^{l-1} \\ &= \delta_j^l a_k^{l-1}. \end{aligned} \quad (BP3)$$

Moreover, we can give the matrix-based form of Eq. BP3 by

$$\frac{\partial \mathcal{L}}{\partial \vec{w^l}} = \vec{\delta^l} (\vec{a^{l-1}})^T. \quad (3.6)$$

We can obtain the rate of change of the cost with respect to any bias in the network. In particular,

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial b_j^l} &= \frac{\partial \mathcal{L}}{\partial z_j^l} \frac{\partial z_j^l}{\partial b_j^l} = \frac{\partial \mathcal{L}}{\partial z_j^l} \frac{\partial \sum_k w_{jk}^l a_k^{l-1} + b_j^l}{\partial b_j^l} \\ &= \frac{\partial \mathcal{L}}{\partial z_j^l} \\ &= \delta_j^l. \end{aligned} \quad (BP4)$$

Moreover, we can give the matrix-based form of Eq. BP4 by

$$\frac{\partial \mathcal{L}}{\partial \vec{b^l}} = \vec{\delta^l}. \quad (3.7)$$

3.2 Why does batch normalization work?

3.2.1 Dataset shift and covariate shift

There exists a inherent hypothesis in supervised learning: The training data in source domain and testing data in target domain are independent and identically drawn from the same distribution.

Usually, the distributions of training data and testing data are not exactly the same, especiallt when the training data is insufficient. This phenomenon is named as “dataset shift”.

In analysis of covariance (ANOVA), convate is any continuous variable that is usually not controlled during data collection. (referring to <https://support.minitab.com/en-us/minitab/18/help-and-how-to/modeling-statistics/anova/supporting-topics/anova-models/understanding-covariates/>).

For the input variable x (also named as explanatory variable or covariate in scene of batch normalizatio), y is the corresponding label, which is regarded as response variable. If the conditional probabilities of training data and testing data are the same but the marinal probabilities are different, *i.e.*

$$\begin{aligned} P_{tr}(y|x) &= P_{te}(y|x) \\ P_{tr}(x) &\neq p_{te}(y) \end{aligned} \quad (3.8)$$

we call this phenomenon as “covariate shift”. Covariate shift is a special kind of dataset shift.

Moreover, we can deduce the Eq. 3.8 as follows:

$$\begin{aligned} P_{tr}(x, y) &= P_{tr}(y|x)P_{tr}(x) \\ P_{te}(x, y) &= P_{te}(y|x)P_{te}(x) \end{aligned} \quad (3.9)$$

and we know the joint probabilities of traning data and testing data are different, *i.e.*,

$$P_{tr}(x, y) \neq P_{te}(x, y). \quad (3.10)$$

. From Eq. 3.10, this scene violates the independence hypothesis of supervised learning, so the prediction is inaccurate.

A more precise definition of **covariate** can’t be found.

3.2.2 Internal covariate shift

Internal covariate shift is defined as the change in the distribution of network activations due to the change in network parameters during training.

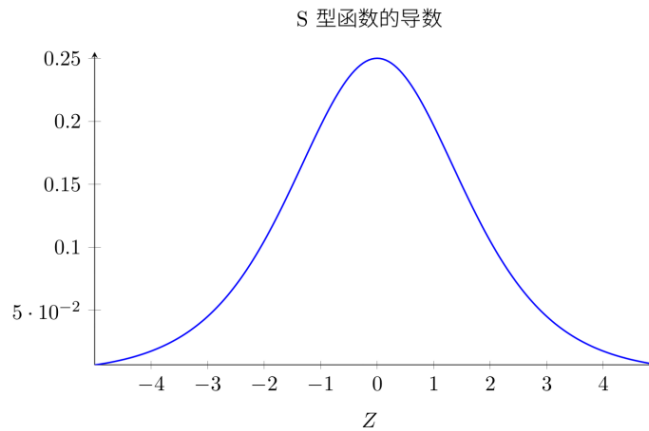


Figure 3.2: Derivative of sigmoid function.

Consequences of internal covariate shift: (1) Intuitively, the layers need to continuously adapt to the new distribution, which lead to a decrease in training speed. (2) Objectively, the changes of W and b during training will likely move many dimensions of x into the saturated regime of the nonlinearity (*e.g.* sigmoid and tanh) and slow down the convergence. The saturation problem is usually addressed by using ReLU, careful initialization and small learning rates. Moreover, if the distribution of nonlinearity inputs remains more stable at the sensitive regions in nonlinearity as the network trains, the optimizer would be less likely to get stuck in the saturated regime, and the training would accelerate.

PCA whitening and ZCA Whitening.

3.2.3 Batch normalization

Problem: because whitening is expensive and maybe change the representation ability of a network by discarding the absolute scale of activations, batch normalization emerges.

There are two treatments for batchnorm to address above two problems. (1) The first is instead of whitening the features in layer inputs and outputs jointly, batchnorm normalizes each scalar feature independently, by making it have the mean of 0 and the variance of 1. For a layer with d -dimensional input $x = [x^{(1)}, x^{(2)}, \dots, x^{(k)}]$, each dimension is normalized as follows

$$\hat{x}^{(k)} = \frac{x^{(k)} - E[x^{(k)}]}{\sqrt{\text{Var}(x^{(k)})}}. \quad (3.11)$$

(2) Simply normalizing each input of a layer may reduce what the layer can represent. For example, normalizing the inputs of a sigmoid would constrain them to the linear regime of the nonlinearity. In this way, batchnorm turns the multi-layer

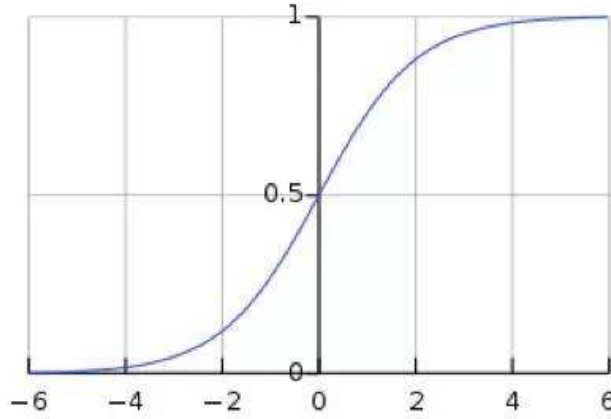


Figure 3.3: Sigmoid function.

nonlinear neural network into a multi-layer linear neural network, which is equivalent to the representation ability of a single layer linear network. This reduces the representation ability of the network. So, introduce $\gamma^{(k)}, \beta^{(k)}$ to scale and shift the normalized value:

$$y^{(k)} = \gamma^{(k)} \hat{x}^{(k)} + \beta^{(k)}, \quad (3.12)$$

which guarantee the representation ability of the network.

Why is multi-layer linear network equivalent to the representation ability of a single layer linear network.

In addition, use mini-batch in stochastic gradient training, each mini-batch produces estimates of the mean and variance of each activation. And obtain the forward algorithm of batch normalization as follows:

Algorithm 1 Forward algorithm of batch normalization.

Input: values of x over a mini-batch: $\mathcal{B} = \{x_1, x_2, \dots, x_m\}$, parameters to be learned: γ, β

Output: $y_i = BN_{\gamma, \beta}(x_i)$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta = BN_{\gamma, \beta}(x_i)$$

In this manner, we know

$$\begin{aligned}
E[\hat{x}] &= E\left[\frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}\right] \\
&= \frac{1}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}(E[x_i] - E[\mu_{\mathcal{B}}]) \\
&= 0 \\
Var[\hat{x}] &= E[\hat{x}^2] - (E[\hat{x}])^2 \\
&= E\left[\left(\frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}\right)^2\right] + 0 \\
&= \frac{1}{\sigma_{\mathcal{B}}^2 + \epsilon} E[(x_i - \mu_{\mathcal{B}})^2] \\
&= \frac{1}{\sigma_{\mathcal{B}}^2 + \epsilon} \sigma_{\mathcal{B}}^2 \\
&\approx 1
\end{aligned} \tag{3.13}$$

The distribution of values of any \hat{x} has the expected value of 0 and the invariance of 1, as long as the elements of each mini-batch are sampled from the same distribution, and if we neglect ϵ . This distribution can be seen as an approximate standard normal distribution $N(0, 1)$.

We compute the gradients with respect to the parameters of the BN transform of loss \mathcal{L} for backpropagation. The chain rule is used as follows:

$$\begin{aligned}
\frac{\partial \mathcal{L}}{\partial \gamma} &= \sum_i \frac{\partial \mathcal{L}}{\partial y_i} \frac{\partial y_i}{\partial \gamma} = \sum_i \frac{\partial \mathcal{L}}{\partial y_i} \hat{x}_i, \\
\frac{\partial \mathcal{L}}{\partial \beta} &= \sum_i \frac{\partial \mathcal{L}}{\partial y_i} \frac{\partial y_i}{\partial \beta} = \sum_i \frac{\partial \mathcal{L}}{\partial y_i}.
\end{aligned} \tag{3.14}$$

There exists \sum_i for y_i because gradient from each y_i will backpropagate for γ and the effect of summation is achieved.

For the process of **backpropagation**, we also need to compute $\frac{\partial \mathcal{L}}{\partial x_i}$, which is got by:

$$\begin{aligned}
\frac{\partial \mathcal{L}}{\partial \hat{x}_i} &= \frac{\partial \mathcal{L}}{\partial y_i} \frac{\partial y_i}{\partial \hat{x}_i} \\
&= \frac{\partial \mathcal{L}}{\partial y_i} \gamma \\
\frac{\partial \mathcal{L}}{\partial \sigma_{\mathcal{B}}^2} &= \sum_i \frac{\partial \mathcal{L}}{\partial \hat{x}_i} \frac{\partial \hat{x}_i}{\partial \sigma_{\mathcal{B}}^2} \\
&= \sum_i \frac{\partial \mathcal{L}}{\partial \hat{x}_i} (x_i - \mu_{\mathcal{B}}) \left(-\frac{1}{2}\right) (\sigma_{\mathcal{B}}^2 + \epsilon)^{-\frac{3}{2}} \\
&= -\frac{1}{2} \sum_i \frac{\partial \mathcal{L}}{\partial \hat{x}_i} (x_i - \mu_{\mathcal{B}}) (\sigma_{\mathcal{B}}^2 + \epsilon)^{-\frac{3}{2}} \\
\frac{\partial \mathcal{L}}{\partial \mu_{\mathcal{B}}} &= \frac{\partial \mathcal{L}}{\partial \hat{x}_i} \frac{\partial \hat{x}_i}{\partial \mu_{\mathcal{B}}} + \frac{\partial \mathcal{L}}{\partial \sigma_{\mathcal{B}}^2} \frac{\partial \sigma_{\mathcal{B}}^2}{\partial \mu_{\mathcal{B}}} \\
&= \sum_i \frac{\partial \mathcal{L}}{\partial \hat{x}_i} (-1) (\sigma_{\mathcal{B}}^2 + \epsilon)^{-\frac{1}{2}} + \frac{\partial \mathcal{L}}{\partial \sigma_{\mathcal{B}}^2} \frac{1}{m} \sum_i (-2) (x_i - \mu_{\mathcal{B}}) \\
&= -\sum_i \frac{\partial \mathcal{L}}{\partial \hat{x}_i} \frac{1}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} - \frac{\partial \mathcal{L}}{\partial \sigma_{\mathcal{B}}^2} \frac{2}{m} \sum_i (x_i - \mu_{\mathcal{B}}) \\
\frac{\partial \mathcal{L}}{\partial x_i} &= \frac{\partial \mathcal{L}}{\partial \hat{x}_i} \frac{\partial \hat{x}_i}{\partial x_i} + \frac{\partial \mathcal{L}}{\partial \sigma_{\mathcal{B}}^2} \frac{\partial \sigma_{\mathcal{B}}^2}{\partial x_i} + \frac{\partial \mathcal{L}}{\partial \mu_{\mathcal{B}}} \frac{\partial \mu_{\mathcal{B}}}{\partial x_i} \\
&= \frac{\partial \mathcal{L}}{\partial \hat{x}_i} \frac{1}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} + \frac{2}{m} \frac{\partial \mathcal{L}}{\partial \sigma_{\mathcal{B}}^2} (x_i - \mu_{\mathcal{B}}) + \frac{1}{m} \frac{\partial \mathcal{L}}{\partial \mu_{\mathcal{B}}}
\end{aligned} \tag{3.15}$$

For the **inference** of the network, sometimes we want the output to depend only on the few or one input. At this time, μ and σ^2 are computed by each mini-batch μ_B and σ_B^2 by their unbiased estimations as proved in 1.17:

$$\begin{aligned}\mu &= E(\mu_B), \\ \sigma^2 &= \frac{n}{n-1} E(\sigma_B^2), \\ BN(X) &= \gamma \frac{X - \mu}{\sqrt{\sigma^2 + \epsilon}} + \beta.\end{aligned}\tag{3.16}$$

The advantages of batch normalization is as follows: (1) batch normalization ensure the stable distributions of input in each layer, which reduces the learning pressure at the later layer and improves the learning speed; (2) batch normalization makes the model less sensitive to the parameters in the network, which simplifies the parameter adjustment precess, and makes the learning of this network more stable; (3) batch normalization can make the input of activation fall into the unsaturated region through the normalization operation, thus alleviating the problem of gradient vanishing.

3.3 Why does residual learning work?

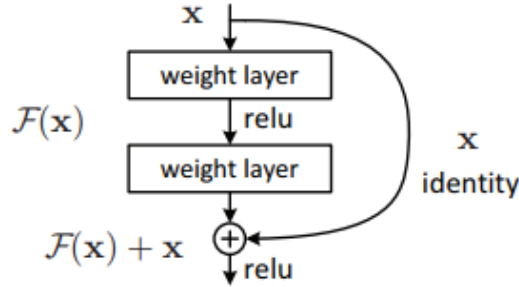


Figure 3.4: Residual Learning.

In the process of optimization in deep network k, the input and output are usually close in the latter layers. In some latter layers, we define the input is x and the output is $\mathcal{H}(x)$. From above experience, we assume x and $\mathcal{H}(x)$ are close. Thus, the function of these layers is mapping from x to $\mathcal{H}(x)$. But usually this process is difficult because the relative gap between x and $\mathcal{H}(x)$ is small. Toward this end, we construct the residual as $\mathcal{F}(x) = \mathcal{H}(x) - x$ to learning the gap. We can hypothesize that it is easier to optimize the residual mapping than to optimize the original mapping. To the extreme, if the mapping from input to output is identity, it would be easier to push the residual to zero than to directly fit an identity mapping.

The real purpose of residual learning is that even if the network deepens, the performance of this network will not degenerate, thus ensuring the leaning of deeper network (even 1000 layers).

Chapter 4

Contents specifically referenced

- the chapter of (2)Perceptron in [2];
- BatchNorm from [1];

Bibliography

- [1] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In ICML, 2015.
- [2] Hang Li. Statistical learning method. Tsinghua University Press, 2012.