

# Learning Machine Learning

Pt

October 18, 2020

## Contents

<b>1</b>	<b>Glossary</b>	<b>2</b>
1.1	Odds and Logit . . . . .	2
<b>2</b>	<b>Loss Functions and Objective Functions</b>	<b>2</b>
2.1	Minimum Square Loss (MSE) . . . . .	2
2.2	Cross-Entropy/Log Loss . . . . .	2
2.3	Maximum (Log-)Likelihood . . . . .	2
<b>3</b>	<b>Probability</b>	<b>3</b>
3.1	Basis . . . . .	3
3.2	Conditional Independent . . . . .	3
3.3	Distributions . . . . .	3
3.3.1	Bernoulli distribution . . . . .	3
<b>4</b>	<b>Matrix</b>	<b>3</b>
4.1	Differentiate/Derivation . . . . .	3
<b>5</b>	<b>Normalization</b>	<b>3</b>
5.0.1	Norm . . . . .	3
5.1	Scale . . . . .	3

# 1 Glossary

## 1.1 Odds and Logit

In *Binary Classification Problem*, the probability of  $label = 1$  divided by the probability of  $label = 0$  is called **Odds**.

$$y = P(label = 1)$$
$$odds = \frac{y}{1 - y}$$

Further, take the logarithm of both sides, we got **Log Odds/Logit**:

$$logit = \ln \frac{y}{1 - y}$$

## 2 Loss Functions and Objective Functions

### 2.1 Minimum Square Loss (MSE)

### 2.2 Cross-Entropy/Log Loss

For *Binary Classification Problem*. Given that  $x_i$  is one of the sample training data,  $y_i$  is the corresponding label, then

$$\hat{y}_i = \sigma(h(x_i|\theta)) \in \mathbb{R}$$
$$\mathcal{L}(y_i, \hat{y}_i | \theta) = \begin{cases} -\log(\hat{y}_i) & y_i = 1 \\ -\log(1 - \hat{y}_i) & y_i = 0 \end{cases}$$

Compress into one equation, then

$$\mathcal{L}(y_i, \hat{y}_i | \theta) = -[y_i * \log(\hat{y}_i) + (1 - y_i) * \log(1 - \hat{y}_i)]$$

More generally, for *Multi-Classification Problem*, given that

- $x_i$  is one of the sample training data, which will be classified into one of  $k$  categories,
- $y_i \in \mathbb{R}^k$  is the **One-Hot Representation** of the corresponding label

the **Cross-Entropy Loss** is

$$\hat{y}_i = \text{softmax}(h(x_i|\theta)) \in \mathbb{R}^k$$
$$\mathcal{L}(y_i, \hat{y}_i | \theta) = - \sum_j^n y_i[j] * \log_2(\hat{y}_i[j])$$

### 2.3 Maximum (Log-)Likelihood

Given the probability of a series of accidents  $A_i, i \in 1, 2, \dots, k$  is  $P(A_i)$ , then we want to **maximize** the probability that all of these accidents happen, thus

$$\max\{\prod_{i=0}^k P(A_i)\}$$

To simplify, we can take the logarithm of both sides, then

$$\max\{\ln \sum_{i=0}^k P(A_i)\}$$

which is namely **Maximum (Log) Likelihood**. While the **Loss Function** derived from above is naturally:

$$\mathcal{L}(y_i, \hat{y} | \theta) = -\ln \sum_{i=0}^k P(A_i)$$

which we want to **Minimize**.

## 3 Probability

### 3.1 Basis

- Priori Probability: the probability which can be empirically inferred
- Posterior Probability: after A happening, sought the probability of the reason of A
- Bayesian Equation

### 3.2 Conditional Independent

$$p(A | C) * p(B | C) = p(AB | C)$$

### 3.3 Distributions

#### 3.3.1 Bernoulli distribution

Defined as is the discrete probability distribution of a random variable which takes the value 1 with probability  $p$  and the value 0 with probability  $q = 1 - p$ . Denote  $B(1, p)$  as *Bernoulli distribution*, then

$$\begin{aligned} \text{Given } X &\sim B(1, p) \\ E(x) &= p; \\ D(x) &= p(1 - p); \end{aligned}$$

## 4 Matrix

### 4.1 Differentiate/Derivation

$$\begin{aligned} Y &= A \cdot X \cdot B \\ \frac{\partial Y}{\partial X} &= A^T \cdot B^T \\ \frac{\partial Y}{\partial X^T} &= B \cdot A \end{aligned}$$

Another scenario,

$$\begin{aligned} Y &= X^T \cdot A \cdot X \\ \frac{\partial Y}{\partial X} &= (A + A^T) \cdot X \end{aligned}$$

## 5 Normalization

### 5.0.1 Norm

Given  $x \in \mathbb{R}^d$ ,

$$||x||_2 = \sqrt{x_1^2 + \dots + x_d^2} \quad (1)$$

### 5.1 Scale

Given  $p \in \mathbb{R}^k$ , where  $p$  is the result of *Dot-Product* in *Dot-Product Attention Mechanism*, in order to counteract gradient vanishing in *softmax*, scale  $p$  by

$$p_{norm} = \frac{p}{\sqrt{k}}$$

Given that  $x \in \mathbb{R}^k$  is one of the record of the datasets, which is a **real-valued input vector**, in order to make *Gradient Descent* faster and more efficient, we'd better scale the value of input features to  $-1 \leq x_i \leq 1$ , or at least the scale of different input features is similar.

$$\bar{x}_i = \frac{x_i - \mu}{\max(x_i) - \min(x_i)}$$

where  $\mu$  is the mean value of  $x_i$  over the datasets;