

# 分类问题-Logistic Regression

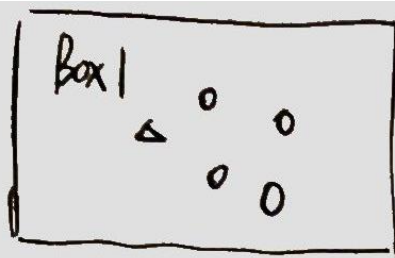
---

汇报人：王茂南

时间：2019年4月29日

# 例子引入

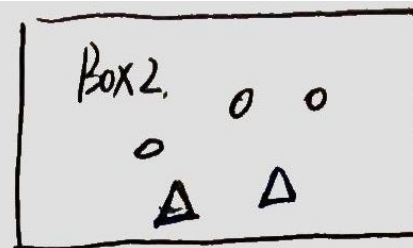
- 现在有两个Box, 里面装有一些○和▲ (个数如下图所示);
- 抽到Box1的概率为2/3, 抽到Box2的概率为1/3;



共5个,  $\begin{cases} 1 \uparrow \Delta \\ 4 \uparrow \circ \end{cases}$

从Box1抽到概率  $P(B_1) = \frac{2}{3}$

$$\begin{cases} P(\Delta|B_1) = \frac{1}{5} \\ P(\circ|B_1) = \frac{4}{5} \end{cases}$$



共5个,  $\begin{cases} 2 \uparrow \Delta \\ 3 \uparrow \circ \end{cases}$

从Box2抽到概率  $P(B_2) = \frac{1}{3}$

$$\begin{cases} P(\Delta|B_2) = \frac{2}{5} \\ P(\circ|B_2) = \frac{3}{5} \end{cases}$$

# 例子引入

- 现在的问题是, 如果抽到了一个▲, 问是来自Box1还是Box2;
- 我们可以使用贝叶斯公式进行求解;

$$P(A|B) = \frac{P(B|A)P(A)}{P(B|A)P(A) + P(B|A^C)P(A^C)}$$

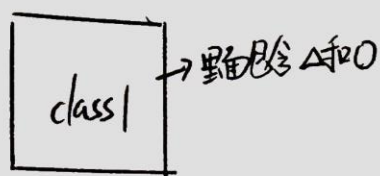
若现在抽-个方△, 问来自B<sub>1</sub>还是B<sub>2</sub>.

即求  $P(B_1|\triangle)$  → 发生的(观察到的现象)

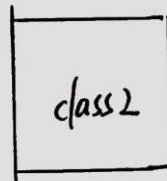
$$P(B_1|\triangle) = \frac{P(\triangle|B_1) \cdot P(B_1)}{P(\triangle|B_1)P(B_1) + P(\triangle|B_2)P(B_2)}$$

# 转换为分类问题

- 如果在分类的问题里面，上面的Box1和Box2代表class1和class2；
- 如果我们抽到的是三角，我们只需要计算 $P(\text{class1} | \text{三角})$ 和 $P(\text{class2} | \text{三角})$ 的概率大小进行比较即可。



已知  $P(\text{class1})$   
 $P(\Delta | \text{class1})$   
 $P(0 | \text{class1})$



已知  $P(\text{class2})$   
 $P(\Delta | \text{class2})$   
 $P(0 | \text{class2})$

**问题**

现在有样本(观测值)  $\Delta$ ，问  $\Delta$  属于 class1 or class2

**解法**

求  $P(\text{class1} | \Delta)$ ，与  $P(\text{class2} | \Delta)$

# 数学式子求解

- 于是, 求解分类问题, 转换为求解上面的贝叶斯式子, 根据解决的方法, 可以分为两个类别:
- Generative Model :  $P(x) = P(x|c_1) \cdot P(c_1) + P(x|c_2) \cdot P(c_2)$ , 我们可以求出 $x$ 的概率分布, 从而模拟 $x$ 的生成;
  - Logistic Model : 与线性回归思想类似;

$$P(c_1|x) = \frac{P(x|c_1) \cdot P(c_1)}{P(x|c_1) \cdot P(c_1) + P(x|c_2) \cdot P(c_2)}$$

↓

{ Generative Model. (计算  $P(x|c_1) \cdot P(c_1)$ )

Logistic Regression.

# Generative Model

□ Generative Model的想法是求出表达式中的所有概率;

- $P(c1), P(c2)$
- $P(x|c1), P(x|c2)$

$$P(c1|x) = \frac{P(x|c1) \cdot P(c1)}{P(x|c1) \cdot P(c1) + P(x|c2) \cdot P(c2)}$$

现在需要从样本中计算

┌	$P(c1), P(c2)$
	$P(x c1), P(x c2)$

①  $P(c1), P(c2)$  直接计算

$$P(c1) = \frac{\text{样本中 } c1 \text{ 的个数}}{\text{样本总数}}$$

# Generative Model

□ Generative Model 的想法是求出表达式中的所有概率;

■  $P(c1), P(c2)$

■  $P(x|c1), P(x|c2)$

□ 离散值

□ 连续值

②  $P(x|c1), P(x|c2)$

- 离散: 可直接进行估计
- 连续: 如  $x$  取值  $[0, 100]$ , 不是  $[0, 100]$  中每值都能取到.

假设  $c1, c2$  中的变量服从多元高斯分布

$f_{\mu_1, \Sigma_1}(x)$   $f_{\mu_2, \Sigma_2}(x)$  [一般假设  $\Sigma_1 = \Sigma_2$ ]

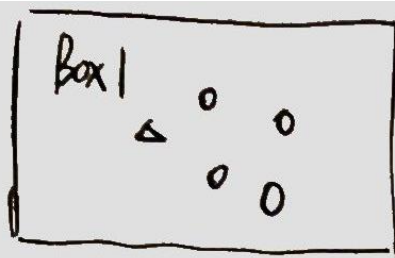
如何求  $\mu_1, \mu_2, \Sigma$   $\Rightarrow$  使用极大似然估计

$$\text{Max } L(\mu_1, \mu_2, \Sigma) = f_{\mu_1, \Sigma}(x_1) f_{\mu_2, \Sigma}(x_2) \dots$$

有了  $\mu_1, \mu_2, \Sigma$ , 即可求出  $P(x|c1), P(x|c2)$

# 离散变量解释

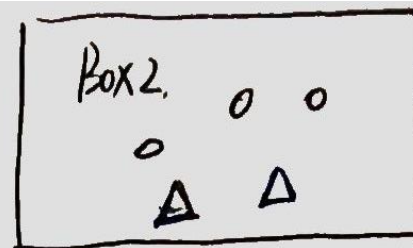
- 现在有两个Box, 里面装有一些○和▲ (个数如下图所示) ;
- 抽到Box1的概率为2/3, 抽到Box2的概率为1/3;



共5个,  $\begin{cases} 1 \text{ 个 } \Delta \\ 4 \text{ 个 } \circ \end{cases}$

从B<sub>1</sub>抽到概率  $P(B_1) = \frac{2}{3}$

$$\begin{cases} P(\Delta|B_1) = \frac{1}{5} \\ P(\circ|B_1) = \frac{4}{5} \end{cases}$$



共5个,  $\begin{cases} 2 \text{ 个 } \Delta \\ 3 \text{ 个 } \circ \end{cases}$

从B<sub>2</sub>抽到概率  $P(B_2) = \frac{1}{3}$

$$\begin{cases} P(\Delta|B_2) = \frac{2}{5} \\ P(\circ|B_2) = \frac{3}{5} \end{cases}$$



# Logistic Model

□ 我们从贝叶斯的表达式出发进行化简;

$$P(C_1|X) = \frac{P(X|C_1) \cdot P(C_1)}{P(X|C_1) \cdot P(C_1) + P(X|C_2) \cdot P(C_2)}$$

同除

$$\frac{1}{1 + \frac{P(X|C_2) \cdot P(C_2)}{P(X|C_1) \cdot P(C_1)}}$$

令  $z = z_n \frac{P(X|C_1) \cdot P(C_1)}{P(X|C_2) \cdot P(C_2)}$       则  $P(C_1|X) = \frac{1}{1 + e^{-z}}$

# Logistic Model

□ 我们从贝叶斯的表达式出发进行化简;

- 注意: 这是其中一种的解释方式, 我们可以从别的解释同样推导出logistic model的表达式;
- 如: 我们可以认为Logistic Model需要输出一个概率, 于是我们使用Sigmoid函数做压缩, 把(负无穷, 正无穷)的值压缩到(0,1)

如何用高斯分布  $f_{\mu_1, \Sigma}(x)$ ,  $f_{\mu_2, \Sigma}(x)$  代入  $\lambda z$ ,  $z$  以化简的表达式.

$$z = \underbrace{(\mu_1 - \mu_2)^T \Sigma^{-1} X}_{w} - \underbrace{\frac{1}{2} \mu_1^T \Sigma^{-1} \mu_2 + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_1 + \ln \frac{N_1}{N_2}}_b$$

这是两个常数!

$$\therefore P(C_1|x) = \frac{1}{1 + e^{wX+b}} = b(z), \quad \text{其中} \begin{cases} b(x) = \frac{1}{1 + e^{-x}} \\ z = wX + b \end{cases}$$

# 详细推导过程

对  $z$  进行化简

$$z = \ln \frac{p(x|c_1) \cdot p(c_1)}{p(x|c_2) \cdot p(c_2)}$$

$$= \ln \frac{p(x|c_1)}{p(x|c_2)} + \ln \frac{p(c_1)}{p(c_2)} \rightarrow \frac{M}{M+N_2}$$

$$= \ln \frac{p(x|c_1)}{p(x|c_2)} + \ln \frac{N_1}{N_2}$$

将  $p(x|c_1) = f_{\mu_1, \Sigma_1}(x)$ ,  $p(x|c_2) = f_{\mu_2, \Sigma_2}(x)$  代入

$$= \ln \frac{f_{\mu_1, \Sigma_1}(x)}{f_{\mu_2, \Sigma_2}(x)} + \ln \frac{N_1}{N_2}$$

$$= \ln \frac{\frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma_1|^{1/2}} e^{-\frac{1}{2}(x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1)}}{\frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma_2|^{1/2}} e^{-\frac{1}{2}(x-\mu_2)^T \Sigma_2^{-1} (x-\mu_2)}} + \ln \frac{N_1}{N_2}$$

$$= \ln \left[ \frac{|\Sigma_2|^{1/2}}{|\Sigma_1|^{1/2}} e^{-\frac{1}{2}[(x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1) - (x-\mu_2)^T \Sigma_2^{-1} (x-\mu_2)]} \right] + \ln \frac{N_1}{N_2}$$

$$= \ln \frac{|\Sigma_2|^{1/2}}{|\Sigma_1|^{1/2}} - \frac{1}{2} [(x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1) - (x-\mu_2)^T \Sigma_2^{-1} (x-\mu_2)] + \ln \frac{N_1}{N_2}$$

↓ 展开

$$= \ln \frac{|\Sigma_2|^{1/2}}{|\Sigma_1|^{1/2}} - \frac{1}{2} x^T (\Sigma_1^{-1})^T x + (\mu_1)^T \Sigma_1^{-1} x - \frac{1}{2} \mu_1^T \Sigma_1^{-1} \mu_1$$

$$+ \frac{1}{2} x^T (\Sigma_2^{-1})^T x - (\mu_2)^T (\Sigma_2^{-1})^T x + \frac{1}{2} \mu_2^T \Sigma_2^{-1} \mu_2$$

$$+ \ln \frac{N_1}{N_2}$$

由于  $\Sigma_1 = \Sigma_2$ , 可以进一步化简

$$z = \underbrace{(\mu_1 - \mu_2)^T \Sigma^{-1} x}_{b} - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 + \ln \frac{N_1}{N_2}$$

✓

# Logistic Model

□ 接下来需要求解Logistic Model 中的  $w$  和  $b$

关于评价函数

现在问题:  $f_{w,b}(x) = \frac{1}{1+e^{-z}}$  其中  $z=wx+b$  [ $f_{w,b}(x)$  表示  $x$  属于  $C_1$  的概率]

要求  $w, b$ . 可以用极大似然估计.

如. 有数据如下:

<u>data</u>	$x_1$	$x_2$	$x_3$	...
<u>label</u>	$c_1$	$c_2$	$c_1$	...

$$\text{Max } L(w,b) = f_{w,b}(x_1) \cdot [1 - f_{w,b}(x_2)] \cdot f_{w,b}(x_3) \cdot \dots$$

# Logistic Model

接下来需要求解Logistic Model 中的  $w$  和  $b \Rightarrow$  Cross Entropy (交叉熵)

$$\boxed{\text{Max}} \quad L(w, b) = f_{w, b}(x_1) \cdot [1 - f_{w, b}(x_2)] \cdot f_{w, b}(x_3) \cdots$$

↓  
手边转加法, 极大化转小

$$\begin{aligned} \boxed{\text{Min}} \quad -\ln L(w, b) &= -\ln f_{w, b}(x_1) \xleftarrow{\text{等价}} -[y_1 \ln f(x_1) + (1-y_1) \ln (1-f(x_1))] \\ &\quad -\ln (1-f_{w, b}(x_2)) \xleftarrow{\text{等价}} -[y_2 \ln f(x_2) + (1-y_2) \ln (1-f(x_2))] \\ &\quad -\ln f_{w, b}(x_3) \xleftarrow{\text{等价}} -[y_3 \ln f(x_3) + (1-y_3) \ln (1-f(x_3))] \end{aligned}$$

于是, 即要求

$$\boxed{\text{Min} \quad -\ln L(w, b) = \sum_{i=1}^n -y_i \ln f(x_i) - (1-y_i) \ln (1-f(x_i))}$$

Cross Entropy



使用梯度下降法求最优解。

## Logistic Regression

Step 1:  $f_{w,b}(x) = \sigma \left( \sum_i w_i x_i + b \right)$

Output: between 0 and 1

Training data:  $(x^n, \hat{y}^n)$

Step 2:  $\hat{y}^n$ : 1 for class 1, 0 for class 2

$$L(f) = \sum_n C(f(x^n), \hat{y}^n)$$

Step 3:

Logistic regression:  $w_i \leftarrow w_i - \eta \sum_n - \left( \hat{y}^n - f_{w,b}(x^n) \right) x_i^n$

Linear regression:  $w_i \leftarrow w_i - \eta \sum_n - \left( \hat{y}^n - f_{w,b}(x^n) \right) x_i^n$

## Linear Regression

$$f_{w,b}(x) = \sum_i w_i x_i + b$$

Output: any value

Training data:  $(x^n, \hat{y}^n)$

$\hat{y}^n$ : a real number

$$L(f) = \frac{1}{2} \sum_n (f(x^n) - \hat{y}^n)^2$$

# 多分类问题—Softmax介绍

$$C_1: w^1, b_1 \quad z_1 = w^1 \cdot x + b_1$$

$$C_2: w^2, b_2 \quad z_2 = w^2 \cdot x + b_2$$

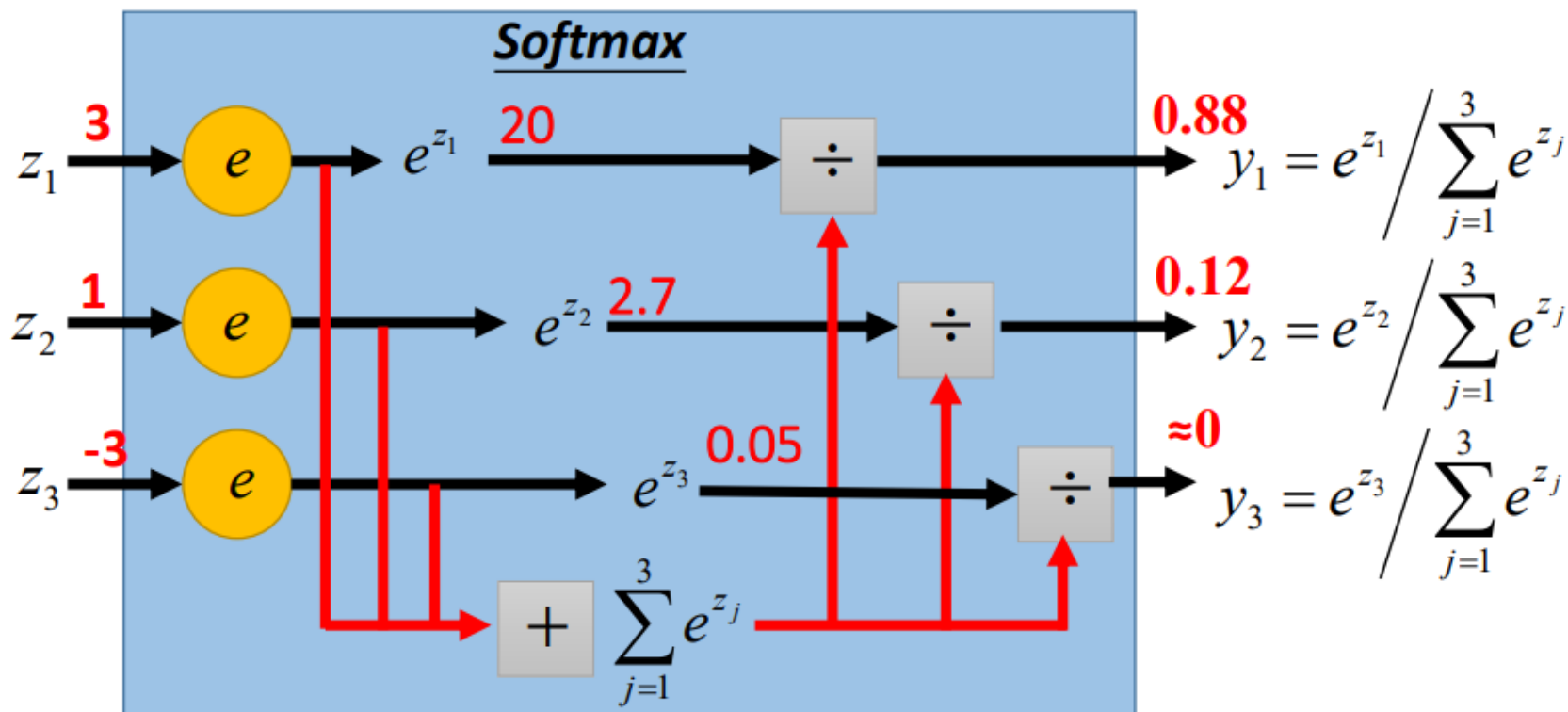
$$C_3: w^3, b_3 \quad z_3 = w^3 \cdot x + b_3$$

**Probability:**

$$\blacksquare 1 > y_i > 0$$

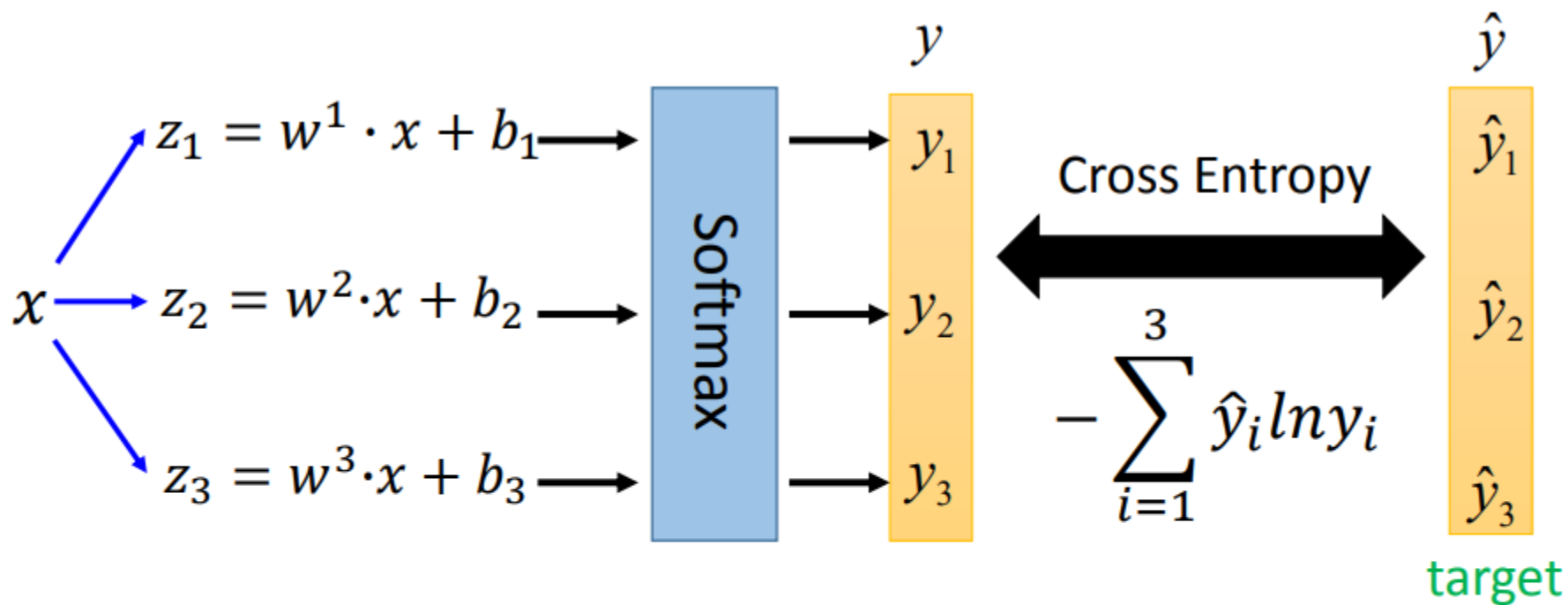
$$\blacksquare \sum_i y_i = 1$$

$$y_i = P(C_i | x)$$





# 多分类问题—Softmax介绍



If  $x \in \text{class 1}$

$$\hat{y} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

If  $x \in \text{class 2}$

$$\hat{y} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

If  $x \in \text{class 3}$

$$\hat{y} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$



# 二分类与多分类关系

$$z_1 = w_1 x + b_1 \longrightarrow y_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2}}$$

$$z_2 = w_2 x + b_2 \longrightarrow y_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2}}$$

在两个变量, 需比较  $z_1, z_2$  大小

$$\frac{z_1}{z_2} > 1 \quad \text{则 class 1}$$

$$z = \frac{z_1}{z_2} = wx + b \longrightarrow y_1 = \frac{e^z}{e^z + 1} \xrightarrow{\text{同除 } e^z} \frac{1}{1 + e^{-z}}$$

$$1 \longrightarrow y_2 = \frac{1}{e^z + 1}$$

⇒ 这样有了  $y$  的表达式.

# 二分类与多分类关系

---

(cross Entropy, 1-样本时)

$$-\sum_{i=1}^2 \hat{y}_i \ln y_i \quad \frac{\text{其中 } \hat{y}_1 + \hat{y}_2 = 1}{x_1 + x_2 = 1} \quad \hat{y} \ln y + (1-\hat{y}) \ln y$$

$$\Rightarrow \begin{cases} \text{其中 } \hat{y} \text{ 为 label} \\ y = f(x) \end{cases}$$

# Logistic Model实验

---

□ 数据集: The NSL-KDD Data Set

□ 训练集 : The full NSL-KDD train set including all difficulty level in CSV format.

- 离散变量使用One-hot编码;
- 离散变量不使用One-hot编码;

□ 测试集

- KDDTest : The full NSL-KDD test set including attack-type labels and difficulty level in CSV format.
- KDDTest-21 : A subset of the KDDTest+.txt file which does not include records with difficulty level of 21 out of 21

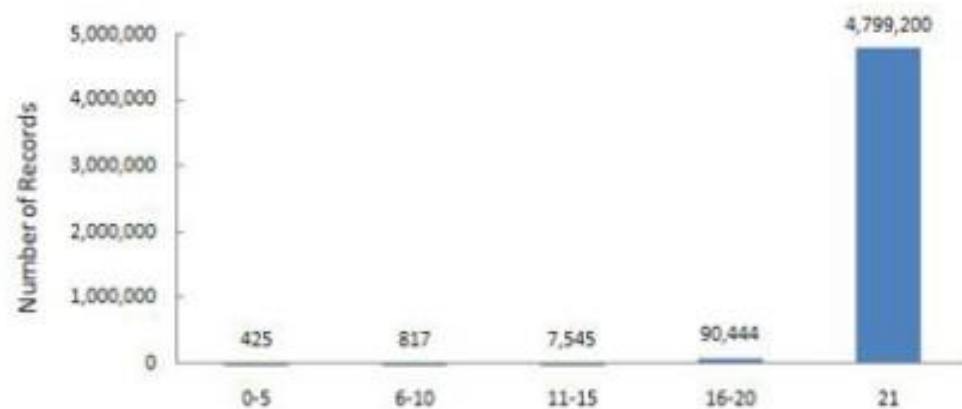
□ 关于 difficulty level 介绍

- 作者在NSL-KDD的时候, 对每一条数据集跑了21个算法;
  - difficulty level 个数表示分类成功的次数;
  - 即difficulty level 越小, 这条数据被分对的次数越小(也就是越难);
-

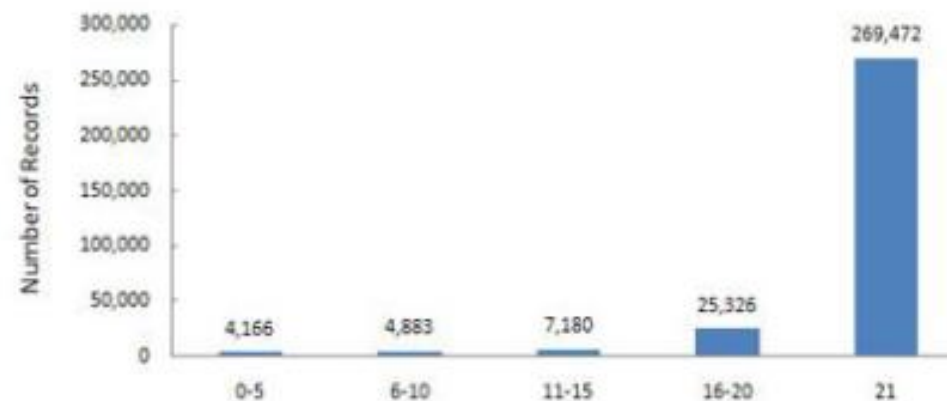
# Logistic Model实验

## □ 关于 difficulty level 介绍

- 作者在 NSL-KDD 的时候, 对每一条数据集跑了 21 个算法;
- difficulty level 个数表示分类成功的次数;
- 即 difficulty level 越小, 这条数据被分对的次数越小 (也就是越难);



The distribution of *#successfulPrediction* values for the KDD train set records



The distribution of *#successfulPrediction* values for the KDD test set records

# 实验结果

	KDDTest-21	KDDTest
使用One-hot编码	<pre>***** Training complete in 0m 31s Best val Acc: 0.974200 ***** 在测试集上进行测试. Accuracy of the network on the 11850 test Data: 57.6456 %</pre>	<pre>***** Training complete in 0m 31s Best val Acc: 0.972930 ***** 在测试集上进行测试. Accuracy of the network on the 22544 test Data: 79.1608 %</pre>
不使用One-hot编码	<pre>***** Training complete in 0m 32s Best val Acc: 0.954116 ***** 在测试集上进行测试. Accuracy of the network on the 11850 test Data: 56.4895 %</pre>	<pre>***** Training complete in 0m 37s Best val Acc: 0.954989 ***** 在测试集上进行测试. Accuracy of the network on the 22544 test Data: 78.3357 %</pre>

**局限性** : Logistic Model是只能进行线性分割的.

# 实验结果—疑问

- ❑ 使用Full的训练集，使用One-hot编码；
- ❑ Hidden-layer加到5层数；
- ❑ Dropout(0.9), 使用BN;

\*\*\*\*\*

Training complete in 4m 52s

Best val Acc: 0.978011

\*\*\*\*\*

在测试集上进行测试.

Accuracy of the Best\_model network on the 22544 test Data: 78.0075 %

Accuracy of the Final\_model network on the 22544 test Data: 77.5949 %

```
NeuralNet(  
  (inLayer): Linear(in_features=122, out_features=100, bias=True)  
  (relu): ReLU()  
  (hiddenLayer): Sequential(  
    (0): Linear(in_features=100, out_features=100, bias=True)  
    (1): BatchNorm1d(100, eps=1e-05, momentum=0.5, affine=True, track_running_stats=True)  
    (2): Dropout(p=0.9)  
    (3): ReLU()  
    (4): Linear(in_features=100, out_features=100, bias=True)  
    (5): BatchNorm1d(100, eps=1e-05, momentum=0.5, affine=True, track_running_stats=True)  
    (6): Dropout(p=0.9)  
    (7): ReLU()  
    (8): Linear(in_features=100, out_features=100, bias=True)  
    (9): BatchNorm1d(100, eps=1e-05, momentum=0.5, affine=True, track_running_stats=True)  
    (10): Dropout(p=0.9)  
    (11): ReLU()  
    (12): Linear(in_features=100, out_features=100, bias=True)  
    (13): BatchNorm1d(100, eps=1e-05, momentum=0.5, affine=True, track_running_stats=True)  
    (14): Dropout(p=0.9)  
    (15): ReLU()  
    (16): Linear(in_features=100, out_features=100, bias=True)  
    (17): BatchNorm1d(100, eps=1e-05, momentum=0.5, affine=True, track_running_stats=True)  
    (18): Dropout(p=0.9)  
    (19): ReLU()  
  )  
  (outLayer): Linear(in_features=100, out_features=2, bias=True)  
)
```

# 实验结果—疑问

---

- 用测试集做训练，训练集做测试；

```
*****
Training complete in 4m 52s
Best val Acc: 0.978011
*****
在测试集上进行测试.
Accuracy of the Best_model network on the 22544 test Data: 78.0075 %
Accuracy of the Final_model network on the 22544 test Data: 77.5949 %
```

```
*****
Training complete in 0m 54s
Best val Acc: 0.891748
*****
在测试集上进行测试.
Accuracy of the Best_model network on the 125973 test Data: 88.6420 %
Accuracy of the Final_model network on the 125973 test Data: 85.7652 %
```

- 我自己感觉test中出现了train中没有的数据；
  - 在KDD99中，测试集中出现了训练集中没有的攻击；
  - 对于NSL-KDD，还没有验证！！！！
-