

# Machine Learning Course (Li Hongyi 2023)

---

## 1.ChatGPT Introduction

---

### Possible Method :

Pre-Train 预训练 => Self-Supervised 自监督式学习 => Supervised learning 监督式学习 => Intensive training 强化训练

### Foreground:

Prompting 工程, Neural Editing, AI 检测, Machine Unlearning

---

## 2.Regression

---

### Step1 Linear Model

$$y = b + \omega * x_{cp} \Rightarrow y = b + \sum w_i * x_i$$

$x_i$  : feature,  $w_i$  : weight,  $b$  : bias

### Step2 Goodness of function

Loss Function L

$$L(f) = L(\omega, b)$$

Normal

$$L(f) = \sum_{n=1}^{10} (\hat{y}^n - (b + \omega * x_{cp}^n))^2$$

### Step3 Gradient Descent

#### For one parameter

$$\omega^* = \operatorname{argmin}_{\omega} L(\omega)$$

1.pick initial value  $\omega^0$

$$2.\omega^1 \leftarrow \omega^0 - \eta * \left. \frac{dL}{d\omega} \right|_{\omega=\omega^0}$$

$$3.\omega^2 \leftarrow \omega^1 - \eta * \left. \frac{dL}{d\omega} \right|_{\omega=\omega^1}$$

.... => Local optimal (not global)

## For two parameters

$$1. \omega^1 \leftarrow \omega^0 - \eta * \frac{dL}{d\omega} \Big|_{\omega=\omega^0, b=b^0}$$

$$b^1 \leftarrow b^0 - \eta * \frac{dL}{db} \Big|_{\omega=\omega^0, b=b^0}$$

$$2. \omega^2 \leftarrow \omega^1 - \eta * \frac{dL}{d\omega} \Big|_{\omega=\omega^1, b=b^1}$$

$$b^2 \leftarrow b^1 - \eta * \frac{dL}{db} \Big|_{\omega=\omega^1, b=b^1}$$

.....

## For Many Types

### back to design model

$$y = b_1 * \delta(x_s = \text{pidgey}) + \omega_1 * \delta(x_s = \text{pidgey}) * x_{cp} + \dots + b_4 * \delta(x_s = \text{Eevee}) + \omega_4 * \delta(x_s = \text{Eevee})$$

$$\delta = \begin{cases} 1, & x_s = \text{type} \\ 0, & x_s \neq \text{type} \end{cases}$$

### Regularization

$$L = \sum_n (\hat{y} - (b + \sum \omega_i x_i))^2 + \lambda \sum w_i^2$$

$\lambda$ 越大，找到的越平滑

---

## 3.Classification

$$x \rightarrow \text{Function} \rightarrow \text{Class}N$$

### Ideal Alternatives

#### Function(Model)

$$x \Rightarrow \begin{cases} g(x) > 0 & | \text{Output} = \text{class1} \\ \text{else} & | \text{Output} = \text{class2} \end{cases}$$

#### Loss Function

$$L(f) = \sum_n \delta(f(x^n) \neq \hat{y}^n)$$

Number of times f get incorrect results on training data.

### Find best function

Perceptron, SVM

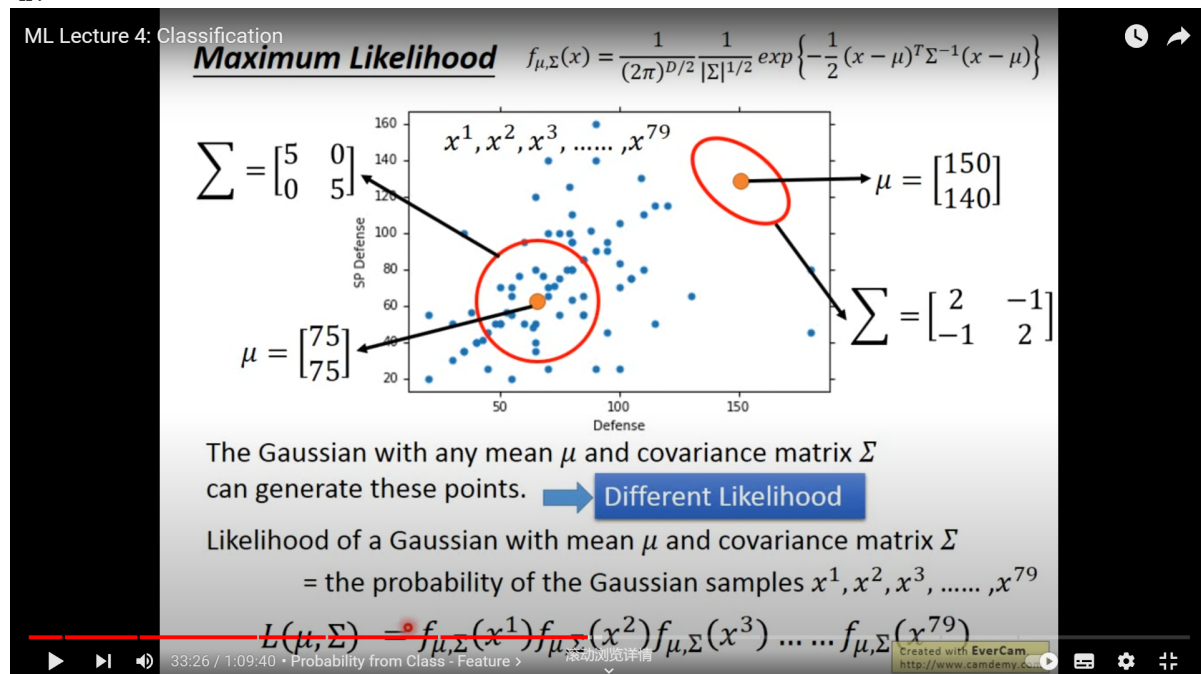
## Generative Model(Gaussian Distribution)

$$f_{\mu, \Sigma}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right\}$$

Determined by Mean  $\mu$ , covariance matrix  $\Sigma$

### Maximum Likelihood

![]()



$$L(\mu, \Sigma) = f_{\mu, \Sigma}(x^1) * f_{\mu, \Sigma}(x^2) * \dots * f_{\mu, \Sigma}(x^{79})$$

$$\Rightarrow (\mu^*, \Sigma^*) = \arg \max_{\mu, \Sigma} L(\mu, \Sigma)$$

$$\mu^* = \frac{1}{79} \sum_{n=1}^{79} x^n$$

$$\Sigma^* = \frac{1}{79} \sum_{n=1}^{79} (x^n - \mu^*)(x^n - \mu^*)^T$$

### Back To Classification

$$f_{\mu^1, \Sigma^1}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^1|^{1/2}} \exp\left\{-\frac{1}{2}(x - \mu^1)^T (\Sigma^1)^{-1}(x - \mu^1)\right\}$$

$\mu^1 = \begin{bmatrix} 75.0 \\ 71.3 \end{bmatrix}$     $\Sigma^1 = \begin{bmatrix} 874 & 327 \\ 327 & 929 \end{bmatrix}$

$P(C_1) = 79 / (79 + 61) = 0.56$

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

$$f_{\mu^2, \Sigma^2}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^2|^{1/2}} \exp\left\{-\frac{1}{2}(x - \mu^2)^T (\Sigma^2)^{-1}(x - \mu^2)\right\}$$

$\mu^2 = \begin{bmatrix} 55.6 \\ 59.8 \end{bmatrix}$     $\Sigma^2 = \begin{bmatrix} 847 & 422 \\ 422 & 685 \end{bmatrix}$

$P(C_2) = 61 / (79 + 61) = 0.44$

$$P(C_1|x) > 0.5 \Rightarrow x \in Class1$$

从多维空间来看增加更多参数容易导致overfitting

## Resolution

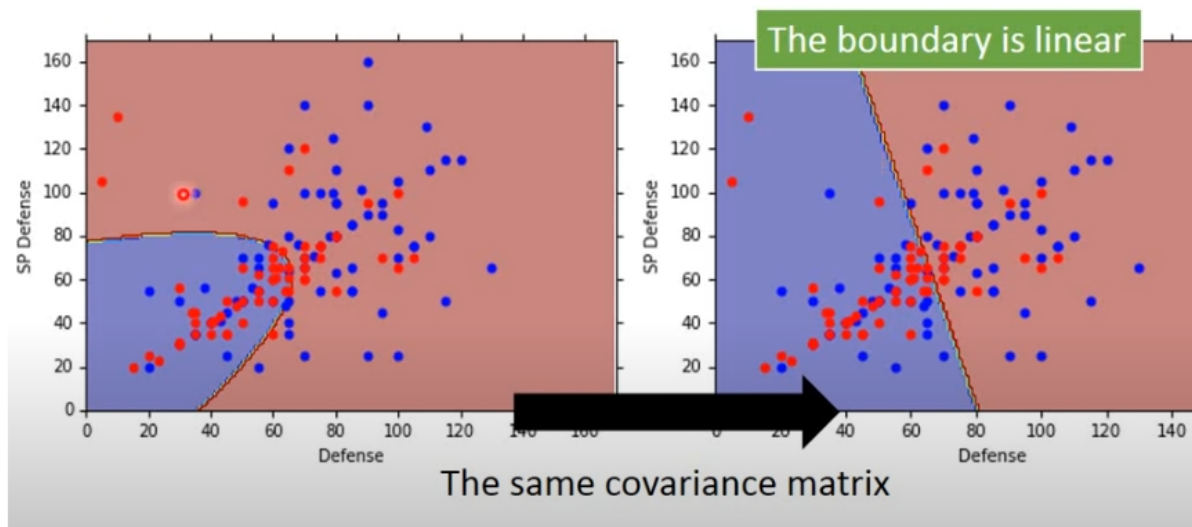
给两边分类相同的 $\Sigma$

$$L(\mu_1, \mu_2, \Sigma) = f_{\mu_1, \Sigma}(x^1) * f_{\mu_1, \Sigma}(x^2) \dots f_{\mu_1, \Sigma}(x^{79}) f_{\mu_2, \Sigma}(x^{80}) \dots f_{\mu_2, \Sigma}(x^{140})$$

$$\Rightarrow \mu_1 = \mu_2$$

$$\Rightarrow \Sigma = \frac{79}{140} \Sigma^1 + \frac{61}{140} \Sigma^2$$

## Modifying Model



## Three Steps

### Function Set(Model)

$$x \Rightarrow P(C_1|x) = \frac{P(x|C_1) * P(C_1)}{P(x|C_1) * P(C_1) + P(x|C_2) * P(C_2)} \Rightarrow \begin{cases} P(C_1|x) > 0.5 \rightarrow class1 \\ P(C_1|x) < 0.5 \rightarrow class2 \end{cases}$$

### Goodness of a function

mean  $\mu$ , covariance  $\Sigma$  Maximizing the likelihood

### Transformation

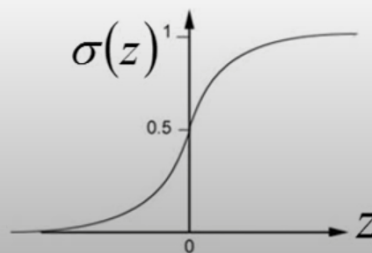
# Posterior Probability

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

$$= \frac{1}{1 + \frac{P(x|C_2)P(C_2)}{P(x|C_1)P(C_1)}} = \frac{1}{1 + \exp(-z)} = \sigma(z)$$

Sigmoid function

$$z = \ln \frac{P(x|C_1)P(C_1)}{P(x|C_2)P(C_2)}$$



Classification

$$z = \ln \frac{P(x|C_1)}{P(x|C_2)} + \ln \frac{P(C_1)}{P(C_2)} = \frac{N_1}{N_2}$$

$$P(x|C_1) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^1|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1) \right\}$$

$$P(x|C_2) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^2|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2) \right\}$$

$$\ln \frac{\frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^1|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1) \right\}}{\frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^2|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2) \right\}}$$

$$\ln \frac{|\Sigma^2|^{1/2}}{|\Sigma^1|^{1/2}} \exp \left\{ -\frac{1}{2} [(x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1) - (x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2)] \right\}$$

If  $\Sigma_1 = \Sigma_2 = \Sigma$

$$z = (\mu^1 - \mu^2)^T \Sigma^{-1} x - \frac{1}{2} (\mu^1)^T (\Sigma^{-1}) \mu^1 + \frac{1}{2} (\mu^2)^T (\Sigma^{-1}) \mu^2 + \ln \frac{N_1}{N_2}$$

$$z = w^T x - b$$

$$\Rightarrow P(C_1|x) = \sigma(w^T x + b)$$

## 4. Logistic Regression

### Comparison with Linear Regression

#### Step1

For Logistic Regression

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

Output: between 0 and 1

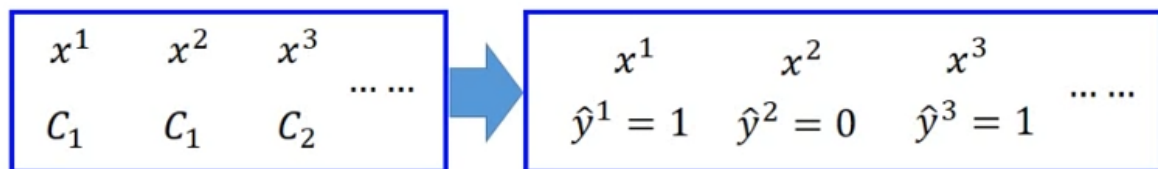
For Linear Regression

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

Output: any value

#### Step2 Goodness of a Function

Training Data :



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

$$L(w, b) = f_{w,b}(x^1) f_{w,b}(x^2) (1 - f_{w,b}(x^3)) \dots f_{w,b}(x^N)$$

The most likely  $w^*$  and  $b^*$  is the one with the largest  $L(w, b)$

$$w^*, b^* = \arg \max_{w, b} L(w, b)$$

$$\Rightarrow w^*, b^* = \arg \max_{w, b} -\ln L(w, b)$$

$$-\ln f_{w,b}(x^1) \Rightarrow -[\hat{y}^1 \ln f(x^1) + (1 - \hat{y}^1) \ln(1 - f(x^1))]$$

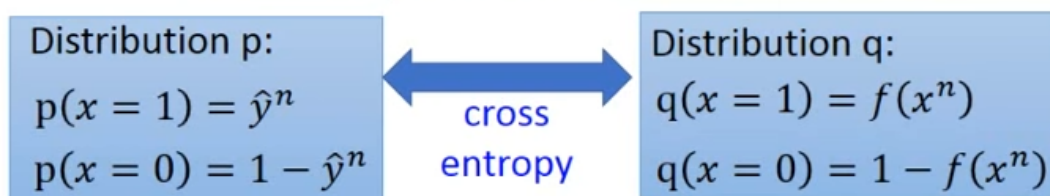
$$-\ln f_{w,b}(x^2) \Rightarrow -[\hat{y}^2 \ln f(x^2) + (1 - \hat{y}^2) \ln(1 - f(x^2))]$$

$$-\ln f_{w,b}(x^3) \Rightarrow -[\hat{y}^3 \ln f(x^3) + (1 - \hat{y}^3) \ln(1 - f(x^3))]$$

...

$$\Rightarrow -\ln L(w, b) = \sum_n -[\hat{y}^n \ln f_{w,b}(x^n) + (1 - \hat{y}^n) \ln(1 - f_{w,b}(x^n))]$$

Cross entropy between two Bernoulli distribution



$$H(p, q) = -\sum_x p(x) \ln(q(x))$$

## For Logistic Regression

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

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

## For Linear Regression

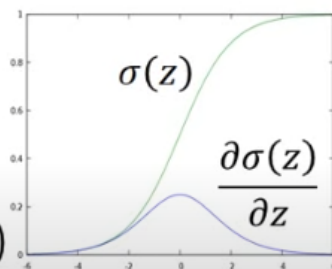
$$L(f) = \frac{1}{2} \sum_n (f(x^n) - \hat{y}^n)^2$$

## Step3 Find the best function

$$\frac{-\ln L(w, b)}{\partial w_i} = \sum_n - \left[ \hat{y}^n \frac{\partial \ln f_{w,b}(x^n)}{\partial w_i} + (1 - \hat{y}^n) \frac{\partial \ln (1 - f_{w,b}(x^n))}{\partial w_i} \right]$$

$$\frac{\partial \ln f_{w,b}(x)}{\partial w_i} = \frac{\partial \ln f_{w,b}(x)}{\partial z} \frac{\partial z}{\partial w_i} \quad \frac{\partial z}{\partial w_i} = x_i$$

$$\frac{\partial \ln \sigma(z)}{\partial z} = \frac{1}{\sigma(z)} \frac{\partial \sigma(z)}{\partial z} = \frac{1}{\sigma(z)} \sigma(z)(1 - \sigma(z))$$



$$\frac{\partial \ln (1 - f_{w,b}(x))}{\partial w_i} = \frac{\partial \ln (1 - f_{w,b}(x))}{\partial z} \frac{\partial z}{\partial w_i} \quad \frac{\partial z}{\partial w_i} = x_i$$

$$\frac{\partial \ln (1 - \sigma(z))}{\partial z} = - \frac{1}{1 - \sigma(z)} \frac{\partial \sigma(z)}{\partial z} = - \frac{1}{1 - \sigma(z)} \sigma(z)(1 - \sigma(z))$$

$$\frac{-\ln L(w, b)}{\partial w_i} = \sum_n - \left[ \hat{y}^n \frac{\partial \ln f_{w,b}(x^n)}{\partial w_i} + (1 - \hat{y}^n) \frac{\partial \ln (1 - f_{w,b}(x^n))}{\partial w_i} \right]$$

$$= \sum_n - \left[ \hat{y}^n \frac{\partial \ln f_{w,b}(x^n)}{\partial w_i} + (1 - \hat{y}^n) \frac{\partial \ln (1 - f_{w,b}(x^n))}{\partial w_i} \right]$$

$$= \sum_n - \left[ \hat{y}^n \frac{\partial \ln f_{w,b}(x^n)}{\partial w_i} + (1 - \hat{y}^n) \frac{\partial \ln (1 - f_{w,b}(x^n))}{\partial w_i} \right]$$

$$= \sum_n - \left( \hat{y}^n - f_{w,b}(x^n) \right) x_i^n$$

Larger difference,  
larger update

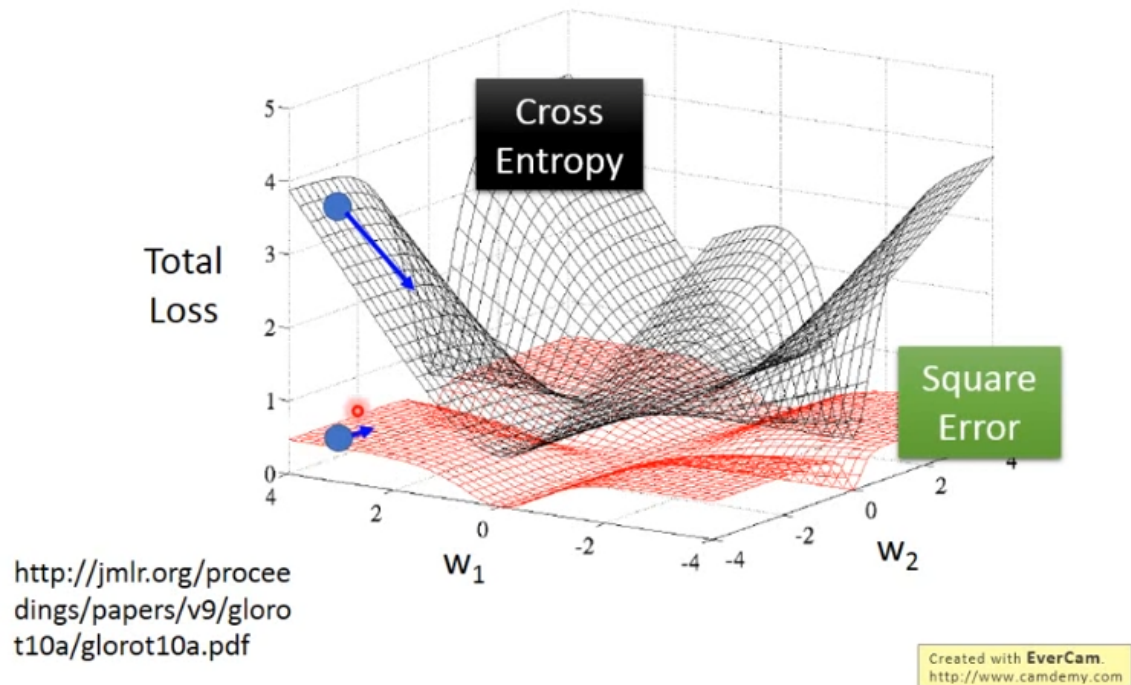
$$w_i \leftarrow w_i - \eta \sum_n - \left( \hat{y}^n - f_{w,b}(x^n) \right) x_i^n$$

## For Logistic Regression and Linear Regression

The same

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

## If Use Logistic Regression with square error



用Square error 离目标很远时趋势也很小