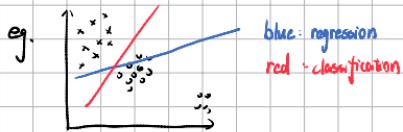


# Classification

## 1. Regression 和 Classification



## 2. Ideal Alternatives

Function (Model)

$$x \rightarrow \begin{cases} g(x) > 0 & \text{output = class 1} \\ \text{else} & \text{output = class 2} \end{cases}$$

Loss function

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

Find the best function

e.g. SVM, Perception

## 3. Generative Model

Bayes Formula

Class 1  Class 2 

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

$$P(x) = P(x|C_1)P(C_1) + P(x|C_2)P(C_2)$$

假设样本各 feature 取值服从 Gaussian Distribution

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

\* 使用 Maximum Likelihood 求解参数

推导:  $L(\mu, \Sigma) = \left( \frac{1}{(2\pi)^D \det(\Sigma)} \right)^{\frac{n}{2}} \exp \left\{ -\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right\}$

$$\ln L(\mu, \Sigma) = -\frac{D}{2} \ln (2\pi)^D \det(\Sigma) - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

$$\frac{\partial \ln L}{\partial \mu} = \sum_{i=1}^n (x_i - \mu) = 0 \Rightarrow \hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\frac{\partial \ln L}{\partial \Sigma} = -\frac{n}{2} \sum_{i=1}^n \sum_{j=1}^n (x_i - \mu_j) (x_i - \mu_j)^T \Sigma \Rightarrow \Sigma^* = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu}) (x_i - \hat{\mu})^T$$

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

其中  $x|C_i \sim N(\mu, \Sigma)$  用上述 ML 估计求  $\mu, \Sigma$

$$P(C_i) = \frac{\text{num}(C_i)}{\sum \text{num}(C_i)}$$

改进：① 减少参数，以避免 overfitting

假设不同 class 分布  $\mu$  不同，但  $\Sigma$  相同

find  $\mu^1, \mu^2, \Sigma$

$$L(\mu^1, \mu^2, \Sigma) = \left( \frac{1}{(2\pi)^n \det(\Sigma)} \right)^{\frac{n(m+n)}{2}} \exp \left[ -\frac{1}{2} \sum_{i=1}^m (x_i^1 - \mu^1)^T \Sigma^{-1} (x_i^1 - \mu^1) + \frac{1}{2} \sum_{i=m+1}^{m+n} (x_i^2 - \mu^2)^T \Sigma^{-1} (x_i^2 - \mu^2) \right]$$

$$L(\mu^1, \mu^2, \Sigma) = -\sum_{i=1}^m \ln \frac{1}{(2\pi)^n \det(\Sigma)} - \frac{1}{2} \left[ \sum_{i=1}^m (x_i^1 - \mu^1)^T \Sigma^{-1} (x_i^1 - \mu^1) + \sum_{i=m+1}^{m+n} (x_i^2 - \mu^2)^T \Sigma^{-1} (x_i^2 - \mu^2) \right]$$

$$\frac{\partial L}{\partial \mu^1} = \frac{m}{2} \sum_{i=1}^m (x_i^1 - \mu^1) \Rightarrow \mu^1 = \frac{1}{m} \sum_{i=1}^m x_i^1 \quad \frac{\partial L}{\partial \mu^2} = \frac{n}{2} \sum_{i=m+1}^{m+n} (x_i^2 - \mu^2) \Rightarrow \mu^2 = \frac{1}{n} \sum_{i=m+1}^{m+n} x_i^2$$

$$\frac{\partial L}{\partial \Sigma} = -\frac{m+n}{2} \Sigma + \frac{1}{2} \sum_{i=1}^m (x_i^1 - \mu^1)(x_i^1 - \mu^1)^T + \frac{1}{2} \sum_{i=m+1}^{m+n} (x_i^2 - \mu^2)(x_i^2 - \mu^2)^T \Rightarrow \Sigma^* = \frac{1}{m+n} \left( \sum_{i=1}^m (x_i^1 - \mu^1)(x_i^1 - \mu^1)^T + \sum_{i=m+1}^{m+n} (x_i^2 - \mu^2)(x_i^2 - \mu^2)^T \right)$$

## ② Naive Bayes Classifier

假设不同 feature 间独立

## ③ 线性的关系

$$P(C_i|x) = \frac{P(x|C_i)P(C_i)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)} = \frac{1}{1 + \frac{P(x|C_2)P(C_2)}{P(x|C_1)P(C_1)}}$$

$$\text{令 } z = \ln \frac{P(x|C_1)P(C_1)}{P(x|C_2)P(C_2)} = \frac{1}{1 + \exp(-z)} = \sigma(z) \quad \text{Sigmoid}$$

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

$$P(x|C_1) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \exp \left\{ -\frac{1}{2} (x - \mu^1)^T \Sigma^{-1} (x - \mu^1) \right\} \quad P(x|C_2) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \exp \left\{ -\frac{1}{2} (x - \mu^2)^T \Sigma^{-1} (x - \mu^2) \right\}$$

$$\therefore \ln \frac{P(x|C_1)}{P(x|C_2)} = \ln \sqrt{\frac{|\Sigma|}{|\Sigma|}} \exp \left\{ -\frac{1}{2} [(x - \mu^1)^T \Sigma^{-1} (x - \mu^1) - (x - \mu^2)^T \Sigma^{-1} (x - \mu^2)] \right\}$$

$$= \ln \sqrt{\frac{|\Sigma|}{|\Sigma|}} - \frac{1}{2} (x - \mu^1)^T \Sigma^{-1} (x - \mu^1) + \frac{1}{2} (x - \mu^2)^T \Sigma^{-1} (x - \mu^2)$$

$$\text{其中 } (x - \mu^1)^T \Sigma^{-1} (x - \mu^1) = x^T (\Sigma^{-1})^T x - x^T (\Sigma^{-1})^T \mu^1 - (\mu^1)^T (\Sigma^{-1})^T x + (\mu^1)^T (\Sigma^{-1})^T \mu^1$$

$$\Sigma' \text{ 对称} \Rightarrow \mathbf{x}'(\Sigma')\mathbf{x} = (\mathbf{y}')^T(\Sigma')^{-1}\mathbf{x} + (\mathbf{y}')^T(\Sigma')^{-1}\mathbf{y}$$

$$\therefore z = \frac{1}{2} \ln \left| \frac{\Sigma'}{\Sigma} \right| - \frac{1}{2} \mathbf{x}'(\Sigma')^{-1}\mathbf{x} + \frac{1}{2} \mathbf{x}'(\Sigma'^{-1}\mathbf{x} + (\mathbf{y}')^T(\Sigma')^{-1}\mathbf{x} - (\mathbf{y}')^T(\Sigma'^{-1}\mathbf{x}) - \frac{1}{2} (\mathbf{y}')^T(\Sigma')^{-1}\mathbf{y}^2 + \frac{1}{2} (\mathbf{y}')^T(\Sigma'^{-1})\mathbf{y}^2 + \ln \frac{M}{N}$$

一般全  $\Sigma' = \Sigma^2 = \Sigma$

$$z = (\mathbf{y}'\mathbf{y}^T)^T \Sigma^{-1} \mathbf{x} - \frac{1}{2} (\mathbf{y}')^T \Sigma^{-1} \mathbf{y}^2 + \frac{1}{2} (\mathbf{y}')^T \Sigma^{-1} \mathbf{y}^2 + \ln \frac{M}{N}$$

$$\text{全 } (\mathbf{y}'\mathbf{y}^T)^T \Sigma^{-1} = \mathbf{w}^T - \frac{1}{2} (\mathbf{y}')^T \Sigma^{-1} \mathbf{y}^2 + \frac{1}{2} (\mathbf{y}')^T \Sigma^{-1} \mathbf{y}^2 + \ln \frac{M}{N} = b$$

$$P(y|C_1) = \sigma(w \cdot x + b)$$

## 4. Logistic Regression

(1) Function Set

$$f_{w,b}(x) = P(C_1|x) = \sigma(w \cdot x + b) \in (0, 1)$$

(2) Goodness of a function

$$\begin{matrix} \text{Training Data:} & x^1 & x^2 & x^3 & \cdots & x^N \\ & C_1 & C_1 & C_2 & \cdots & C_1 \end{matrix}$$

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

$$w^*, b^* = \underset{w, b}{\operatorname{argmax}} L(w, b)$$

$$\text{转化为 } w^*, b^* = \underset{w, b}{\operatorname{argmin}} -\ln L(w, b)$$

$$= -\ln f_{w,b}(x_1) - \ln f_{w,b}(x_2) - \ln (1 - f_{w,b}(x_3)) - \cdots - \ln f_{w,b}(x^N)$$

$$\text{对 } i \text{ 来说 } x^i: \hat{y}^i = 1, x^2: \hat{y}^2 = 1, x^3: \hat{y}^3 = 0, \cdots$$

$$\text{每一项变为 } -[\hat{y}^i \ln f(x^i) + (1 - \hat{y}^i) \ln (1 - f(x^i))]$$

$$= -\sum_{i=1}^N [\hat{y}^i \ln f(x^i) + (1 - \hat{y}^i) \ln (1 - f(x^i))] \quad \text{Cross Entropy}$$

$$\begin{aligned} P(x=1) &= \hat{y}^n \\ P(x=0) &= 1 - \hat{y}^n \end{aligned}$$

$$\text{Cross Entropy} = -\sum_{i=1}^N p_i \ln q_i$$

$$\begin{aligned} q(x=1) &= f(x^n) \\ q(x=0) &= 1 - f(x^n) \end{aligned}$$

### 13) Optimize

$$-\frac{\partial \ln(f_{w,b})}{\partial w} = \sum_n -\hat{y}^n \frac{\partial \ln(f_{w,b})}{\partial w} + (1-\hat{y}^n) \frac{\partial \ln(1-f_{w,b})}{\partial w}$$

$$\frac{\partial \ln(f_{w,b})}{\partial w} = \frac{\partial \ln(\sigma(w \cdot x + b))}{\partial w} = \frac{1}{\sigma(w \cdot x + b)} \frac{\partial \sigma(w \cdot x + b)}{\partial w} = \frac{1}{\sigma(w \cdot x + b)} \cdot \sigma'(w \cdot x + b) \cdot x = [1 - \sigma(w \cdot x + b)] \cdot x$$

$$\frac{\partial \ln(1-f_{w,b})}{\partial w} = -\frac{1}{1-f_{w,b}} \cdot f_{w,b}[1-f_{w,b}] \cdot x = -f_{w,b}x$$

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

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

$$w^{(t+1)} \leftarrow w^{(t)} + \eta \sum_n (\hat{y}^n - f_{w,b}(x^n)) x^n$$

hint: ① 不用 square loss 原因

$$L(f) = \frac{1}{N} \sum_n (f_{w,b}(x^n) - \hat{y}^n)^2$$

$$\frac{\partial L(f)}{\partial w} = 2(f_{w,b}(x^n) - \hat{y}^n) f_{w,b}(x^n) (1 - f_{w,b}(x^n)) x^n$$

$$\text{if } \hat{y}^n = 0 \quad \text{if } f_{w,b}(x^n) = 1 \quad \text{update } \frac{\partial L}{\partial w} = 0$$

② 与 Generative model 关系.

$$\text{都为 } P(C_i|x) = \sigma(w \cdot x + b) \quad (\text{shared } \Sigma)$$

discriminative model: directly find w, b

generative model: find  $\mu^1, \mu^2, \Sigma^{-1}$ , then calculate w, b

w, b 不同: generative model: 带有一层假设

discriminative model: 可能局部最优

一代 discriminative model 效果好一些

### ③ multi-class classification

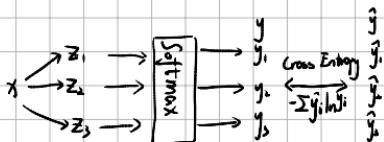
原理同上

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

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

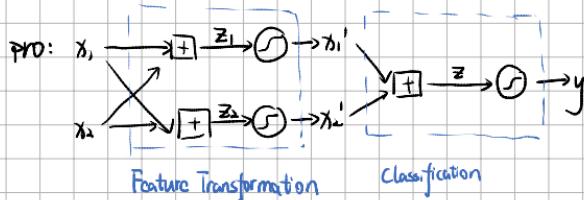
$$C_3: w^3, b^3 \quad z_3 = w^3 \cdot x + b$$

$$\text{Softmax} \quad y_t = \frac{e^{z_t}}{\sum_i e^{z_i}}$$



class 1	class 2	class 3
$\hat{y} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 1 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 1 \end{pmatrix}$

\* boundary 为一条直线



∴ 每一个 Neuron 是一个 Logistic Regression