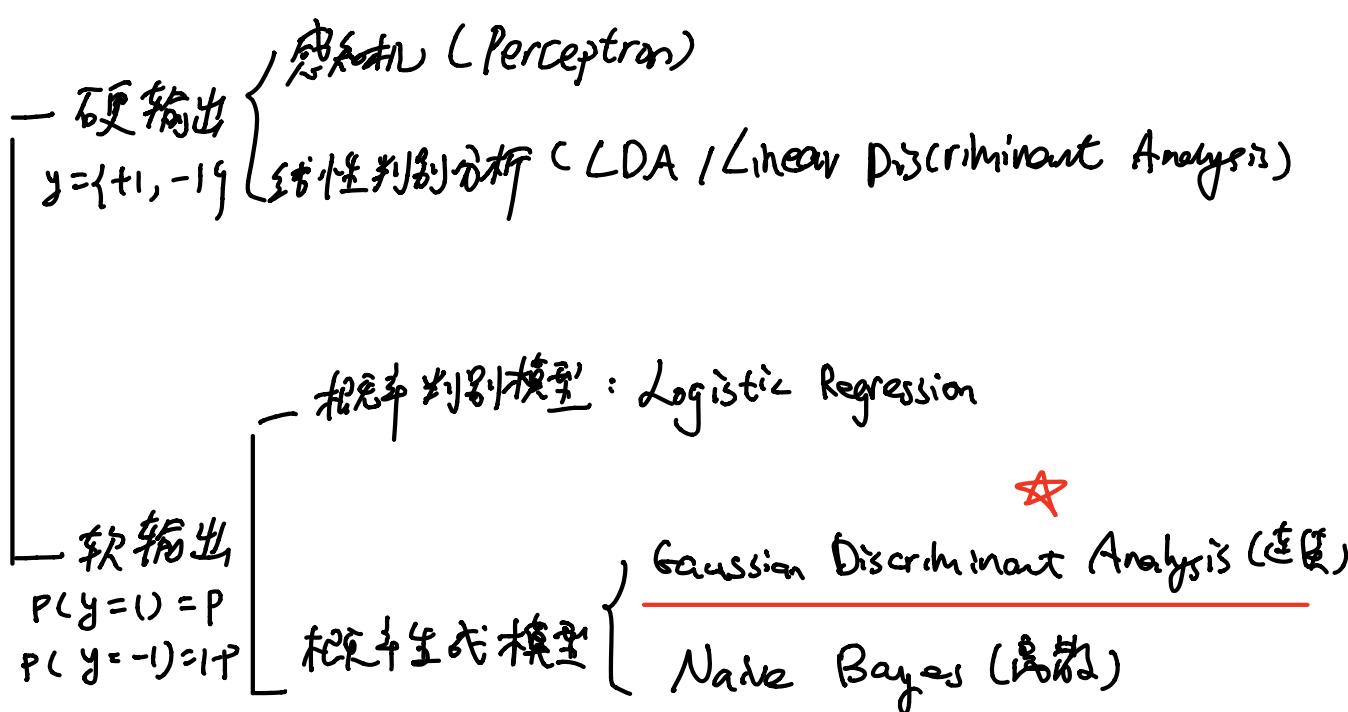
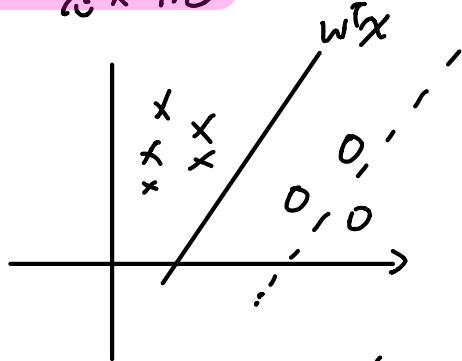


# 线性分类



## 一. 感知机



思想: 错误驱动

模型:  $f(x) = \text{sign}(w^T x)$ ,  $x \in \mathbb{R}^p, w \in \mathbb{R}^p$

$$\text{sign}(a) = \begin{cases} +1 & a > 0 \\ -1 & a \leq 0 \end{cases}$$

D. 《被错误分类的样本》

样本集:  $\{(x_i, y_i)\}_{i=1}^N$

策略: Loss Function:

$$L(w) = \sum_{i=1}^N \{y_i; w^T x_i \leq 0\} \quad \text{不可导}$$

$$\hookrightarrow L(w) = \sum_{x_i \in D} -y_i w^T x_i$$

Loss Function: 将随机事件或其有关随机变量的取值映射为非负数以表示该随机事件的“风险”或“损失”的函数.

求解感知机, 可以用 SGD 算法. (见统计学习方法第 2 版)

## 二. 线性判别分析 (LDA)

思想 ① 相同类内部的点距样本距离越小越好  
 ② 不同类之间距离越大越好.

要点: 找一个合适的投影方向

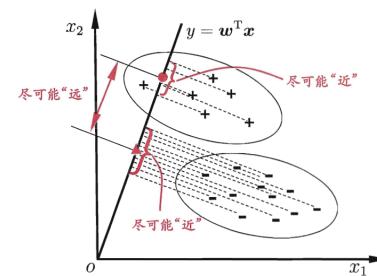
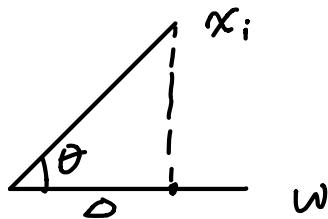


图 3.3 LDA 的二维示意图. “+”、“-”分别代表正例和反例, 椭圆表示数据簇的外轮廓, 虚线表示投影, 红色实心圆和实心三角形分别表示两类样本投影后的中心点.

$$S_w^{-1} (\bar{x}_{c_1} - \bar{x}_{c_2})$$



投影长度为  $\Delta = x_i \cos \theta$

当  $\|w\|=1$  时,  $w^T x_i = |x_i| \cdot \|w\| \cdot \cos \theta = |x_i| \cdot \cos \theta$

即  $\|w\|=1$  时,  $w^T x_i$  为投影长度.

如何表示同种类样本类间距离最小?  $\rightarrow$  用方差表示.

由上可知投影长度为  $z_i = w^T x_i$

则平均投影长度为  $\bar{z} = \frac{1}{N} \sum_{i=1}^N w^T x_i$

则方差为  $S = \frac{1}{N} \sum_{i=1}^N (z_i - \bar{z})(z_i - \bar{z})^T$

$$= \frac{1}{N} \sum_{i=1}^N (w^T x_i - \bar{z})(w^T x_i - \bar{z})^T$$

设分为  $C_1$  和  $C_2$  两组,  $\bar{z}_1, S_1$  为组  $C_1$  的平均值和方差,

$\bar{z}_2, S_2 \dots C_2 \dots \dots \dots$

$$\text{有 } C_1 \text{ 组: } \bar{z}_1 = \frac{1}{N_1} \sum_{i=1}^{N_1} w^T x_i \quad ①$$

$$S_1 = \frac{1}{N_1} \sum_{i=1}^{N_1} (w^T x_i - \bar{z}_1)(w^T x_i - \bar{z}_1)^T \quad ②$$

联立 ① ②:

$$S_1 = \frac{1}{N_1} \sum_{i=1}^{N_1} (w^T x_i - \frac{1}{N_1} \sum_{j=1}^{N_1} w^T x_j)(w^T x_i - \frac{1}{N_1} \sum_{j=1}^{N_1} w^T x_j)^T$$

$$= w^T \frac{1}{N_1} \sum_{i=1}^{N_1} (x_i - \frac{1}{N_1} \sum_{j=1}^{N_1} x_j) (x_i - \frac{1}{N_1} \sum_{j=1}^{N_1} x_j)^T$$

$$= w^T \left[ \frac{1}{N_1} \sum_{i=1}^{N_1} (x_i - \bar{x}_{c_1})(x_i - \bar{x}_{c_1})^T \right] w$$

$\bar{x}_{c_1}$  的方差

$$S_1 = W^T S_{C1} W$$

同理,  $C_2$  组 算得  $S_2 = W^T S_{C2} W$

两个类的类内分散度之和.

则类内距离可以表示为  $S_1 + S_2 = W^T (S_{C1} + S_{C2}) W$

- 如何表示不同样本类间距离最大?  $\rightarrow$  用平均值表示、

$$\begin{aligned} d &= (\bar{x}_1 - \bar{x}_2)^2 \\ &= \left( \frac{1}{N_1} \sum_{i=1}^{N_1} W^T x_i - \frac{1}{N_2} \sum_{i=1}^{N_2} W^T x_i \right)^2 \\ &= \left( W^T \left( \frac{1}{N_1} \sum_{i=1}^{N_1} x_i - \frac{1}{N_2} \sum_{i=1}^{N_2} x_i \right) \right)^2 \\ &= [W^T (\bar{x}_{C1} - \bar{x}_{C2})]^2 \end{aligned}$$

$$d = W^T (\bar{x}_{C1} - \bar{x}_{C2}) (\bar{x}_{C1} - \bar{x}_{C2}) W$$

- 综合类间大, 类内小思想. 可得目标函数

$$\begin{aligned} J(W) &= \frac{(Z_1 - Z_2)^2}{S_1 + S_2} \\ &= \frac{W^T (\bar{x}_{C1} - \bar{x}_{C2}) (\bar{x}_{C1} - \bar{x}_{C2}) W}{W^T (S_{C1} + S_{C2}) W} \end{aligned}$$

$$\hat{W} = \arg \max_W J(W)$$

$$\therefore S_b = (\bar{x}_{C1} - \bar{x}_{C2})(\bar{x}_{C1} - \bar{x}_{C2})^T$$

$$S_w = S_{C1} + S_{C2}$$

$$\therefore \hat{W} = \arg \max_W \frac{W^T S_b W}{W^T S_w W}$$

$$\Rightarrow \underline{\partial J(W)} = 0$$

矩阵求导公式  

$$\frac{\partial x^T A x}{\partial x} = 2 A x$$

$\omega(\omega)$

**方法一**

$$\Rightarrow \frac{\partial}{\partial \omega} \omega^T S_b \omega \cdot (\omega^T S_w \omega)^{-1} = 0$$

$$\cancel{S_b \omega (\omega^T S_w \omega)^{-1}} - \cancel{\omega^T S_b \omega (\omega^T S_w \omega)^{-2} \cdot S_w \omega} = 0$$

$$\Rightarrow S_b \omega (\omega^T S_w \omega)^{-1} = \omega^T S_b \omega \underline{(\omega^T S_w \omega)^{-2} S_w \omega}$$

$$\Rightarrow \frac{\omega^T S_b \omega}{R} S_w \omega = S_b \omega \frac{\omega^T S_w \omega}{R}$$

两边同乘

$$\Rightarrow S_w \omega = \frac{\omega^T S_w \omega}{\omega^T S_b \omega} \cdot S_b \omega$$

$$\omega = \frac{\omega^T S_w \omega}{\omega^T S_b \omega} \cdot S_w^{-1} \cdot S_b \omega$$

$$\propto S_w^{-1} \frac{S_b \cdot \omega}{\downarrow}$$

$$(\overline{x_{c_1}} - \overline{x_{c_2}}) \frac{(\overline{x_{c_1}} - \overline{x_{c_2}})^T \cdot \omega}{R}$$

$$\Rightarrow \propto \boxed{S_w^{-1} (\overline{x_{c_1}} - \overline{x_{c_2}})}$$

↑  
就是  $\omega$  的方向

## 方法二(西瓜书)

$$\hat{f} = \frac{\omega^T S_b \omega}{\omega^T S_w \omega}$$

因为  $f$  的值与  $\omega$  的长度无关，只与它的方向有关。令  $\omega^T S_w \omega = 1$   
 故问题转化为最优化问题

$$\begin{aligned} & \min -\omega^T S_b \omega \\ \text{s.t. } & \omega^T S_w \omega = 1 \end{aligned}$$

$$\frac{\partial x^T A x}{\partial x} = (A + A^T) x$$

$$\text{对称: } \frac{\partial x^T A x}{\partial x} = 2Ax$$

$$\Rightarrow L(\omega, \lambda) = -\omega^T S_b \omega + \lambda (\omega^T S_w \omega - 1)$$

$$\frac{\partial L(w, \lambda)}{\partial w} = - (S_b + S_b^T) w + \lambda (S_w + S_w^T) w$$

$\therefore S_b = S_b^T, S_w = S_w^T$  (对称)

$$\therefore \frac{\partial L(w, \lambda)}{\partial w} = 2S_b w + 2\lambda S_w w = 0$$

$$\Rightarrow S_b w = \lambda S_w w$$

$$(\bar{x}_{c_1} - \bar{x}_{c_2})(\bar{x}_{c_1} - \bar{x}_{c_2})^T w = \lambda S_w w$$

令  $(\bar{x}_{c_1} - \bar{x}_{c_2})^T w = \varphi$

有  $\varphi(\bar{x}_{c_1} - \bar{x}_{c_2}) = \lambda S_w w$

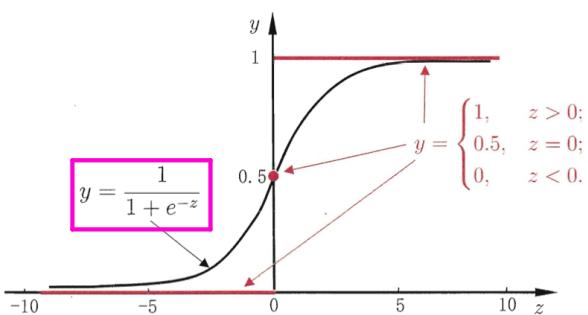
$$w = \frac{\varphi}{\lambda} S_w^{-1} (\bar{x}_{c_1} - \bar{x}_{c_2})$$

$\because w$ 不必关心其值. 取  $\frac{\varphi}{\lambda} = 1$ . 得  $w = S_w^{-1} (\bar{x}_{c_1} - \bar{x}_{c_2})$

### 三. logistics Regression

• 原理: 线性回归  $\xrightarrow{\text{激活函数}}$  线性分类  
 $w^T x$  ( $0, 1$ )

• logistics Regression 使用的激活函数为 sigmoid function:  $\sigma(z) = \frac{1}{1+e^{-z}}$



Sigmoid Function

Data:  $\{(x_i, y_i)\}_{i=1}^n$   
 $x_i \in \mathbb{R}^p, y_i \in \{0, 1\}$

$P_j = P(y=1|x) = \sigma(w^T x) = \frac{1}{1+e^{-w^T x}}, j=1$

$$P_0 = P(y=0|x) = 1 - P_1 = 1 - \frac{1}{1+e^{-w^T x}} = \frac{e^{-w^T x}}{1+e^{-w^T x}}$$

$$\rightarrow P(y|x) = P_1 P_0^{1-y}$$

BP 可用 MLE 估计  $w$  的最大值.

$$\text{MLE: } \arg \max_w I_n P(y|x)$$

$$= \arg \max_w I_n \prod_{i=1}^N P(y_i|x_i)$$

$$= \arg \max_w \prod_{i=1}^N I_n(P_1^{y_i} P_0^{1-y_i})$$

$$= \arg \max_w \sum_{i=1}^N (y_i \log P_1 + (1-y_i) \log P_0) \quad \text{Cross Entropy}$$

$$= \arg \max_w \sum_{i=1}^N (y_i (\ln P_1 + \ln P_0) - y_i \ln P_0)$$

$$= \arg \max_w \sum_{i=1}^N (y_i (\ln \frac{P_1}{P_0}) + \ln P_0)$$

$$\Rightarrow P_1 = \frac{1}{1+e^{-w^T x}}, \quad P_0 = \frac{e^{-w^T x}}{1+e^{-w^T x}} \text{ 似然函数.}$$

$$\text{解: } \arg \max_w \sum_{i=1}^N (y_i \ln (e^{w^T x}) + \ln (\frac{1}{1+e^{w^T x}}))$$

$$= \arg \max_w \sum_{i=1}^N y_i w^T x - \ln (1+e^{w^T x})$$

$$\Rightarrow \arg \min_w \sum_{i=1}^N [y_i w^T x + \ln (1+e^{w^T x})]$$

求解此问题. 可以用 Gradient descent method 或 Newton method

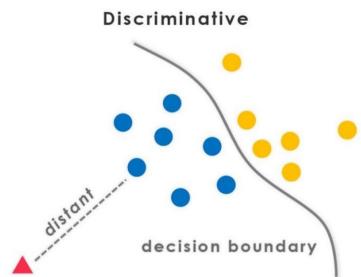
## ④. 高斯判别分析 (GDA)

生成模型

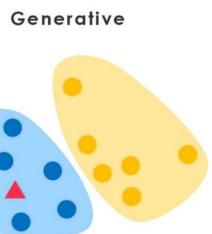
· 判别模型和生成模型的区别

判别模型: 求得  $P(Y|X)$ , 对未知示例  $x$ , 根据  $P(Y|X)$  可以获得标记  $Y$  (对二分类).

## Discriminative vs. Generative



- Only care about estimating the conditional probabilities
- Very good when underlying distribution of data is really complicated (e.g. texts, images, movies)



- Model observations  $(x, y)$  first, then infer  $p(y|x)$
- Good for missing variables, better diagnostics
- Easy to add prior knowledge about data

其实是得到一个 score. 考 score >

threshold, 则为正类). 常见的模型有 logistic Regression, SVM.

生成模型：求得  $p(y|x)$ . 对未见示例  $x$ , 你要求

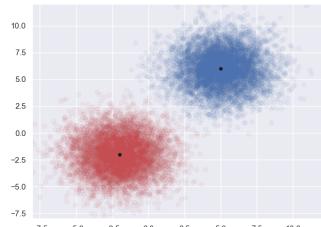
$x$  与不同标记之间的联合概率分布，然后大的那个获胜. 常见模型：GDA,

HMM, Naive Bayes

$$P(y|x) \propto \frac{P(x|y) \cdot P(y)}{\text{后验} \quad \text{似然} \quad \text{先验}} \rightarrow p(x, y)$$

$$\hat{y} = \arg \max_{y \in \{0, 1\}} P(y|x) = \arg \max_y P(y) \cdot P(x|y)$$

最大后验概率 (MAP)



GDA 模型定义

$$\text{定义 } y \in \text{Bernoulli } (\phi) \Rightarrow \phi^y (1-\phi)^{1-y}$$

$$x|y=1 \sim N(\mu_1, \Sigma) \quad \int \Rightarrow N(\mu_1, \Sigma)^y \cdot N(\mu_2, \Sigma)^{1-y}$$

$$x|y=0 \sim N(\mu_2, \Sigma) \quad \int \rightarrow \text{强假设: 基于不同分类的条件概率满足高斯分布. 它们拥有不同的均值, 但它们的方差是相同的. (} P(x|y=C_i) \text{ 是一样的!})$$

$$\log - \text{likelihood: } L(\theta) = \log \prod_{i=1}^N P(x_i, y_i) \quad \text{值相加但协方差相同的矩阵!}$$

$$= \prod_{i=1}^N \log [P(y_i) P(x_i|y_i)]$$

$$= \prod_{i=1}^N \log P(y_i) + \log P(x_i|y_i)$$

$$= \prod_{i=1}^N \log (\phi^{y_i} (1-\phi)^{1-y_i}) + \log (N(\mu_1, \Sigma)^{y_i} \cdot N(\mu_2, \Sigma)^{1-y_i})$$

$$L(\theta) = \prod_{i=1}^N (\log \phi^{y_i} + \log (1-\phi)^{1-y_i} + \log (N(\mu_1, \Sigma)^{y_i}) + \log (N(\mu_2, \Sigma)^{1-y_i}))$$

$$y_i \sim B(\phi)$$

$$\Theta = (\mu_1, \mu_2, \Sigma, \phi)$$

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmax}} L(\Theta)$$

$$\text{设 } y=1 \Rightarrow N_1.$$

$$y=0 \Rightarrow N_2$$

$$N = N_1 + N_2$$

• GDA 模型的求解

• 首先是对于  $\phi$  求解，有  $\phi$  的只有前两项，对其进行求偏导。

$$\text{即 } \frac{\partial}{\partial \phi} \left[ \sum_{i=1}^N y_i \log \phi + (1-y_i) \log (1-\phi) \right]$$

$$= \sum_{i=1}^N \left( \frac{y_i}{\phi} + \frac{1-y_i}{1-\phi} \right) = 0$$

$$\Rightarrow \sum_{i=1}^N [(1-\phi)y_i - \phi(1-y_i)] = 0$$

$$\Rightarrow \sum_{i=1}^N (y_i - \phi) = 0$$

$$\Rightarrow \sum_{i=1}^N y_i - N\phi = 0$$

$$\Rightarrow \phi = \frac{\sum_{i=1}^N y_i}{N} = \frac{N_1}{N}$$

• 然后对  $\mu_1$  求解。

$$\sum_{i=1}^N y_i \log N(\mu_1, \Sigma)$$

$$\sum_{i=1}^N y_i \log \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left( -\frac{1}{2} (x_i - \mu_1)^T \Sigma^{-1} (x_i - \mu_1) \right)$$

无关

$$\sum_{i=1}^N y_i - \frac{1}{2} (x_i - \mu_1)^T \Sigma^{-1} (x_i - \mu_1)$$

$$\text{求解得 } \hat{\mu}_1 = \underset{\mu_1}{\operatorname{argmax}} \sum_{i=1}^N y_i - \frac{1}{2} (x_i - \mu_1)^T \Sigma^{-1} (x_i - \mu_1) \Rightarrow \Delta$$

$$\Delta = -\frac{1}{2} \sum_{i=1}^N (x_i - \mu_1)^T \Sigma^{-1} (x_i - \mu_1)$$

$$(A+B)^T = A^T + B^T$$

$$= -\frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i^\top \Sigma^{-1} - \boldsymbol{\mu}_1^\top \Sigma^{-1}) (\mathbf{x}_i - \boldsymbol{\mu}_1) \quad (\text{实数})^T = \text{实数}$$

$$= -\frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i - \underline{\mathbf{x}_i^\top \Sigma^{-1} \boldsymbol{\mu}_1} - \underline{\boldsymbol{\mu}_1^\top \Sigma^{-1} \mathbf{x}_i} + \underline{\boldsymbol{\mu}_1^\top \Sigma^{-1} \boldsymbol{\mu}_1})$$

(x P x P x P x P x I)  
↓ 实数

$$= -\frac{1}{2} \sum_{i=1}^N (\underline{\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i} - 2 \boldsymbol{\mu}_1^\top \Sigma^{-1} \mathbf{x}_i + \underline{\boldsymbol{\mu}_1^\top \Sigma^{-1} \boldsymbol{\mu}_1})$$

(x P x P x P x P x I)  
R

$$\frac{\partial D}{\partial \boldsymbol{\mu}_1} = -\frac{1}{2} \sum_{i=1}^N (-2 \underline{\Sigma^{-1} \mathbf{x}_i} - 2 \underline{\Sigma^{-1} \boldsymbol{\mu}_1}) = 0$$

$$= \sum_{i=1}^N y_i (\underline{\Sigma^{-1} \mathbf{x}_i} + \underline{\Sigma^{-1} \boldsymbol{\mu}_1}) = 0$$

$\frac{\partial \mathbf{x}^\top A}{\partial \mathbf{x}} = A \quad \textcircled{1}$   
 $\frac{\partial \mathbf{x}^\top A \mathbf{x}}{\partial \mathbf{x}} = 2A\mathbf{x},$   
当  $A$  is symmetric

$$= \sum_{i=1}^N y_i (\mathbf{x}_i - \boldsymbol{\mu}_1) = 0$$

$$\Rightarrow \sum_{i=1}^N y_i \mathbf{x}_i = \sum_{i=1}^N y_i \boldsymbol{\mu}_1$$

$$\Rightarrow \boldsymbol{\mu}_1 = \frac{\sum_{i=1}^N \mathbf{x}_i y_i}{\sum_{i=1}^N y_i} = \frac{\sum_{i=1}^N \mathbf{x}_i y_i}{N_1}$$

由对称性 .  $\boldsymbol{\mu}_2 = \frac{\sum_{i=1}^N (1-y_i) \mathbf{x}_i}{N_2}$

最后对  $\Sigma$  求解 :

$$L(\theta) = \prod_{i=1}^n (\log \phi^{y_i} + \log (1-\phi)^{1-y_i}) + \frac{\log(N(\boldsymbol{\mu}_1, \Sigma))}{\textcircled{1}} + \frac{\log(N(\boldsymbol{\mu}_2, \Sigma))}{\textcircled{2}}$$

有的只有后两项.  $\hat{\Sigma} = \operatorname{argmax} \textcircled{1} + \textcircled{2}$

$$= \operatorname{argmax} \sum_i y_i \log N(\boldsymbol{\mu}_1, \Sigma) + (1-y_i) \log N(\boldsymbol{\mu}_2, \Sigma)$$

$$C_1 = \{\mathbf{x}_i | y_i = 1, i = 1, \dots, N\}$$

$$C_2 = \{\mathbf{x}_i | y_i = 0, i = 1, \dots, N\}$$

$$C_1 = |N_1|, C_2 = |N_2|, N_1 + N_2 = N$$

$$\text{则有 } \textcircled{1} + \textcircled{2} = \frac{\sum_{x \in C_1} \log N(M_1, \Sigma) + \sum_{x \in C_2} \log N(M_2, \Sigma)}{\textcircled{1}}$$

$$\text{先看 } \textcircled{1}: \sum_{i=1}^N \log N(M_1, \Sigma)$$

$$= \sum_{i=1}^N \log \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x_i - M_1)^T \Sigma^{-1} (x_i - M_1) \right\}$$

$$= \sum_{i=1}^N \frac{\log \frac{1}{(2\pi)^{\frac{p}{2}}}}{C} - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x_i - M_1)^T \Sigma^{-1} (x_i - M_1)$$

$$= C - N \frac{1}{2} \log |\Sigma| - \sum_{i=1}^N \frac{1}{2} (x_i - M_1)^T \Sigma^{-1} (x_i - M_1) \quad \text{exp} \quad \text{pxp} \quad \text{pxp} \Rightarrow R$$

$$\therefore \sum_{i=1}^N \text{tr}[(x_i - M_1)^T \Sigma^{-1} (x_i - M_1)]$$

$$\text{tr}(ABC) = \text{tr}(CAB) \\ = \text{tr}C \text{tr}B \text{tr}A$$

$$= \sum_{i=1}^N \text{tr}[(x_i - M_1)(x_i - M_1)^T \Sigma^{-1}]$$

$$S = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T$$

$$= \text{tr} \left[ \sum_{i=1}^N (x_i - M_1)(x_i - M_1)^T \Sigma^{-1} \right] \quad NS$$

$$\therefore \hat{f}_2(x) = -\frac{1}{2} N \log |\Sigma| - \frac{1}{2} N \cdot \text{tr}(S \cdot \Sigma^{-1}) + C$$

$$\textcircled{1} + \textcircled{2} \Rightarrow -\frac{1}{2} N_1 \log |\Sigma_1| - \frac{1}{2} N_1 \cdot \text{tr}(S_1 \cdot \Sigma_1^{-1}) + C$$

$$-\frac{1}{2} N_2 \log |\Sigma_2| - \frac{1}{2} N_2 \cdot \text{tr}(S_2 \cdot \Sigma_2^{-1}) + C$$

$$L(\Sigma) = -\frac{1}{2} N \log |\Sigma| - \frac{1}{2} N_1 \text{tr}(S_1 \cdot \Sigma^{-1}) - \frac{1}{2} N_2 \text{tr}(S_2 \cdot \Sigma^{-1}) + C$$

$$= -\frac{1}{2} (N \log |\Sigma| + N_1 \text{tr}(S_1 \cdot \Sigma^{-1}) + \frac{1}{2} N_2 \text{tr}(S_2 \cdot \Sigma^{-1})) + C$$

$$\frac{\partial L(\Sigma)}{\partial \Sigma} = N \Sigma^{-1} - N_1 S_1^T \Sigma^{-2} - N_2 S_2^T \Sigma^{-2} = 0 \quad \frac{\partial \text{tr}(AB)}{\partial A} = B^T$$

$$\frac{\partial |A|}{\partial A} = |A| \cdot A^{-1}$$