

1. 高斯分布

大釋然函數

最 MLE: $\theta_{MLE} = \arg \max_{\theta} P(X|\theta)$, \Rightarrow 令 $P=1$ (等角), $\cdot \theta = (\mu, \sigma^2)$

$$\log P(X|\theta) = \log \prod_{i=1}^N P(X_i|\theta) = \sum_{i=1}^N \log P(X_i|\theta)$$

Data: $x = (x_1, x_2, \dots, x_N)^T = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{pmatrix}$ $N \times p$.

$$\mu_{MLE} = \frac{1}{N} \sum_{i=1}^N x_i \quad (\text{均值})$$

$$\theta = (\mu, \Sigma) = (\mu, \sigma^2)$$

均值

$$E[\mu_{MLE}] = \frac{1}{N} \sum_{i=1}^N E[x_i] = \frac{1}{N} \sum_{i=1}^N \mu = \mu \Rightarrow \text{无偏的}$$

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

有偏估计

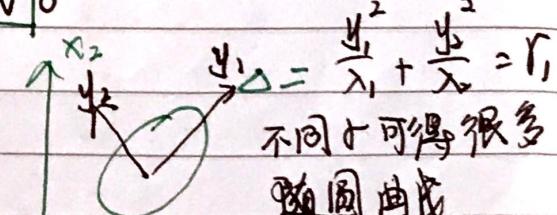
$$\text{真正的无偏: } \hat{\sigma}^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \mu_{MLE})^2$$

期望计算后等于本身 \Rightarrow 无偏估计

$$\begin{aligned} &= \frac{1}{N} \sum_{i=1}^N (x_i^2 - 2x_i \mu_{MLE} + \mu_{MLE}^2) \\ &= \frac{1}{N} \sum_{i=1}^N x_i^2 - \mu_{MLE}^2 \end{aligned}$$

$$E(\sigma_{MLE}^2) = \sigma^2 - \frac{1}{N} \sigma^2 = \frac{N-1}{N} \sigma^2 \rightarrow \text{所以有偏}$$

$$x \sim N(\mu, \Sigma) = \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp(-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu))$$



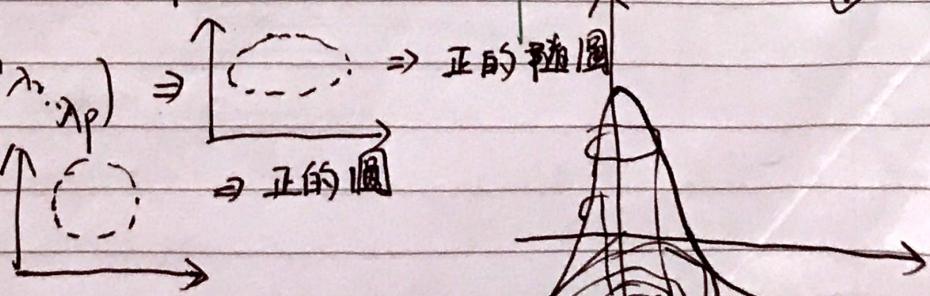
不同 Σ 可得很多椭圆曲线

Σ 是 $p \times p$ 維的

$$\Sigma_{p \times p} \rightarrow \frac{p^2 + p}{2} = O(p^2)$$

① $\Sigma \Rightarrow$ 对角矩阵 $(\lambda_1, \lambda_2, \dots, \lambda_p) \Rightarrow$ 正的椭圆

② 若 $\lambda_1 = \lambda_2 = \lambda_p$ \Rightarrow 正的圆



2. 高斯分布的局限性：有的数据不能统计为一个

3. 联合概率分布

$$X = \begin{pmatrix} X_a \\ X_b \end{pmatrix} \xrightarrow{\text{联合分布}} n \text{ 维} \quad m+n = p \quad \mu = \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix}$$

$$\Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}$$

① 边缘概率分布 $P(X_b | X_a)$
条件概率分布 $P(X_a | X_b)$

② 定理 已知 $X \sim N(\mu, \Sigma)$

$$Y = AX + B$$

$$\therefore Y \sim N(A\mu + B, A\Sigma A^T)$$

$$E[Y] = E[AX + B] = A\mu + B$$

$$\text{Var}[Y] = A\Sigma A^T \quad (\text{多维的})$$

③ $E[X_a] = (I_m \ 0) \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix} = \mu_a$

$$\text{Var}[X_a] = \Sigma_{aa}$$

$\therefore X_a \sim N(\mu_a, \Sigma_{aa})$

a 和 b 的关系

$$X_b | X_a = X_b - \Sigma_{ba} \Sigma_{aa}^{-1} X_a$$

$$\mu_{b|a} = \mu_b - \Sigma_{ba} \cdot \Sigma_{aa}^{-1} \mu_a$$

$$\Sigma_{b|a} = \Sigma_{bb} - \Sigma_{ba} \cdot \Sigma_{aa}^{-1} \Sigma_{ab}$$

$$\Rightarrow X_{b|a} = (-\Sigma_{ba} \Sigma_{aa}^{-1} \cdot I_n) \begin{pmatrix} X_a \\ X_b \end{pmatrix}$$

A

$$E[X_{b|a}] = \mu_b - \Sigma_{ba} \cdot \Sigma_{aa}^{-1} \mu_a = \mu_{b|a}$$

$$\text{Var}[X_{b|a}] = (-\Sigma_{ba} I_{aa} + I_n) \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix} \begin{pmatrix} -\Sigma_{aa}^{-1} \Sigma_{ba} \\ I \end{pmatrix}$$

$$= \Sigma_{b|a}$$

$$X_b = X_{b|a} + \Sigma_{ba} \Sigma_{aa}^{-1} X_a$$

$$E[X_b | X_a] = \mu_{b|a} + \Sigma_{ba} \cdot \Sigma_{aa}^{-1} X_a$$

$$\text{Var}[X_b | X_a] = \text{Var}[X_{b|a}] = \Sigma_{b|a}$$

$X_b | X_a \sim N(\mu_{b|a} + \Sigma_{ba} \Sigma_{aa}^{-1} X_a, \Sigma_{b|a})$

已知 $P(x) = N(x|\mu, \Sigma^{-1})$ 精度矩阵
 $P(y|x) = N(y|Ax+b, L^{-1})$

求: $P(y)$ $P(x|y)$

解: $y = Ax + b + \epsilon$

$\epsilon \sim N(0, L^{-1})$

$E[y] = E[Ax+b+\epsilon]$

$= A\mu + b + 0$

方差: $\text{Var}[y] = \text{Var}[Ax+b+\epsilon]$

$= A \cdot L^{-1} A^T + L^{-1}$

$\therefore y \sim N(A\mu+b; L^{-1} + A\Delta A^T)$

设: $z = \begin{pmatrix} x \\ y \end{pmatrix} \sim N\left(\begin{bmatrix} \mu \\ A\mu+b \end{bmatrix}, \begin{bmatrix} \Sigma & \Delta \\ \Delta & L^{-1} + A\Delta A^T \end{bmatrix}\right)$

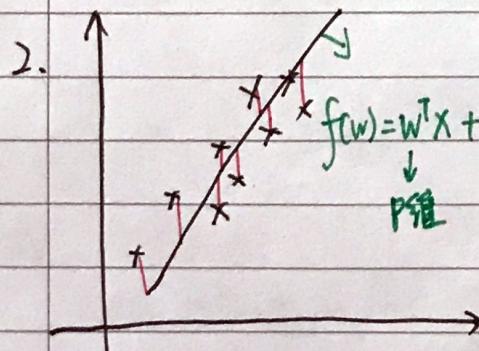
$\therefore P(x|y) \sim$ 代入求得.

前面求过的: $P(x_b|x_a) = N(\mu_{ba}; \Sigma_{bb|a})$

(三) 线性回归

1. 最小二乘法 (几何意义: 误差的和 // 投影
 概率: 噪声为高斯的 MLE)

正则化: $\begin{cases} L_1 \rightarrow \text{Lasso} \\ L_2 \rightarrow \text{Ridge, 岭回归} \end{cases}$



Data: $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$

$x_i \in \mathbb{R}^p, y_i \in \mathbb{R}, i=1, 2, \dots, N$

$$X = (x_1, x_2, \dots, x_N)^T = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{pmatrix} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{Np} \end{pmatrix}$$

$$Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}_{N \times 1}$$

最小二乘法:

$$L(w) = \sum_{i=1}^N \|w^T x_i - y_i\|^2 \rightarrow \text{Loss function}$$

$$= \sum_{i=1}^N (w^T x_i - y_i)^2 \Rightarrow \text{写成矩阵的形式}$$

向量 \times 向量的转置.

$$= (w^T x^T - y^T)(x w - y)$$

$$= w^T x^T x w - 2 w^T x^T y + y^T y$$

$$\hat{w} = \arg \min L(w)$$

$$\frac{\partial L(w)}{\partial w} = 2 x^T x w - 2 x^T y = 0 \text{ 对 } w \text{ 求导}$$

\therefore

$$x^T x w = x^T y$$

$$w = (x^T x)^{-1} x^T y$$

逆

② 概率视角: $LSE = MLE$ 前提
(噪声是高斯分布)

设 $\Sigma \sim N(0, \sigma^2)$

$$y = f(w) + \varepsilon$$

$$f(w) = w^T x \rightarrow \text{拟合}$$

$$p(y|x; w) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{(y - w^T x)^2}{2\sigma^2} \right\}$$

$$y = w^T x + \varepsilon$$

$\therefore y|x, w \sim N(w^T x, \sigma^2)$ 仍然符合高斯分布

MLE = 因为每个点都是独立的.

$$S(w) = \log P(Y|X, w) = \log \prod_{i=1}^N P(y_i|x_i; w)$$

= $\sum_{i=1}^N \log \dots$

$$= \sum_{i=1}^N \left(\log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} (y_i - w^T x_i)^2 \right)$$

$$\hat{w} = \arg \max_w S(w)$$

$$= \arg \max_w -\frac{1}{2\sigma^2} (y_i - w^T x_i)^2$$

$$= \arg \min_w (y_i - w^T x_i)^2 \rightarrow \text{就是最小二乘法.}$$

3. 正则化:

$$\text{Loss Function: } L(w) = \sum_{i=1}^N \| (w^T x_i - y_i) \|^2 \quad \hat{w} = (x^T x)^{-1} x^T y$$

$X_{N \times p}$ N 个样本 p 维 理论上 $N \gg p$

若 $N < p$, \Rightarrow 过拟合 \Rightarrow 可以有很多拟合方法

① 框架: 对 $L(w)$ 加约束

$$\arg \min_w [L(w) + \lambda P(w)] \quad L_1: \text{Lasso}, \quad P(w) = \|w\|$$

$L_2: \text{Ridge 岭回归} \quad P(w) = \|w\|^2 = w^T w$

$$\sum_{i=1}^N \|w^T x_i - y_i\|^2 + \lambda w^T w$$

$$= (w^T x^T - y^T)(w - y) + \lambda w^T w$$

$$= w^T (x^T x + \lambda I) w - 2 w^T x^T y + y^T y$$

$$\hat{w} = \arg \min_w J(w)$$

$$\frac{\partial J(w)}{\partial w} = 2(x^T x + \lambda I)w - 2x^T y = 0$$

$$\hat{w} = (x^T x + \lambda I)^{-1} x^T y$$

贝叶斯计算后: $P(w|y) = \frac{P(y|w) P(w)}{P(y)}$

$$\hat{w} = \arg \max_w P(w|y)$$

$$= \arg \min_w \sum_{i=1}^N (y_i - w^T x_i)^2 + \left(\frac{6^2}{G^2} \right) \|w\|^2$$

极大似然估计

一样

LSE \Leftrightarrow MLE (噪声为高斯分布)

Regularized LSE \Leftrightarrow MAP (噪声为高斯分布) 试验也是高斯分布
 ↓ 最大后验估计

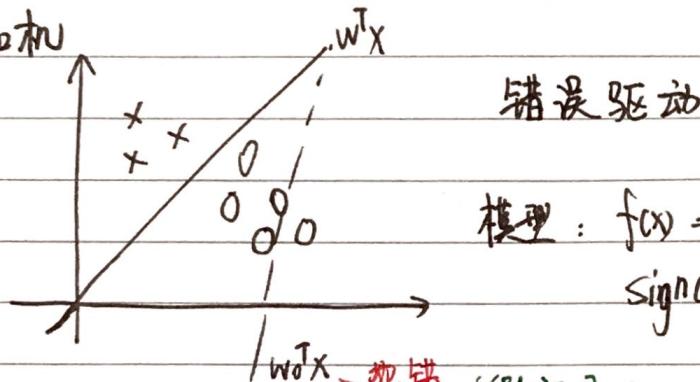
(四) 线性分类

硬分类 $\{y \in \{0, 1\}\}$
线性判别 fisher
感知机

软分类 $\{y \in [0, 1]\}$
判别式: Logistic Regression

线性回归 激活函数 $y = f(w^T x + b)$ f^{-1} : link function
降维 $\{0, 1\} \rightarrow w^T x + b$

① 感知机



模型: $f(x) = \text{sign}(w^T x)$ $x \in \mathbb{R}^P$, $w \in \mathbb{R}^R$
 $\text{sign}(a) = \begin{cases} +1 & a \geq 0 \\ -1 & a < 0 \end{cases}$

犯错 (假定是线性可分的)
向正确的移动

策略: Loss function

$$L(w) = \sum_{i=1}^N I\{y_i w^T x_i \leq 0\}$$

$w^T x_i > 0 \quad y_i = +1 \quad \left. \right\} \rightarrow \text{分类正确}$
 $w^T x_i < 0 \quad y_i = -1 \quad \left. \right\} \quad y_i w^T x_i > 0$

\Rightarrow 小于 0 则错误分类

$w \rightarrow w + \Delta w$ 整体的值无限趋近于 0

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

类间: 间:

$$(\bar{x}_1 - \bar{x}_2)^2$$

$$\text{梯度: } \nabla_w L = -y_i x_i;$$

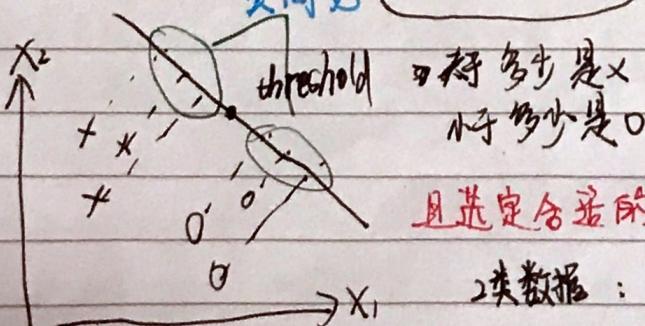
类内:

$$(S_1 + S_2)$$

$$w^{(t+1)} \leftarrow w^{(t)} - \lambda \nabla_w L$$

类内小
类间宽

线性判别分析



$$\text{投影: } w^T x = z; \\ \bar{z} = \frac{1}{N} \sum_{i=1}^N w^T x_i;$$

且选定合适的轴

$$2 \text{类数据: } C_1: \bar{z}_1 = \frac{1}{N_1} \sum_{i=1}^{N_1} w^T x_i;$$

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

$$C_2: \bar{z}_2 = \frac{1}{N_2} \sum_{i=1}^{N_2} w^T x_i;$$

$$S_2 = \frac{1}{N_2} \sum_{i=1}^{N_2} (w^T x_i - \bar{z}_2)(w^T x_i - \bar{z}_2)^T$$

类内小类间大：

$$\text{目标函数: } J(w) = \frac{(\bar{z}_1 - \bar{z}_2)^T w}{S_b + S_w} \uparrow \text{大} = \frac{[w^T (\bar{x}_{c_1} - \bar{x}_{c_2})]^T}{w^T (S_{c_1} + S_{c_2}) w}$$

$$\hat{w} = \operatorname{argmax}_w J(w)$$

$$= \frac{w^T (\bar{x}_{c_1} - \bar{x}_{c_2}) (\bar{x}_{c_1} - \bar{x}_{c_2})^T w}{w^T (S_{c_1} + S_{c_2}) w}$$

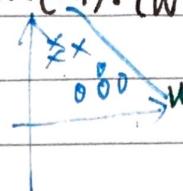
$$S_b: \text{between-class: 类间方差} \quad S_w: \text{within-class: 类内方差}$$

$$= \frac{w^T S_b w}{w^T S_w w}$$

求最优解：

$$\frac{\partial J(w)}{\partial w} = 2S_b w (w^T S_w w)^{-1} + w^T S_b w \cdot (-1) \cdot (w^T S_w w)^{-2} \cdot 2S_w \cdot w = 0$$

$$\therefore S_w \cdot w = \frac{w^T S_w w}{w^T S_b w} \cdot S_b \cdot w \quad \text{实数}$$



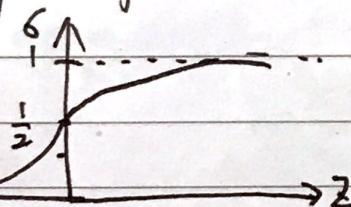
不关心 w 的大小，只关心方向

$$w = \frac{w^T S_w w}{w^T S_b w} S_b^{-1} \cdot S_b \cdot w \quad \propto S_b^{-1} (\bar{x}_{c_1} - \bar{x}_{c_2})$$

只有这部分与方向有关
如果 S_w^{-1} 是对角. 各向同性. 即 $S_w^{-1} \propto I$
则 $w \propto (\bar{x}_{c_1} - \bar{x}_{c_2})$

② 软输出：概率判别模型

$$\text{sigmoid function: } \sigma(z) = \frac{1}{1+e^{-z}}$$



$$\left. \begin{array}{l} z \rightarrow \infty, \lim \sigma(z) \xrightarrow{=} 1 \\ z \rightarrow 0 \quad \sigma(z) = \frac{1}{2} \\ z \rightarrow -\infty, \lim \sigma(z) \xrightarrow{=} 0 \end{array} \right\}$$

$\sigma: \mathbb{R} \rightarrow (0,1)$ 实数映射到 $(0,1)$ 区间

若 $w^T x \mapsto p$

$$p_1 = P(y=1|x) = \sigma(w^T x) = \frac{1}{1+e^{-w^T x}}$$

$$p_0 = P(y=0|x) = 1 - P(y=1|x) = \frac{e^{-w^T x}}{1+e^{-w^T x}}, \quad y=0$$

$$p(y|x) = p_1^y p_0^{1-y}$$

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

交叉熵: cross Entropy

$$\begin{aligned} \text{MLE: } \hat{w} &= \operatorname{argmax}_w \log p(y|x) = \operatorname{argmax}_w \log \prod_{i=1}^N p(y_i|x_i) = \operatorname{argmax}_w \log \prod_{i=1}^N p(y_i|x_i) \\ &= \operatorname{argmax}_w \sum_{i=1}^N \log p(y_i|x_i) = \operatorname{argmax}_w \sum_{i=1}^N (y_i \log p_i + (1-y_i) \log p_0) \\ &= \operatorname{argmax}_w \sum_{i=1}^N y_i \log \sigma(w^T x_i) + (1-y_i) \log (1-\sigma(w^T x_i)) \end{aligned}$$

-Cross Entropy

看哪个大
就属于哪
一类

Gaussian Discriminant Analysis (高斯)

④ 软输出：概率生成模型

Naive Bayes (高斯)

判别： $\hat{y} = \arg\max_{y \in \{0, 1\}} P(y|x)$ 目的 $P(y=0|x) \geq P(y=1|x)$
谁大就属于哪一类

生成： $P(y|x) = \frac{P(x|y) P(y)}{P(x)}$ (贝叶斯模型)

~~正比~~ $P(y=0)$

$P(y|x) \propto P(x|y) P(y)$
posterior likelihood prior

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

$y \sim \text{Bernoulli}(\phi) \Rightarrow \begin{array}{c|cc} y & 0 & 1 \\ \hline \phi & \phi & 1-\phi \end{array} \quad \begin{cases} \phi^y (1-\phi)^{1-y} & y=1 \\ (1-\phi)^y \phi^{1-y} & y=0 \end{cases}$

$x|y=1 \sim N(\mu_1, \Sigma)$
 $x|y=0 \sim N(\mu_2, \Sigma)$

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

$$\begin{aligned} L(\theta) &= \sum_{i=1}^N \log [P(x_i|y_i) \cdot P(y_i)] \\ &= \sum_{i=1}^N [\log P(x_i|y_i) + \log P(y_i)] \\ &= \sum_{i=1}^N [\log N(\mu_1, \Sigma)^{y_i} \cdot N(\mu_2, \Sigma)^{1-y_i} + \log \phi^{y_i} (1-\phi)^{1-y_i}] \\ &= \sum_{i=1}^N \underbrace{\log N(\mu_1, \Sigma)^{y_i}}_{\textcircled{1}} + \underbrace{\log N(\mu_2, \Sigma)^{1-y_i}}_{\textcircled{2}} + \underbrace{\log \phi^{y_i} (1-\phi)^{1-y_i}}_{\textcircled{3}} \end{aligned}$$

$y=1$	N_1
$y=0$	N_2
$N = N_1 + N_2$	

求 ϕ : $\textcircled{3} = \sum_{i=1}^N \log \phi^{y_i} (1-\phi)^{1-y_i}$

$\frac{\partial \textcircled{3}}{\partial \phi} = \sum_{i=1}^N y_i \cdot \frac{1}{\phi} - (1-y_i) \cdot \frac{1}{1-\phi} = 0$

$\Rightarrow \phi = \frac{1}{N} \sum_{i=1}^N y_i = \frac{N_1}{N}$ ($y_i=0$ 时约掉, 因为连加)

求 μ_1 : $\textcircled{1} = \sum_{i=1}^N \log N(\mu_1, \Sigma)^{y_i} = \sum_{i=1}^N y_i \log \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} \exp(-\frac{1}{2}(x_i - \mu_1)^T \Sigma^{-1} (x_i - \mu_1))$

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

$\frac{\partial \textcircled{1}}{\partial \mu_1} = \frac{\sum_{i=1}^N y_i x_i}{\sum_{i=1}^N y_i} = \frac{\sum_{i=1}^N y_i x_i}{N_1}$

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

$$\hat{\Sigma} = \arg \max_{\Sigma} \text{①} + \text{②}$$

$\sum_{i=1}^N y_i$

$$C_1 = \{x_i | y_i = 1, i=1 \dots N\}$$

$$C_2 = \{x_i | y_i = 0, i=1 \dots N\}$$

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

$$\text{①} + \text{②} = \sum_{x_i \in C_1} \log N(\mu_1, \Sigma) + \sum_{x_i \in C_2} \log N(\mu_2, \Sigma)$$

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

引入概念:

$$\text{tr}(AB) = \text{tr}(BA)$$

$$= \log \frac{1}{(2\pi)^{\frac{D}{2}}} + \log \frac{1}{|\Sigma|^{\frac{1}{2}}} + (-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu))$$

$$= \sum_{i=1}^N [C - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)]$$

\Rightarrow 实数

$$\text{tr}(ABC) = \text{tr}(ACB) = \text{tr}(BCA)$$

可以任意交换次序

$$\frac{\partial \text{tr}(AB)}{\partial A} = B^T$$

$$\frac{\partial A}{\partial A} = I A A^T$$

$$= C - \frac{1}{2} N \log |\Sigma| - \frac{1}{2} \text{tr}[(x_i - \mu) \cdot (x_i - \mu)^T \Sigma^{-1}]$$

$$= C - \frac{1}{2} N \log |\Sigma| - \frac{1}{2} \text{tr}(NS \cdot \Sigma^{-1})$$

$$S = \frac{1}{N} \sum (x_i - \mu)(x_i - \mu)^T$$

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

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

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

$$\therefore \frac{\partial \text{①} + \text{②}}{\partial \Sigma} = -\frac{1}{2} (N \Sigma^{-1} - N_1 S_1 \Sigma^{-2} N_2 S_2 \Sigma^{-2})$$

$$= 0$$

$$\therefore \hat{\Sigma} = \frac{1}{N} (N_1 S_1 + N_2 S_2)$$

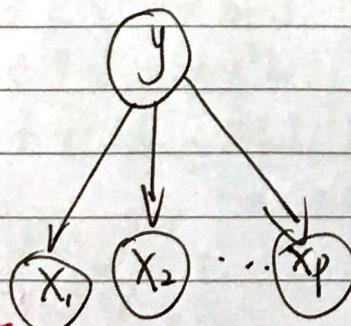
⑤ Naive Bayes 朴素贝叶斯：最简单的概率图模型

$$P(y|x) = \frac{P(x,y)}{P(x)}$$

$$= P(y) P(x|y)$$

$$P(x)$$

$$\propto P(y) P(x|y)$$



x_i 与 x_j 之间被 y 阻断

$x_i | x_j | y$ ($i \neq j$) 独立的表达

\times 表示 $x_j \sim \text{Categorical}$ 分布

$P(x_i|y) = \prod_{j=1}^D P(x_j|y)$ \times 表示 $x_j \sim N(\mu_j, \sigma_j^2)$

二分类 (0/1 分类) = y 属于 Bernoulli Dist

多分类 0, 1, ..., k = y 属于 Categorical Dist
做 n 次 \rightarrow 二项式 = Binomial

做 n 次 \rightarrow 多项式分布 Multinomial

目的：简化运算

$$\text{③ } x \in \mathbb{R}^p, y \in \{0, 1\}$$

给定 x , $y \in \{0, 1\}$

$$\hat{y} = \arg \max_y P(y|x)$$

$$= \arg \max_{y \in \{0, 1\}} \frac{P(x,y)}{P(x)}$$

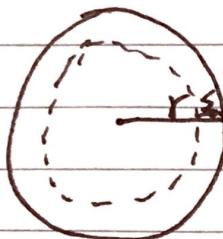
$$= \arg \max_y P(y) \cdot P(x|y)$$

(3) 降维

①

过拟合 $\left\{ \begin{array}{l} \text{正则化 (加约束 } L(W) \text{ 项)} \\ \text{降维. } \end{array} \right.$

直接降维: 特征选择
本质降维: PCA, 多维空间缩放
难度灾难: 非线性降维: 高维



$$V_{\text{球体}} = K \cdot r^D \quad D \text{ 维}$$

$$V_{\text{外}} = K \cdot r^D = K$$

$$\cancel{V_{\text{环}}} = V_{\text{环}} = K - K \cdot (r - \varepsilon)^D$$

$$\frac{V_{\text{环}}}{V_{\text{外}}} = \frac{K - K(r - \varepsilon)^D}{K} = 1 - (1 - \varepsilon)^D \quad 0 < \varepsilon < 1$$

$$\lim_{D \rightarrow \infty} (1 - \varepsilon)^D \rightarrow 0 \quad \therefore \text{维度越大, 比例} \rightarrow 0$$

②

$$\text{Data: } X = (X_1 \ X_2 \ \dots \ X_N)^T_{n \times p} = \begin{pmatrix} X_1^T \\ X_2^T \\ \vdots \\ X_N^T \end{pmatrix}_{n \times p} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & & & \\ x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix}_{n \times p}$$

$$\text{样本均值: } \bar{X}_{p,n} = \frac{1}{n} \sum_{i=1}^n X_i$$

$$\text{样本方差: } S_{p,p} = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})^T$$

$$\boxed{\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i} = \frac{1}{n} (X_1 \ X_2 \ \dots \ X_N) \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}_{n \times 1} = \boxed{\frac{1}{n} X^T I_n}$$

$$\begin{aligned} S &= \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})^T \\ &= \frac{1}{n} [X^T - \bar{X} I_n^T] \cdot (X - \bar{X})^T \\ &= \frac{1}{n} [X^T - \frac{1}{n} X^T I_n I_n^T] \cdot (I_n - \frac{1}{n} I_n I_n^T)^T \cdot X \\ &= \frac{1}{n} X^T (I_n - \frac{1}{n} I_n I_n^T) \cdot (I_n - \frac{1}{n} I_n I_n^T)^T \cdot X \\ &\quad H_N \rightarrow \text{centering matrix} \end{aligned}$$

$$= \frac{1}{n} X^T H \cdot H^T \cdot X$$

$$= \frac{1}{n} X^T H \cdot X$$

$$H = I_n - \frac{1}{n} I_n \cdot I_n^T$$

$$H^T = H \quad (\text{对称矩阵})$$

$$H^2 = H \cdot H = H$$

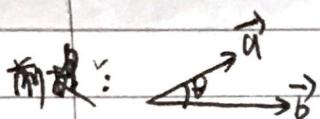
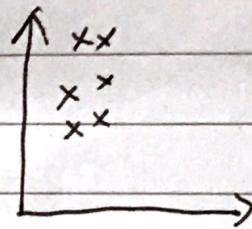
$$H^n = H$$

③ 一个中心: 原始特征空间的重构

相关 \rightarrow 无关

两个基本点: $\left\{ \begin{array}{l} \text{最大 投影方差} \rightarrow \text{使投影到线上的点尽量散开} \\ \text{最小重构距离} \end{array} \right.$

对数据的处理：① 中心化 $x_i - \bar{x}$



先移动

$$\vec{a} \cdot \vec{b} = |\vec{a}| |\vec{b}| \cos \theta \quad \text{且} |\vec{b}|=1$$

$$\text{则 } \vec{a} \text{ 的投影 } |\vec{a}| \cdot \cos \theta = \vec{a} \cdot \vec{b}$$

对应到矩阵为 $\vec{a}^T b$ () ()

均值

由于中心化 \rightarrow 前提 \Rightarrow 实数可以任意交换位置

$$\begin{aligned} \therefore \text{方差} &= \frac{1}{N} \sum_{i=1}^N ((x_i - \bar{x})^T u_1)^2 \quad \text{设 } u_1^T \cdot u_1 = 1 \\ &= \sum_{i=1}^N \frac{1}{N} u_1^T (x_i - \bar{x}) \cdot (x_i - \bar{x})^T u_1 \\ &= u_1^T \left(\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x}) \cdot (x_i - \bar{x})^T \right) u_1 \\ &= u_1^T S \cdot u_1 \end{aligned}$$

\therefore 最大投影方差 $\hat{u}_1 = \arg \max u_1^T S \cdot u_1$, 前提 $u_1^T \cdot u_1 = 1$

求解时, 设 $S(u_1, \lambda) = u_1^T S u_1 + \lambda (u_1^T u_1 - 1)$ 拉格朗日求解法

$$\frac{\partial S}{\partial u_1} = S \cdot 2u_1 - \lambda \cdot 2u_1 = 0$$

方差矩阵 特征值

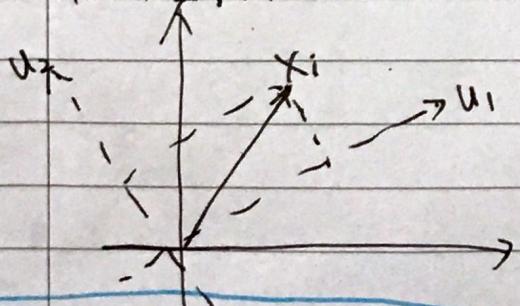
$$\therefore S u_1 = \lambda u_1$$

特征向量

④.

对坐标空间的重构：相当于重新定义坐标轴位置

最小重构代价（2个基本点的第2点）



在u轴上的坐标值

↑ 单位向量

$$x_i = (x_i^T \cdot u_1) \cdot u_1 + (x_i^T \cdot u_2) \cdot u_2 \quad (2\text{维})$$

总的就是个横向量

$$x_i = \sum_{k=1}^P (x_i^T \cdot u_k) \cdot u_k \quad (P\text{维})$$

$$\Rightarrow u_k = \arg \min_{k=1}^P u_k^T \cdot S \cdot u_k \text{ 最小重构代价}$$

s.t. $u_1^T u_k = 1$ (前提)

$$\text{降重后: } \hat{x}_i = \sum_{k=1}^q (x_i^T \cdot u_k) \cdot u_k$$

$$\begin{aligned} J &= \sum_{i=1}^N \|x_i - \hat{x}_i\|^2 \\ &= N \sum_{i=1}^N \left\| \sum_{k=q+1}^P (x_i^T \cdot u_k) \cdot u_k \right\|^2 \\ &= \frac{1}{N} \sum_{i=1}^N \sum_{k=q+1}^P (x_i^T \cdot u_k)^2 \quad \text{实数} \end{aligned}$$

考虑没有中心化

$$\hat{x}_i = \frac{1}{N} \sum_{i=1}^N \sum_{k=q+1}^P ((x_i - \bar{x})^T \cdot u_k)^2$$

$$= \sum_{k=q+1}^P u_k^T S u_k$$

主成分分析 ① SVD 角度

中心矩阵 $H \in \mathbb{R}^{n \times p}$ 奇异值分解
 $H = U \Sigma V^T \rightarrow SVD$ $\begin{cases} U^T U = I \\ V^T V = VV^T = I \end{cases} \rightarrow 正交的$
 $S = \frac{1}{\sqrt{n}} X^T H X$ Σ 对角

$$= X^T H \cdot H X = V \Sigma U^T \cdot U \Sigma V^T = V \Sigma^2 V^T$$

固定的 G, k

$$\text{之前有: } S = G K G^T \rightarrow G = V \quad K = \Sigma^2$$

$$\text{设 } T = H X X^T H = U \Sigma V^T \cdot V \Sigma U^T = V \Sigma^2 V^T$$

且 T 和 S 有相同的 特征 (eigen value)

且 $T = U \Sigma^2 V^T$ \rightarrow 特征向量组成的矩阵

$S = \text{特征分解: 得到方向 (主成分) } H \cdot X \cdot V \rightarrow \text{坐标}$

$T = \text{特征分解: 直接得到坐标. } \leftarrow$

特征值矩阵

主坐标分析 (principle coordinate analysis)

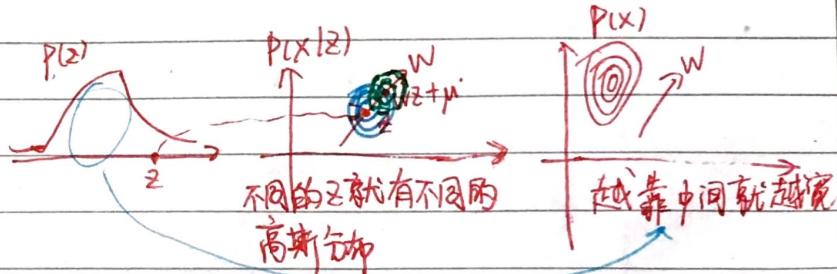
PCA

⑥ 极率角度. P-PCA

$$x \in \mathbb{R}^p \quad z \in \mathbb{R}^q \quad q < p$$

latent variable

observable data



$$z \sim N(0_q, I_q)$$

$$E[z] = E[wz + \mu + \Sigma] = w\bar{z} + \mu$$

$$X = wz + \mu + \Sigma$$

$$\text{Var}[z|z] = \text{Var}[wz + \mu + \Sigma] = w^2 I$$

$$\Sigma \sim N(0, \sigma^2 I_p) \text{ 噪声, 对角矩阵. } X|z = N(wz + \mu, \sigma^2 I)$$

$$\Sigma \perp z \text{ 独立的}$$

$$E[X] = E[wz + \mu + \Sigma] = \mu$$

$$\text{Inference: } P(z|x)$$

$$\text{Var}[X] = \text{Var}[wz + \mu + \Sigma] = \text{Var}[wz] + \text{Var}[\Sigma]$$

$$= w \cdot I \cdot w^T + \sigma^2 I = ww^T + \sigma^2 I$$

P-PCA

$$\text{Learning: } w, \mu, \sigma^2 \rightarrow \text{EM 算法} / \boxed{X \sim N(\mu, ww^T + \sigma^2 I)}$$

(求参数)

$$z \sim X|z. X \sim z|X$$

$$z|x = \begin{pmatrix} x_a \\ x_b \end{pmatrix} \sim N\left(\begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix}, \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}\right)$$

$$\Sigma_{b-a} = \Sigma_{bb} - \Sigma_{ba} \Sigma_{aa}^{-1} \Sigma_{ab}$$

$$\mu_{b-a} = \mu_b - \Sigma_{ba} \Sigma_{aa}^{-1} \mu_a$$

$$\Sigma_{bba} = \Sigma_{bb} - \Sigma_{ba} \Sigma_{aa}^{-1} \Sigma_{ab}$$

$$x_{b-a} \sim N(\mu_{b-a}, \Sigma_{bba})$$

$$x = x_{b-a} + \Sigma_{ba} \Sigma_{aa}^{-1} x_a$$

求 $z|x$ (高斯第6节内容)

$$\text{cov}(x, z) = E[(x - \mu)(z - \mu)^T] = E[wz + \mu z^T]$$

$$= E[wz z^T] = w E(z z^T) = w$$

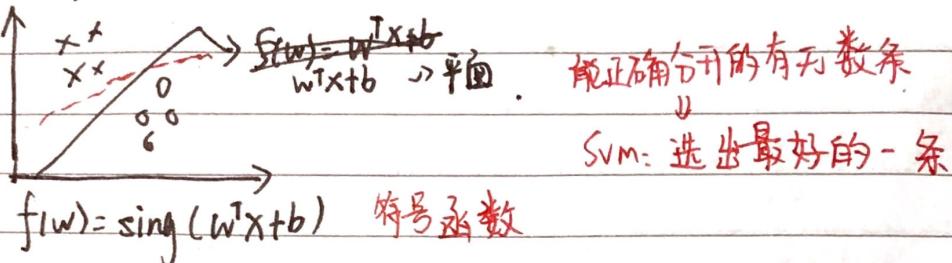
$$x_b | x_a \sim N(\mu_b + \dots, \Sigma_{bb|a})$$

$$\left. \begin{aligned} E[x_b | x_a] &= E[x_{b-a}] + \Sigma_{ba} \cdot \Sigma_{aa}^{-1} x_a = \mu_{b-a} + \Sigma_{ba} \cdot \Sigma_{aa}^{-1} x_a \\ &= \mu_b + \Sigma_{ba} \Sigma_{aa}^{-1} (\mu_a - \mu_b) \\ \text{Var}[x_b | x_a] &= \text{Var}[x_{b-a}] = \Sigma_{bb|a} \end{aligned} \right\}$$

(第六) 支持向量机 SVM 有三宝：间隔，对偶，核技巧

SVM $\left\{ \begin{array}{ll} \text{hard-margin SVM} & \text{硬间隔} \\ \text{soft-margin SVM} & \text{软间隔} \\ \text{kernel SVM} & \text{约束优化} \end{array} \right.$

SVM 最初是为了解决 2 分类问题



① 最大间隔分类器 $\max_{w,b} \text{margin}(w,b)$

$$\text{s.t. } \begin{cases} w^T x_i + b \geq 0 & y_i = +1 \\ w^T x_i + b \leq 0 & y_i = -1 \end{cases} \Rightarrow \begin{cases} y_i (w^T x_i + b) \geq 0 & i=1,2,\dots,N \\ y_i (w^T x_i + b) \leq 0 & \text{同上} \end{cases}$$

定 $\text{margin}(w,b) = \min_{x_i} \text{distance}(w,b, x_i) = \min_{x_i} \frac{1}{\|w\|} |w^T x_i + b|$
 (点到直线的距离)

$\text{distance} = \frac{1}{\|w\|} \cdot |w^T x_i + b|$

$$\therefore \max_{w,b} \text{margin}(w,b) = \max_{w,b} \min_{x_i} \frac{1}{\|w\|} |w^T x_i + b|$$

st. $y_i (w^T x_i + b) \geq 0$ 去掉绝对值: $y_i (w^T x_i + b)$

$$\rightarrow \exists r > 0, \quad \min_{x_i} y_i (w^T x_i + b) = r$$

$$= \max_{w,b} \frac{1}{\|w\|} \left[\min_{x_i} y_i (w^T x_i + b) \right]$$

响 → 变成 → 规化问题

$$\downarrow = \max_{w,b} \frac{1}{\|w\|} = \min_{w,b} \|w\| \quad \text{因为 } w^T x_i + b \text{ 直线可以以任意比例缩放}$$

st. $\min_{x_i} y_i (w^T x_i + b) = 1$
 可写成 $y_i (w^T x_i + b) \geq 1$

$$\Rightarrow \left\{ \begin{array}{l} \min_{w,b} \frac{1}{2} w^T w \rightarrow \text{二次} \\ \text{st. } y_i (w^T x_i + b) \geq 1 \text{ for all } i=1\dots N \end{array} \right.$$

N 个约束

接上面：

写成拉格朗日：

$$\mathcal{L}(w, b, \lambda) = \frac{1}{2} w^T w + \sum_{i=1}^N \lambda_i (1 - y_i (w^T x_i + b))$$

≥ 0 ≤ 0

问题转化为：
 (从带约束
无约束)

$$\begin{cases} \min_{w, b} \max_{\lambda} \mathcal{L}(w, b, \lambda) \\ \text{s.t. } \lambda_i \geq 0 \end{cases}$$

若： $1 - y_i (w^T x_i + b) > \max \rightarrow \infty$

若 $1 - y_i (w^T x_i + b) < \max \lambda_i$ 一定存在
 $\max = \frac{1}{2} w^T w + 0$

$$\min_{w, b} \left(\frac{1}{2} w^T w, \infty \right) = \min_{w, b} \frac{1}{2} w^T w$$

$\min \max \mathcal{L} \geq \max \min \mathcal{L}$ (宁为凤尾不为鸡头)

弱对偶关系

$$\min \max \mathcal{L} = \max \min \mathcal{L}$$

强对偶关系

对于 $\min \mathcal{L}(w, b, \lambda)$

$$\frac{\partial f}{\partial b} \triangleq 0 \Rightarrow \sum_{i=1}^N \lambda_i y_i = 0 \quad \text{将其代入 } \mathcal{L}(w, b, \lambda)$$

$$\mathcal{L}(w, b, \lambda) = \frac{1}{2} w^T w + \sum_{i=1}^N \lambda_i - \sum_{i=1}^N \lambda_i y_i w^T x_i$$

$$\frac{\partial f}{\partial w} = \frac{1}{2} \cdot 2 \cdot w - \sum_{i=1}^N \lambda_i y_i x_i \triangleq 0 \Rightarrow w = \sum_{i=1}^N \lambda_i y_i x_i \quad \text{将其代入 } \mathcal{L}(w, b, \lambda)$$

$$\text{有: } \mathcal{L}(w, b, \lambda) = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_i^T x_j + \sum_{i=1}^N \lambda_i$$

相当于求 \min (因为是最小值的最优解的值)

$$\text{再求: } \max -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_i^T x_j + \sum_{i=1}^N \lambda_i$$

$$\text{s.t. } \lambda_i \geq 0$$

$$\sum_{i=1}^N \lambda_i y_i = 0 \quad \text{也是一个向量}$$

$$\text{KKT 条件: } \frac{\partial \mathcal{L}}{\partial w} = 0 \quad \frac{\partial \mathcal{L}}{\partial b} = 0$$

原, 对偶问题具有强对偶关系

$$\text{是数据 data 相关 } \lambda_i (1 - y_i (w^T x_i + b)) = 0$$

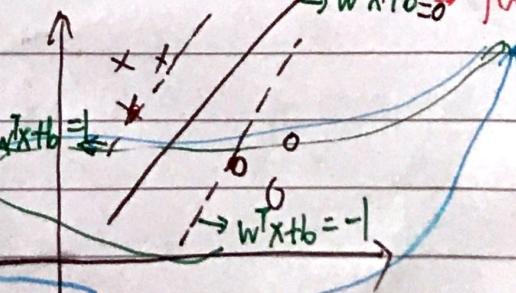
\Leftrightarrow 满足 KKT 条件 $w^T x + b = 0 \Rightarrow f(x) = \text{sign}(w^T x + b)$

不在 2 条线上
的点对 w 的贡献
都为 0

$$w^* = \sum_{i=1}^N \lambda_i y_i x_i$$

$$1 - y_i (w^* x_i + b) \leq 0$$

$$\text{计算 } b: \exists (x_k, y_k), \text{s.t. } 1 - y_k (w^T x_k + b) = 0$$



$$b^* = y_k - w^* x_k = y_k - \sum_{i=1}^N \lambda_i y_i x_i^T x_k$$

超平面

② soft-margin SVM 软间隔分类

$$\begin{cases} \min_{w,b} \frac{1}{2} w^T w \\ \text{s.t. } y_i(w^T x_i + b) \geq 1 \end{cases}$$

hard 是基于样本本身就是可分的，soft：数据不可分/或存在噪声
soft思想：允许一点错误

$$\min \frac{1}{2} w^T w + \text{loss}$$

$$\text{① loss} = \sum_{i=1}^N \{ y_i(w^T x_i + b) < 1 \}$$

但此时 w 不连续

$$\text{令 } z = y_i(w^T x_i + b)$$

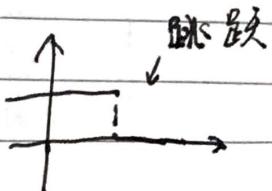
$$\text{② loss: 距离} \rightarrow \text{hinge loss}$$

$$\text{若 } y_i(w^T x_i + b) \geq 1, \text{ loss} = 0$$

$$\text{若 } \dots < 1, \text{ loss} = 1 - y_i(w^T x_i + b)$$

$$\text{loss} = \max \{ 0, 1 - y_i(w^T x_i + b) \} \rightarrow \text{函数本身是连续的}$$

$$\text{loss}_0/1 = \begin{cases} 1 & z < 1 \\ 0 & \text{otherwise} \end{cases}$$

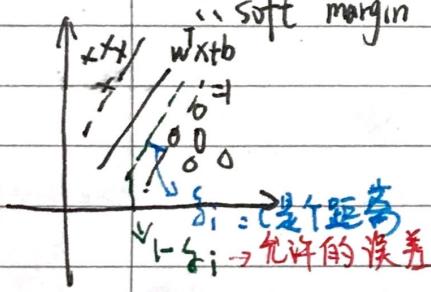


$$\therefore \text{soft margin: } \min_{w,b} \frac{1}{2} w^T w + C \sum_{i=1}^N \max \{ 0, 1 - y_i(w^T x_i + b) \}$$

一般不用 max:

$$\exists i \quad \delta_i = 1 - y_i(w^T x_i + b), \quad \delta_i \geq 0$$

$$\boxed{\begin{cases} \min_{w,b} \frac{1}{2} w^T w + C \sum_{i=1}^N \delta_i \\ \text{s.t. } y_i(w^T x_i + b) \geq 1 - \delta_i \end{cases}}$$



③ 约束优化问题

等价的

$$\text{拉格朗日函数: } \mathcal{L}(x, \lambda, \eta) = f(x) + \sum_{i=1}^m \lambda_i m_i + \sum_{j=1}^n \eta_j n_j$$

$$\left\{ \begin{array}{l} \min_x \max_{\eta} \mathcal{L}(x, \lambda, \eta) \Rightarrow \text{自动排除了不好的 } x \text{ 的集合} \quad (\text{即: } m_i(x) < 0) \end{array} \right.$$

$$\text{SVM: } \left\{ \begin{array}{l} \min_x \max_{\eta} \mathcal{L}(x, \lambda, \eta) \\ \text{s.t. } \lambda_i \geq 0 \end{array} \right.$$

(原问题的无约束

形式)

(原是关于 x 的函数)

$$\text{对偶问题: } \left\{ \begin{array}{l} \max_{\lambda, \eta} \min_x \mathcal{L}(x, \lambda, \eta) \\ \text{s.t. } \lambda_i \geq 0 \end{array} \right.$$

(对...是关于 λ, η 的

函数)

证明弱对偶性: $\max \min \mathcal{L} \leq \min \max \mathcal{L}$

$$\min_x \mathcal{L} \leq \mathcal{L}(x, \lambda, \eta) \leq \max_{\lambda, \eta} \mathcal{L}(x, \lambda, \eta)$$

$$A(\lambda, \eta) \leq B(x)$$

$$\max_{\lambda, \eta} A(\lambda, \eta) \leq \min_x B(x)$$

$$\text{BP: } \max_{\lambda, \eta} \min_x \mathcal{L} \leq \min_x \max_{\lambda, \eta} \mathcal{L}$$

④ 对偶性的几何解释.

简化问题 $\begin{cases} \min_{x \in D} f(x) \\ \text{s.t. } m_i(x) \leq 0 \end{cases}$ D: 定义域 $x \in D$

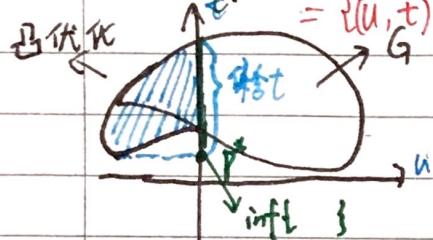
拉格朗日函数 $L(x, \lambda) = f(x) + \lambda m_i(x), \lambda \geq 0$

$p^* = \min f(x)$ 原问题最优解

$d^* = \max_{\lambda} \min_x f(x, \lambda)$ 对偶最优解

引入集合的定义:

$$G = \{(m_i(x), f(x)) \mid x \in D\} \quad \text{二维坐标系的一个区域}$$

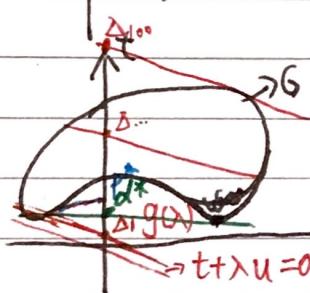


此时: $\begin{cases} p^* = \min f(x) = \min_t \\ \rightarrow p^* = \inf \{t \mid (u, t) \in G\}, u \leq 0 \end{cases}$

$d^* = \max_{\lambda} \min_x (t + \lambda u)$ \rightarrow 下确界(集合中) = min 的意思
 $\rightarrow g(\lambda) = \max_t (t + \lambda u)$

$= \max_{\lambda} g(\lambda)$

其中 $\boxed{g(\lambda) = \inf \{t + \lambda u \mid (u, t) \in G\}}$

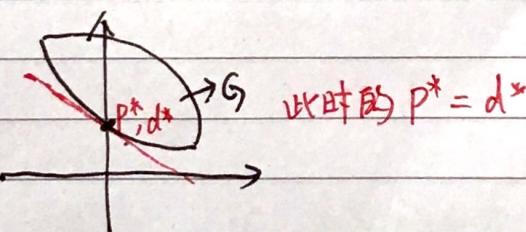


$d^* = \max_{\lambda} g(\lambda)$: 变动斜率, 使同时切到两边

此时的 $g(\lambda)$ 是 $\max g(x)$

有: $d^* \leq p^*$

若: G的集合为:



凸优化 + slater 条件 $\Rightarrow d^* = p^*$ (强对偶)

并不是充要. (充分非必要, 还是有其他情况推出强对偶)

⑤ slater 条件:

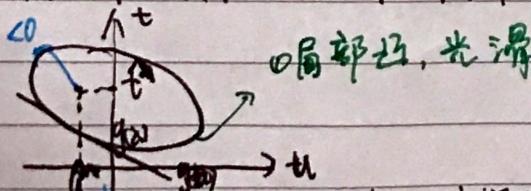
$\exists \hat{x} \in \text{relint } D$

relint: 相对内部 例: (D) 去掉边界的 D

s.t. $\forall i = 1 \dots m, m_i(\hat{x}) < 0$

I. 对于大多数凸优化, slater 成立

II. 放松 slater: 若 M 中有 k 个仿射函数, $m - k$



凸二次规划问题天然满足 II 成立

② 至少有元素在 $U \subset D \Rightarrow$ 使得有 $g(x)$ 这样的构成.

⑩ KKT 条件

$$p^* \rightarrow x^*$$

$$d^* \rightarrow \lambda^*, \eta^*$$

KKT 条件给出了 x^*, λ^*, η^* 的一种关系

可行条件: $\begin{cases} m_i(x^*) \leq 0 \\ n_j(x^*) = 0 \\ \lambda^* \geq 0 \end{cases}$ 必须满足

KKT: $\begin{cases} \text{互补松弛: } \lambda_i^* m_i = 0 \quad \forall i = 1, 2, \dots, n \\ \text{梯度为0: } \frac{\partial f(x, x^*, \eta^*)}{\partial x} \Big|_{x=x^*} = 0 \end{cases}$

II. 互补松弛:

$$d^* = \max_{\lambda, \eta} g(\lambda, \eta) = g(\lambda^*, \eta^*)$$

$$= \min_x S(x, x^*, \eta^*) \leq S(x, x^*, \eta^*) \quad \forall x \text{ 成立. 即 } x^* \text{ 也}$$

$$= f(x^*) + \sum_{i=1}^n \lambda_i^* m_i + \sum_{j=1}^m \eta_j^* n_j \quad \eta_j = 0 \rightarrow \text{只能取等}$$

$$\leq f(x^*)$$

$$= p^* \quad : \sum_i \lambda_i^* m_i = 0 \quad \forall i = 1, 2, \dots, n$$

$$\text{不等号} \leftarrow \sum_i \lambda_i^* \geq 0 \rightarrow m_i \leq 0$$

III. 梯度为0:

$$(x = x^* \text{ 时取得最小值})$$

$$\min S(x, \lambda^*, \eta^*) = S(x^*, \lambda^*, \eta^*)$$

$$\frac{\partial S}{\partial x} \Big|_{x=x^*} = 0$$

(1) 核方法 Kernel Function

线性可分 | 一点错误

PLA

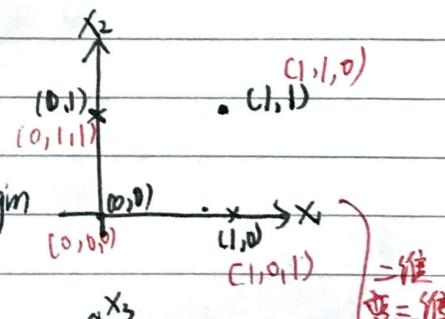
严格非线性

严格非线性

$$\phi(x) + PLA$$

Hard-Margin SVM | Soft-Margin SVM

$\phi(x) + \text{Hard-Margin}$



I. 高数据不能线性可分 (例如 螺圆等)

I. PLA \rightarrow 多层感知机 (神经网络) \rightarrow deep learning

II. 非线性可分问题 \rightarrow 转化为 线性可分问题

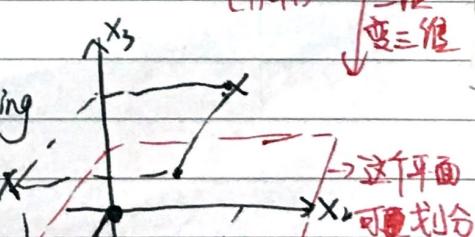
$\phi(x)$ = 非线性转换

$$\text{例: } X \xrightarrow{\phi(x)}$$

input space feature space

$$X = (x_1, x_2) = \text{二维}$$

$$(\phi(x))$$



Cover theorem: 高维比低维更易线性可分 $= (x_1, x_2, (x_1 - x_2)^2)$ 三维 只要满足需求

如何划分即

II. 对偶问题：对偶表示带来内积

Primal problem:

$$\begin{cases} \min_{w \in \mathbb{R}^N} \frac{1}{2} w^T w \\ \text{s.t. } y_i (w^T x_i + b) \geq 1 \end{cases} \quad (\text{N个约束})$$

Dual Problem

$$\begin{cases} \max_{\lambda \in \mathbb{R}^N} \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_i^T x_j - \frac{1}{2} \sum_i \lambda_i \\ \text{s.t. } \lambda_i \geq 0 \quad \text{for all } i=1, 2, \dots, N \\ \sum_{i=1}^N \lambda_i y_i = 0 \end{cases}$$

内积 变成 $\phi(x_i)^T \phi(x_j)$, 但 $\phi(x_i)$ 在现实情况下维度可能会非常高.

Kernel Function:

$$k(x, x') = \phi(x)^T \phi(x') = \langle \phi(x), \phi(x') \rangle$$

input space feature space

$\forall x, x' \in X$, $\exists \Phi: X \rightarrow Z$

$$\text{s.t. } k(x, x') = \phi(x)^T \phi(x')$$

则称 $k(x, x')$ 是一个核函数

例: 假设已经找到了一个核函数 $k(x, x') = \exp(-\frac{(x-x')^2}{2\sigma^2})$

则不用求 $\phi(x)$, 则可直接得到内积

② 正定核的定义

通常情况下的核函数为正定核函数

核: $K: X \times X \rightarrow \mathbb{R}$

$\forall x, z \in X$, 则称 $k(x, z)$ 为核函数

正定核: $K: X \times X \rightarrow \mathbb{R}$

$\forall x, z \in X$, 有 $k(x, z)$

若 $\exists \Phi: X \rightarrow \mathbb{R}$, $\Phi \in \mathcal{H}$

$$\text{s.t. } K(x, z) = \langle \Phi(x), \Phi(z) \rangle$$

那么则称 $K(x, z)$ 为正定核函数

正定核: $K: X \times X \rightarrow \mathbb{R}$

$\forall x, z \in X$, 有 $k(x, z)$

若 $K(x, z)$ 满足如下两条性质:

I. 对称性

II. 正定性

则称 $K(x, z)$ 为正定核

① 对称性 $\Leftrightarrow K(x, z) = K(z, x)$

② 正定性: \Leftrightarrow 任取 N 个元素

$x_1, x_2, \dots, x_N \in X$

对应的 Gram matrix 是半正定的

要证: $K(x, z) = \langle \Phi(x), \Phi(z) \rangle$

\Leftrightarrow Gram matrix 半正定

Hilbert Space: 完备的, 可能是无限维的, 被赋予内积的线性空间

向量空间(加法和数乘)

对称性, 正定性, 线性

$\langle fg \rangle = \langle g f \rangle$ $\langle f, f \rangle \geq 0$, $\langle f, f \rangle = 0 \Leftrightarrow f = 0$

$\langle r_1 f_1 + r_2 f_2, g \rangle = r_1 \langle f_1, g \rangle + r_2 \langle f_2, g \rangle$

$$\lim_{n \rightarrow \infty} K_n = K \quad \exists H$$

对极限是封闭的

$\lim_{n \rightarrow \infty} K_n = K$

$\exists H$

$\lim_{n \rightarrow \infty} K_n = K$

$\exists H$

$\lim_{n \rightarrow \infty} K_n = K$

$\exists H$

④ 正定核 - 必要性证明. (\Rightarrow)

已知 $K(x, z) = \langle \phi(x), \phi(z) \rangle$, 证: Gram matrix 半正定, 且 $K(x, z)$ 对称
 证: $K(x, z) = \langle \phi(x), \phi(z) \rangle$

$$K(z, x) = \langle \phi(z), \phi(x) \rangle$$

\because 内积运算本身具有对称性, 即 $\langle \phi(x), \phi(z) \rangle = \langle \phi(z), \phi(x) \rangle$

$$\therefore K(x, z) = K(z, x)$$

$\boxed{\therefore K(x, z) \text{ 满足对称性}}$

欲证 Gram Matrix 半正定. \rightarrow 方法: ① 特征值 ≥ 0

$$\text{Gram matrix: } K = [K(x_i, x_j)]_{N \times N}$$

$$\text{② } \forall \alpha \in \mathbb{R}^N, \alpha^T K \alpha \geq 0$$

即证. $\forall \forall \alpha \in \mathbb{R}^N, \alpha^T K \alpha \geq 0$

$$\alpha^T K \alpha = (\alpha_1, \alpha_2, \dots, \alpha_N) \begin{pmatrix} K_{11} & K_{12} & \dots & K_{1N} \\ K_{21} & K_{22} & \dots & K_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ K_{N1} & K_{N2} & \dots & K_{NN} \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_N \end{pmatrix}$$

$$= \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j K_{ij}$$

$$= \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j k(x_i, x_j) \quad \text{就是内积}$$

$$= \left(\sum_{i=1}^N \alpha_i \phi(x_i) \right)^T \cdot \left(\sum_{j=1}^N \alpha_j \phi(x_j) \right)$$

$$= \left[\sum_{i=1}^N \alpha_i \phi(x_i) \right]^T \cdot \left[\sum_{j=1}^N \alpha_j \phi(x_j) \right]$$

$$= \langle \sum_{i=1}^N \alpha_i \phi(x_i), \sum_{j=1}^N \alpha_j \phi(x_j) \rangle$$

$$= \left\| \sum_{i=1}^N \alpha_i \phi(x_i) \right\|^2 \geq 0$$

$\boxed{\therefore K \text{ 是半正定的}}$

\therefore 必要性得证

(1) 指数族分布

$$① P(x|\eta) = h(x) \exp(\eta^T \phi(x) - A(\eta))$$

↑ 成性

η : 参数(向量), $x \in \mathbb{R}^P$

$A(\eta) = \log \text{partition function}$

$\Rightarrow \log \text{配分函数}$

例如样本方差

$$\text{例: } P(x|\theta) = \frac{1}{Z} \hat{P}(x|\theta)$$

↑ 1-化因子

$$\int P(x|\theta) dx = \int \frac{1}{Z} \hat{P}(x|\theta) dx$$

$$= \frac{1}{Z} \int \hat{P}(x|\theta) dx$$

$$\therefore Z = \int \hat{P}(x|\theta) dx$$

$$P(x|\eta) = h(x) \cdot \exp(\eta^T \phi(x)) \exp(-A(\eta))$$

$$= \frac{\exp(A\eta)}{h(x)} h(x) \cdot \exp(\eta^T \phi(x))$$

$$= \frac{1}{Z} \cdot \hat{P}(x|\eta)$$

↑ 充分统计量

共轭

最大熵(无信息先验)

$$\therefore Z = \exp(A\eta) = Z$$

$$A(\eta) = \log Z$$

↑ 指数族分布

↑ 互推断

↑ 线性组合 $w^T x$ 义线性模型

概率图模型

link function \rightarrow (激活函数)

无向图 = RBM

II. 指数族分布: $y|x \sim$ 指数分布 link function \rightarrow (激活函数)

线性回归: $y|x \sim N(\mu, \sigma^2)$ 因为它是指数族, 所以有这个性质

共轭:

$y|x \sim \text{Bernoulli}$

$$P(z|x) \propto P(x|z) \cdot P(z)$$

等等

后验

共轭的效果: 先验和后验有相同的分布形式

$$\text{例: } \text{Beta}$$

$$\theta = \text{参数}$$

$$\text{Beta}$$

prior

① 共轭 \rightarrow 计算方便

② 取大商 \rightarrow 无信息先验

③ Jeffif

EM 算法：EM 期望最大

MLE : $P(x|\theta)$

$$\hat{\theta}_{MLE} = \arg \max_{\theta} P(x|\theta)$$

$$\rightarrow \arg \max_{\theta} \log P(x|\theta) \quad \text{log likelihood}$$

简化

① $\theta^{(t+1)} = \arg \max_{\theta} \int_Z \underbrace{\log P(x, z|\theta)}_{\text{完整数据}} \cdot \underbrace{P(z|x, \theta^{(t)})}_{\text{后验}} dz$ 证明是 θ^{t+1} 与 θ^t 之间的关系

② 会有 $\theta^{(t)} \rightarrow \theta^{(t+1)}$ $\rightarrow E_{z|x, \theta^{(t)}} [\log P(x, z|\theta)]$

$(\log P(x|\theta^{(t)}) \leq \log P(x|\theta^{(t+1)})$ 收敛性 证明这个等式：

✓ $\log P(x|\theta) = \log P(x, z|\theta) - \log P(z|x, \theta)$

同时求期望：左边 = $\int_Z P(z|x, \theta^{(t)}) \log P(x|\theta) dz$

= $\log P(x|\theta) \int_Z P(z|x, \theta^{(t)}) dz \rightarrow$ 只跟 θ 相关，所以积分=1

= $\log P(x|\theta)$

右边 = $\int_Z \underbrace{P(z|x, \theta^{(t)})}_{Q(\theta, \theta^{(t)})} \cdot \log P(x, z|\theta) dz - \int_Z P(z|x, \theta^{(t)}) \cdot \log P(z|x, \theta) dz$ H(\theta, \theta^{(t)})

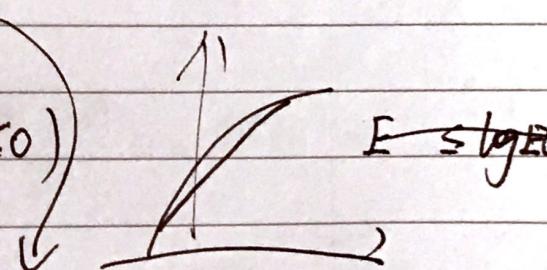
要证： $\left\{ \begin{array}{l} Q(\theta^{(t+1)}, \theta^{(t)}) \geq Q(\theta, \theta^{(t)}) \\ H(\theta^{(t+1)}, \theta^{(t)}) \leq H(\theta^{(t)}, \theta^{(t)}) \end{array} \right.$

$H(\theta^{(t+1)}, \theta^{(t)}) - H(\theta^{(t)}, \theta^{(t)})$

= $\int_Z P(z|x, \theta^{(t)}) \log P(z|x, \theta^{(t+1)}) dz - \int_Z P(z|x, \theta^{(t)}) \log P(z|x, \theta^{(t)}) dz$

= $\int_Z P(z|x, \theta^{(t)}) \log \frac{P(z|x, \theta^{(t+1)})}{P(z|x, \theta^{(t)})} dz$

$(= -KL(P(z|x, \theta^{(t)}) || P(z|x, \theta^{(t+1)})) \leq 0)$



$E[\log x] \leq \log E(x)$

$\therefore \log \int_Z P(z|x, \theta^{(t+1)}) dz = \log$

$= 0$

证明收敛性

② 证明公式从何

基础篇

② 证明公式从何而来：

$$\begin{aligned}\log P(x|\theta) &= \log P(x, z|\theta) - \log P(z|x, \theta) \\ &= \log \frac{P(x, z|\theta)}{q(z)} - \log \frac{P(z|x, \theta)}{q(z)}\end{aligned}$$

(log的一减一加后恒等)

乘两边求积分：

$$\text{左边} = \int_Z q(z) \cdot \log P(x|\theta) dz = \log P(x|\theta) \int_Z q(z) \cdot dz = \log P(x|\theta)$$

$$\text{右边} = \int_Z q(z) \log \frac{P(x, z|\theta)}{q(z)} dz = \int_Z q(z) \cdot \log \frac{P(z|x, \theta)}{q(z)} dz$$

$\underbrace{\text{ELBO}}_{\text{Evidence lower bound}} \leq \text{对数熵} : \text{KL}(q(z)||P(z|x, \theta))$

$$\log P(x|\theta) = \text{ELBO} + \text{KL}(q||P)$$

后验
 ≥ 0

$\therefore \log P(x|\theta) \geq \text{ELBO}$ [当 q 与 P 相等时，取等号]

$$\hat{\theta} = \arg \max_{\theta} \text{ELBO}$$

$$= \arg \max_{\theta} \int_Z q(z) \log \frac{P(x, z|\theta)}{q(z)} dz \quad (\text{取等号时}, q(z) = P(z|x, \theta^{(+)})$$

$$= \arg \max_{\theta} \int_Z q(z|x, \theta^{(+)}) \log \frac{P(x, z|\theta)}{P(z|x, \theta^{(+)})} dz$$

\rightarrow 与无关 ($\theta^{(+)}$ 是常数)

$$= \arg \max_{\theta} \int_Z P(z|x, \theta^{(+)}) \log P(x, z|\theta) dz \rightarrow \text{EM 公式}$$

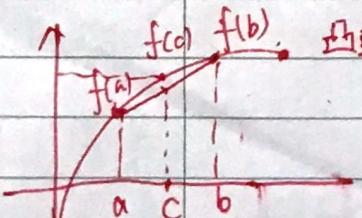
③ 从另一个角度推公式

X : observed data $\cdot X = \{x_1, x_2, \dots, x_N\}$

$$\log P(x|\theta) = \log \int_Z P(x, z|\theta) dz \xrightarrow{\text{直接引出}, \text{所以求积分}}$$

$$= \log \int_Z \frac{P(x, z|\theta)}{q(z)} \cdot q(z) dz \quad (\text{第一讲: 积分可理解为求期望})$$

Jensen 不等式:



$$= \log \mathbb{E}_{q(z)} \left[\frac{P(x, z|\theta)}{q(z)} \right]$$

$\geq \mathbb{E}_{q(z)} \left[\log \frac{P(x, z|\theta)}{q(z)} \right]$ 在 $\frac{P(x, z|\theta)}{q(z)} = C$ 时取“=”

ELBO

$$q(z) = \frac{1}{C} P(x, z|\theta)$$

$$= \int_Z \frac{1}{C} P(x, z|\theta) dz$$

$$\text{因为是 } \frac{1}{C} \text{ 的 } \int_Z \text{ 的话} = \frac{1}{C} \int_Z P(x, z|\theta) dz$$

$$= \frac{1}{C} P(x|\theta) = 1$$

$$P(x|\theta) = C$$

$$\therefore q(z) = \frac{1}{P(x|\theta)} \cdot P(x, z|\theta) = P(z|x, \theta)$$

后验

HMM
GMM \rightarrow 模型
EM \rightarrow 算法，解决估计问题
 \rightarrow 一个迭代算法

① 从狭义 EM \rightarrow 广义 EM

② ~~key~~ 广义 EM 是 EM 的一种

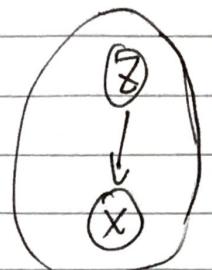
③ 真正的 EM, EM 的变种, 问题等

设：模型：概率生成模型

Data: X

Latent Variable $\rightarrow Z$

θ : model parameter



$P(X|\theta)$

EM是用于估计参数 θ 的

$$\hat{\theta} = \arg \max \theta P(X|\theta) \quad P(X) \text{ 未知, 所以才有归一化步骤}$$

$$= \arg \max \log P(X|\theta) \quad \rightarrow \text{引入 } P(X, Z)$$

为了方便起见

$$P(X) = \int_Z P(X, Z) dZ \quad \text{互积律}$$

$$P(X) = \frac{P(X, Z)}{P(Z, X)}$$

$$\text{其中 } \log P(X|\theta) = \text{ELBO} + \text{KL}(q||p) \geq S(q, \theta)$$

$$\text{ELBO} = \mathbb{E}_{q(z)} [\log \frac{P(X, z|\theta)}{q(z)}] \quad \log P(X) = \log \frac{P(X, z)}{P(z, X)}$$

$$\text{KL}(q||p) = \int q(z) \cdot \log \frac{q(z)}{p(z|x, \theta)} dz \quad \text{取等号} \quad = \text{KL}(z||p) = \log P(x, z) - \log P(z, x)$$

但有时 $P(z|x, \theta)$ 有时无法计算出

广义 EM:

E-step 固定时, $q \rightarrow p$, KL 越小

左右两边同时求期望:

\Leftrightarrow ELBO 越大

$$\text{左边: } \mathbb{E}_{q(z)} [\log P(x)] = \int q(z) \log P(x) dz$$

$$= \log P(x)$$

M-step 固定 q 时, $\theta = \arg \max_{\theta} S(q, \theta)$

$$\text{右边: } \int q(z) \cdot \log \frac{P(x, z)}{q(z)} dz - \int q(z) \log \frac{P(z, x)}{q(z)} dz$$

ELBO

$\text{KL}(q||p)$

更规范化:

广义 EM:

$$\text{E-step: } q^{(t+1)} = \arg \max_q S(q, \theta^{(t)})$$

$$\text{M-step: } \theta^{(t+1)} = \arg \max_{\theta} S(q^{(t+1)}, \theta)$$

$$S(q, \theta) = \mathbb{E}_q [\log P(x, z) - \log q]$$

$$\text{简化为: } \mathbb{E}_q [\log P(x, z)] - \mathbb{E}_q [\log q]$$

$$= \mathbb{E}_q [\log P(x, z|\theta)] + H(q)$$

坐标上升法:

\Downarrow M-E 的本质无所谓

$$\text{熵: } \int q(z) \cdot \frac{1}{\log q(z)} dz$$

梯度上升法

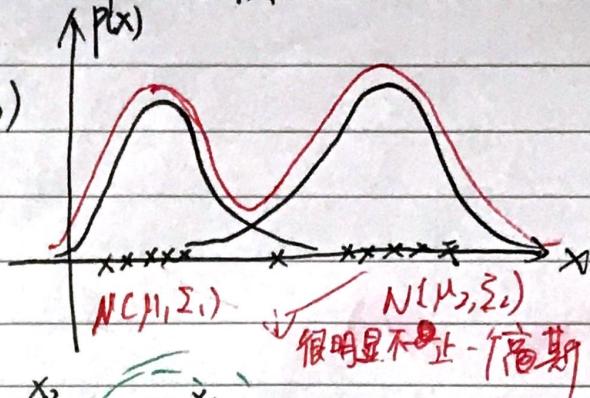
(+) GMM : Gaussian mixture Model 高斯混合模型.

① I. 从几何角度来看:

关于多个高斯的加权平均 (叠加)

$$p(x) = \sum_{k=1}^K \omega_k N(\mu_k, \Sigma_k), \sum_{k=1}^K \omega_k = 1$$

权重



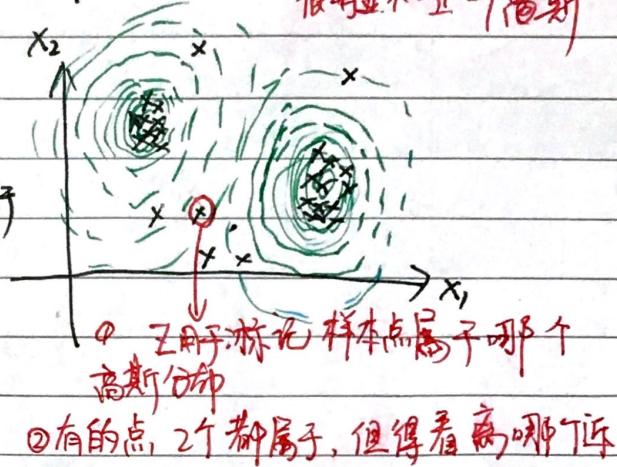
② II. 从混合模型角度来看: (生成模型)

x : observed variable

引进 z : latent variable: 对应的样本 x 是属于
离散的随机变量 | 那一个高斯分布.

z	1	2	\cdots	K
-----	---	---	----------	-----

$$P(z) | p_1, p_2, \dots, p_K \quad \sum_{k=1}^K p_k = 1$$



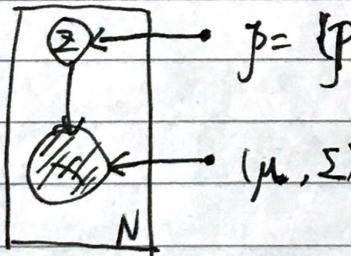
③ I. 假设有一个骰子, 有六个面 (例子) 每多个面都有概率, 概率更大的更重

$x \sim C_1$ (概率更大)

C_2 (概率更小)

$x \sim Z$.

II. 概率图:



$$p = (p_1, p_2, \dots, p_K) \quad (p \text{ 是 } z \text{ 的参数})$$

$$\begin{aligned} p(x) &= \sum_{k=1}^K p(x|z) \\ &= \sum_{k=1}^K p(x, z=k) \\ &= \sum_{k=1}^K p(z=k) \cdot p(x|z=k) \end{aligned}$$

$$= \sum_{k=1}^K p_k \cdot N(x|\mu_k, \Sigma_k)$$

↑ 概率

$$\theta = \text{parameter} \rightarrow \theta = \{p_1, p_2, \dots, p_N, \mu_1, \mu_2, \dots, \mu_N, \Sigma_1, \Sigma_2, \dots, \Sigma_K\}$$

$$\hat{\theta}_{MLE} = \arg \max_{\theta} \log P(x) \rightarrow \text{极大似然估计方法求得}$$

② 用EM求解GMM

$$EM: \theta^{(t+1)} = \arg \max_{\theta} E_{z|x, \theta^{(t)}} [\log P(x, z|\theta)]$$

$\checkmark Q\text{函数 } Q(\theta, \theta^{(t)})$

$$Q(\theta, \theta^{(t)}) = \int_z \log P(x, z|\theta) \cdot P(z|x, \theta^{(t)}) dz$$

这是高斯的

E-step

$$= \sum_i^N \underbrace{\log \prod_{j=1}^N P(x_j, z_j|\theta)}_{P(x, z|\theta)} \cdot \prod_{j=1}^N P(z_j|x_j, \theta^{(t)})$$

$$= \sum_{z_1, z_2, \dots, z_N} \log P(x_1, z_1|\theta) \prod_{i=1}^N P(z_i|x_i, \theta^{(t)})$$

$$= \sum_{z_1, z_2, \dots, z_N} [\log P(x_1, z_1|\theta) + \log P(x_2, z_2|\theta) + \dots + \log P(x_N, z_N|\theta)] \prod_{i=1}^N P(z_i|x_i, \theta^{(t)})$$

$$\textcircled{2} \quad \therefore \sum_{z_1, z_2, \dots, z_N} \log P(x_1, z_1|\theta) \prod_{i=1}^N P(z_i|x_i, \theta^{(t)})$$

$$= \sum_{z_1} \log P(x_1, z_1|\theta) P(z_1|x_1, \theta^{(t)})$$

$$\therefore \text{原式} = \sum_{z_1} \log P(x_1, z_1|\theta) \cdot P(z_1|x_1, \theta^{(t)}) + \dots + \sum_{z_N} \log P(x_N, z_N|\theta) P(z_N|x_N, \theta^{(t)})$$

$$= \sum_{k=1}^K \sum_{z_i} \log P(x_i, z_i|\theta) P(z_i|x_i, \theta^{(t)})$$

$$P(x, z) = P(z) P(x|z)$$

$$= P_z \cdot N(x|\mu_z, \Sigma_z)$$

$$P(z|x) = \frac{P(x, z)}{P(x)} = \frac{P_z N(x|\mu_z, \Sigma_z)}{\sum_{k=1}^K P_k N(x|\mu_k, \Sigma_k)}$$

$$\begin{aligned} &= \sum_{i=1}^N \sum_{z_i} \log P_{z_i} N(x_i|\mu_{z_i}, \Sigma_{z_i}) \\ &\cdot P_{z_i} N(x|\mu_{z_i}, \Sigma_{z_i}) \\ &\cdot \sum_{k=1}^K P_k N(x_i|\mu_k, \Sigma_k) \\ &= \sum_{k=1}^K \sum_{i=1}^N [\log P_k + \log N(x_i|\mu_k, \Sigma_k)] P_{z_i} = C_k | x_i, \theta^{(t)} \end{aligned}$$

$$\textcircled{3} \quad M\text{-step} \quad \theta^{(t+1)} = \arg \max_{\theta} Q(\theta, \theta^{(t)})$$

$$\text{求 } P_k^{(t+1)} = (P_1^{(t+1)}, P_2^{(t+1)}, \dots, P_K^{(t+1)})$$

$$P_k^{(t+1)} = \arg \max_{P_k} \sum_{k=1}^K \sum_{i=1}^N \log P_k \cdot P(z_i=c_k|x_i, \theta^{(t)}) \quad \text{s.t. } \sum_{k=1}^K P_k = 1$$

$$\sum_{k=1}^K \sum_{i=1}^N \log P_k \cdot P(z_i=c_k|x_i, \theta^{(t)}) + \lambda (\sum_{k=1}^K P_k - 1)$$

$$\frac{\partial \mathcal{L}}{\partial P_k} = \sum_{i=1}^N \frac{1}{P_k} \cdot P(z_i=c_k|x_i, \theta^{(t)}) + \lambda \triangleq 0$$

$$\Rightarrow \sum_{i=1}^N P(z_i=c_k|x_i, \theta^{(t)}) + P_k \cdot \lambda = 0$$

$$\Rightarrow \sum_{i=1}^N \sum_{k=1}^K P(z_i=c_k|x_i, \theta^{(t)}) + \sum_{k=1}^K P_k \cdot \lambda = 0$$

$$\Rightarrow N + \lambda = 0 \quad \therefore \lambda = -N$$

$$\therefore P_k^{(t+1)} = \frac{1}{N} \sum_{i=1}^N P(z_i=c_k|x_i, \theta^{(t)})$$

$$P^{(t+1)} = (P_1^{(t+1)}, P_2^{(t+1)}, \dots, P_N^{(t+1)})$$