

HW9+10

证

HW9

$$1. \text{err}^*(x) = 1 - \max_{c \in Y} P(c|x)$$

$$\text{err}(x) = 1 - \sum_c P(c|x) P(c|\bar{x})$$

$$|Y| \times \max_{c \in Y} P(c|x) \geq \sum_c P(c|x)$$

$$\therefore \sum_c P(c|\bar{x}) \leq 1 \quad \therefore \text{err}^*(x) \leq \text{err}(x)$$

$$\text{下面证明: } \text{err}(x) \leq \text{err}^*(x) \left( 2 - \frac{|Y|}{|Y|-1} \times \text{err}^*(x) \right)$$

$$\text{err}(x) = 1 - \sum_c P(c|x) P(c|\bar{x}) = 1 - \sum_c P^2(c|x)$$

$$\cdot \text{由于 } \left( \sum_{c \in Y} P(c|x) \right)^2 \leq |Y| \left( \sum_{c \in Y} P^2(c|x) \right)$$

$$\text{即有 } \text{err}(x) = 1 - \sum_{c \in Y} P^2(c|x) \leq 1 - \frac{1}{|Y|} \left( \sum_{c \in Y} P(c|x) \right)^2$$

$$\therefore \max_{c \in Y} P(c|x) \geq \frac{1}{|Y|} \sum_{c \in Y} P(c|x)$$

$$\therefore \text{err}^*(x) \leq 1 - \frac{1}{|Y|-1} \sum_{c \in Y} P(c|x)$$

$$\text{err}^*(x) \left[ 2 - \frac{|Y|}{|Y|-1} \text{err}^*(x) \right] = 2 - 2 \max_{c \in Y} P(c|x) - \frac{|Y|}{|Y|-1} [\text{err}^*(x)]^2$$

$$\geq 2 - 2 \max_{c \in Y} P(c|x) - \frac{|Y|}{|Y|-1} + \frac{2|Y|}{|Y|-1} \max_{c \in Y} P(c|x) - \frac{1}{|Y|(|Y|-1)} \left( \sum_{c \in Y} P(c|x) \right)^2$$

$$= \frac{|Y|-2}{|Y|-1} - \frac{1}{|Y|(|Y|-1)} \left( \sum_{c \in Y} P(c|x) \right)^2 + \frac{2}{|Y|-1} \max_{c \in Y} P(c|x)$$

$$\begin{aligned}
 & \text{即有 } \text{err}^*(x) \left[ 2 - \frac{|y|}{|y|-1} \text{err}^*(x) \right] \\
 & \geq \frac{|y|-2}{|y|-1} - \frac{1}{|y|(|y|-1)} \left( \sum_{c \in y} P(c|x) \right)^2 + \frac{2}{|y|-1} \max_{c \in y} P(c|x) \\
 & \geq \frac{|y|-2}{|y|-1} - \frac{1}{|y|(|y|-1)} + \frac{2}{|y|-1} \times \frac{1}{|y|} \\
 & = \frac{|y|-2}{|y|-1} + \frac{1}{|y|(|y|-1)} \\
 & \geq 1 - \frac{1}{|y|} \left( \sum_{c \in y} P(c|x) \right)^2 \geq \text{err}(x) \\
 & \text{综上, 证得 } \text{err}^*(x) \leq \text{err}(x) \leq \text{err}^*(x) \left( 2 - \frac{|y|}{|y|-1} \text{err}^*(x) \right)
 \end{aligned}$$

10.4. ~~因为奇异值分解后直接就得到了最大的特征值,~~  
~~不需要先求所有特征值再筛选。~~

$$\begin{aligned}
 & \text{对 } X \text{ 作奇异值分解得 } X = U \Sigma V^T, \quad XX^T = U \Sigma V^T (U \Sigma V^T)^T \\
 & \quad \quad \quad = U \Sigma V^T V \Sigma^T U^T \\
 & \quad \quad \quad = U \Sigma \Sigma^T U^T = U \Lambda U^T
 \end{aligned}$$

$X$  的奇异值分解与  $XX^T$  特征分解的效果等价。

但对于存储高维数据的  $X$ , 如  $X \in \mathbb{R}^{10000 \times 100}$

对  $XX^T \in \mathbb{R}^{10000 \times 10000}$  特征分解的计算成本显然更高。

且  $X$  的奇异值分解的计算精度更高。



3. 求解:  $\max_W \text{tr}(W^T X X^T W)$  s.t.  $W^T W = I_{d'}$

$$\Rightarrow \min_W -\text{tr}(W^T X X^T W) \quad \text{s.t. } W^T W = I_{d'}$$

其中  $X = (x_1, \dots, x_m) \in \mathbb{R}^{d \times m}$   $W = (w_1, w_2, \dots, w_{d'}) \in \mathbb{R}^{d \times d'}$   $I_{d'} \in \mathbb{R}^{d' \times d'}$

令  $\theta \in \mathbb{R}^{d' \times d'}$  为拉格朗日乘子矩阵.

$$L(W, \theta) = -\text{tr}(W^T X X^T W) + \langle \theta, W^T W - I \rangle$$

$$= -\text{tr}(W^T X X^T W) + \text{tr}(\theta(W^T W - I))$$

若仅考虑约束  $w_i^T w_i = 1 \quad (i=1, 2, \dots, d')$

此时  $\theta$  为对角矩阵, 令新的拉格朗日乘子矩阵为  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_{d'}) \in \mathbb{R}^{d' \times d'}$

$$L(W, \Lambda) = -\text{tr}(W^T X X^T W) + \text{tr}(\Lambda^T (W^T W - I))$$

$$\frac{\partial L(W, \Lambda)}{\partial W} = -\frac{\partial}{\partial W} \text{tr}(W^T X X^T W) + \frac{\partial}{\partial W} \text{tr}(\Lambda^T (W^T W - I))$$

$$= -X X^T W + X X^T W + W \Lambda + W \Lambda^T$$

$$= -2X X^T W + 2W \Lambda$$

$$\text{令 } \frac{\partial L(W, \Lambda)}{\partial W} = 0 \text{ 得 } X X^T W = W \Lambda$$

$$\text{显然 } \lambda_1, \dots, \lambda_{d'} \text{ 为 } X X^T \text{ 的特征值,}$$

$w_1, \dots, w_{d'}$  为  $\lambda_1, \dots, \lambda_{d'}$  对应单位特征向量

又  $w_i^T w_j = 0 \quad (i \neq j)$  为原问题的约束.

且当  $\lambda_i \neq \lambda_j$  时,  $w_i$  与  $w_j$  正交, 即

满足原约束条件.

从  $XX^T$  的  $d$  个特征向量中找出  $d'$  个能使得目标函数达到最优值的特征向量作为最优解。

将  $XX^T w_i = \lambda_i w_i$  代入  $\min_W -\text{tr}(W^T XX^T W)$

$$\min_W -\text{tr}(W^T XX^T W) = \max_W \text{tr}(W^T XX^T W)$$

$$= \max_W \sum_{i=1}^{d'} w_i^T XX^T w_i$$

$$= \max_W \sum_{i=1}^{d'} w_i^T \lambda_i w_i$$

$$= \max_W \sum_{i=1}^{d'} \lambda_i w_i^T w_i$$

$$= \max_W \sum_{i=1}^{d'} \lambda_i$$

此时只需令  $\lambda_1, \dots, \lambda_{d'}$  为  $XX^T$  的前  $d'$  个最大特征值,

$w_1, \dots, w_{d'}$  为  $\lambda_1, \dots, \lambda_{d'}$  分别对应的单位特征向量

就可使目标函数取到最优值。

HW10

11.5, 11.7

11.5  $L_1$ 正则化可以产生稀疏解, 是因为平方误差项等值线与 $L_1$ 等值线的第一个交点位于坐标轴上。当平方误差项等值线的曲率较大时, 会导致其与 $L_1$ 等值线的第一个交点不再位于坐标轴上, 此时无法产生稀疏解。

11.7  $L_0$ 范数是不连续的, 且是非凸函数, 无法通过优化直接求解, 必须采用遍历的方式, 从而导致这个问题是NP难问题。

3. ① 证明回归:  $f(x) = wx + b$

$$L(w, b) = \frac{1}{2n} \sum_{i=1}^n (f(x_i) - y_i)^2$$

$$\nabla L(w, b) = \left( \frac{\partial L}{\partial w}, \frac{\partial L}{\partial b} \right) \quad \frac{\partial L}{\partial w} = \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i) \cdot x_i$$

$$\frac{\partial L}{\partial b} = \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)$$

$$|f(x) - f(y)| \leq L \|x - y\|$$

$$L = \frac{1}{n} \sum_{i=1}^n x_i^2$$

② 对数几率回归:  $f(w, b) = -\frac{1}{n} \sum_{i=1}^n [y_i \log(f(x_i)) + (1 - y_i) \log(1 - f(x_i))]$

$$f(x_i) = \frac{1}{1 + e^{-x_i}} \quad \frac{\partial L}{\partial w} = \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i) \cdot x_i$$

$$\frac{\partial L}{\partial b} = \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)$$

~~$f(x) = f(y) \in [0, 1]$~~

$$\frac{\partial L}{\partial w} = \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i) \cdot x_i \quad L = \frac{1}{4} \sigma_{\max}(X^T X) = \frac{1}{4} \max_i \|x_i\|^2$$