

• 记  $\text{err}^*(x) = 1 - \max_{c \in Y} P(c|x)$ ,  $\text{err}(x) = 1 - \sum_c P(c|x)P(c|z)$ , 其中  $z$  为  $x$  的最近邻, 试证明在样本无穷多时

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

提示: 柯西-施瓦兹不等式  $(\sum_i a_i)^2 \leq n(\sum_i a_i^2)$

设  $\max_{c \in Y} P(c|x) = P(c^*|x)$

$$\sum_c P(c|x)P(c|z) \leq \sum_c P(c^*|x)P(c|z) = P(c^*|x) \sum_c P(c|z) = P(c^*|x)$$

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

$$\text{err}(x) = 1 - \sum_c P(c|x)P(c|z) \geq 1 - P(c^*|x)$$

右)  $\text{err}^*(x) \leq \text{err}(x)$ , 不等式左边得证.

在样本无穷多时,

$$\begin{aligned} \text{err}(x) &= 1 - \sum_c P(c|x)P(c|z) \approx 1 - \sum_c P^2(c|x) = 1 - P^2(c^*|x) - \sum_{c \neq c^*} P^2(c|x) \\ &\leq 1 - P^2(c^*|x) - \frac{1}{|Y|-1} \left( \sum_{c \neq c^*} P(c|x) \right)^2 = 1 - P^2(c^*|x) - \frac{1}{|Y|-1} (1 - P(c^*|x))^2 \\ &= (1 - P(c^*|x)) \left[ 1 + P(c^*|x) - \frac{1}{|Y|-1} (1 - P(c^*|x)) \right] \\ &= \text{err}^*(x) \left[ -(1 - P(c^*|x)) + 2 - \frac{1}{|Y|-1} (1 - P(c^*|x)) \right] \\ &= \text{err}^*(x) \left( 2 - \frac{1}{|Y|-1} \text{err}^*(x) \right) \end{aligned}$$

故  $\text{err}(x) \leq \text{err}^*(x) \left( 2 - \frac{1}{|Y|-1} \text{err}^*(x) \right)$ , 不等式右边得证.

10.4 在实践中, 协方差矩阵  $\mathbf{X}\mathbf{X}^T$  的特征值分解常由中心化后的样本矩阵  $\mathbf{X}$  的奇异值分解代替, 试述其原因.

对任意大小为  $m \times n$  的矩阵  $X$ ,  $X^T X v = \lambda v$

令  $Xv = \sigma u$ , 那么  $X^T \sigma u = \lambda v$ ,

分别左乘  $A$  得  $XX^T u = \frac{\lambda}{\sigma} Xv = \lambda u$ ,  $u$  对应  $XX^T$  的特征值为  $\lambda$  的特征向量,

$$\left. \begin{aligned} \text{在 } Xv = \sigma u \text{ 两边分别乘 } u^T, \text{ 得 } u^T Xv &= \sigma \\ \text{在 } X^T \sigma u = \lambda v \text{ 两边分别乘 } v^T, \text{ 得 } v^T X^T u &= \frac{\lambda}{\sigma} \end{aligned} \right\} \Rightarrow \sigma^2 = \lambda$$

将  $Xv = \sigma u$  写成矩阵形式为  $XV = U\Sigma$ , 其中  $\Sigma = \text{diag}[\sigma_1, \dots, \sigma_r, 0, \dots, 0]$ ,  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ ,  $r = \text{rank}(X)$

$\therefore V$  为正交矩阵  $\therefore X$  可分解为  $X = XVV^T = U\Sigma V^T$ ,  $U \in R^{m \times m}$ ,  $\Sigma \in R^{m \times n}$ ,  $V \in R^{n \times n}$

其中  $U$  的列向量是  $XX^T$  的特征向量,  $V$  的列向量是  $X^T X$  的特征向量,

$\Sigma$  的非零对角元的平方即为  $XX^T$  和  $X^T X$  的共同非零特征值.

故两种分解具有等价性, 而奇异值分解所需的计算和存储的成本更小. 中心化消除量纲对数据的影响, 解决模型不稳定等, 故在实践中使用中心化后的样本矩阵的奇异值分解代替.

- 求解20页的优化问题

## 主成分分析—求解



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

- 使用拉格朗日乘子法可得

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

只需对协方差矩阵  $X X^T$  进行特征值分解，并将求得特征值排序： $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$ ，再取前  $d'$  个特征值对应的特征向量构成  $W = (w_1, w_2, \dots, w_{d'})$

这就是主成分分析的解。

$$\begin{aligned} \max_W \text{tr}(W^T X X^T W) \quad \text{s.t. } W^T W &= I_{d'} \\ \Leftrightarrow \min -\text{tr}(W^T X X^T W) \quad \text{s.t. } W^T W &= I_{d'} \end{aligned}$$

其拉格朗日主数为  $L(W, \Lambda) = -\text{tr}(W^T X X^T W) + (W^T W - I) \Lambda$

其中  $\Lambda = \begin{bmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots \\ & & & \lambda_{d'} \end{bmatrix} \in \mathbb{R}^{d' \times d'}$ ,  $I = \begin{bmatrix} 1 & & \\ & 1 & \\ & & \ddots \\ & & & 1 \end{bmatrix} \in \mathbb{R}^{d' \times d'}$ ,  $W \in \mathbb{R}^{d \times d'}$

$$\begin{aligned} \frac{\partial L(W, \Lambda)}{\partial W} &= -\frac{\partial \text{tr}(W^T X X^T W)}{\partial W} + \frac{\partial (W^T W - I)}{\partial W} \Lambda \\ &= -X X^T W - (X X^T)^T W + 2W \Lambda = -2X X^T W + 2W \Lambda \end{aligned}$$

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

$X X^T \in \mathbb{R}^{d \times d}$  有  $d$  个特征向量，我们只需要其中  $d'$  个。

对  $X X^T W = W \Lambda$  两边同时左乘  $W^T$ , 得  $W^T X X^T W = W^T W \Lambda = \Lambda$

我们的优化目标  $\text{tr}(W^T X X^T W) = \text{tr}(\Lambda) = \sum_{i=1}^{d'} \lambda_i$

$\therefore$  最大化迹即选择最大的  $d'$  个  $\lambda_i$ ，和它们对应的特征向量  $w_i$  组成矩阵  $W = (w_1, w_2, \dots, w_{d'})$

这就是主成分分析的解。

附加题：令  $M = PP^T$ ，那么下列问题还是凸优化问题吗？试证明之。

$$\min_P \sum_{(x_i, x_j) \in \mathcal{M}} \|x_i - x_j\|_M^2 \quad \text{s. t.} \quad \sum_{(x_i, x_j) \in \mathcal{C}} \|x_i - x_j\|_M^2 \geq 1$$

$$\text{目标函数: } f_0 = \sum \|x_i - x_j\|_M^2 = \sum (x_i - x_j)^T P P^T (x_i - x_j)$$

$$\text{不等式约束: } f_1 = 1 - \sum \|x_i - x_j\|_M^2 = 1 - \sum (x_i - x_j)^T P P^T (x_i - x_j)$$

$$\frac{\partial f_0}{\partial P} = \sum [(x_i - x_j)(x_i - x_j)^T P + (x_i - x_j)(x_i - x_j)^T P] = \sum [2(x_i - x_j)(x_i - x_j)^T P]$$

$$\frac{\partial f_0}{\partial P \partial P^T} = \sum 2(x_i - x_j)(x_i - x_j)^T, \quad \text{元素不全为 0}$$

$$\frac{\partial f_1}{\partial P} = - \sum [2(x_i - x_j)(x_i - x_j)^T P]$$

$$\frac{\partial f_1}{\partial P \partial P^T} = - \sum 2(x_i - x_j)(x_i - x_j)^T$$

$\frac{\partial f_0}{\partial P \partial P^T}$  与  $\frac{\partial f_1}{\partial P \partial P^T}$  异号，故  $f_0$  与  $f_1$  不可能同时为凸函数，故该问题不是凸优化问题。

### 11.5 结合图 11.2，试举例说明 $L_1$ 正则化在何种情形下不能产生稀疏解。

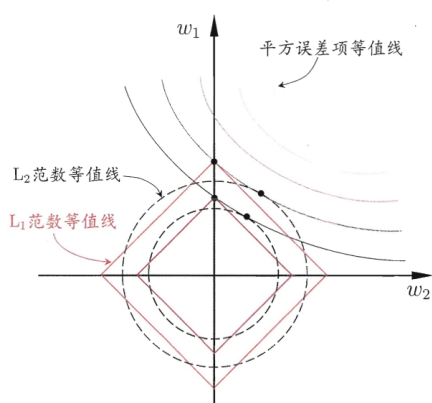


图 11.2  $L_1$  正则化比  $L_2$  正则化更易于得到稀疏解

对于  $L_1$  范数而言，约束为红色菱形框及其内，我们要在这之中找到一点尽可能使得平方误差项更小。当平方误差等值线在第一、三象限存在斜率为 -1 的点，或在第二、四象限存在斜率为 1 的点时，其与  $L_1$  等值线的第一个交点不在坐标轴上，此时，无法产生稀疏解。

### 11.7 试述直接求解 $L_0$ 范数正则化会遇到的困难。

$L_0$  范数是非零元素的个数，是很好的“稀疏约束”，

但  $L_0$  范数不连续，非凸，难以优化求解，

只能采用遍历的方法，设特征数为  $N$ ，则要在  $2^N$  的情况下分别求解问题，因此比较困难。

- PPT 20页: 证明回归和对率回归的损失函数的梯度是否满足L-Lipschitz条件, 并求出L

## L<sub>1</sub>正则化问题的求解(1)



- 可使用近端梯度下降(Proximal Gradient Descend, 简称PGD) 求解  
[Boyd and Vandenberghe, 2004]
- 考虑更一般的问题

$$\min_w f(w) + \lambda \|w\|_1$$

- 假设 $f(w)$ 为凸函数, 且 $\nabla f(w)$ 满足L-Lipschitz条件, 即存在常数 $L > 0$ , 使得

$$\forall w, w' \quad \|\nabla f(w') - \nabla f(w)\| \leq L \|w' - w\|$$

线性回归的损失函数:  $f(w) = (y - Xw)^T (y - Xw)$

$$\nabla f(w) = 2X^T(Xw - y)$$

$$\|\nabla f(w') - \nabla f(w)\| = \|2X^T(Xw' - y) - 2X^T(Xw - y)\| = \|2X^T X(w' - w)\| \leq 2\|X^T X\| \|w' - w\|$$

取  $L \geq 2\|X^T X\|$ , 即满足  $\forall w, w' \quad \|\nabla f(w') - \nabla f(w)\| \leq L\|w' - w\|$

对率回归的损失函数:  $f(w) = \sum_{i=1}^m (-y_i w^T x_i + \ln(1 + e^{w^T x_i}))$

$$\nabla f(w) = \sum_{i=1}^m (-y_i + \frac{1}{1 + e^{-w^T x_i}}) x_i$$

$$\|\nabla f(w') - \nabla f(w)\| = \left\| \sum_{i=1}^m x_i \left( \frac{1}{1 + e^{-w'^T x_i}} - \frac{1}{1 + e^{-w^T x_i}} \right) \right\|$$

$$\text{设 } g(w) = \frac{1}{1 + e^{-w^T x}}, \quad \frac{\partial g(w)}{\partial w} = -\frac{e^{-w^T x} \cdot (-x)}{(1 + e^{-w^T x})^2} = \frac{x e^{-w^T x}}{(1 + e^{-w^T x})^2}$$

$$\text{设 } a = e^{-w^T x}, \text{ 则 } a > 0 \quad \frac{e^{-w^T x}}{(1 + e^{-w^T x})^2} = \frac{a}{(1+a)^2} = \frac{a}{a^2 + 2a + 1} = \frac{1}{a + \frac{1}{a} + 2} \leq \frac{1}{4}$$

$$\text{故 } \frac{\partial g(w)}{\partial w} = \frac{x e^{-w^T x}}{(1 + e^{-w^T x})^2} \leq \frac{1}{4} x$$

$$\text{由微分中值定理: } \frac{g(w') - g(w)}{w' - w} = g'(\xi), \quad \xi \in (w', w) \text{ 或 } (w, w')$$

$$\leq \frac{1}{4} x$$

$$\text{故 } \|\nabla f(w') - \nabla f(w)\| = \left\| \sum_{i=1}^m x_i \left( \frac{1}{1 + e^{-w'^T x_i}} - \frac{1}{1 + e^{-w^T x_i}} \right) \right\| \leq \left\| \sum_{i=1}^m x_i \cdot \frac{1}{4} x_i (w' - w) \right\|$$

$$\leq \frac{1}{4} \left\| \sum_{i=1}^m x_i^2 \right\| \|w' - w\|$$

取  $L \geq \frac{1}{4} \|X^T X\|$ , 即满足  $\forall w, w' \quad \|\nabla f(w') - \nabla f(w)\| \leq L\|w' - w\|$