

# 一. 复习题

1. (1) 对原式进行变形. 同时左乘与右乘  $(XX^T + \lambda I_m)$  和  $(X^T X + \lambda I_d)$  有:

$$X(X^T X + \lambda I_d) = (X X^T + \lambda I_m) X$$

进行分配得:  $XX^T X + \lambda X I_d = XX^T X + \lambda I_m X$

即:  $XX^T X + \lambda X = XX^T X + \lambda X$  故式(4)均成立.

(2) 优化目标函数为:  $\frac{1}{2} \|Xw + mb - y\|_2^2 + \lambda \|w\|_2^2$

先优化  $b$ . 且显然有  $1 \cdot m^T / m = m$

由于  $\lambda \|w\|_2^2$  不包含  $b$ . 故这里对其不必关心.

对  $\frac{1}{2} (Xw + mb - y)^T (Xw + mb - y)$  中的  $b$  求导

即  $mb - 1m^T y + 1m^T Xw$ . 令其等于 0.

$$\text{可解得 } b = \frac{1m^T (y - Xw^*)}{m}$$

再优化  $w$ . 对  $\frac{1}{2} \|Xw + mb - y\|_2^2 + \lambda \|w\|_2^2$  中的  $w$  求导

得  $X^T (Xw + mb - y) + 2\lambda w = 0$  由于  $\lambda$  为常数. 故可解得:

$$w_R^* = (X^T X + \lambda I_d)^{-1} X^T (y - 1mb), \quad b_R^* = \frac{1m^T (y - Xw_R^*)}{m}$$

$$\text{而 } w_{LS}^* = (X^T X)^{-1} X^T (y - 1mb), \quad b_{LS}^* = \frac{1m^T (y - Xw_{LS}^*)}{m}$$

区别在于岭回归引入了  $\lambda I_d$  这一项. 以保证  $w_R^*$  为满秩矩阵. 避免了  $WX^T X$  无法求逆的问题. 同时减小了模型复杂度. 降低了过拟合风险.

2. (1) ①证明  $0 \leq \text{Ent}(D)$

由定义知  $\text{Ent}(D) = -\sum_{k=1}^{|Y|} P_k \log_2 P_k$  因为  $P_k$  为样本所占比例, 故  $0 \leq P_k \leq 1$ .

所以,  $\log_2 P_k \leq 0, -\log_2 P_k \geq 0 \therefore \text{Ent}(D) = \sum_{k=1}^{|Y|} -P_k \log_2 P_k \geq 0$

在  $P_k=1$  即  $|Y|=1$ , 仅有一类样本时取到等号.  $\text{Ent}(D)=0$

②证明  $\text{Ent}(D) \leq \log_2 |Y|$ .

设  $f(x) = -x \log x, (0 \leq x \leq 1)$ .  $f'(x) = -1 - \ln x$   $f''(x) = -\frac{1}{x} < 0 \therefore f(x)$  在  $(0, 1)$  上为上凸函数

由 Jensen 不等式得:  $-\frac{\sum_{k=1}^{|Y|} P_k \ln P_k}{|Y|} \leq f\left(\frac{\sum_{k=1}^{|Y|} P_k}{|Y|}\right) = -\frac{\sum_{k=1}^{|Y|} P_k}{|Y|} \ln\left(\frac{\sum_{k=1}^{|Y|} P_k}{|Y|}\right) = \frac{1}{|Y|} \ln |Y|$

故  $\text{Ent}(D) = -\sum_{k=1}^{|Y|} P_k \log_2 P_k \leq \ln |Y|$  由  $\sum_{k=1}^{|Y|} P_k = 1$  和  $f(x) = -x \log x$  得:

当  $P_1 = P_2 = \dots = P_k = \frac{1}{|Y|}$  时,  $\text{Ent}(D)$  取得最大值  $\text{Ent}(D) = \log_2 |Y|$ .

2. (2) 由 (1) 可知, 在二分类问题中,  $|Y|=2$ .

① 信息熵:  $\text{Ent}(D) = -(p \log_2 p + (1-p) \log_2 (1-p)) = -p \log_2 p - (1-p) \log_2 (1-p)$

② 基尼指数:  $1 - [p^2 + (1-p)^2] = 2p - 2p^2$

③ 设分类错误率:  $1 - \max_k P_k = 1 - \max\{p, 1-p\}$ .

故  $p \geq 0.5$  时, 选择  $P$ , 即为  $1-p$ .  $p < 0.5$  时选择  $1-p$ , 即为  $p$ .

$\therefore \text{Error}(D) = \begin{cases} 1-p, & p \geq 0.5 \\ p, & p < 0.5 \end{cases}$

(2) 由 (1) 可知,  $\text{Ent}(D)$  为一个上凸函数, 且  $\sum_{v=1}^V \frac{|D^v|}{|D|} = 1$ . 数据条  $D$  为  $D^v, v \in [1, V]$  之和.

故由 Jensen 不等式有:  $\text{Ent}(D) \geq \sum_{v=1}^V \frac{|D^v|}{|D|} \text{Ent}(D^v)$

故  $\text{Gain}(D, a) = \text{Ent}(D) - \sum_{v=1}^V \frac{|D^v|}{|D|} \text{Ent}(D^v) \geq 0$

更正: 一. 2 (1) 中信息熵下界的证明:

设  $f(p) = -p \log_2 p$ , 则有  $f'(p) = -\log_2 p - \frac{1}{\ln 2}$   $f''(p) = -\frac{1}{p \ln 2}$

$\therefore p \in [0, 1]$  所以  $f'(p) < 0$ , 说明  $f(p)$  单调递减, 且  $f'(1/e) = 0$

当  $p=0$  和  $1$  时  $f(p) = 0 \therefore f(p)$  在  $[0, 1]$  上最小值为  $0$ . 即对任意  $-p \log_2 p$

都有  $-p \log_2 p \geq 0 \therefore \text{Ent}(D) \geq 0$ .

3. (1) 由链式法则:  $\frac{\partial L(y, \hat{y}_1)}{\partial \beta_1} = \frac{\partial L(y, \hat{y}_1)}{\partial \hat{y}_1} \cdot \frac{\partial \hat{y}_1}{\partial \beta_1}$

对  $\frac{\partial L(y, \hat{y}_1)}{\partial \hat{y}_1} = -[y \cdot \frac{1}{\hat{y}_1} + (1-y) \cdot (-\frac{1}{1-\hat{y}_1})] = \frac{1-y}{1-\hat{y}_1} - \frac{y}{\hat{y}_1}$

由 Sigmoid 函数性质:  $f(x) = f(x) \cdot (1-f(x))$   $\frac{\partial \hat{y}_1}{\partial \beta_1} = \hat{y}_1 \cdot (1-\hat{y}_1)$

代入可得  $\frac{\partial L(y, \hat{y}_1)}{\partial \beta_1} = (\frac{1-y}{1-\hat{y}_1} - \frac{y}{\hat{y}_1}) \cdot (\hat{y}_1 \cdot (1-\hat{y}_1)) = (1-y)\hat{y}_1 - y(1-\hat{y}_1) = \hat{y}_1 - y$

3. (2) 对于  $i=j$  时:  $\frac{\partial L(y, \hat{y}_j)}{\partial \hat{y}_j} \cdot \frac{\partial \hat{y}_j}{\partial \beta_j} = (-\frac{y_j}{\hat{y}_j}) \cdot \hat{y}_j (1-\hat{y}_j) = \hat{y}_j - y_j$

对于  $i \neq j$  时:  $\frac{\partial L(y, \hat{y}_j)}{\partial \hat{y}_j} \cdot \frac{\partial \hat{y}_j}{\partial \beta_j} = (-\frac{y_j}{\hat{y}_j}) \cdot (-\hat{y}_i \hat{y}_j) = \hat{y}_i \hat{y}_j$

且对于  $i=j$  时:  $y_i=0$  故  $\frac{\partial L(y, \hat{y}_j)}{\partial \hat{y}_j} \cdot \frac{\partial \hat{y}_j}{\partial \beta_j} = 0$

所以  $\frac{\partial L(y, \hat{y}_j)}{\partial \beta_j} = \frac{\partial L(y, \hat{y}_j)}{\partial \hat{y}_j} \cdot \frac{\partial \hat{y}_j}{\partial \beta_j} = \hat{y}_j - y_j$

$\therefore \frac{\partial L(y, \hat{y})}{\partial \beta} = \hat{y} - y$  其中  $\hat{y}$  为 softmax 输出的概率向量  
 $y$  为真实类别的标签向量

一. 3. (3) 对于二分类问题而言, 两者无本质区别. 下面以输入  $\beta_1$  为例:

Sigmoid 函数:  $\hat{y}_1 = \frac{1}{1+e^{-\beta_1}}$

Softmax 函数:  $\hat{y}_1 = \frac{e^{\beta_1}}{e^{\beta_1} + e^{\beta_2}} = \frac{1}{1+e^{-(\beta_1-\beta_2)}}$

由于  $(x_1, x_2)$  可以用  $z_1$  代替, 即可写为:  $\hat{y}_1(z_1) = \frac{1}{1+e^{-z_1}}$

这与 Sigmoid 形式完全相同. 所以理论上两者没有区别. 但是对于实际的输入输出来说, Sigmoid 输出一个值  $\hat{y}_1$ , 表示正例的概率, 反例的概率则为  $1-\hat{y}_1$ .  
Softmax 在二分类中输出两个值  $\hat{y}_1$  和  $\hat{y}_2$ , 分别对应正例和反例的概率, 且满足  $\hat{y}_1 + \hat{y}_2 = 1$ .

## 二. PCA降维

1. ①对任意非零向量  $v \in \mathbb{R}^d$ ,  $v^T \hat{X}^T \hat{X} v = (\hat{X} v)^T (\hat{X} v) = \|\hat{X} v\|^2 \geq 0$

故  $\hat{X}^T \hat{X}$  是半正定矩阵. 同理任意非零向量  $u \in \mathbb{R}^m$

$u^T \hat{X} \cdot \hat{X}^T u = (\hat{X}^T u)^T (\hat{X}^T u) = \|\hat{X}^T u\|^2 \geq 0$  故  $\hat{X} \hat{X}^T$  也是半正定矩阵.

②对  $\hat{X}$  进行奇异值分解, 有  $\hat{X} = U \Sigma V^T$ , 故  $\hat{X}^T \hat{X} = V \Sigma^T U^T U \Sigma V^T$

由于  $U$  是正交矩阵, 故  $U^T U = I$  所以  $\hat{X}^T \hat{X} = V \Sigma^T \Sigma V^T$

故其特征值为  $\Sigma^T \Sigma$  的对角元素, 即奇异值的平方.

同理  $\hat{X} \hat{X}^T = U \Sigma \Sigma^T U^T$ , 特征值为  $\Sigma \Sigma^T$  的对角元素, 也是奇异值的平方.

故  $\hat{X}^T \hat{X}$  和  $\hat{X} \hat{X}^T$  的非零特征值相同.

③为了高效, 由于  $\hat{X} \hat{X}^T$  的维度较小, 可以先用其推导出  $\hat{X}^T \hat{X}$  的特征值, 再对其进行特征值分解.

## 二. 2.

①关系: 奇异值是  $\hat{X}^T \hat{X}$  或  $\hat{X} \hat{X}^T$  的特征值的平方根

即  $\sigma_i = \sqrt{\lambda_i}$ , 其中  $\sigma_i$  为  $\hat{X}$  的奇异值,  $\lambda_i$  为  $\hat{X}^T \hat{X}$  或  $\hat{X} \hat{X}^T$  的特征值.

②若求得  $\hat{X}$  的奇异值分解, 即有  $\hat{X} = U \Sigma V^T$

样本的协方差矩阵  $\hat{X} \cdot \hat{X}^T = U \Sigma V^T \cdot V \Sigma^T U^T = U \Sigma \Sigma^T U^T$  (忽略常数因子).

由于矩阵  $U$  正交,  $U U^T = I$ ,  $U U^T = I$  故  $U^T = U^{-1}$

故  $\hat{X} \hat{X}^T = U \Sigma \Sigma^T U^T$  为特征值分解.

$U$  中对应的列向量,  $u_i$  即为对应的特征向量.

这里取要求的前  $K$  个向量, 即  $u_1, u_2, \dots, u_K$  组成的投影矩阵  $W^*$

③优势: 1) 适用于高维数据, 计算开销少; 当  $d \gg m$  时, 特征值的分解

成本过高, 而 SVD 可以直接分解.

2) SVD 数值稳定而特征值分解不行.

3) 对于某些稀疏数据, 比 SVD 的计算量和内存开销更小.

二. 3. 中心化将数据平移至均值为零的状态, 这意味着去除了数据的偏移, 使得所有主成分的计算基于数据的内部结构和相对分布, 而不是受到原始数据均值位置的影响. 如果不进行中心化, PCA 可能会错误地将均值位置作为主成分的一个方向, 从而偏离对数据方差的真实描述.



二.4. 首先对数据进行中心化处理:

$$\bar{x}_1 = \frac{3+4+4+6+3}{5} = 4$$

$$\bar{x}_2 = \frac{2+3+2+3+0}{5} = 2$$

计算中心化后的数据矩阵  $\hat{X}^T = \begin{pmatrix} -1 & 0 & 0 & 2 & -1 \\ 0 & 1 & 0 & 1 & -2 \end{pmatrix}$

$$\hat{X}^T \hat{X} = \begin{pmatrix} -1 & 0 & 0 & 2 & -1 \\ 0 & 1 & 0 & 1 & -2 \end{pmatrix} \begin{pmatrix} -1 \\ 0 \\ 0 \\ 1 \\ -2 \end{pmatrix} = \begin{pmatrix} 6 & 4 \\ 4 & 6 \end{pmatrix}$$

协方差矩阵 (为:  $C = \frac{1}{n-1} \hat{X}^T \hat{X} = \begin{pmatrix} \frac{3}{2} & \frac{1}{2} \\ 1 & \frac{3}{2} \end{pmatrix}$ )

求其特征值与特征向量:  $|C - \lambda I| = 0$  即  $\begin{vmatrix} \frac{3}{2} - \lambda & \frac{1}{2} \\ 1 & \frac{3}{2} - \lambda \end{vmatrix} = 0$

$$\text{即 } (\frac{3}{2} - \lambda)^2 - 1 = 0 \quad \therefore \lambda = \frac{5}{2}, \frac{1}{2}$$

$$(C - \frac{5}{2}I) v_1 = 0 \Rightarrow \begin{pmatrix} -1 & 1 \\ 1 & -1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = 0$$

解得: 特征向量  $v_1 = (1, 1)^T$ . 对应特征值  $\lambda_1 = \frac{5}{2}$

$$(C - \frac{1}{2}I) v_2 = 0 \Rightarrow \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = 0$$

解得: 特征向量  $v_2 = (1, -1)^T$ . 对应特征值  $\lambda_2 = \frac{1}{2}$

因此投影矩阵为:  $W = (v_1, v_2) = \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$

$\therefore$  降维后的数据为  $Z = W^T \hat{X}^T = \begin{pmatrix} -1 & 1 & 0 & 3 & -3 \\ -1 & -1 & 0 & 1 & 1 \end{pmatrix}$

### 三. 度量学习应用:

1. (1) 对标准马氏距离的  $M$  进行讨论, 对协方差矩阵  $\Sigma$  进行特征值分解.

$\Sigma = V \Lambda V^T$  由于  $V$  是一个正交矩阵, 有  $V^{-1} = V^T$ .

故可写成  $\Sigma = V \Lambda V^T$   $\Sigma^{-1} = V \Lambda^{-1} V^T$ .  $\Lambda$  的对角元素为特征值, 对应主成分的方差.

重马氏距离:  $\text{dist}^2_{\text{mah}}(X_i, X_j) = (X_i - X_j)^T V \Lambda^{-1} V^T (X_i - X_j)$

在这种表示中, 引入  $\Lambda^{-1}$ , 使得每个特征贡献大都会被其方差标准化.

消去了变量之间的相关性.

由于  $\Sigma$  对角元素表示各个特征值的方差, 故使用  $\Lambda^{-1} (\Sigma^{-1})$  对每个特征标准化.

将每个特征转换为相同的度量单位, 消去了数据的量纲影响.

(2) 是的, 在数据维度大于样本数量或特征存在共线性时, 或者数据中的某些特征的方差为 0 时, 就会出现协方差矩阵不可逆. 应该这样应对:

① 正则化协方差矩阵, 添加一个小的正则项, 即  $\Sigma_{\text{reg}} = \Sigma + \lambda I$

② 使用 PCA 降维, 挑选其主要成分, 去除冗余或低方差的特征再计算.

一. (3)

① 非负性: 由于  $M$  是半正定矩阵,  $M \geq 0$ , 故  $(X_i - X_j)^T M (X_i - X_j) \geq 0$

故  $\text{dist}^2_{\text{mah}}(X_i, X_j) \geq 0$  成立.

② 同一性: 当  $X_i = X_j$  时, 显然  $(X_i - X_j)^T M (X_i - X_j) = 0$   $\text{dist}^2_{\text{mah}}(X_i, X_j) = 0$

当  $\text{dist}^2_{\text{mah}}(X_i, X_j) = 0$  时, 即  $(X_i - X_j)^T M (X_i - X_j) = 0$

由于  $M$  为半正定矩阵, 所以  $X_i - X_j = 0$  即  $X_i = X_j$  成立.

③ 对称性: 由于  $\|X_i - X_j\|_2 = \|X_j - X_i\|_2$

故显然有:  $(X_i - X_j)^T M (X_i - X_j) = (X_j - X_i)^T M (X_j - X_i)$

即  $\text{dist}^2_{\text{mah}}(X_i - X_j) = \text{dist}^2_{\text{mah}}(X_j - X_i)$  故满足对称性

④ 直递性: 将  $M$  分解成  $V \Lambda^{-1} V^T$ . 由于  $\Lambda^{-1}$  为对角阵, 故将  $V^T$  分解成  $PP^T$

故  $\text{dist}^2_{\text{mah}}(X_i, X_j) = (X_i - X_j)^T V P P^T V^T (X_i - X_j)$

$= (V^T X_i - V^T X_j)^T \cdot P \cdot P^T (V^T X_i - V^T X_j)$

$= (P^T V^T X_i - P^T V^T X_j)^T (P^T V^T X_i - P^T V^T X_j)$

$= \text{dist}^2_L(P^T V^T X_i, P^T V^T X_j)$

同理有:  $\text{dist}^2_{\text{mah}}(X_i, X_k) = \text{dist}^2_L(P^T V^T X_i, P^T V^T X_k)$

$\text{dist}^2_{\text{mah}}(X_k, X_j) = \text{dist}^2_L(P^T V^T X_k, P^T V^T X_j)$

由欧氏距离性质有:  $\text{dist}_L(P^T V^T X_i, P^T V^T X_j) \leq \text{dist}_L(P^T V^T X_i, P^T V^T X_k) + \text{dist}_L(P^T V^T X_k, P^T V^T X_j)$

故有:  $\text{dist}^2_{\text{mah}}(X_i, X_j) \leq \text{dist}^2_{\text{mah}}(X_i, X_k) + \text{dist}^2_{\text{mah}}(X_k, X_j)$ .

### 三.1.4

对于 LDA 优化目标为:

$$J(W) = \frac{\text{tr}(W^\top S_b W)}{\text{tr}(W^\top S_w W)}$$

其中,  $W$  是降维投影矩阵,  $S_b$  为类间散度矩阵,  $S_w$  为类内散度矩阵.

在马氏距离的度量中, LDA 可以看作是通过选择一个  $M$  矩阵来衡量样本之间的距离, 对于 LDA, 降维后的马氏距离中的矩阵  $M$  可以写为,

$$M_{\text{LDA}} = S_w^{-1}$$

PCA 的目标函数可以表示为:

$$J(W) = \text{tr}(W^\top S_t W)$$

其中,  $S_t$  是数据的总散度矩阵 (也就是协方差矩阵),  $W$  是投影矩阵.

在马氏距离的度量中, PCA 可以看作是通过选择一个矩阵  $M$  来衡量样本之间的距离, 对于 PCA, 降维后的马氏距离中的矩阵  $M$  可以写为:

$$M_{\text{PCA}} = S_t^{-1}$$

LDA 和 PCA 两者都是线性降维的方法, 并且在度量学习中也都可以用马氏距离来解释, 都是通过选择不同的矩阵  $M$  来实现数据的投影与度量, 最终在降维后的空间中计算样本之间的欧氏距离, 不同点在于 LDA 是有监督的降维方法, 而 PCA 是无监督的降维方法. LDA 对应的矩阵  $M$  是类内散度矩阵的逆, 即  $M_{\text{LDA}} = S_w^{-1}$ , 表示样本所属类别的关系, 特别是类内的相关性, PCA 对应的矩阵  $M$  是协方差矩阵的逆, 即  $M_{\text{PCA}} = S_t^{-1}$ , 只关注数据在各个方向上的总方差, 最大化数据的方差, 以保留尽可能多的信息量.

2.

题1: LMNN的损失函数可以表示为:  $\mathcal{E}(L) = (1-\mu) \mathcal{E}_{pull}(L) + \mu \mathcal{E}_{push}(L)$

$$\text{其中: } \mathcal{E}_{pull}(L) = \sum_{i,j} \|L(\vec{x}_i - \vec{x}_j)\|^2 \quad \mathcal{E}_{push}(L) = \sum_{i,j} (1-y_{ij}) [1 + \|L(\vec{x}_i - \vec{x}_j)\|^2 - \|L(\vec{x}_i - \vec{x}_l)\|^2]_+$$

这其中  $j \sim i$  表示样本  $j$  是样本  $i$  的邻居,  $y_{ij}$  是一个指示函数, 用于区分是否为不同类别,  $[z]_+ = \max(0, z)$  表示 hinge 损失.

① 类内紧密性部分的梯度:

令马氏距离矩阵  $M = L^T L$ , 则这部分可以表示为:

$$\mathcal{E}_{pull}(M) = \sum_{i,j} (\vec{x}_i - \vec{x}_j)^T M (\vec{x}_i - \vec{x}_j)$$

$$\frac{\partial \mathcal{E}_{pull}(M)}{\partial M} = \sum_{i,j} (\vec{x}_i - \vec{x}_j) (\vec{x}_i - \vec{x}_j)^T$$

② 类间分离性部分的梯度: 这部分包含 hinge 损失, 只有在连接约束的情况下

才会对梯度有贡献. 损失形式为:  $\mathcal{E}_{push}(M) = \sum_{i,j} \sum_l (1-y_{ij}) [1 + (\vec{x}_i - \vec{x}_j)^T M (\vec{x}_i - \vec{x}_j) - (\vec{x}_i - \vec{x}_l)^T M (\vec{x}_i - \vec{x}_l)]_+$

$$\frac{\partial \mathcal{E}_{push}(M)}{\partial M} = \sum_{i,j} \sum_l (1-y_{ij}) I_{[z]_+ > 0} [(\vec{x}_i - \vec{x}_j) (\vec{x}_i - \vec{x}_j)^T - (\vec{x}_i - \vec{x}_l) (\vec{x}_i - \vec{x}_l)^T]$$

其中  $I_{[z]_+ > 0}$  是指示函数, 表示当 hinge 损失大于 0 时, 该项对梯度有贡献.

因此 LMNN 损失函数对于  $M$  的梯度为:

$$\begin{aligned} \frac{\partial \mathcal{E}(M)}{\partial M} &= (1-\mu) \sum_{i,j} (\vec{x}_i - \vec{x}_j) (\vec{x}_i - \vec{x}_j)^T + \mu \sum_{i,j} \sum_l (1-y_{ij}) I_{[z]_+ > 0} [(\vec{x}_i - \vec{x}_j) (\vec{x}_i - \vec{x}_j)^T - (\vec{x}_i - \vec{x}_l) (\vec{x}_i - \vec{x}_l)^T] \\ &= \sum_{i,j} \left[ (1-\mu) (\vec{x}_i - \vec{x}_j) (\vec{x}_i - \vec{x}_j)^T + \mu \sum_l (1-y_{ij}) I_{[z]_+ > 0} [(\vec{x}_i - \vec{x}_j) (\vec{x}_i - \vec{x}_j)^T - (\vec{x}_i - \vec{x}_l) (\vec{x}_i - \vec{x}_l)^T] \right] \end{aligned}$$

在优化过程中, 每次更新  $M$  后, 可以通过以下操作保持其对称性.

$$M = \frac{1}{2} (M + M^T)$$

这样在梯度下降的过程中, 任何对  $M$  的非对称扰动都被对称化.



题2: 可以将目标函数写为  $\| \hat{M} - M \|^F = \text{tr}(\hat{M} - M)^T(\hat{M} - M) = \text{tr}(\hat{M}^2) - 2\text{tr}(\hat{M}M) + \text{tr}(M^2)$

由于  $\text{tr}(M^2)$  是常数, 只需最小化  $f(\hat{M}) = \text{tr}(\hat{M}^2) - 2\text{tr}(\hat{M}M)$

$\because \hat{M}$  是半正定对称矩阵, 对  $\hat{M}$  进行特征值分解:  $\hat{M} = P\Lambda P^T$

其中:  $P$  是由  $\hat{M}$  的特征向量组成的正交矩阵.

$\Lambda = \text{diag}(\hat{\lambda}_1, \dots, \hat{\lambda}_d)$  是由  $\hat{M}$  的特征值组成的对角矩阵, 且  $\hat{\lambda}_i \geq 0$

由于  $\text{tr}(AB) = \text{tr}(BA)$  且  $P^T P = I$ , 因此有:  $\text{tr}(\hat{M}^2) = \text{tr}(P\Lambda P^T P\Lambda P^T) = \text{tr}(P\Lambda^2 P^T) = \text{tr}(\Lambda^2)$

$$\text{tr}(\hat{M}M) = \text{tr}(P\Lambda P^T Q\Lambda Q^T) = \text{tr}(\Lambda^T P^T Q \Lambda Q^T P)$$

令  $W = P^T Q$ , 则  $W$  是正交矩阵, 且  $\text{tr}(\hat{M}M) = \text{tr}(\Lambda W \Lambda W^T)$

因此目标函数变为:  $f(\hat{M}) = \text{tr}(\Lambda^2) - 2\text{tr}(\Lambda W \Lambda W^T)$

展开  $\text{tr}(\Lambda W \Lambda W^T)$ , 其中  $w_{ij}$  是矩阵  $W$  的元素:  $\text{tr}(\Lambda W \Lambda W^T) = \sum_{i=1}^d \sum_{j=1}^d \hat{\lambda}_i \lambda_j (w_{ij})^2$

由于  $\hat{\lambda}_i \geq 0$ , 且  $W$  为正交矩阵, 因此  $\sum_{j=1}^d (w_{ij})^2 = 1$ , 而对于  $w_{ij}$ , 为了最小化  $f(\hat{M})$ ,

需要最大化  $\sum_{j=1}^d \hat{\lambda}_i \lambda_j (w_{ij})^2$ , 即寻找  $w_{ij}$ , 使得  $\sum_{j=1}^d \hat{\lambda}_i \lambda_j (w_{ij})^2$  最大化.

注意到当  $W = I$  时  $(w_{ij})^2 = \delta_{ij}$  (克罗内克函数), 即  $w_{ij} = \delta_{ij}$ , 此时有:

$$\sum_{j=1}^d \hat{\lambda}_i \lambda_j (w_{ij})^2 = \sum_{j=1}^d \hat{\lambda}_i \lambda_j \text{ 这能取到的最大值. 对于任意正交矩阵 } W.$$

根据 Schur 凸性的性质:  $\sum_{j=1}^d \hat{\lambda}_i \lambda_j \geq \sum_{j=1}^d \hat{\lambda}_i \lambda_j (w_{ij})^2$

因此, 为了最小化  $f(\hat{M})$  应选择  $P = Q$  即:  $W = I$ .

$$\text{当 } P=Q \text{ 时, 目标函数变为: } f(\hat{M}) = \sum_{i=1}^d \hat{\lambda}_i^2 - 2 \sum_{i=1}^d \hat{\lambda}_i \lambda_i = \sum_{i=1}^d (\hat{\lambda}_i^2 - 2\hat{\lambda}_i \lambda_i) = \sum_{i=1}^d (\hat{\lambda}_i - \lambda_i)^2$$

由于  $\lambda_i$  是常数, 所以只需最小化  $g(\hat{\lambda}_i) = \sum_{i=1}^d (\hat{\lambda}_i - \lambda_i)^2$

需要指明的是, 此时并未假设  $\hat{\lambda}$  是特定的对角矩阵, 而是通过优化过程得出了  
应选择  $P=Q$  从而使得  $\hat{M}$  与  $M$  在同一基下进行比较. 现在需要在约束  $\hat{\lambda}_i \geq 0$  下,

最小化  $g(\hat{\lambda}_i)$ . 对于每一个  $i$  这是一个简单的凸优化问题:  $\min_{\hat{\lambda}_i \geq 0} (\hat{\lambda}_i - \lambda_i)^2$

解为: 若  $\lambda_i \geq 0$ , 则最优解为  $\hat{\lambda}_i = \lambda_i$ . 因此最优  $\hat{\lambda}_i$  为:  $\hat{\lambda}_i = \max(\lambda_i, 0)$

若  $\lambda_i < 0$ , 则最优解为  $\hat{\lambda}_i = 0$  代入最优特征值.

$$\hat{\Lambda} = \hat{\Lambda}^+ = \text{diag}(\hat{\lambda}_1, \dots, \hat{\lambda}_d) = \text{diag}(\max(\lambda_1, 0), \max(\lambda_2, 0), \dots, \max(\lambda_d, 0)).$$

$\therefore P=Q$ ,  $\therefore$  最优半正定矩阵为:  $\hat{M} = P\hat{\Lambda}P^T = Q\hat{\Lambda}^+Q^T$ .

题3: 因为  $M = P^T P$  ∴ 距离度量可表示为:

$$\text{dist}^2_{\text{mah}}(X_i, X_j) = (X_i - X_j)^T P^T P (X_i - X_j) = \|P(X_i - X_j)\|^2$$

$L_{MNN}$  可以转化为关于  $P$  的无约束优化问题:

$$L(P) = \sum_{i,j} \|P^T(X_i - X_j)\|^2 + c \sum_{i,j,k} [1 + \|P^T(X_i - X_j)\|^2 - \|P^T(X_i - X_k)\|^2]$$

其中  $[\cdot]_+$  表示取非负部分,  $c = \frac{\alpha}{1-\alpha}$  代表近邻保持项和排斥项之间的权重比例

首先考虑第一项,  $L_1(P) = \sum_{i,j} \|P^T(X_i - X_j)\|^2$  对  $P$  求导:

$$\nabla L_1(P) = \frac{\partial}{\partial P} \|P^T(X_i - X_j)\|^2 = 2P(X_i - X_j)(X_i - X_j)^T$$

再考虑排斥项:  $L_2(P) = c \sum_{i,j,k} [1 + \|P^T(X_i - X_j)\|^2 - \|P^T(X_i - X_k)\|^2]$

该项对  $\|P^T(X_i - X_k)\|^2$  和  $\|P^T(X_i - X_j)\|^2$  求导, 梯度推导同上.

$$\therefore \nabla L_2(P) = 2c \sum_{i,j,k} I[1 + \|P^T(X_i - X_j)\|^2 - \|P^T(X_i - X_k)\|^2 > 0]$$

$$\cdot P((X_i - X_j)(X_i - X_j)^T - (X_i - X_k)(X_i - X_k)^T)$$

两部分梯度相加最终梯度为:

$$\nabla L(P) = 2 \sum_{i,j} P(X_i - X_j)(X_i - X_j)^T + 2c \sum_{i,j,k} I[1 + \|P^T(X_i - X_j)\|^2 - \|P^T(X_i - X_k)\|^2 > 0] \cdot P((X_i - X_j)(X_i - X_j)^T - (X_i - X_k)(X_i - X_k)^T)$$

由于  $\nabla L(P)$  中包含指示符函数  $I$  的项使得梯度不再是线性的

使得该问题关于  $P$  是非凸的.