

Penalized Mahalanobis Distance

Theoretical Justification

Let \mathbf{X} be $n \times p$ dimensional matrix with $n > p$. The primary interest is the distances between the columns of the matrix \mathbf{X} using the correlation structure between the rows of the matrix \mathbf{X} . The singular value decomposition (SVD) for \mathbf{X} has the form \mathbf{USV}^T . Here $\mathbf{UU}^T = \mathbf{I}$ and $\mathbf{VV}^T = \mathbf{I}$ are the orthogonal matrices of dimension $n \times n$ and $p \times p$ respectively. The matrix \mathbf{S} has dimension $n \times p$ and has non-zero elements only on diagonal when row indexes are equal to column indexes. The rank of \mathbf{S} is equal to the rank of \mathbf{X} . Also let $\mathbf{XX}^T = \mathbf{UDU}^T$ be a spectral decomposition (SD) of \mathbf{XX}^T . The matrix \mathbf{D} is a square diagonal matrix with dimension $n \times n$. The relationship between two decompositions is the following:

$$\mathbf{XX}^T = [\mathbf{USV}^T][\mathbf{USV}^T]^T = \mathbf{USV}^T\mathbf{VS}^T\mathbf{U}^T = \mathbf{USS}^T\mathbf{U}^T = \mathbf{UDU}^T \quad (1)$$

Since $n > p$, the number of non-zero diagonal elements can be at most p . The relationship between two decompositions imply $\mathbf{SS}^T = \mathbf{D}$. The diagonal elements in \mathbf{S} are not defined uniquely since the sign in SVD can be arbitrary. The diagonal elements of \mathbf{D} however are defined uniquely and are all non-negative for non-negative definite matrices such as \mathbf{XX}^T . We use the notations d_i and s_i^2 where $i = 1, 2, \dots, n$ for diagonal elements of \mathbf{D} and \mathbf{SS}^T respectively. The values $|s_i|$ are called singular values of the matrix \mathbf{X} and they are related to eigenvalues d_i of matrix \mathbf{XX}^T using relationship $d_i = s_i^2$.

If we assume that matrix \mathbf{XX}^T is invertible (i.e. $n < p$ and row space has a full rank, which is *not* the case for our matrix) then all d_i -s are positive and the inverse matrix $[\mathbf{XX}^T]^{-1}$ has the form

$$[\mathbf{XX}^T]^{-1} = \mathbf{UD}^{-1}\mathbf{U}^T = \mathbf{U}[\mathbf{SS}^T]^{-1}\mathbf{U}^T \quad (2)$$

where

$$\begin{aligned} \mathbf{D}^{-1} &= \text{diag} \left[\frac{1}{d_1}, \frac{1}{d_2}, \dots, \frac{1}{d_n} \right] \\ [\mathbf{S}\mathbf{S}^T]^{-1} &= \text{diag} \left[\frac{1}{s_1^2}, \frac{1}{s_2^2}, \dots, \frac{1}{s_n^2} \right] \end{aligned}$$

For the considered matrix \mathbf{X} with dimensions $n \times p$ with $n > p$ the matrix $\mathbf{X}\mathbf{X}^T$ does not have a full rank and is *not* invertible. In this case the SVD decomposition of the inverse (2) does not make sense.

To fix that the penalty is introduced. The idea is to add penalty to the original matrix $\mathbf{X}\mathbf{X}^T$ to make it invertible. The suggested matrix is $\mathbf{X}\mathbf{X}^T + \lambda\mathbf{I}$ and the corresponding inverse has the form

$$\begin{aligned} [\mathbf{X}\mathbf{X}^T + \lambda\mathbf{I}]^{-1} &= [\mathbf{U}\mathbf{D}\mathbf{U}^T + \mathbf{U}(\lambda\mathbf{I})\mathbf{U}^T]^{-1} \\ &= \mathbf{U} \left[\text{diag} \left[\frac{1}{d_1 + \lambda}, \frac{1}{d_2 + \lambda}, \dots, \frac{1}{d_n + \lambda} \right] \right] \mathbf{U}^T \\ &= \mathbf{U} \left[\text{diag} \left[\frac{1}{s_1^2 + \lambda}, \frac{1}{s_2^2 + \lambda}, \dots, \frac{1}{s_n^2 + \lambda} \right] \right] \mathbf{U}^T \end{aligned}$$

The matrix $\mathbf{X}\mathbf{X}^T + \lambda\mathbf{I}$ is invertible and the inverse $[\mathbf{X}\mathbf{X}^T + \lambda\mathbf{I}]^{-1}$ exists provided $\lambda \neq 0$. The question is how to select appropriate penalty λ ? Intuitively, the smallest eigenvalues d_i are the least important ones. Those values experience maximal shrinkage d_i and the largest values of $\frac{1}{d_i}$ which degenerates to ∞ when $d_i = 0$. To follow the logic the reasonable choose is to shrink based on the smallest d_i in the dataset which is *not* equal to 0. That implies

$$\lambda := \min \{d_i : d_i \neq 0\}.$$

The approach tends to produce the results when all distances become very similar i.e. the penalty is too severe to see the difference b/w the samples. The more general alternative will be to use some rescaling

$$\lambda := \min \{d_i\gamma : d_i \neq 0\}$$

where scaling constant $\gamma \in (0; 1]$. The suggested choice is to use the median \tilde{d} of all non-zero eigenvalues as a penalty:

$$\lambda := \tilde{d} = \text{median} \{d_i : d_i \neq 0\}.$$

Suggested Computation Algorithm

The algorithm has to be outlined as follows.

1. Standardize the matrix \mathbf{X} for variance computation i.e. compute the differences $[\mathbf{X} - \bar{\mathbf{X}}_R]$ where $\bar{\mathbf{X}}_R$ is $n \times p$ dimensional matrix where every element is replaced by row average of the original matrix \mathbf{X} .
2. Compute singular value decomposition of $[\mathbf{X} - \bar{\mathbf{X}}_R]$. The decomposition has the form

$$[\mathbf{X} - \bar{\mathbf{X}}_R] = \mathbf{U}_c \mathbf{S}_c \mathbf{V}_c^T \quad (3)$$

where subscript c emphasizes that it is for the *centered* matrix $[\mathbf{X} - \bar{\mathbf{X}}_R]$. The variance estimate has the form

$$\hat{\Sigma} = \widehat{\text{Var}(\mathbf{X})} = \frac{1}{p-1} [\mathbf{X} - \bar{\mathbf{X}}_R][\mathbf{X} - \bar{\mathbf{X}}_R]^T$$

3. Define the penalty $\lambda := \min\{d_i \gamma : d_i \neq 0\}$ where γ is a user-defined penalty parameter and d_i -s are diagonal elements of $\mathbf{D}_c = \mathbf{S}_c \mathbf{S}_c^T$.
4. Compute the estimate of the penalized sigma $\hat{\Sigma}_\lambda^{-1}$

$$\begin{aligned} \hat{\Sigma}_\lambda^{-1} &= (p-1) [([\mathbf{X} - \bar{\mathbf{X}}_R][\mathbf{X} - \bar{\mathbf{X}}_R]^T + \lambda \mathbf{I})^{-1}] \\ &= (p-1) \mathbf{U}_c \left[\text{diag} \left[\frac{1}{d_1 + \lambda}, \frac{1}{d_2 + \lambda}, \dots, \frac{1}{d_n + \lambda} \right] \right] \mathbf{U}_c^T \end{aligned} \quad (4)$$

5. Compute the Mahalanobis distance using the penalized estimate (4). The Mahalanobis distance for individual columns \mathbf{X}_1 and \mathbf{X}_2 of \mathbf{X} is denoted as $\hat{d}_M(\mathbf{X}_1, \mathbf{X}_2)$ and has

the form

$$\hat{d}_M(\mathbf{X}_1, \mathbf{X}_2) = \sqrt{[\mathbf{X}_1 - \mathbf{X}_2]^T \hat{\Sigma}_\lambda^{-1} [\mathbf{X}_1 - \mathbf{X}_2]} \quad (5)$$