

# Feature Selection for Latent Factor Models

Rittwika Kansabanik, Adrian Barbu  
Department of Statistics, Florida State University

## 1. Extended Literature Review

Class-specific feature selection and SNR-based methods have been explored in recent years. Cluster-based pattern discrimination [7] relies on classifiers to compute the similarity between unknown patterns and clusters, which is crucial for feature selection, while our approach operates independently of classifiers. [9] proposes a one-vs-all binary classification technique with traditional feature selectors for class-specific feature selection. SMBA-CSFS[8] employs a sparse model with a two-step strategy. It builds classifier ensembles and incurs higher computational costs. Authors in [6] leverage fuzzy entropy and mutual information to evaluate feature relevance and redundancy. These methods lack scalability and solid theoretical foundations. Our approach stands out using low-rank generative models and SNR, ensuring strong theoretical support for feature recovery and improved scalability using independent class models. [1] integrates Independent Components Analysis and SNR for feature selection via a linear transformation of all features for classification. In contrast, our method calculates SNRs on original features and uses only those selected for classification.

## 2. Theorems and Proofs

### 2.1. ELF Parameter Estimation

**Theorem 1** *The ELF objective is:*

$$(\hat{\mathbf{W}}_{ELF}, \hat{\mathbf{\Gamma}}_{ELF}) = \underset{(\mathbf{\Gamma}, \mathbf{W}), \mathbf{\Gamma}^T \mathbf{\Gamma} = \mathbf{I}_r}{\operatorname{argmin}} \|(\mathbf{X} - \mathbf{\Gamma} \mathbf{W}^T) \mathbf{\Psi}^{-\frac{1}{2}}\|_F^2, \quad (1)$$

without the constraint  $\mathbf{\Gamma}^T \mathbf{\Gamma} = \mathbf{I}_r$ , (1) is minimized w.r.t  $\mathbf{\Gamma}$  and  $\mathbf{W}$  by

$$\hat{\mathbf{\Gamma}} = \mathbf{X} \mathbf{\Psi}^{-1} \mathbf{W} (\mathbf{W}^T \mathbf{\Psi}^{-1} \mathbf{W})^{-1} \text{ and } \hat{\mathbf{W}} = \mathbf{X}^T \mathbf{\Gamma} (\mathbf{\Gamma}^T \mathbf{\Gamma})^{-1}. \quad (2)$$

**Proof.** Let  $l(\mathbf{\Gamma}) = \|(\mathbf{X} - \mathbf{\Gamma} \mathbf{W}^T) \sqrt{\mathbf{\Psi}^{-1}}\|_F^2$ . Then

$$\begin{aligned} \arg \min_{\mathbf{\Gamma}} l(\mathbf{\Gamma}) &= \arg \min_{\mathbf{\Gamma}} \|\mathbf{X} \sqrt{\mathbf{\Psi}^{-1}} - \mathbf{\Gamma} \mathbf{W}^T \sqrt{\mathbf{\Psi}^{-1}}\|_F^2 \\ &= \arg \min_{\mathbf{\Gamma}} \operatorname{Tr}((\mathbf{X} \sqrt{\mathbf{\Psi}^{-1}} - \mathbf{\Gamma} \mathbf{W}^T \sqrt{\mathbf{\Psi}^{-1}})^T (\mathbf{X} \sqrt{\mathbf{\Psi}^{-1}} - \mathbf{\Gamma} \mathbf{W}^T \sqrt{\mathbf{\Psi}^{-1}})) \\ &= \arg \min_{\mathbf{\Gamma}} \operatorname{Tr}(\sqrt{\mathbf{\Psi}^{-1}} \mathbf{X}^T \mathbf{X} \sqrt{\mathbf{\Psi}^{-1}} - 2 \sqrt{\mathbf{\Psi}^{-1}} \mathbf{X}^T \mathbf{\Gamma} \mathbf{W}^T \sqrt{\mathbf{\Psi}^{-1}} + \sqrt{\mathbf{\Psi}^{-1}} \mathbf{W} \mathbf{\Gamma}^T \mathbf{\Gamma} \mathbf{W}^T \sqrt{\mathbf{\Psi}^{-1}}) \\ &= \arg \min_{\mathbf{\Gamma}} \operatorname{Tr}(\mathbf{\Psi}^{-1} \mathbf{X}^T \mathbf{X} - 2 \mathbf{\Psi}^{-1} \mathbf{X}^T \mathbf{\Gamma} \mathbf{W}^T + \mathbf{\Psi}^{-1} \mathbf{W} \mathbf{\Gamma}^T \mathbf{\Gamma} \mathbf{W}^T) \\ &= \arg \min_{\mathbf{\Gamma}} \operatorname{Tr}(-2 \mathbf{\Psi}^{-1} \mathbf{X}^T \mathbf{\Gamma} \mathbf{W}^T + \mathbf{\Psi}^{-1} \mathbf{W} \mathbf{\Gamma}^T \mathbf{\Gamma} \mathbf{W}^T) \\ &= \arg \min_{\mathbf{\Gamma}} \operatorname{Tr}(-2 \mathbf{W}^T \mathbf{\Psi}^{-1} \mathbf{X}^T \mathbf{\Gamma} + \mathbf{\Gamma}^T \mathbf{\Psi}^{-1} \mathbf{W} \mathbf{\Gamma}^T) \end{aligned}$$

$$\frac{\partial l(\mathbf{\Gamma})}{\partial \mathbf{\Gamma}} = \frac{\partial}{\partial \mathbf{\Gamma}} \operatorname{Tr}(-2 \mathbf{W}^T \mathbf{\Psi}^{-1} \mathbf{X}^T \mathbf{\Gamma} + \mathbf{\Gamma}^T \mathbf{\Psi}^{-1} \mathbf{W} \mathbf{\Gamma}^T) = -2 \mathbf{X} \mathbf{\Psi}^{-1} \mathbf{W} + 2 \mathbf{\Gamma} \mathbf{W}^T \mathbf{\Psi}^{-1} \mathbf{W}$$

$$\frac{\partial l(\mathbf{\Gamma})}{\partial \mathbf{\Gamma}} = 0 \implies \mathbf{\Gamma} = \mathbf{X} \mathbf{\Psi}^{-1} \mathbf{W} (\mathbf{W}^T \mathbf{\Psi}^{-1} \mathbf{W})^{-1}.$$

$$\hat{\mathbf{W}} = \underset{\mathbf{W}}{\operatorname{argmin}} \sum_{j=1}^d \frac{\|\mathbf{X}_{\cdot j} - \mathbf{\Gamma} \mathbf{W}_{\cdot j}^T\|^2}{\sigma_j^2}$$

which is minimized individually for each  $\mathbf{W}_j$  as  $\mathbf{W}_j^T = (\mathbf{\Gamma}^T \mathbf{\Gamma})^{-1} \mathbf{\Gamma}^T \mathbf{X}_{.j}$ , which gives the result.  $\square$

**Proposition 1** If  $\mathbf{U}\mathbf{D}\mathbf{V}^T = \mathbf{\Gamma}$  is the SVD of  $\mathbf{\Gamma}$ , then  $\mathbf{\Gamma}_1 = \mathbf{U}$ , and  $\mathbf{W}_1 = \mathbf{W}\mathbf{V}\mathbf{D}$  satisfy  $\mathbf{\Gamma}_1 \mathbf{W}_1^T = \mathbf{\Gamma} \mathbf{W}^T$  along with  $\mathbf{\Gamma}_1^T \mathbf{\Gamma}_1 = \mathbf{I}_r$ .

**Proof.** It is easy to verify that  $\mathbf{\Gamma}_1^T \mathbf{\Gamma}_1 = \mathbf{I}_r$ .  $\square$

## 2.2. True Feature Recovery Guarantees

We have considered a model that aims to find a relationship between the observed  $\mathbf{x} \in \mathbb{R}^d$  and a hidden set of variables (latent variables)  $\boldsymbol{\gamma} \in \mathbb{R}^r$  with  $r < d$  and assumes the latent factors and noise variables are independent of each other. It is as follows:

$$\mathbf{x} = \boldsymbol{\mu} + \mathbf{W}\boldsymbol{\gamma} + \boldsymbol{\epsilon}, \quad \text{with } E(\boldsymbol{\epsilon}) = \mathbf{0} \text{ and } \text{var}(\boldsymbol{\epsilon}) = \boldsymbol{\Psi}. \quad (3)$$

We have also considered the following assumptions:

- (A1) The observations,  $(\mathbf{x}_i, i = 1, 2, \dots, n)$  are independently generated from the LFA model (3) with  $\boldsymbol{\mu} = \mathbf{0}$  and  $\boldsymbol{\epsilon} \stackrel{i.i.d}{\sim} \mathcal{N}(\mathbf{0}, \boldsymbol{\Psi})$ .
- (A2) Denote  $\boldsymbol{\Gamma} = (\gamma_1, \gamma_2, \dots, \gamma_n)^T$ .  $\boldsymbol{\Gamma}_{ij}$  are i.i.d random variables with  $E(\boldsymbol{\Gamma}_{ij}) = 0$ ,  $\text{Var}(\boldsymbol{\Gamma}_{ij}) = 1$ ,  $E(\boldsymbol{\Gamma}_{ij}^4) < \infty$ , for all  $(i, j) \in \{1, 2, \dots, n\} \times \{1, 2, \dots, r\}$ .
- (A3) There are  $m$  true features with indices  $S = \{s_1, s_2, \dots, s_m\}$  and  $(d - m)$  noisy features in our data, which satisfy, for some positive constant  $\gamma > 0$ :

$$\min\{\mathbf{S}\mathbf{N}\mathbf{R}_i^*, i \in S\} \geq \max\{\mathbf{S}\mathbf{N}\mathbf{R}_i^*, i \notin S\} + \gamma. \quad (4)$$

To prove the following proposition, we first introduce the spiked covariance model[5]:

**Definition 1 (Spiked Covariance model [due to [5]]):** Under this model, the data matrix  $\mathbf{X}$ , can be viewed as  $\mathbf{X}^T = \mathbf{E}\mathbf{\Lambda}^{\frac{1}{2}}\mathbf{Z}$ , where  $\mathbf{E} = [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_d]$  is a  $d \times d$  orthogonal matrix,  $\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_d)$  with  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$  and  $\mathbf{Z}$  is a  $d \times n$  matrix constructed with iid random variables  $\mathbf{Z}_{ij}$  with  $E(\mathbf{Z}_{ij}) = 0$ ,  $E(\mathbf{Z}_{ij}^2) = 1$  and  $E(\mathbf{Z}_{ij}^4) \leq \infty$ . The population covariance matrix is  $\boldsymbol{\Sigma} = \mathbf{E}\mathbf{\Lambda}\mathbf{E}^T$ . Here,  $\lambda_k$ 's are assumed to follow a specific structure,  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > \lambda_{r+1} = \dots = \lambda_d = 1$ .

For the following section, we assume that  $\lim_{n \rightarrow \infty} \frac{d}{n} = \delta$ . Also there are  $k$  population eigenvalues such that  $\lambda_i > 1 + \sqrt{\delta}$ , for  $i \leq k$ . The following result is due to [3].

**Theorem 2 (due to [3])** For  $\delta \in (0, 1)$ , the following holds:

$$s_i \xrightarrow{a.s.} \begin{cases} \rho(\lambda_i), & \text{if } i \leq k \\ (1 + \sqrt{\delta})^2 & \text{otherwise,} \end{cases}$$

where  $\rho(x) = x(1 + \frac{\delta}{x-1})$ .

Although consistency could not be proved for  $\delta > 0$ , [5] proved consistency for  $\delta = 0$ .

**Lemma 1 (due to [5])** If  $\lim_{n \rightarrow \infty} \frac{d}{n} = \delta = 0$ , then,

$$s_i \xrightarrow{a.s.} \begin{cases} \lambda_i & \text{if } i \leq m \\ 1 & \text{otherwise.} \end{cases}$$

**Proposition 2** Under the assumptions (A1, A2), if  $\frac{d}{n} \rightarrow \delta \in (0, 1)$  as  $n \rightarrow \infty$ , and  $\boldsymbol{\Psi} = \sigma^{*2} \mathbf{I}_d$  then  $\sigma_{ML}^2 \xrightarrow{a.s.} \sigma^{*2}(1 + \sqrt{\delta})^2$ . If  $\delta = 0$  then  $\sigma_{ML}^2 \xrightarrow{a.s.} \sigma^{*2}$ .

**Proof.** The covariance matrix under the PPCA model,  $\boldsymbol{\Sigma}^* = \mathbf{W}^* \mathbf{W}^{*T} + \sigma^{*2} \mathbf{I}_d = \mathbf{E}^* \mathbf{L}^* \mathbf{E}^{*T} + \sigma^{*2} \mathbf{I}_d$  can be viewed as a spiked covariance model.

Let  $\mathbf{l}_i^*$  be the  $i^{th}$  diagonal element of  $\mathbf{L}^*$ . As  $\mathbf{W}^* \in \mathbb{R}^{d \times r}$  is a tall skinny matrix, the first  $r$  diagonal elements of  $\mathbf{L}$  are greater than 0, and the rest of the diagonal elements are equal to 0.

Therefore, the SVD on  $\Sigma^*$  provides  $\Sigma^* = \mathbf{E}^* \Lambda^* \mathbf{E}^{*T}$ , where  $\lambda_i^*$ , the  $i^{th}$  element of  $\Lambda^*$  is:

$$\lambda_i^* = \begin{cases} l_i^* + \sigma^{*2} & \text{if } i \leq r \\ \sigma^{*2} & \text{otherwise.} \end{cases} \quad (5)$$

The ML estimates of  $\sigma^{*2}$  from [10] are given below:

$$\sigma_{ML}^2 = \frac{1}{d-r} \sum_{j=r+1}^d s_j \quad (6)$$

where  $s_i$  is the  $i^{th}$  largest eigen value of  $\hat{\Sigma}$ .

From (5), we have  $\lambda_{r+1}^* = \lambda_{r+2}^* = \dots = \lambda_d^* = \sigma^{*2}$ . Let,  $\Sigma_0^* = \Sigma^* / \sigma^{*2}$ . Therefore the estimated  $\hat{\Sigma}_0 = \hat{\Sigma} / \sigma^2$ . Let  $\lambda_i^{*0}$  and  $s_i^0$  be the  $i^{th}$  largest eigenvalues of  $\Sigma_0^*$  and  $\hat{\Sigma}_0$  respectively.

It is easy to verify that the set of principal eigenvectors for  $\Sigma_0^*$  and  $\Sigma^*$  are the same and  $\lambda_i^{*0} = \frac{\lambda_i^*}{\sigma^{*2}}$ . A similar logic holds for  $\hat{\Sigma}_0$  and  $\hat{\Sigma}$  as well, i.e.  $s_i^0 = \frac{s_i}{\sigma^{*2}}$ .

Therefore,  $\Lambda_0^* = \text{diag}(\lambda_1^{*0}, \dots, \lambda_d^{*0})$  and  $\lambda_1^{*0} \geq \lambda_2^{*0} \geq \dots \geq \lambda_r^{*0} > \lambda_{r+1}^{*0} = \lambda_{r+2}^{*0} = \dots = \lambda_d^{*0} = 1$

It can be easily verified that  $\tilde{\mathbf{X}}_{n \times d} = \frac{1}{\sigma^*} \mathbf{X}$  follows a spiked covariance model. As,  $\tilde{\mathbf{X}}^T = \Sigma_0^{* \frac{1}{2}} \mathbf{Z}_{d \times n}$ . From the assumptions (A1,A2),  $\mathbf{Z}$  can be viewed as a matrix constructed with iid random variables  $Z_{ij}$  with  $E(Z_{ij}) = 0$ ,  $E(Z_{ij}^2) = 1$  and  $E(Z_{ij}^4) \leq \infty$ .

Therefore, from the Theorem 2 we get, for  $\delta \in (0, 1)$ ,

$$s_i^0 \xrightarrow{a.s.} (1 + \sqrt{\delta})^2, \forall i > r.$$

Now  $s_i^0 = s_i / \sigma^{*2}$ , therefore  $s_i \xrightarrow{a.s.} \sigma^{*2} (1 + \sqrt{\delta})^2, \forall i > r$ .

Let us denote  $\mathbf{s}_{d-r} = \{s_{r+1}, s_{r+2}, \dots, s_d\}$ . It can be easily seen that  $\sigma_{ML}^2$  is a continuous transformation of  $\mathbf{s}_{d-r}$ , which can be defined as:  $\sigma_{ML}^2 = h(\mathbf{s}_{d-r}) = (\mathbf{s}_{d-r}^T \mathbf{1}_{d-r}) / (d-r)$ .

Then  $h(\mathbf{s}_{d-r}) \xrightarrow{a.s.} h((\sigma^{*2} (1 + \sqrt{\delta})^2) \mathbf{1}_{d-r}) \implies \sigma_{ML}^2 \xrightarrow{a.s.} \sigma^{*2} (1 + \sqrt{\delta})^2$ .

Now if  $\delta \rightarrow 0$ , it is evident from the Lemma 1 that  $\sigma_{ML}^2 \xrightarrow{a.s.} \sigma^{*2}$ .  $\square$

**Theorem 3** Let  $d$  be fixed, and let  $n \rightarrow \infty$ :

(C1) Under assumptions (A1,A2), if  $\Psi^* = \sigma^{*2} \mathbf{I}_d$ , then

$$\hat{\mathbf{S}} \hat{\mathbf{N}} \mathbf{R}_i^{PPCA} \xrightarrow{p} \mathbf{S} \mathbf{N} \mathbf{R}_i^*,$$

for all  $i \in \{1, 2, \dots, d\}$ .

(C2) Under the assumption (A1), if  $\Psi^* = \text{diag}(\sigma_1^{*2}, \sigma_2^{*2}, \dots, \sigma_d^{*2})$  and  $\gamma_i \stackrel{i.i.d}{\sim} \mathcal{N}(0, \mathbf{I}_r)$  then

$$\hat{\mathbf{S}} \hat{\mathbf{N}} \mathbf{R}_i^{LFA} \xrightarrow{p} \mathbf{S} \mathbf{N} \mathbf{R}_i^*,$$

for all  $i \in \{1, 2, \dots, d\}$ .

Furthermore, under the assumption (A3), the probability that the  $m$  features with the highest SNRs are the true features converges to 1 as  $n \rightarrow \infty$  for both (C1) and (C2).

Here,  $\hat{\mathbf{S}} \hat{\mathbf{N}} \mathbf{R}^{PPCA}$  and  $\hat{\mathbf{S}} \hat{\mathbf{N}} \mathbf{R}^{LFA}$  denote the estimated SNRs from PPCA and LFA, respectively.

**Proof.** We have defined SNRs as:

$$\mathbf{S} \mathbf{N} \mathbf{R}_i = \frac{\sum_{j=1}^r \mathbf{W}_{ij}^2}{\sigma_i^2} = \frac{(\mathbf{W} \mathbf{W}^T)_{ii}}{\sigma_i^2}, i \in \{1, 2, \dots, d\}. \quad (7)$$

Under the assumptions (A1,A2), for a particular dimension ( $i$ ), the  $n$  observations corresponding to dimension  $i$ , denoted as  $\{\mathbf{X}_{ji}, j = 1, 2, \dots, n\}$  are i.i.d single valued random variables with

$$E(\mathbf{X}_{ji}) = 0 \quad (8)$$

$$\text{Var}(\mathbf{X}_{ji}) = \sum_{k=1}^d \lambda_k^* \mathbf{E}_{ik}^{*2} \quad (9)$$

Where  $Cov(\mathbf{X}) = \Sigma^* = \mathbf{E}^* \mathbf{\Lambda}^* \mathbf{E}^{*T}$ . Therefore, when  $\delta = 0$  (that is,  $d$  is fixed and  $n \rightarrow \infty$ ), by the law of large numbers, the corresponding sample variance, which is also the element  $i^{th}$  of the principal diagonal of  $\hat{\Sigma}$  (denoted as  $\hat{\Sigma}_{ii}$ ) converges to  $\Sigma_{ii}^*$  in probability. i.e.

$$\hat{\Sigma}_{ii} \xrightarrow{p} \Sigma_{ii}^* \quad (10)$$

When  $\Psi^* = \sigma^{*2} \mathbf{I}_d$ ,

$$\Sigma_{ii}^* = (\mathbf{W}^* \mathbf{W}^{*T})_{ii} + \sigma^{*2} \quad (11)$$

$$\hat{\Sigma}_{ii} = (\hat{\mathbf{W}} \hat{\mathbf{W}}^T)_{ii} + \hat{\sigma}^2 \quad (12)$$

From Proposition 2, we have  $\hat{\sigma}^2 \xrightarrow{a.s.} \sigma^{*2}$  for  $\delta = 0$ . Therefore combining this result with (10), from (11), we get:

$$(\hat{\mathbf{W}} \hat{\mathbf{W}}^T)_{ii} \xrightarrow{p} (\mathbf{W}^* \mathbf{W}^{*T})_{ii} \quad (13)$$

From the definition of SNR (7), we get that the estimated SNR is a continuous transformation of  $((\hat{\mathbf{W}} \hat{\mathbf{W}}^T)_{ii}, \hat{\sigma}^2)$ ,  $i = 1, 2, \dots, d$  and provided  $\hat{\sigma}^2 > 0$ . Therefore, under condition (C1), we get,  $\hat{SNR}_i \xrightarrow{PPCA} SNR_i^*$ .

Under condition (C2), when  $d$  is fixed and  $n \rightarrow \infty$ , the ML estimate of the noise covariance matrix,  $\hat{\Psi}$ , we get from [4], is consistent[2], i.e.

$$\hat{\sigma}_i^2 \xrightarrow{p} \sigma_i^{*2}, i = 1, 2, \dots, d \quad (14)$$

Also, the  $n$  observations corresponding to dimension  $i$ , denoted as  $\{\mathbf{X}_{ji}, j = 1, 2, \dots, n\}$  are i.i.d single valued random variables with  $\mathbf{X}_{ji} \stackrel{i.i.d}{\sim} \mathcal{N}(0, \sum_{k=1}^d \lambda_k^* \mathbf{E}_{ik}^{*2})$ . Therefore, by the law of large numbers, we get  $\hat{\Sigma}_{ii} \xrightarrow{p} \Sigma_{ii}^*$  and using similar logic from the condition (C1), we get,

$$(\hat{\mathbf{W}} \hat{\mathbf{W}}^T)_{ii} \xrightarrow{p} (\mathbf{W}^* \mathbf{W}^{*T})_{ii}, i = 1, 2, \dots, d \quad (15)$$

Under the condition (C2), the definition of SNR (7) suggests that the estimated SNR is a continuous transformation of  $((\hat{\mathbf{W}} \hat{\mathbf{W}}^T)_{ii}, \hat{\sigma}_i^2)$ , provided  $\min_{i \in \{1, 2, \dots, d\}} \hat{\sigma}_i^2 > 0$ . Therefore, combining the results of (14) and (15), we get:  $\hat{SNR}_i \xrightarrow{LFA} SNR_i^*$ .

Now, we prove the last part of the theorem. Let  $\hat{SNR}_i$  be the estimate of true  $SNR_i^*$ , for  $i = 1, 2, \dots, d$ . Under condition (C1 and C2), we get,  $\hat{SNR}_i \xrightarrow{p} SNR_i^*$ . Therefore, for any  $\epsilon > 0$ , there exists  $n_0$  such that for any  $n \geq n_0$

$$P(|\hat{SNR}_i - SNR_i^*| < \gamma/2) > 1 - \epsilon.$$

Therefore, with at least  $(1 - \epsilon)$  probability,

$$\begin{aligned} & SNR_i^* - \gamma/2 < \hat{SNR}_i < SNR_i^* + \gamma/2, \forall i \\ & \implies \min_{i \in S} \hat{SNR}_i > \min_{i \in S} SNR_i^* - \gamma/2, \text{ and } \max_{j \notin S} SNR_j^* + \gamma/2 > \max_{j \notin S} \hat{SNR}_j, \\ & \text{(by A3)} \implies \min_{i \in S} \hat{SNR}_i > \min_{i \in S} SNR_i^* - \gamma/2 > \max_{j \notin S} SNR_j^* + (\gamma - \gamma/2) > \max_{j \notin S} \hat{SNR}_j, \end{aligned}$$

which proves that with probability  $1 - \epsilon$ , the  $m$  features with the highest estimated SNRs are the true features  $S$  for  $n \geq n_0$ .  $\square$

We have used Proposition 3 from [11], to prove the theorem related to the generalized score  $r$ — (Theorem 4).

**Proposition 3 (due to [11])** *If  $\Sigma$  admits a rank- $r$  eigendecomposition of the form:*

$$\Sigma = \mathbf{L} \mathbf{D} \mathbf{L}^T + \lambda \mathbf{I}_m, \quad (16)$$

*with  $\mathbf{L} \in \mathbb{R}^{m \times r}$ , diagonal  $\mathbf{D} = \text{diag}(\mathbf{d}) \in \mathbb{R}^{r \times r}$  with positive entries, and  $\lambda > 0$ , the Mahalanobis distance can be computed as:*

$$MD(\mathbf{x}, \boldsymbol{\mu}, \Sigma) = r(\mathbf{x}; \boldsymbol{\mu}, \mathbf{L}, \mathbf{D}, \lambda), \quad (17)$$

*where*

$$r(\mathbf{x}; \boldsymbol{\mu}, \mathbf{L}, \mathbf{D}, \lambda) = \|\mathbf{x} - \boldsymbol{\mu}\|_2^2 / \lambda - \|\mathbf{u}(\mathbf{x})\|_2^2 / \lambda, \quad (18)$$

*with  $\mathbf{u}(\mathbf{x}) = \text{diag}(\frac{\sqrt{\mathbf{d}}}{\sqrt{\mathbf{d} + \lambda \mathbf{1}_r}}) \mathbf{L}^T (\mathbf{x} - \boldsymbol{\mu})$ , and  $\mathbf{1}_r = (1, 1, \dots, 1)^T$ . The operation  $\frac{\sqrt{\mathbf{d}}}{\sqrt{\mathbf{d} + \lambda \mathbf{1}_r}}$  is performed element-wise.*

**Theorem 4** If

$$\Sigma = \mathbf{L}\mathbf{D}\mathbf{L}^T + \Psi, \quad (19)$$

with  $\mathbf{L} \in \mathbb{R}^{m \times r}$ , the diagonal matrices  $\mathbf{D} \in \mathbb{R}^{r \times r}$  and  $\Psi \in \mathbb{R}^{m \times m}$  with positive entries, the Mahalanobis distance can be computed as:

$$MD(\mathbf{x}, \mu, \Sigma) = r(\Psi^{-\frac{1}{2}}\mathbf{x}; \Psi^{-\frac{1}{2}}\mu, \mathbf{L}', \mathbf{D}', 1) \quad (20)$$

where  $r(\mathbf{x}; \mu, \mathbf{L}, \mathbf{D}, \lambda)$  is defined in Eq. (18), and  $\mathbf{L}'$  and  $\mathbf{D}'$  are obtained by SVD on  $\Sigma' = \Psi^{-\frac{1}{2}}\Sigma\Psi^{-\frac{1}{2}}$ .

**Proof.** We consider the following transformation:

$$\begin{aligned} \mathbf{x}' &= \Psi^{-\frac{1}{2}}\mathbf{x}, \\ \mu' &= \Psi^{-\frac{1}{2}}\mu, \\ \Sigma' &= (\Psi^{-\frac{1}{2}}\mathbf{W})(\Psi^{-\frac{1}{2}}\mathbf{W}^T) + \mathbf{I}_m. \end{aligned}$$

$\Sigma'$  looks similar to (16), with  $\lambda = 1$ . Therefore, using the Proposition 3, we get:  $MD(\mathbf{x}', \mu', \Sigma') = r(\mathbf{x}'; \mu', \mathbf{L}', \mathbf{D}', 1)$ . Here,  $\Sigma' = \Psi^{-\frac{1}{2}}\Sigma\Psi^{-\frac{1}{2}} = \mathbf{L}'\mathbf{D}'\mathbf{L}'^T$ . Also,

$$\begin{aligned} MD(\mathbf{x}, \mu, \Sigma) &= (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) \\ &= (\mathbf{x} - \mu)^T (\Psi^{-\frac{1}{2}})(\Psi^{\frac{1}{2}})\Sigma^{-1}(\Psi^{\frac{1}{2}})(\Psi^{-\frac{1}{2}})(\mathbf{x} - \mu) \\ &= (\mathbf{x}' - \mu')^T (\Psi^{-\frac{1}{2}}\Sigma\Psi^{-\frac{1}{2}})^{-1} (\mathbf{x}' - \mu') \\ &= (\mathbf{x}' - \mu')^T \Sigma'^{-1} (\mathbf{x}' - \mu') \\ &= MD(\mathbf{x}', \mu', \Sigma') \\ &= r(\mathbf{x}'; \mu', \mathbf{L}', \mathbf{D}', 1) \\ &= r(\Psi^{-\frac{1}{2}}\mathbf{x}; \Psi^{-\frac{1}{2}}\mu, \mathbf{L}', \mathbf{D}', 1). \quad \square \end{aligned}$$

## References

- [1] R. Aziz, C. Verma, and N. Srivastava. A weighted-snr feature selection from independent component subspace for nb classification of microarray data. *Int J Adv Biotech Res*, 6:245–255, 2015. [1](#)
- [2] Jushan Bai and Kunpeng Li. Statistical analysis of factor models of high dimension. *Ann. Stat.*, pages 436–465, 2012. [4](#)
- [3] Jinho Baik and Jack W Silverstein. Eigenvalues of large sample covariance matrices of spiked population models. *Journal of multivariate analysis*, 97(6):1382–1408, 2006. [2](#)
- [4] Zoubin Ghahramani, Geoffrey E Hinton, et al. The em algorithm for mixtures of factor analyzers. Technical report, CRG-TR-96-1, University of Toronto, 1996. [4](#)
- [5] Seunggeun Lee, Fei Zou, and Fred A Wright. Convergence and prediction of principal component scores in high-dimensional settings. *Ann. Stat.*, 38(6):3605, 2010. [2](#)
- [6] Xi-Ao Ma, Hao Xu, Yi Liu, and Justin Zuopeng Zhang. Class-specific feature selection using fuzzy information-theoretic metrics. *Engineering Apps. of AI*, 136:109035, 2024. [1](#)
- [7] L. Nanni. Cluster-based pattern discrimination: A novel technique for feature selection. *Patt. Rec. Lett.*, 27(6), 2006. [1](#)
- [8] D. Nardone, A. Ciaramella, and A. Staiano. A sparse-modeling based approach for class specific feature selection. *PeerJ Computer Science*, 5:e237, 2019. [1](#)
- [9] B. Pineda-Bautista, J. Carrasco-Ochoa, and J. Martinez-Trinidad. General framework for class-specific feature selection. *Expert Sys. with Apps.*, 38(8):10018–10024, 2011. [1](#)
- [10] Michael E Tipping and Christopher M Bishop. Probabilistic principal component analysis. *JRSS B*, 61(3):611–622, 1999. [3](#)
- [11] Boshi Wang and Adrian Barbu. Scalable learning with incremental probabilistic pca. In *IEEE Int. Conf. on Big Data*, pages 5615–5622, 2022. [4](#)