

Supplementary Material for Unsupervised Discriminative Feature Selection With $\ell_{2,0}$ -Norm Constrained Sparse Projection

Xia Dong, Feiping Nie*, *Senior Member, IEEE*, Lai Tian, Rong Wang, and Xuelong Li, *Fellow, IEEE*

I. NOTATIONS

TABLE I: Summary of Notations

Notations	Descriptions
n	Number of samples
d	Number of features
c	Number of clusters
m	Reduced dimensionality
k	Number of selected features
$\mathbf{1}_n$	Vector with all n elements as one
$\mathbf{I}_{n \times n}$	Identity matrix with size $n \times n$
\mathbb{R}	Set of real numbers
\mathbb{Z}^+	Set of positive integers
$\text{Tr}(\mathbf{X})$	Trace of square matrix $\mathbf{X} \in \mathbb{R}^{n \times n}$
$\text{rank}(\mathbf{X})$	Rank of matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$
\mathbf{x}_i	The i -th column of matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$
x_{ij}	The (ij) -th element of matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$
\mathbf{X}^\top	Transpose of matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$
\mathbf{X}^\dagger	Moore-Penrose inverse of matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$
$\{\lambda_i(\mathbf{X})\}_{i=1}^n$	Eigenvalues of \mathbf{X} , ordered in descending order
$\ \mathbf{X}\ _F = \sqrt{\text{Tr}(\mathbf{X}^\top \mathbf{X})}$	Frobenius norm of matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$
$\ \mathbf{X}\ _{p,q} = \left(\sum_{i=1}^n \ \mathbf{x}_i\ _p^q \right)^{1/q}$	$\ell_{p,q}$ -norm of matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$
$\ \mathbf{X}\ _1 = \max_{j \in [1,d]} \sum_{i=1}^n x_{ij} $	1-norm of matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$
$\ \mathbf{X}\ _\infty = \max_{i \in [1,n]} \sum_{j=1}^d x_{ij} $	Infinity norm of matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$

II. AN EXAMPLE OF MATRIX A

To clarify the description of matrix $\mathbf{A} \in \{0, 1\}^{d \times k}$ in Section 4.1.1, we provide an example. Suppose there are $d = 6$ inputs, and we select $k = 3$ with row indices $\mathbf{q} = [2, 4, 5]$. According to the definition of the operator $\Omega_d^k(\mathbf{q})$, the corresponding row-selection matrix \mathbf{A} is:

$$\mathbf{A} = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}.$$

From this example, we see that \mathbf{A} is a sparse matrix with k columns, each containing exactly one 1 at the row index specified by \mathbf{q} , which implies that $\mathbf{A}^\top \mathbf{1}_d = \mathbf{1}_k$.

III. THE PROOF OF THEOREM 1

Proof. Suppose \mathbf{x}_* is the globally optimal solution to problem (7), with the corresponding globally minimal objective value α_* . This implies that $\frac{h(\mathbf{x}_*)}{p(\mathbf{x}_*)} = \alpha_*$. Consequently, $\forall \mathbf{x} \in \mathcal{C}$, we have $\frac{h(\mathbf{x})}{p(\mathbf{x})} \geq \alpha_*$. Since $p(\mathbf{x}) > 0$, it follows that $h(\mathbf{x}) - \alpha_* p(\mathbf{x}) \geq 0$. Moreover, noting that $h(\mathbf{x}_*) - \alpha_* p(\mathbf{x}_*) = 0$, we conclude that $\min_{\mathbf{x} \in \mathcal{C}} (h(\mathbf{x}) - \alpha_* p(\mathbf{x})) = 0$. Now, define the function $f(\alpha) = \min_{\mathbf{x} \in \mathcal{C}} (h(\mathbf{x}) - \alpha p(\mathbf{x}))$. Then, we have $f(\alpha_*) = 0$. This completes the proof. \square

IV. THE PROOF OF THEOREM 2

Proof. In Algorithm 1, we observe from lines 1–2 that $h(\mathbf{x}_t) - \alpha_t p(\mathbf{x}_t) = 0$ and $h(\mathbf{x}_{t+1}) - \alpha_t p(\mathbf{x}_{t+1}) \leq h(\mathbf{x}_t) - \alpha_t p(\mathbf{x}_t)$. Accordingly, it follows that $h(\mathbf{x}_{t+1}) - \alpha_t p(\mathbf{x}_{t+1}) \leq 0$, which implies $\frac{h(\mathbf{x}_{t+1})}{p(\mathbf{x}_{t+1})} \leq \alpha_t = \frac{h(\mathbf{x}_t)}{p(\mathbf{x}_t)}$. This indicates that Algorithm 1 guarantees the objective function of problem (7) is non-increasing at each iteration until convergence.

According to Theorem 1, the global minimum of the objective in problem (7) corresponds to the root of the function $f(\alpha)$. It is well known that Newton's method is widely regarded as an effective algorithm for root-finding under standard regularity conditions. According to line 2 of Algorithm 1, let $f(\alpha_t) = h(\mathbf{x}_{t+1}) - \alpha_t p(\mathbf{x}_{t+1})$, then the derivative is $f'(\alpha_t) = -p(\mathbf{x}_{t+1})$. Applying the Newton's update rule, we obtain

$$\alpha_{t+1} = \alpha_t - \frac{f(\alpha_t)}{f'(\alpha_t)} = \alpha_t - \frac{h(\mathbf{x}_{t+1}) - \alpha_t p(\mathbf{x}_{t+1})}{-p(\mathbf{x}_{t+1})} = \frac{h(\mathbf{x}_{t+1})}{p(\mathbf{x}_{t+1})},$$

which coincides with line 1 of Algorithm 1. Therefore, the iterative scheme in Algorithm 1 is equivalent to applying Newton's method to find the root of $f(\alpha)$. According to [1], Newton's method enjoys a quadratic convergence rate under standard regularity conditions. This completes the proof. \square

V. THE PROOF OF REMARK 1

Proof. According to problems (16) and (22), we have

$$f(\mathbf{W}_t) = \text{Tr}(\mathbf{W}_t^\top \mathbf{S}_d \mathbf{W}_t), \quad (1)$$

$$g(\mathbf{W}_t | \mathbf{W}_t) = \text{Tr} \left(\mathbf{W}_t^\top \left(\mathbf{S}_d \mathbf{W}_t (\mathbf{W}_t^\top \mathbf{S}_d \mathbf{W}_t)^\dagger \mathbf{W}_t^\top \mathbf{S}_d \right) \mathbf{W}_t \right). \quad (2)$$

It is straightforward to verify that $f(\mathbf{W}_t) = g(\mathbf{W}_t | \mathbf{W}_t)$ since $\mathbf{P} = \mathbf{P} \mathbf{P}^\dagger \mathbf{P}$ for any matrix \mathbf{P} .

Since S_d is positive semi-definite, it admits a factorization $S_d = QQ^\top$. Denote the following matrices:

$$\Upsilon = Q^\top W_t (W_t^\top S_d W_t)^\dagger W_t^\top Q, \quad (3)$$

$$\Psi = Q^\top W W^\top Q. \quad (4)$$

Then, the function $g(W_t|W_t)$ can be rewritten as

$$\begin{aligned} g(W|W_t) &= \text{Tr} \left(W^\top \left(S_d W_t (W_t^\top S_d W_t)^\dagger W_t^\top S_d \right) W \right) \\ &= \text{Tr}(\Upsilon \Psi). \end{aligned} \quad (5)$$

According to Theorems 4.3.53 and 1.3.22 [2], and noting that $\lambda_i(\Psi) \geq 0$ for all $i \in [1, m]$, we obtain

$$\text{Tr}(\Upsilon \Psi) \leq \sum_{i=1}^d \lambda_i(\Upsilon) \lambda_i(\Psi) \leq \sum_{i=1}^m \lambda_i(\Psi). \quad (6)$$

Since $\text{rank}(\Psi) \leq \text{rank}(W) = m$, we have $\sum_{i=1}^m \lambda_i(\Psi) = \text{Tr}(\Psi)$. That is, $\text{Tr}(\Upsilon \Psi) \leq \text{Tr}(\Psi) = \text{Tr}(W^\top S_d W) = f(W)$. In summary, we have $g(W|W_t) \leq f(W)$. This completes the proof. \square

VI. THE PROOF OF THEOREM 3

According to [3], we provide the proof of Theorem 3 below.

Proof. Recall that Remark 1 demonstrates that the surrogate problem (22) for optimizing W meets the condition (20) required by the majorize-minimization (MM) framework [4], [5]. Let $\tilde{W}_{t+1} = \arg \max_W g(W|W_t)$, according to Eq. (21), the following inequality holds:

$$f(\tilde{W}_{t+1}) \geq g(\tilde{W}_{t+1}|W_t) \geq g(W_t|W_t) = f(W_t). \quad (7)$$

According to Eq. (1) and inequality (7), we have

$$\text{Tr}(W_t^\top S_d W_t) \leq \text{Tr}(\tilde{W}_{t+1}^\top S_d \tilde{W}_{t+1}). \quad (8)$$

Given $\tilde{W}_{t+1} = A_{t+1} \tilde{B}_{t+1}$ and $W_{t+1} = A_{t+1} B_{t+1}$. According to problem (23), B_{t+1} maximizes its objective in the $(t+1)$ -th iteration, then we have

$$\begin{aligned} \text{Tr}(\tilde{W}_{t+1}^\top S_d \tilde{W}_{t+1}) &= \text{Tr}(\tilde{B}_{t+1}^\top A_{t+1}^\top S_d A_{t+1} \tilde{B}_{t+1}) \\ &\leq \text{Tr}(B_{t+1}^\top A_{t+1}^\top S_d A_{t+1} B_{t+1}) \\ &= \text{Tr}(W_{t+1}^\top S_d W_{t+1}). \end{aligned} \quad (9)$$

According to inequalities (8) and (9), we have

$$\text{Tr}(W_t^\top S_d W_t) \leq \text{Tr}(W_{t+1}^\top S_d W_{t+1}). \quad (10)$$

This indicates that Algorithm 2 ensures the objective of problem (16) remains non-decreasing with each iteration. Then we aim to prove that if $W_t \neq W_{t+1}$, then $\text{Tr}(W_t^\top S_d W_t) \neq \text{Tr}(W_{t+1}^\top S_d W_{t+1})$. This result demonstrates the ascent property of Algorithm 2, namely, $\text{Tr}(W_t^\top S_d W_t) < \text{Tr}(W_{t+1}^\top S_d W_{t+1})$.

Note that if $W_t \neq W_{t+1}$, then $A_t \neq A_{t+1}$, since $W = AB$ and B is formed by the leading m eigenvectors of $(A^\top S_d A)$. Therefore, suppose that there exists $A_t \neq A_{t+1}$ such that $\text{Tr}(W_t^\top S_d W_t) = \text{Tr}(W_{t+1}^\top S_d W_{t+1})$. Then the equality in inequality (6) holds. According to the equality condition in Theorem 4.3.53 [2], the matrices

$Q^\top W_t (W_t^\top S_d W_t)^\dagger W_t^\top Q$ and $Q^\top W_{t+1} W_{t+1}^\top Q$ are simultaneously diagonalizable. Assuming that S_d is full rank, we have that $\Omega_t = W_t^\top S_d W_t$ is diagonal. Define $\Phi_t = Q^\top W_t \Omega_t^{-1/2}$, then

$$Q^\top W_t (W_t^\top S_d W_t)^\dagger W_t^\top Q = Q^\top W_t \Omega_t^{-1} W_t^\top Q = \Phi_t \Phi_t^\top, \quad (11)$$

$$\Phi_t^\top \Phi_t = \Omega_t^{-1/2} W_t^\top S_d W_t \Omega_t^{-1/2} = \Omega_t^{-1/2} \Omega_t \Omega_t^{-1/2} = I_{m \times m}. \quad (12)$$

From the simultaneously diagonalizable property and Theorem 1.3.22 [2], it follows that

$$\begin{aligned} Q^\top W_{t+1} W_{t+1}^\top Q &= \Phi_t \Omega_{t+1} \Phi_t^\top \\ &= Q^\top W_t \Omega_t^{-1/2} \Omega_{t+1} \Omega_t^{-1/2} W_t^\top Q. \end{aligned} \quad (13)$$

Based on Eq. (13), we have

$$\begin{aligned} &W_t^\top S_d W_{t+1} W_{t+1}^\top S_d W_t \\ &= W_t^\top Q (Q^\top W_{t+1} W_{t+1}^\top Q) Q^\top W_t \\ &= W_t^\top Q Q^\top W_t \Omega_t^{-1/2} \Omega_{t+1} \Omega_t^{-1/2} W_t^\top Q Q^\top W_t \\ &= \Omega_t \Omega_{t+1}. \end{aligned} \quad (14)$$

We now consider the objective of the surrogate problem (22), leading to

$$\begin{aligned} &\text{Tr}(W_t^\top S_d W_{t+1} (W_{t+1}^\top S_d W_{t+1})^\dagger W_{t+1}^\top S_d W_t) \\ &= \text{Tr}(W_t^\top S_d W_{t+1} \Omega_{t+1}^{-1} W_{t+1}^\top S_d W_t) \\ &= \text{Tr}(W_t^\top Q Q^\top W_{t+1} \Omega_{t+1}^{-1} W_{t+1}^\top Q Q^\top W_t) \\ &= \text{Tr}(Q^\top W_{t+1} \Omega_{t+1}^{-1} W_{t+1}^\top Q Q^\top W_t W_t^\top Q) \\ &= \text{Tr}(\Upsilon_{t+1} \Psi_t). \end{aligned} \quad (15)$$

From Eq. (15), inequality (6), and Theorem 4.3.53 [2], we obtain

$$\begin{aligned} &\text{Tr}(W_t^\top S_d W_{t+1} \Omega_{t+1}^{-1} W_{t+1}^\top S_d W_t) \\ &\leq \text{Tr}(W_t^\top S_d W_t) = \text{Tr}(\Omega_t). \end{aligned} \quad (16)$$

Let $\Gamma = W_{t+1}^\top S_d W_t W_t^\top S_d W_{t+1} \in \mathbb{R}^{m \times m}$ and $\mathbf{U} = \Omega_{t+1}^{-1} \in \mathbb{R}^{m \times m}$, then based on Theorem 4.3.53 [2], we have

$$\begin{aligned} &\text{Tr}(W_t^\top S_d W_{t+1} \Omega_{t+1}^{-1} W_{t+1}^\top S_d W_t) \\ &= \text{Tr}(\Gamma \mathbf{U}) \geq \sum_{i=1}^m \lambda_i(\Gamma) \lambda_{m-i+1}(\mathbf{U}). \end{aligned} \quad (17)$$

Note that $\lambda_{m-i+1}(\mathbf{U}) = \lambda_i(\Omega_{t+1})^{-1}$, and by Eq. (14) and Theorem 1.3.22 [2], we get

$$\sum_{i=1}^m \lambda_i(\Gamma) \lambda_{m-i+1}(\mathbf{U}) = \sum_{i=1}^m \frac{\lambda_i(\Omega_t \Omega_{t+1})}{\lambda_i(\Omega_{t+1})} = \text{Tr}(\Omega_t). \quad (18)$$

Combining inequalities (16), (17) and Eq. (18), we conclude that $\text{Tr}(W_t^\top S_d W_{t+1} \Omega_{t+1}^{-1} W_{t+1}^\top S_d W_t) = \text{Tr}(\Omega_t)$. According to Theorem 4.3.53 [2], the equality in inequality (17) implies that Γ and \mathbf{U} are simultaneously diagonalizable. From Eq. (14), we know that $\Gamma = W_{t+1}^\top S_d W_t W_t^\top S_d W_{t+1} = \Omega_t \Omega_{t+1}$, which implies that $W_{t+1}^\top S_d W_t = \Omega_t^{1/2} \Omega_{t+1}^{1/2}$ is diagonal. Recall that

$\Phi_t = Q^\top W_t \Omega_t^{-1/2}$ and $\Phi_{t+1} = Q^\top W_{t+1} \Omega_{t+1}^{-1/2}$, then we have

$$\begin{aligned}\Phi_t^\top \Phi_t &= \Omega_t^{-1/2} W_t^\top S_d W_t \Omega_t^{-1/2} = I_{m \times m}, \\ \Phi_{t+1}^\top \Phi_{t+1} &= \Omega_{t+1}^{-1/2} W_{t+1}^\top S_d W_{t+1} \Omega_{t+1}^{-1/2} = I_{m \times m}, \\ \Phi_{t+1}^\top \Phi_t &= \Omega_{t+1}^{-1/2} W_{t+1}^\top S_d W_t \Omega_t^{-1/2} = I_{m \times m}.\end{aligned}\quad (19)$$

Thus, we conclude that $\Phi_t = \Phi_{t+1}$, which implies that $S_d W_t \Omega_t^{-1/2} = S_d W_{t+1} \Omega_{t+1}^{-1/2}$. Since S_d is full rank, it follows that $W_t \Omega_t^{-1/2} = W_{t+1} \Omega_{t+1}^{-1/2}$. Therefore, we conclude that $A_t = A_{t+1}$, since the operation $\Omega^{-1/2}$ does not affect the sparsity pattern of W . This leads to a contradiction with our initial assumption.

The above analysis establishes the non-decreasing property of the sequence $\{\text{Tr}(W_t^\top S_d W_t)\}_{t \in \mathbb{Z}^+}$. Note that $\text{Tr}(W_t^\top S_d W_t)$ is upper bounded by $\text{Tr}(S_d)$. Therefore, the sequence will eventually converge after a finite number of iterations. Moreover, we have shown that if $W_t \neq W_{t+1}$, then $\text{Tr}(W_t^\top S_d W_t) < \text{Tr}(W_{t+1}^\top S_d W_{t+1})$. Therefore, by the contrapositive, if the objective value converges, i.e., $\text{Tr}(W_t^\top S_d W_t) = \text{Tr}(W_{t+1}^\top S_d W_{t+1})$, then it must be that $W_t = W_{t+1}$ and $A_t = A_{t+1}$. Consequently, the sequence $\{W_t\}_{t \in \mathbb{Z}^+}$ converges to a fixed point \widehat{W} , and as $W_t \rightarrow \widehat{W}$, we have $\text{Tr}(W_t^\top S_d W_t) \rightarrow \text{Tr}(\widehat{W}^\top S_d \widehat{W})$. This completes the proof. \square

VII. THE PROOF OF THEOREM 4

Proof. In Algorithm 3, the optimization strategy for problem (6) consists of an outer loop and a two-layer inner loop. The two-layer inner loop involves three optimization variables: W , M , and Y . For M and Y , we can get their closed-form solutions through Eqs. (12) and (26), respectively. For W , when $\text{rank}(S_d) \leq m$, we can get the globally optimal solution through Algorithm 2; When $\text{rank}(S_d) > m$, Theorem 3 guarantees that Algorithm 2 achieves convergence in both the objective and the iterates. For the outer loop, Theorem 2 shows that the quadratic convergence of the objective in the ratio minimization problem (7) ensures the convergence of the objective in problem (6). Specifically, assume that the optimization variables in the t -th iteration are W_t , M_t , and Y_t , and in the $(t+1)$ -th iteration are W_{t+1} , M_{t+1} , and Y_{t+1} . Let $\mathcal{L}(W, M, Y) = \sum_{i=1}^n \sum_{j=1}^c y_{ij}^r \|W^\top x_i - m_j\|_2^2$ and $\mathcal{L}(W) = \text{Tr}(W^\top S_t W)$. Then, we have

$$\begin{aligned}\mathcal{L}(W_{t+1}, M_{t+1}, Y_{t+1}) - \alpha_t \mathcal{L}(W_{t+1}) \\ \leq \mathcal{L}(W_t, M_t, Y_t) - \alpha_t \mathcal{L}(W_t) = 0,\end{aligned}$$

and

$$\frac{\mathcal{L}(W_{t+1}, M_{t+1}, Y_{t+1})}{\mathcal{L}(W_{t+1})} \leq \alpha_t = \frac{\mathcal{L}(W_t, M_t, Y_t)}{\mathcal{L}(W_t)}.$$

This indicates that Algorithm 3 ensures the objective of problem (6) remains non-increasing with each iteration and ultimately converges at the objective level. This completes the proof. \square

VIII. THE PROOF OF LEMMA 1

According to [3], we provide the proof of Lemma 1 below.

Proof. Let $S_d = S_d^m + S_d^c$. We then have the following derivation:

$$\begin{aligned}\text{Tr}(W_*^\top S_d W_*) &= \max_{W^\top W = I_{m \times m}, \|W\|_{2,0}=k} \text{Tr}(W^\top (S_d^m + S_d^c) W) \\ &\leq \max_{\substack{W^\top W = I_{m \times m}, \\ \|W\|_{2,0}=k}} \text{Tr}(W^\top S_d^m W) + \max_{\substack{W^\top W = I_{m \times m}, \\ \|W\|_{2,0}=k}} \text{Tr}(W^\top S_d^c W) \\ &\leq \text{Tr}(W_m^\top S_d^m W_m) + \max_{W^\top W = I_{m \times m}} \text{Tr}(W^\top S_d^c W) \\ &\leq \text{Tr}(W_m^\top S_d^m W_m) + \sum_{i=m+1}^{2m} \lambda_i(S_d),\end{aligned}\quad (20)$$

from which it follows that

$$\frac{\text{Tr}(W_m^\top S_d^m W_m)}{\text{Tr}(W_*^\top S_d W_*)} \geq 1 - \frac{\sum_{i=m+1}^{2m} \lambda_i(S_d)}{\text{Tr}(W_*^\top S_d W_*)}.\quad (21)$$

Furthermore, we have the following two observations:

$$\begin{aligned}\text{Tr}(W_*^\top S_d W_*) &\stackrel{(a)}{\geq} \max_{W^\top W = I_{m \times m}, \|W\|_{2,0}=m} \text{Tr}(W^\top S_d W) \\ &\stackrel{(b)}{\geq} \frac{m}{d} \text{Tr}(S_d) = \frac{m}{d} \sum_{i=1}^d \lambda_i(S_d),\end{aligned}\quad (22)$$

$$\begin{aligned}\text{Tr}(W_*^\top S_d W_*) &\stackrel{(c)}{\geq} \max_{W^\top W = I_{m \times m}, \|W\|_{2,0}=k} \text{Tr}(W^\top S_d^m W) \\ &\stackrel{(d)}{\geq} \frac{k}{d} \sum_{i=1}^d \lambda_i(S_d^m) = \frac{k}{d} \sum_{i=1}^m \lambda_i(S_d),\end{aligned}\quad (23)$$

where (a) holds because $m \leq k$, and (c) holds since S_d is positive semi-definite. Inequality (b) holds since $k = m$, implying that the problem on the left side of (b) achieves its global optimum via Algorithm 2. Inequality (d) holds because $\text{rank}(S_d^m) = m$, and thus the problem on the left side of (d) also achieves its global optimum via Algorithm 2.

Let $z = \min\{\text{rank}(S_d), 2m\}$, $c_1 = \frac{\sum_{i=m+1}^z \lambda_i(S_d)}{\sum_{i=1}^m \lambda_i(S_d)}$, and $c_2 = \frac{\sum_{i=m+1}^z \lambda_i(S_d)}{\sum_{i=1}^d \lambda_i(S_d)}$. Combining inequalities (21), (22), (23), we obtain

$$1 \geq \frac{\text{Tr}(W_m^\top S_d^m W_m)}{\text{Tr}(W_*^\top S_d W_*)} \geq 1 - \min\left\{\frac{d \cdot c_1}{k}, \frac{d \cdot c_2}{m}\right\}.\quad (24)$$

Since S_d is positive semi-definite, we have

$$\begin{aligned}\text{Tr}(W_m^\top S_d W_m) &= \text{Tr}(W_m^\top S_d^m W_m) + \text{Tr}(W_m^\top S_d^c W_m) \\ &\geq \text{Tr}(W_m^\top S_d^m W_m).\end{aligned}\quad (25)$$

Combining inequalities (24) and (25), we obtain $\varepsilon \leq \min\{(d \cdot c_1/k), (d \cdot c_2/m)\}$. Moreover, according to Theorem 4.3.53 [2], we have

$$\text{Tr}(W_m^\top S_d W_m) \geq \sum_{i=d-m+1}^d \lambda_i(S_d) \geq m \cdot \lambda_d(S_d).\quad (26)$$

Let $\kappa = \lambda_1(\mathbf{S}_d)/\lambda_d(\mathbf{S}_d)$, then by Ky Fan's Theorem [6], we obtain

$$\begin{aligned} \text{Tr}(\mathbf{W}_*^\top \mathbf{S}_d \mathbf{W}_*) &\leq \max_{\mathbf{W}^\top \mathbf{W} = \mathbf{I}_{m \times m}} \text{Tr}(\mathbf{W}^\top \mathbf{S}_d \mathbf{W}) = \sum_{i=1}^m \lambda_i(\mathbf{S}_d) \\ &\leq m \cdot \lambda_1(\mathbf{S}_d) = m \cdot \kappa \cdot \lambda_d(\mathbf{S}_d). \end{aligned} \quad (27)$$

Combing inequalities (26) and (27), we obtain $\varepsilon \leq 1 - \kappa^{-1}$. Furthermore, we have

$$\begin{aligned} \text{Tr}(\mathbf{W}_m^\top \mathbf{S}_d \mathbf{W}_m) &\geq \text{Tr}(\mathbf{W}_m^\top \mathbf{S}_d^m \mathbf{W}_m) \\ &\geq \frac{k}{d} \text{Tr}(\mathbf{S}_d^m) \geq \frac{k}{d} \sum_{i=1}^m \lambda_i(\mathbf{S}_d). \end{aligned} \quad (28)$$

According to inequality (27), we have

$$\text{Tr}(\mathbf{W}_*^\top \mathbf{S}_d \mathbf{W}_*) \leq \sum_{i=1}^m \lambda_i(\mathbf{S}_d). \quad (29)$$

Combing inequalities (28) and (29), we obtain $\varepsilon \leq 1 - k/d$. In summary, we conclude that $\varepsilon \leq \min\{(d \cdot c_1/k), (d \cdot c_2/m), 1 - \kappa^{-1}, 1 - k/d\}$. This completes the proof. \square

REFERENCES

- [1] H. W. Kuhn and A. W. Tucker, "Nonlinear programming," in *Traces and Emergence of Nonlinear Programming*. Springer, 2014, pp. 247–258.
- [2] R. A. Horn and C. R. Johnson, *Matrix analysis*, 2nd ed. Cambridge; New York: Cambridge University Press, 2013.
- [3] L. Tian, F. Nie, and X. Li, "Learning feature sparse principal components," *arXiv preprint arXiv:1904.10155*, 2019.
- [4] Y. Sun, P. Babu, and D. P. Palomar, "Majorization-minimization algorithms in signal processing, communications, and machine learning," *IEEE Transactions on Signal Processing*, vol. 65, no. 3, pp. 794–816, 2016.
- [5] F. Nie, D. Wu, R. Wang, and X. Li, "Truncated robust principle component analysis with a general optimization framework," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 2, pp. 1081–1097, 2022.
- [6] K. Fan, "On a theorem of weyl concerning eigenvalues of linear transformations: II," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 36, no. 1, p. 31, 1950.