

Introduction to Machine Learning
Fall 2019
University of Science and Technology of China

Lecturer: Jie Wang
Posted: Nov. 27, 2019
Name: San Zhang

Homework 6
Due: Dec. 11, 2019
ID: PBXXXXXXXX

Notice, to get the full credits, please show your solutions step by step.

Exercise 1: 10pts

Show that $(1 - \epsilon)^m \leq e^{-m\epsilon}$, where $m \in \mathbb{N}^+$ and $0 \leq \epsilon < 1$.

Solution: Let $f(x) = e^{-x}$. The first-order Taylor expansion of f at 0 with lagrangian remainder is

$$e^{-x} = 1 - x + \frac{1}{2!}e^{-\xi}x^2 \text{ (assume } x \geq 0 \text{)}$$

for some $\xi \in [0, x]$. Then for all $\epsilon \in [0, 1]$,

$$e^{-\epsilon} \geq 1 - \epsilon > 0.$$

Therefore, for all $m \in \mathbb{N}^+$, $\epsilon \in [0, 1]$,

$$(1 - \epsilon)^m \leq e^{-m\epsilon}.$$

■

Exercise 2: Markov inequality 10pts

Let X be a nonnegative random variable on \mathbb{R} . Then, for all $t > 0$, show that

$$\mathbf{P}(X \geq t) \leq \frac{\mathbf{E}[X]}{t}.$$

You can assume that X is a continuous random variable.

Solution: Let $f_X(x)$ be the probability density function of X . We have

$$\begin{aligned}\mathbf{E}[X] &= \int_{-\infty}^{+\infty} x f_X(x) dx \\ &= \int_0^{+\infty} x f_X(x) dx \\ &= \int_0^t x f_X(x) dx + \int_t^{+\infty} x f_X(x) dx \\ &\geq \int_t^{+\infty} x f_X(x) dx \\ &\geq t \int_t^{+\infty} f_X(x) dx \\ &= t \mathbf{P}(X \geq t).\end{aligned}$$

Thus,

$$\mathbf{P}(X \geq t) \leq \frac{\mathbf{E}[X]}{t}.$$

■

Exercise 3: VC-dimension 10pts

Assume that the instance space $X = \mathbb{R}^2$ and the hypothesis space H be the set of all linear threshold functions defined on \mathbb{R}^2 . Find $VC(H)$ and prove it.

Solution: We show that $VC(H) = 3$.

For convenience, we assume that the labels are -1 and 1 , and $\text{sign}(0) \triangleq 1$.

Suppose that x is a sample, y is the corresponding label, and \hat{y} is the predicted label. WLOG, we assume that x is a column vector in \mathbb{R}^2 .

A linear classifier can be formulated as $\hat{y} = \text{sign}(w^T x + b)$, where $w \in \mathbb{R}^{2 \times 1}$ and $b \in \mathbb{R}$. We can rewrite it as $\hat{y} = \text{sign}(\tilde{w}^T \tilde{x})$, where $\tilde{w} = [w; b]$ and $\tilde{x} = [x; 1]$.

Let x_1, x_2, x_3 be three different samples, then we can get the predicted labels

$$[\hat{y}_1, \hat{y}_2, \hat{y}_3] = \text{sign}(\tilde{w}^T [\tilde{x}_1, \tilde{x}_2, \tilde{x}_3]).$$

Since $X_1 \triangleq [\tilde{x}_1, \tilde{x}_2, \tilde{x}_3]$ is a 3×3 matrix, if $\text{rank}(X_1) = 3$, then for any possible labels $[y_1, y_2, y_3]$, there exists a \tilde{w}_0 , such that

$$[y_1, y_2, y_3] = \text{sign}(\tilde{w}_0^T [\tilde{x}_1, \tilde{x}_2, \tilde{x}_3]) = [\hat{y}_1, \hat{y}_2, \hat{y}_3].$$

That is, $VC(H) \geq 3$.

Let x_1, x_2, x_3, x_4 be four (not necessarily different) samples, then we can get the predicted labels

$$[\hat{y}_1, \hat{y}_2, \hat{y}_3, \hat{y}_4] = \text{sign}(\tilde{w}^T [\tilde{x}_1, \tilde{x}_2, \tilde{x}_3, \tilde{x}_4]).$$

Let $X_2 \triangleq [\tilde{x}_1, \tilde{x}_2, \tilde{x}_3, \tilde{x}_4]$. Since X_2 is a 3×4 matrix, $\text{rank}(X_2) \leq 3 < 4$. We can write at least one of \tilde{x}_i ($1 \leq i \leq 4$) as a linear combination of the other vectors. WLOG, we assume that $\tilde{x}_4 = \alpha_1 \tilde{x}_1 + \alpha_2 \tilde{x}_2 + \alpha_3 \tilde{x}_3$, where not all of α_i are zero. Then we have

$$\begin{aligned} \tilde{w}^T \tilde{x}_4 &= \tilde{w}^T (\alpha_1 \tilde{x}_1 + \alpha_2 \tilde{x}_2 + \alpha_3 \tilde{x}_3) \\ &= \alpha_1 (\tilde{w}^T \tilde{x}_1) + \alpha_2 (\tilde{w}^T \tilde{x}_2) + \alpha_3 (\tilde{w}^T \tilde{x}_3). \end{aligned}$$

Let the true labels

$$[y_1, y_2, y_3] \triangleq [\text{sign}(\alpha_1), \text{sign}(\alpha_2), \text{sign}(\alpha_3)].$$

If the classifier can correctly classify x_1, x_2 and x_3 , then

$$[\text{sign}(\tilde{w}^T \tilde{x}_1), \text{sign}(\tilde{w}^T \tilde{x}_2), \text{sign}(\tilde{w}^T \tilde{x}_3)] = [\text{sign}(\alpha_1), \text{sign}(\alpha_2), \text{sign}(\alpha_3)].$$

Thus,

$$\alpha_i (\tilde{w}^T \tilde{x}_i) \geq 0, \quad i = 1, 2, 3,$$

which implies that

$$\hat{y}_4 = [\text{sign}(\tilde{w}^T \tilde{x}_4)] = 1.$$

Now if we let the true label $y_4 = -1$, then there is a contradictory.

Thus, $VC(H) < 4$.

From the above all, we know that $VC(H) = 3$. ■

Exercise 4: Learning intervals 20pts

Let the target concept class be $C = \{[a, b] : a < b, a, b \in \mathbb{R}\}$ and the hypotheses class $H = C$, and the version space be $VS_{H,D}$. Each $c \in C$ labels the points inside the interval positive and the others negative. A consistent learner will pick a consistent hypothesis—if any— $h \in H$ according to a set of i.i.d. samples $\{(x_1, c(x_1)), (x_2, c(x_2)), \dots, (x_m, c(x_m))\}$ that obey an unknown absolute continuous distribution \mathcal{D} . \mathcal{D} 's p.d.f. is $p(x)$. Please find

$$\mathbf{P}[\exists h \in VS_{H,D} \text{ and } error_{\mathcal{D}}(h) > \epsilon],$$

and the corresponding sample complexity.

Solution: Let \hat{a} and \hat{b} be the estimation of a and b .

We define ϵ_1, ϵ_2 in the following way:

$$\epsilon_1 = \begin{cases} \int_a^{\hat{a}} p(x) dx, & \hat{a} \geq a \\ \int_{\hat{a}}^a p(x) dx, & \hat{a} < a \end{cases}$$

$$\epsilon_2 = \begin{cases} \int_b^{\hat{b}} p(x) dx, & \hat{b} \geq b \\ \int_{\hat{b}}^b p(x) dx, & \hat{b} < b \end{cases}$$

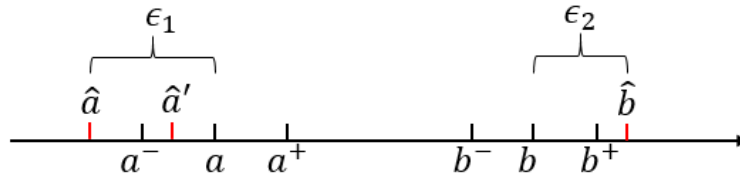
The event “ $\exists h \in VS_{H,D}$ and $error_{\mathcal{D}}(h) > \epsilon$ ” is equivalent to the event “ $\epsilon_1 + \epsilon_2 > \epsilon$ ”. Note that “ $\epsilon_1 + \epsilon_2 > \epsilon$ ” \subset “ $(\epsilon_1 > \frac{\epsilon}{2}) \vee (\epsilon_2 > \frac{\epsilon}{2})$ ” since “ $(\epsilon_1 \leq \frac{\epsilon}{2}) \wedge (\epsilon_2 \leq \frac{\epsilon}{2}) \Rightarrow \epsilon_1 + \epsilon_2 \leq \epsilon$ ”. We have the following inequality:

$$\begin{aligned} \mathbf{P}[\exists h \in VS_{H,D} \text{ and } error_{\mathcal{D}}(h) > \epsilon] &= \mathbf{P}[\epsilon_1 + \epsilon_2 > \epsilon] \\ &\leq \mathbf{P}[\epsilon_1 > \frac{\epsilon}{2}] + \mathbf{P}[\epsilon_2 > \frac{\epsilon}{2}] \end{aligned}$$

We also define the a^+, a^-, b^+, b^- in the following way:

$$a^+ = \inf\{x \geq a : \int_a^x p(t) dt \geq \frac{\epsilon}{2}\}; \quad a^- = \sup\{x \leq a : \int_x^a p(t) dt \geq \frac{\epsilon}{2}\};$$

$$b^+ = \inf\{x \geq b : \int_b^x p(t) dt \geq \frac{\epsilon}{2}\}; \quad b^- = \sup\{x \leq b : \int_x^b p(t) dt \geq \frac{\epsilon}{2}\};$$



We also have the following conclusions, and the picture above may help you understand:

$$“\epsilon_1 > \frac{\epsilon}{2}” \subset “all \text{ samples } \notin [a^-, a^+]”$$

$$“\epsilon_2 > \frac{\epsilon}{2}” \subset “all \text{ samples } \notin [b^-, b^+]”$$

Thus we have the following inequalities:

$$\mathbf{P}[\epsilon_1 > \frac{\epsilon}{2}] \leq \mathbf{P}[\text{all samples} \notin [a^-, a^+]]$$

$$\mathbf{P}[\epsilon_2 > \frac{\epsilon}{2}] \leq \mathbf{P}[\text{all samples} \notin [b^-, b^+]]$$

Then according to the above inequalities, we will have:

$$\begin{aligned} \mathbf{P}[\exists h \in VS_{H,D} \text{ and } error_{\mathcal{D}}(h) > \epsilon] &= \mathbf{P}[\epsilon_1 + \epsilon_2 > \epsilon] \\ &\leq \mathbf{P}[\epsilon_1 > \frac{\epsilon}{2}] + \mathbf{P}[\epsilon_2 > \frac{\epsilon}{2}] \\ &\leq \mathbf{P}[\text{all samples} \notin [a^-, a^+]] + \mathbf{P}[\text{all samples} \notin [b^-, b^+]] \\ &\leq 2e^{-m\frac{\epsilon}{2}} + 2e^{-m\frac{\epsilon}{2}} \\ &= 4e^{-m\frac{\epsilon}{2}}. \end{aligned}$$

If let $4e^{-m\epsilon/2} \leq \delta$, then we have

$$m \geq \frac{2}{\epsilon} \ln \frac{4}{\delta}.$$

■

Exercise 5: Basic Matrix Manipulations 20pts

For an arbitrary matrix M , we denote its i^{th} row, j^{th} column, and $(i, j)^{th}$ entry by $\mathbf{m}_{i,*}$, $\mathbf{m}_{*,j}$, and $m_{i,j}$, respectively.

1. Suppose that $A \in \mathbb{R}^{m \times n}$, $B \in \mathbb{R}^{m \times d}$, $C \in \mathbb{R}^{d \times n}$, and $A = BC$. Show that

$$A = \sum_{\ell=1}^d \mathbf{b}_{*,\ell} \mathbf{c}_{\ell,*}.$$

2. Suppose that $A \in \mathbb{R}^{m \times n}$, $B \in \mathbb{R}^{m \times p}$, $C \in \mathbb{R}^{p \times q}$, $D \in \mathbb{R}^{q \times n}$, and $A = BCD$. Show that

$$A = \sum_{i=1}^p \sum_{j=1}^q c_{i,j} \mathbf{b}_{*,i} \mathbf{d}_{j,*}.$$

Solution: 1. Suppose that

$$\begin{aligned} B &= (\mathbf{b}_{*,1}, \mathbf{b}_{*,2}, \dots, \mathbf{b}_{*,d}), \\ C &= (\mathbf{c}_{1,*}^\top, \mathbf{c}_{2,*}^\top, \dots, \mathbf{c}_{d,*}^\top)^\top. \end{aligned}$$

Suppose that

$$\begin{aligned} A &= BC \\ &= \sum_{\ell=1}^d \mathbf{b}_{*,\ell} \mathbf{c}_{\ell,*}. \end{aligned}$$

2. Suppose that

$$\begin{aligned} B &= (\mathbf{b}_{*,1}, \mathbf{b}_{*,2}, \dots, \mathbf{b}_{*,p}), \\ D &= (\mathbf{d}_{1,*}^\top, \mathbf{d}_{2,*}^\top, \dots, \mathbf{d}_{q,*}^\top)^\top. \end{aligned}$$

Then

$$\begin{aligned}
 A &= BCD \\
 &= (\mathbf{b}_{*,1}, \mathbf{b}_{*,2}, \dots, \mathbf{b}_{*,p}) \begin{bmatrix} c_{1,1} & \dots & c_{1,q} \\ \dots & \dots & \dots \\ c_{p,1} & \dots & c_{p,q} \end{bmatrix} (\mathbf{d}_{1,*}^\top, \mathbf{d}_{2,*}^\top, \dots, \mathbf{d}_{q,*}^\top)^\top \\
 &= \left(\sum_{i=1}^p \mathbf{b}_{*,i} c_{i,1}, \dots, \sum_{i=1}^p \mathbf{b}_{*,i} c_{i,q} \right) (\mathbf{d}_{1,*}^\top, \mathbf{d}_{2,*}^\top, \dots, \mathbf{d}_{q,*}^\top)^\top \\
 &= \sum_{j=1}^q \sum_{i=1}^p \mathbf{b}_{*,i} c_{i,j} \mathbf{d}_{j,*} \\
 &= \sum_{i=1}^p \sum_{j=1}^q c_{i,j} \mathbf{b}_{*,i} \mathbf{d}_{j,*}.
 \end{aligned}$$

■

Exercise 6: Subspace 90pts

The column space of a matrix $A \in \mathbb{R}^{m \times n}$ is the set

$$\mathcal{C}(A) = \{\mathbf{y} \in \mathbb{R}^m : \mathbf{y} = A\mathbf{x}, \mathbf{x} \in \mathbb{R}^n\}. \quad (1)$$

1. Let $A \in \mathbb{R}^{m \times n}$, $B \in \mathbb{R}^{n \times p}$, and $C = AB$.
 - (a) Show that $\mathcal{C}(C) \subseteq \mathcal{C}(A)$.
 - (b) Suppose that B is nonsingular, that is, B is invertible. Show that $\mathcal{C}(C) = \mathcal{C}(A)$.
2. Suppose that $A \in \mathbb{R}^{m \times n}$ has full column rank, that is, the column vectors in A are linearly independent. Let $\mathbf{x} \in \mathbb{R}^m$ and

$$P_{\mathcal{C}(A)}(\mathbf{x}) := \underset{\mathbf{z} \in \mathbb{R}^m}{\operatorname{argmin}} \{\|\mathbf{x} - \mathbf{z}\|_2 : \mathbf{z} \in \mathcal{C}(A)\}. \quad (2)$$

- (a) Is $P_{\mathcal{C}(A)}$ unique? If so, please justify your answer and find $P_{\mathcal{C}(A)}$; otherwise, please find all the projections.
 - (b) What are the coordinates of $P_{\mathcal{C}(A)}$ with respect to the column vectors in A ? Are the coordinates unique?
3. Suppose that the column vectors in $A \in \mathbb{R}^{m \times n}$ are orthonormal.
 - (a) Please answer the questions in 2.
 - (b) Suppose that the column vectors in $\tilde{A} \in \mathbb{R}^{m \times n}$ are also orthonormal, and $\mathcal{C}(A) = \mathcal{C}(\tilde{A})$. Show that $P_{\mathcal{C}(A)}(\mathbf{x}) = P_{\mathcal{C}(\tilde{A})}(\mathbf{x})$ for any $\mathbf{x} \in \mathbb{R}^m$.
4. Suppose that the column vectors in $A \in \mathbb{R}^{m \times n}$ are linearly dependent.
 - (a) Is $P_{\mathcal{C}(A)}$ unique? If so, please justify your answer; otherwise, please find all the projections.
 - (b) Are the coordinates of $P_{\mathcal{C}(A)}$ with respect to the column vectors in A unique? If so, please justify your answer; otherwise, please find all the possible coordinates.

Hint: you may assume that the first r column vectors with $r < n$ are a basis of $\mathcal{C}(A)$.

Solution: 1. (a) Suppose that $A = (\mathbf{a}_1, \dots, \mathbf{a}_n)$ and

$$B = \begin{bmatrix} b_{11} & \dots & b_{1p} \\ \dots & \dots & \dots \\ b_{n1} & \dots & b_{np} \end{bmatrix}.$$

Then

$$\begin{aligned} C &= AB \\ &= \left(\sum_{i=1}^n b_{i1} \mathbf{a}_i, \dots, \sum_{i=1}^n b_{ip} \mathbf{a}_i \right). \end{aligned}$$

For all $\mathbf{v} \in \mathcal{C}(C)$, there exist d_1, \dots, d_p such that

$$\begin{aligned}\mathbf{v} &= \sum_{j=1}^p d_j \left(\sum_{i=1}^n b_{ij} \mathbf{a}_i \right) \\ &= \sum_{i=1}^n \left(\sum_{j=1}^p d_j b_{ij} \right) \mathbf{a}_i.\end{aligned}$$

Thus, $\mathbf{v} \in \mathcal{C}(A)$, i.e. $\mathcal{C}(C) \subseteq \mathcal{C}(A)$.

(b) We have proven that $\mathcal{C}(C) \subseteq \mathcal{C}(A)$.

We next show that $\mathcal{C}(A) \subseteq \mathcal{C}(C)$. As $C = AB$ and B is invertible,

$$A = CB^{-1}.$$

Then it is easy to show that

$$\mathcal{C}(A) \subseteq \mathcal{C}(C).$$

Therefore, $\mathcal{C}(C) = \mathcal{C}(A)$.

2. (a) $P_{\mathcal{C}(A)}$ is unique.

We first show that $P_{\mathcal{C}(A)}$ is unique. Suppose that there exist two different points $\mathbf{z}_1, \mathbf{z}_2 \in \mathcal{C}(A)$ such that

$$\|\mathbf{x} - \mathbf{z}_1\|_2 = \|\mathbf{x} - \mathbf{z}_2\|_2 = \min_{\mathbf{z} \in \mathcal{C}(A)} \|\mathbf{x} - \mathbf{z}\|_2$$

Let $\mathbf{z}_0 = \frac{\mathbf{z}_1 + \mathbf{z}_2}{2}$, then

$$\begin{aligned}\left\| \mathbf{x} - \frac{\mathbf{z}_1 + \mathbf{z}_2}{2} \right\|_2^2 &< \|\mathbf{x} - \mathbf{z}_1\|_2^2 \\ \text{and } \frac{\mathbf{z}_1 + \mathbf{z}_2}{2} &\in \mathcal{C}(A).\end{aligned}$$

which leads to a contradiction. Thus, $P_{\mathcal{C}(A)}$ is unique.

We next find $P_{\mathcal{C}(A)}$.

Assume that $\{\mathbf{a}_1, \dots, \mathbf{a}_n, \mathbf{a}_{n+1}, \dots, \mathbf{a}_m\}$ is a basis of \mathbb{R}^m , where $\langle \mathbf{a}_i, \mathbf{a}_j \rangle = 0$ for all $1 \leq i \leq n$ and $n+1 \leq j \leq m$. For all $\mathbf{y} \in \mathcal{C}(A)$, there exist $\lambda_y^i \in \mathbb{R}$, such that $\mathbf{y} = \sum_{i=1}^n \lambda_y^i \mathbf{a}_i$. For all $\mathbf{x} \in \mathbb{R}^m$, there exist $\lambda_x^i \in \mathbb{R}$, such that $\mathbf{x} = \sum_{i=1}^m \lambda_x^i \mathbf{a}_i$. Let $\boldsymbol{\lambda}_x = (\lambda_x^1, \dots, \lambda_x^m)^\top$. Then, $\mathbf{x} = A\boldsymbol{\lambda}_x$ and $\boldsymbol{\lambda}_x = (A^T A)^{-1} A^T \mathbf{x}$.

Suppose that $\mathbf{z} = P_{C(A)}(\mathbf{x})$, then

$$\begin{aligned}
 \|\mathbf{x} - \mathbf{z}\|_2^2 &= \left\| \sum_{i=1}^m \lambda_x^i \mathbf{a}_i - \sum_{i=1}^n \lambda_z^i \mathbf{a}_i \right\|_2^2 \\
 &= \min_{\lambda_y^i \in \mathbb{R}} \left\| \sum_{i=1}^m \lambda_x^i \mathbf{a}_i - \sum_{i=1}^n \lambda_y^i \mathbf{a}_i \right\|_2^2 \\
 &= \min_{\lambda_y^i \in \mathbb{R}} \left\| \sum_{i=1}^n (\lambda_x^i - \lambda_y^i) \mathbf{a}_i + \sum_{i=n+1}^m \lambda_x^i \mathbf{a}_i \right\|_2^2 \\
 &= \min_{\lambda_y^i \in \mathbb{R}} \left\| \sum_{i=1}^n (\lambda_x^i - \lambda_y^i) \mathbf{a}_i \right\|_2^2 + \left\| \sum_{i=n+1}^m \lambda_x^i \mathbf{a}_i \right\|_2^2
 \end{aligned}$$

where the last equality holds for $\langle \mathbf{a}_i, \mathbf{a}_j \rangle = 0$, $1 \leq i \leq n$ and $n+1 \leq j \leq m$.

It's clear that $\lambda_z^i = \lambda_x^i$, $i = 1, \dots, n$, i.e.,

$$\begin{aligned}
 \mathbf{z} &= A\boldsymbol{\lambda}_z \\
 &= A\boldsymbol{\lambda}_x \\
 &= A(A^T A)^{-1} A^T \mathbf{x}
 \end{aligned}$$

Hence we know that $P_{C(A)} = A(A^T A)^{-1} A^T$.

(b) The coordinates are unique.

For all $\mathbf{x} \in \mathbb{R}^m$,

$$P_{C(A)}(\mathbf{x}) = A(A^T A)^{-1} A^T \mathbf{x},$$

which is a vector in \mathbb{R}^m . Suppose that $\alpha = (\alpha_1, \dots, \alpha_n)^T$ are the coordinates of $P_{C(A)}(\mathbf{x})$ with respect to the column vectors in A , then

$$\begin{aligned}
 A\alpha &= A(A^T A)^{-1} A^T \mathbf{x} \\
 \Rightarrow A(\alpha - (A^T A)^{-1} A^T \mathbf{x}) &= \mathbf{0}
 \end{aligned}$$

Since A has full column rank,

$$\alpha = (A^T A)^{-1} A^T \mathbf{x}.$$

3. (a) It follows from Exercise 6.2 that

$$P_{C(A)} = A(A^T A)^{-1} A^T$$

and

$$\alpha = (A^T A)^{-1} A^T \mathbf{x}.$$

Since the column vectors in A are orthonormal,

$$A^T A = I.$$

Thus,

$$P_{\mathcal{C}(A)} = AA^\top$$

and

$$\alpha = A^\top \mathbf{x}$$

- (b) We can prove that $P_{\mathcal{C}(A)}(x) = P_{\mathcal{C}(\tilde{A})}(x)$ in the same way as 2.(a).

Suppose that there exist two different points $\mathbf{z}_1, \mathbf{z}_2 \in \mathcal{C}(A)$ such that

$$\|\mathbf{x} - \mathbf{z}_1\|_2 = \|\mathbf{x} - \mathbf{z}_2\|_2 = \min_{\mathbf{z} \in \mathcal{C}(A)} \|\mathbf{x} - \mathbf{z}\|_2$$

Let $\mathbf{z}_0 = \frac{\mathbf{z}_1 + \mathbf{z}_2}{2}$, then

$$\begin{aligned} \left\| \mathbf{x} - \frac{\mathbf{z}_1 + \mathbf{z}_2}{2} \right\|_2^2 &< \|\mathbf{x} - \mathbf{z}_1\|_2^2 \\ \text{and } \frac{\mathbf{w}_1 + \mathbf{w}_2}{2} &\in \mathcal{C}(A). \end{aligned}$$

This is a contradiction. Thus, $P_{\mathcal{C}(A)}(x) = P_{\mathcal{C}(\tilde{A})}(x)$.

4. (a) The uniqueness comes from Exercise 6.2.

We next find the projection. Suppose that $\mathbf{rank}(A) = r$ and the first r columns of A are linearly independent. Let $A = (\mathbf{a}_1, \dots, \mathbf{a}_n)$ and $A_r = (\mathbf{a}_1, \dots, \mathbf{a}_r)$. Then

$$\mathcal{C}(A) = \mathcal{C}(A_r)$$

and $A_r \in \mathbb{R}^{m \times r}$ has full column rank. Then

$$\begin{aligned} P_{\mathcal{C}(A)} &= P_{\mathcal{C}(A_r)} \\ &= A_r(A_r^\top A_r)^{-1} A_r^\top. \end{aligned}$$

- (b) They may not be unique.

For example, if $A = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix}$, then for $\mathbf{x} = (3, 3)^\top$, we know that $(3, 3, 0)^\top$ and $(0, 0, 3)^\top$ are two possible coordinates.

■

Exercise 7: SVD 80pts

Let $A \in \mathbb{R}^{m \times n}$, $\mathbf{rank}(A) = r$, its SVD be $A = U\Sigma V^\top$, where we sort the diagonal entries of Σ in the descending order $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$, and

$$U_1 = (\mathbf{u}_{*,1}, \mathbf{u}_{*,2}, \dots, \mathbf{u}_{*,r}), U_2 = (\mathbf{u}_{*,r+1}, \dots, \mathbf{u}_{*,m}), \\ V_1 = (\mathbf{v}_{*,1}, \mathbf{v}_{*,2}, \dots, \mathbf{v}_{*,r}), V_2 = (\mathbf{v}_{*,r+1}, \dots, \mathbf{v}_{*,n}).$$

We define the column space of a matrix A in (3). The null space of A is the set

$$\mathcal{N}(A) = \{\mathbf{y} \in \mathbb{R}^n : A\mathbf{y} = 0\}. \quad (3)$$

1. Show that

- (a) $P_{\mathcal{C}(A)}(\mathbf{x}) = U_1 U_1^\top \mathbf{x}$;
- (b) $P_{\mathcal{N}(A)}(\mathbf{x}) = V_2 V_2^\top \mathbf{x}$;
- (c) $P_{\mathcal{C}(A^\top)}(\mathbf{x}) = V_1 V_1^\top \mathbf{x}$;
- (d) $P_{\mathcal{N}(A^\top)}(\mathbf{x}) = U_2 U_2^\top \mathbf{x}$.

2. The Frobenius norm of A is

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n a_{i,j}^2}.$$

- (a) Show that $\|A\|_F^2 = \mathbf{tr}(A^\top A)$.
- (b) Let $B \in \mathbb{R}^{m \times n}$. Suppose that $\mathcal{C}(A) \perp \mathcal{C}(B)$, that is,

$$\langle \mathbf{a}, \mathbf{b} \rangle = 0, \forall \mathbf{a} \in \mathcal{C}(A), \mathbf{b} \in \mathcal{C}(B).$$

Show that

$$\|A + B\|_F^2 = \|A\|_F^2 + \|B\|_F^2.$$

3. Please solve the problem as follows.

$$\min_{X \in \mathbb{R}^{m \times n}} \{\|A - X\|_F : \mathbf{rank}(X) \leq K\}. \quad (4)$$

For simplicity, you can assume that all singular values of A are different.

4. **Programming Exercise** We provide you a grayscale image (“Alan_Turing.jpg”). Suppose that A is the data matrix of the image. We have $A \in \mathbb{R}^{512 \times 512}$ and $r = \mathbf{rank}(A) = 512$. In this exercise, you are expected to implement an image compression algorithm following the steps below. You can use your favorite programming language.

- (a) Compute the SVD $A = U\Sigma V^\top = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$, where $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ are the diagonal entries of Σ , \mathbf{u}_i is the i th column of U , and \mathbf{v}_i is the i th column of V .

- (b) Use the first k ($k < r$) terms of SVD to approximate the original image A . Then, we get the compressed images, of which the data matrices are $A_k = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$. Compute A_k for $k = 2, 4, 8, 16, 32, 64, 128, 256$.
- (c) Plot A_k as images for all k .

Please put the compressed images and their corresponding k in this file.

Solution: 1. (a) We first show that

$$\mathcal{C}(A) = \mathcal{C}(U_1).$$

Let $\Sigma_r = \text{diag}(\sigma_1, \dots, \sigma_r)$, then we have

$$\begin{aligned} A &= U \Sigma V^\top \\ &= U_1 \Sigma_r V_1^\top \\ \Rightarrow U_1 &= A V_1 \Sigma_r^{-1}. \end{aligned}$$

It follows from $A = U_1 \Sigma_r V_1^\top$ that $\mathcal{C}(A) \subseteq \mathcal{C}(U_1)$. And it follows from $U_1 = A V_1 \Sigma_r^{-1}$ that $\mathcal{C}(U_1) \subseteq \mathcal{C}(A)$.

Thus, $\mathcal{C}(A) = \mathcal{C}(U_1)$. Therefore,

$$\begin{aligned} P_{\mathcal{C}(A)}(\mathbf{x}) &= P_{\mathcal{C}(U_1)}(\mathbf{x}) \\ &= U_1 (U_1^\top U_1)^{-1} U_1^\top \mathbf{x} \\ &= U_1 U_1^\top \mathbf{x}. \end{aligned}$$

- (b) Similarly, we show that

$$\mathcal{N}(A) = \mathcal{C}(V_2).$$

As

$$\begin{aligned} AV &= U \Sigma \\ \Rightarrow AV_2 &= \mathbf{0}, \end{aligned}$$

$\mathbf{v}_{*,i} \in \mathcal{N}(A)$, $i = r+1, \dots, n$. Since $\dim(\mathcal{N}(A)) = n - \text{rank}(A) = n - r$ and $(\mathbf{v}_{*,r+1}, \mathbf{v}_{*,r+2}, \dots, \mathbf{v}_{*,n})$ are linearly independent, $(\mathbf{v}_{*,r+1}, \mathbf{v}_{*,r+2}, \dots, \mathbf{v}_{*,n})$ is a basis of $\mathcal{N}(A)$, i.e. $\mathcal{C}(V_2) = \mathcal{N}(A)$. Thus, it follows from Exercise 6.2 that

$$\begin{aligned} P_{\mathcal{N}(A)}(\mathbf{x}) &= P_{\mathcal{C}(V_2)}(\mathbf{x}) \\ &= V_2 (V_2^\top V_2)^{-1} V_2^\top \mathbf{x} \\ &= V_2 V_2^\top \mathbf{x}. \end{aligned}$$

- (c) Similarly, we can prove that

$$\mathcal{C}(A^\top) = \mathcal{C}(V_1).$$

Thus,

$$\begin{aligned} P_{\mathcal{C}(A^\top)}(\mathbf{x}) &= P_{\mathcal{C}(V_1)}(\mathbf{x}) \\ &= V_1 V_1^\top \mathbf{x}. \end{aligned}$$

(d) Similarly, we can prove that

$$\mathcal{N}(A^\top) = \mathcal{C}(U_2).$$

Thus,

$$\begin{aligned} P_{\mathcal{N}(A^\top)}(\mathbf{x}) &= P_{\mathcal{C}(U_2)}(\mathbf{x}) \\ &= U_2 U_2^\top \mathbf{x}. \end{aligned}$$

2. (a) Letting $A = (\mathbf{a}_{*,1}, \dots, \mathbf{a}_{*,n})$, we have

$$A^\top A = \begin{bmatrix} \mathbf{a}_{*,1}^\top \mathbf{a}_{*,1} & \dots & \mathbf{a}_{*,1}^\top \mathbf{a}_{*,n} \\ \dots & \dots & \dots \\ \mathbf{a}_{*,n}^\top \mathbf{a}_{*,1} & \dots & \mathbf{a}_{*,n}^\top \mathbf{a}_{*,n} \end{bmatrix}.$$

Thus,

$$\begin{aligned} \text{tr}(A^\top A) &= \sum_{j=1}^n \mathbf{a}_{*,j}^\top \mathbf{a}_{*,j} \\ &= \sum_{j=1}^n \sum_{i=1}^m a_{i,j}^2 \\ &= \|A\|_F^2. \end{aligned}$$

(b) As $\mathcal{C}(A) \perp \mathcal{C}(B)$,

$$A^\top B = 0 \text{ and } B^\top A = 0.$$

Thus,

$$\begin{aligned} \|A + B\|_F^2 &= \text{tr}((A + B)^\top (A + B)) \\ &= \text{tr}(A^\top A) + \text{tr}(B^\top B) + \text{tr}(A^\top B) + \text{tr}(B^\top A) \\ &= \text{tr}(A^\top A) + \text{tr}(B^\top B) \\ &= \|A\|_F^2 + \|B\|_F^2. \end{aligned}$$

3. (a) **Solution 1:** Assume that all singular values of A are different. Note that $\text{rank}(A) = r$. Suppose $A^* \in \mathbb{R}^{m \times n}$ is an optimal solution to the problem (4).

If $K \geq r$, it is trivial that $A^* = A$.

If $K < r$, suppose $X = WY$, where $W \in \mathbb{R}^{m \times K}$, $W^\top W = I$, and $Y = (\mathbf{y}_1, \dots, \mathbf{y}_n) \in \mathbb{R}^{K \times n}$. Clearly, $\mathcal{C}(X) \subset \mathcal{C}(W)$.

Let $A = (\mathbf{a}_1, \dots, \mathbf{a}_n)$. The problem (4) can be formulated as

$$\begin{aligned}
& \min_{W, Y} \{ \|A - WY\|_F^2 : W^\top W = I, W \in \mathbb{R}^{m \times K} \} \\
&= \min_{W \in \mathbb{R}^{m \times K}} \min_{Y \in \mathbb{R}^{K \times n}} \{ \|A - WY\|_F^2 : W^\top W = I \} \\
&= \min_{W \in \mathbb{R}^{m \times K}} \min_{\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^K} \left\{ \sum_{i=1}^n \|\mathbf{a}_i - W\mathbf{y}_i\|_2^2 : W^\top W = I \right\} \\
&= \min_{W \in \mathbb{R}^{m \times K}} \sum_{i=1}^n \min_{\mathbf{y}_i \in \mathbb{R}^K} \{ \|\mathbf{a}_i - W\mathbf{y}_i\|_2^2 : W^\top W = I \}
\end{aligned}$$

Let $\mathbf{z}_i = W\mathbf{y}_i$. Note that the solution to the subproblem $\min_{\mathbf{z}_i \in \mathbb{R}^m} \{ \|\mathbf{a}_i - \mathbf{z}_i\|_2^2 : \mathbf{z}_i \in \mathcal{C}(W), W^\top W = I \}$ is the projection of the point \mathbf{a}_i onto $\mathcal{C}(W)$. It follows from Exercise 6.3 that

$$P_{\mathcal{C}(W)}(\mathbf{a}_i) = WW^\top \mathbf{a}_i.$$

Hence the problem (4) becomes

$$\min_W \{ \|A - WW^\top A\|_F^2 : W^\top W = I, W \in \mathbb{R}^{m \times K} \} \quad (5)$$

Note that

$$\begin{aligned}
\|A - WW^\top A\|_F^2 &= \text{tr}((A - WW^\top A)^\top (A - WW^\top A)) \\
&= 2 \text{tr}(A^\top A) - 2 \text{tr}(A^\top WW^\top A) \\
&= 2 \text{tr}(V\Sigma^\top \Sigma V^\top) - 2 \text{tr}(W^\top U\Sigma\Sigma^\top U^\top W).
\end{aligned}$$

Letting $Q = U^\top W = (\mathbf{q}_1, \dots, \mathbf{q}_K)$, the problem (5) becomes

$$\begin{aligned}
& \max_{Q \in \mathbb{R}^{m \times K}} \text{tr}(Q^\top \Sigma \Sigma^\top Q) \\
& \text{s.t. } Q^\top Q = I.
\end{aligned} \quad (6)$$

We have

$$\begin{aligned}
\text{tr}(Q^\top \Sigma \Sigma^\top Q) &= \sum_{j=1}^K \mathbf{q}_j^\top \Sigma \Sigma^\top \mathbf{q}_j \\
&= \sum_{j=1}^K \sum_{i=1}^r \sigma_i^2 q_{i,j}^2 \\
&= \sum_{i=1}^r \sigma_i^2 \sum_{j=1}^K q_{i,j}^2,
\end{aligned}$$

where $q_{i,j}$ is the entry in the i^{th} row, j^{th} column of Q . Denote

$$\alpha_i = \sum_{j=1}^K q_{i,j}^2, i = 1, \dots, m.$$

We can see that

$$\begin{aligned} \alpha_i &\in [0, 1], i = 1, \dots, m, \\ \sum_{i=1}^m \alpha_i &= K. \end{aligned} \tag{7}$$

Thus, we consider the following problem

$$\begin{aligned} \max_{\alpha \in \mathbb{R}^m} \quad & \sum_{i=1}^r \alpha_i \sigma_i^2 \\ \text{s.t.} \quad & \alpha_i \in [0, 1], i = 1, \dots, m, \\ & \sum_{i=1}^m \alpha_i = K. \end{aligned} \tag{8}$$

Let Q^* be a solution to problem (6) and α^* be a solution to problem (8) respectively. Notice that $\text{tr}((Q^*)^\top \Sigma \Sigma^\top Q^*) \leq (\alpha^*)^\top \Sigma \Sigma^\top \alpha^*$.

Since $\sigma_1^2 \geq \sigma_2^2 \geq \dots \geq \sigma_r^2$, $\alpha^* = (\alpha_1^*, \dots, \alpha_m^*)^\top$ is a solution to problem (8) with

$$\alpha_i^* = \begin{cases} 1 & i = 1, \dots, K, \\ 0 & i = K + 1, \dots, m. \end{cases} \tag{9}$$

In view of Eq. (7) and (9), we can see that the last $m - K$ entries of \mathbf{q}_i^* are 0 for all $i = 1, \dots, K$, that is

$$Q^* = \begin{pmatrix} \tilde{Q}^* \\ 0 \end{pmatrix}_{m \times K},$$

where

$$\tilde{Q}^* \in \mathbb{R}^{K \times K} \text{ and } (\tilde{Q}^*)^\top \tilde{Q}^* = I.$$

We can see that $\text{tr}((Q^*)^\top \Sigma \Sigma^\top Q^*) = (\alpha^*)^\top \Sigma \Sigma^\top \alpha^*$.

Therefore, the solution to problem (4) is

$$\begin{aligned} A^* &= W^*(W^*)^\top A \\ &= UQ^*(Q^*)^\top U^\top A \\ &= U\Sigma_k V^\top, \end{aligned}$$

where $\Sigma_k = \text{diag}(\sigma_1, \dots, \sigma_K, 0, \dots, 0)$

- (b) **Solution 2:** Assume that all singular values of A are different, and suppose $S \in \mathbb{R}^{m \times n}$ is an optimal solution. Note that $\mathbf{rank}(A) = r$.

If $K \geq r$, it is trivial that $S = A$.

If $K < r$, let $A' = U\Sigma'V^T$, where $\Sigma' = \begin{bmatrix} \Sigma_K & 0 \\ 0 & 0 \end{bmatrix} \in \mathbb{R}^{m \times n}$, $\Sigma_K = \text{diag}(\sigma_1, \dots, \sigma_K)$.

Then we have

$$\|A - A'\|_F = (\sigma_{K+1}^2 + \dots + \sigma_r^2)^{\frac{1}{2}} \geq \|A - S\|_F$$

Now we prove that $\|A - S\|_F \geq \|A - A'\|_F$:

Suppose that

$$X = Q\Omega P^T,$$

where $QQ^T = I$, $PP^T = I$, $\Omega = \begin{bmatrix} \Omega_K & 0 \\ 0 & 0 \end{bmatrix}$, and $\Omega_K = \text{diag}(\omega_1, \dots, \omega_K)$. Let $B = Q^T A P$. Then $A = QBP^T$. Thus,

$$\|A - S\|_F = \|Q(B - \Omega)P^T\|_F = \|B - \Omega\|_F.$$

Block the matrix

$$B = \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix},$$

where $B_{11} \in \mathbb{R}^{K \times K}$, $B_{12} \in \mathbb{R}^{K \times (n-K)}$, $B_{21} \in \mathbb{R}^{(m-K) \times K}$, $B_{22} \in \mathbb{R}^{(m-K) \times (n-K)}$. Then by the definition of Frobenius norm, we have

$$\|A - S\|_F^2 = \|B_{11} - \Omega_K\|_F^2 + \|B_{12}\|_F^2 + \|B_{21}\|_F^2 + \|B_{22}\|_F^2.$$

Firstly, we prove that $B_{12} = 0$ and $B_{21} = 0$ by contradiction.

Suppose that $B_{12} \neq 0$, then let

$$Y = Q \begin{bmatrix} B_{11} & B_{12} \\ 0 & 0 \end{bmatrix} P^T.$$

Note that $\mathbf{rank}(Y) \leq K$. And

$$\|A - Y\|_F^2 = \|B_{21}\|_F^2 + \|B_{22}\|_F^2 < \|A - S\|_F^2,$$

which leads to a contradiction. Similarly, $B_{21} = 0$.

Next, we show $B_{11} = \Omega_K$. Suppose that $B_{11} \neq \Omega_K$. Let

$$Z = Q \begin{bmatrix} B_{11} & 0 \\ 0 & 0 \end{bmatrix} P^T$$

and note that $\mathbf{rank}(Z) \leq K$. Then

$$\begin{aligned}\|A - Z\|_F^2 &= \|B_{22}\|_F^2 \leq \|B_{11} - \Omega_K\|_F^2 + \|B_{22}\|_F^2 = \|A - S\|_F^2 \\ &\Rightarrow B_{11} = \Omega_K.\end{aligned}$$

Finally, we find B_{22} . Take SVD for B_{22} , i.e.

$$B_{22} = U_1 \Lambda V_1^T.$$

Now we prove the diagonal entries of Λ are A 's singular values. Let

$$\begin{aligned}U_2 &= \begin{bmatrix} I_K & 0 \\ 0 & U_1 \end{bmatrix}, \\ V_2 &= \begin{bmatrix} I_K & 0 \\ 0 & V_1 \end{bmatrix}.\end{aligned}$$

Note that

$$\begin{aligned}U_2^T Q^T A P V_2 &= \begin{bmatrix} \Omega_K & 0 \\ 0 & \Lambda \end{bmatrix}, \\ \Rightarrow A &= (Q U_2) \begin{bmatrix} \Omega_K & 0 \\ 0 & \Lambda \end{bmatrix} (P V_2)^T.\end{aligned}$$

which implies that the diagonal entries of Λ are A 's singular values. Hence

$$\|A - S\|_F = \|\Lambda\|_F \geq (\sigma_{K+1}^2 + \dots + \sigma_r^2)^{\frac{1}{2}}$$

which completes the proof.

Therefore,

$$\min_{X \in \mathbb{R}^{m \times n}} \{\|A - X\|_F : \mathbf{rank}(X) \leq K\} = \begin{cases} 0, & K \geq r \\ (\sigma_{K+1}^2 + \dots + \sigma_r^2)^{\frac{1}{2}}, & K < r. \end{cases}$$

4. ■

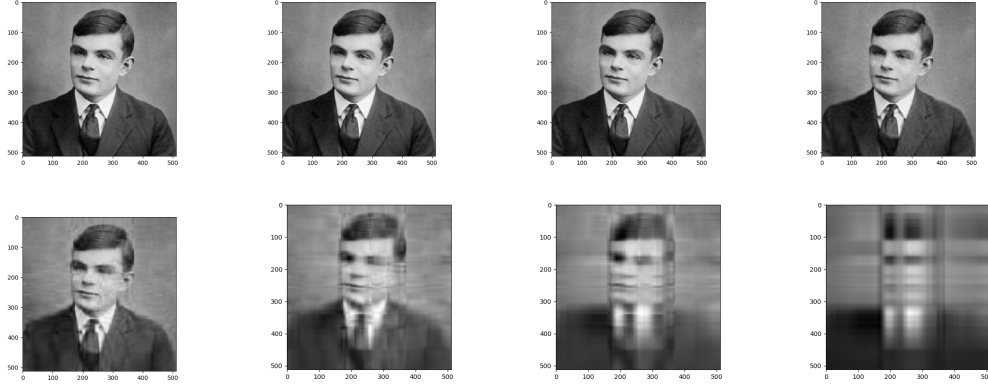


Figure 1: Grayscale images. The number of K decreases from left to right.

Exercise 8: PCA 60pts

Suppose that we have a set of data instances $\{\mathbf{x}_i\}_{i=1}^n \subset \mathbb{R}^d$. Let $\tilde{X} \in \mathbb{R}^{d \times n}$ be the matrix whose i^{th} column is $\mathbf{x}_i - \bar{\mathbf{x}}$, where $\bar{\mathbf{x}}$ is the sample mean, and S be the sample variance matrix.

1. For $G \in \mathbb{R}^{d \times K}$, let us define

$$f(G) = \text{tr}(G^\top S G). \quad (10)$$

Show that $f(GQ) = f(G)$ for any orthogonal matrix $Q \in \mathbb{R}^{K \times K}$.

2. Please find \mathbf{g}_1 defined as follows by the Lagrange multiplier method.

$$\mathbf{g}_1 := \underset{\mathbf{g} \in \mathbb{R}^d}{\text{argmax}} \{f(\mathbf{g}) : \|\mathbf{g}\|_2 = 1\}, \quad (11)$$

where f is defined by (10). Notice that, the vector \mathbf{g}_1 is the first principle component vector of the data.

3. Please find \mathbf{g}_2 defined as follows by the Lagrange multiplier method.

$$\mathbf{g}_2 := \underset{\mathbf{g} \in \mathbb{R}^d}{\text{argmax}} \{f(\mathbf{g}) : \|\mathbf{g}\|_2 = 1, \langle \mathbf{g}, \mathbf{g}_1 \rangle = 0\}, \quad (12)$$

where \mathbf{g}_1 is given by (11). Similar to \mathbf{g}_1 , the vector \mathbf{g}_2 is the second principle component vector of the data.

4. Please derive the first K principle component vectors by repeating the above process.
5. What is $f(\mathbf{g}_k)$, $k = 1, \dots, K$? What about their meaning?
6. When the first K principle component vectors are unique?

Solution: 1. Since $\mathbf{tr}(AB) = \mathbf{tr}(BA)$, we have

$$\begin{aligned} f(GQ) &= \mathbf{tr}(Q^\top G^\top SGQ) \\ &= \mathbf{tr}(QQ^\top G^\top SG) \\ &= \mathbf{tr}(G^\top SG) \\ &= f(G), \end{aligned}$$

which completes the proof.

2. Suppose $\tilde{X} = U\Sigma V^\top$, where $\Sigma = \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 & 0 & \cdots & 0 \\ 0 & 0 & \ddots & \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & \sigma_r & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 & 0 & \cdots & 0 \end{bmatrix} \in \mathbb{R}^{d \times n}$, assume

$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$. Thus

$$S = \frac{1}{n-1} \tilde{X} \tilde{X}^\top = \frac{1}{n-1} U \Sigma \Sigma^\top U^\top = \frac{1}{n-1} U \Sigma_d^2 U^\top,$$

where $\Sigma_d = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r, 0, \dots, 0) \in \mathbb{R}^{d \times d}$.

We have

$$f(G) = \mathbf{tr}\left(\frac{1}{n-1} G^\top U \Sigma_d^2 U^\top G\right).$$

Let $Q = U^\top G$, and we have $Q^\top Q = I$. Hence the problem becomes

$$\begin{aligned} \max_{\mathbf{q} \in \mathbb{R}^d} \quad & \mathbf{tr}(\mathbf{q}^\top \Sigma_d^2 \mathbf{q}), \\ \text{s.t.} \quad & \|\mathbf{q}\|^2 = 1. \end{aligned} \tag{13}$$

Since $\mathbf{tr}(\mathbf{q}^\top \Sigma_d^2 \mathbf{q}) = \mathbf{q}^\top \Sigma_d^2 \mathbf{q}$, the Lagrangian is

$$L(\mathbf{q}, \lambda) = \mathbf{q}^\top \Sigma_d^2 \mathbf{q} - \lambda(\mathbf{q}^\top \mathbf{q} - 1).$$

We have

$$\nabla_{\mathbf{q}} L(\mathbf{q}, \lambda)|_{\mathbf{q}=\mathbf{q}_1} = 2\Sigma_d^2 \mathbf{q}_1 - 2\lambda \mathbf{q}_1 = 0,$$

i.e., $\Sigma_d^2 \mathbf{q}_1 = \lambda \mathbf{q}_1$, which implies that λ is an eigenvalue of Σ_d^2 and \mathbf{q}_1 is the corresponding unit eigenvector.

Then the objective function is $f(\mathbf{q}_1) = \frac{1}{n-1} \mathbf{q}_1^\top \Sigma_d^2 \mathbf{q}_1 = \frac{1}{n-1} \mathbf{q}_1^\top \lambda \mathbf{q}_1 = \frac{1}{n-1} \lambda$. The optimality of $f(\mathbf{q}_1)$ implies that λ is the largest eigenvalue of Σ_d^2 , i.e., $\lambda = \max_i \sigma_i^2$ and \mathbf{q}_1 is the corresponding unit eigenvector.

Since $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$, $\mathbf{q}_1 = (1, 0, \dots, 0)^\top \in \mathbb{R}^d$ is a solution to problem (13). And $\mathbf{g}_1 = U \mathbf{q}_1 \in \mathbf{argmax}_{\mathbf{g} \in \mathbb{R}^d} \{f(\mathbf{g}) : \|\mathbf{g}\|_2 = 1\}$, where $U = (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d)$, i.e., \mathbf{u}_1 is a solution to problem (11).

3. Since $\mathbf{q} = U^\top \mathbf{g}$, we have $\mathbf{q}_i^\top \mathbf{q}_j = 0$ if $\mathbf{g}_i^\top \mathbf{g}_j = 0$. Thus the problem becomes

$$\begin{aligned} \max_{\mathbf{q} \in \mathbb{R}^d} & \mathbf{q}^\top \Sigma_d^2 \mathbf{q}, \\ \text{s.t. } & \|\mathbf{q}\|^2 = 1, \\ & \mathbf{q}^\top \mathbf{q}_1 = 0. \end{aligned} \quad (14)$$

The Lagrangian is

$$L(\mathbf{q}, \lambda, \mu) = \mathbf{q}^\top \Sigma_d^2 \mathbf{q} - \lambda(\mathbf{q}^\top \mathbf{q} - 1) - \mu(\mathbf{q}^\top \mathbf{q}_1).$$

We have

$$\nabla_{\mathbf{q}} L(\mathbf{q}, \lambda, \mu)|_{\mathbf{q}=\mathbf{q}_2} = 2\Sigma_d^2 \mathbf{q}_2 - 2\lambda \mathbf{q}_2 - \mu \mathbf{q}_1 = 0. \quad (15)$$

Note that $\mathbf{q}_2^\top \mathbf{q}_1 = 0$ and $\mathbf{q}_1^\top \Sigma_d^2 \mathbf{q}_2 = \mathbf{q}_2^\top (\Sigma_d^2 \mathbf{q}_1) = \mathbf{q}_2^\top (\sigma_1^2 \mathbf{q}_1) = 0$.

Left-multiplying (15) by \mathbf{q}_1^\top , we have

$$0 = 2\mathbf{q}_1^\top \Sigma_d^2 \mathbf{q}_2 - 2\lambda \mathbf{q}_1^\top \mathbf{q}_2 - \mu \mathbf{q}_1^\top \mathbf{q}_1 = -\mu.$$

Thus $\Sigma_d^2 \mathbf{q}_2 = \lambda \mathbf{q}_2$. Similarly, we derive that $\mathbf{q}_2 = (0, 1, 0, \dots, 0)^\top \in \mathbb{R}^d$ is a solution to problem (14). Thus $\mathbf{g}_2 = U \mathbf{q}_2 = \mathbf{u}_2$ is a solution to problem (12).

4. Let $\mathbf{e}_i \in \mathbb{R}^d$ denote the one-hot vector with i^{th} entry 1 and all other entries 0. Repeating the above process, we derive that $\mathbf{q}_k = \mathbf{e}_k$, $k = 1, \dots, K$. Therefore, $\mathbf{g}_k = \mathbf{u}_k$, $k = 1, \dots, K$ are the first K principle component vectors.
5. Assume $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ and $K \leq r \leq d$. For any $k \in [K]$,

$$f(\mathbf{g}_k) = \text{tr} \left(\frac{1}{n-1} \mathbf{g}_k^\top U \Sigma_d^2 U^\top \mathbf{g}_k \right) = \frac{1}{n-1} \sigma_k^2$$

That is, $f(\mathbf{g}_k)$ corresponds to the square of the k^{th} largest singular value of \tilde{X} .

6. Assume $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ and $r \leq d$. If $K < r$

$$\sigma_1 > \sigma_2 > \dots > \sigma_K > \sigma_{K+1},$$

then the first K principle component vectors $\{\mathbf{g}_i\}_{i=1}^K$ are unique.

If $K = r$

$$\sigma_1 > \sigma_2 > \dots > \sigma_K,$$

then the first K principle component vectors $\{\mathbf{g}_i\}_{i=1}^K$ are unique. ■