

1 Fisher Discriminant Analysis (Cont'd)

Recall Let $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^T$ be an $n \times p$ matrix and $\mathbf{Y} = \mathbf{H}\mathbf{X}\mathbf{G}$ be an $n \times q$ matrix ($q < p$). Our goal is to make the *with-in covariance* \mathbf{S}_w 'small' and the *between covariance* \mathbf{S}_b 'large', i.e.

$$\max_{\mathbf{G}} \text{tr}((\mathbf{G}^T \mathbf{S}_w \mathbf{G})^{-1} (\mathbf{G}^T \mathbf{S}_b \mathbf{G}))$$

In the last lecture, we have derived that

$$\mathbf{S}_b \mathbf{G} = \mathbf{S}_w \mathbf{G} \mathbf{\Lambda}$$

Now we move to solve this augmented eigenvalue problem. Consider each column \mathbf{g}_i of matrix \mathbf{G}

$$\begin{aligned} \mathbf{S}_b \mathbf{g}_i &= \lambda_i \mathbf{S}_w \mathbf{g}_i \\ &= \lambda_i (\mathbf{S}_t - \mathbf{S}_b) \mathbf{g}_i \\ \mathbf{S}_b \mathbf{g}_i &= \frac{\lambda_i}{1 + \lambda_i} \mathbf{S}_t \mathbf{g}_i \end{aligned} \tag{1}$$

If \mathbf{S}_t is invertible, equation (1) can be solved as an ordinary eigenvalue problem

$$\mathbf{S}_t^{-1} \mathbf{S}_b \mathbf{g}_i = \lambda \mathbf{g}_i$$

where $\lambda = \frac{\lambda_i}{1 + \lambda_i}$. But it's often the case that \mathbf{S}_t is not invertible. We give the following theorem to show that equation (1) can be solved using the pseudo-inverse \mathbf{S}_t^\dagger .

Theorem 1.1. Let $\mathbf{\Sigma}_1$ and $\mathbf{\Sigma}_2$ be two $m \times m$ real matrices. Assume $\mathcal{R}(\mathbf{\Sigma}_1) \subseteq \mathcal{R}(\mathbf{\Sigma}_2)$. (Here $\mathcal{R}(\cdot)$ is the range of a matrix) Then if $(\mathbf{\Lambda}, \mathbf{A})$ are the nonzero eigenpairs of $\mathbf{\Sigma}_2^\dagger \mathbf{\Sigma}_1$, we have that $(\mathbf{\Lambda}, \mathbf{A})$ are the nonzero eigenpairs of matrix pencil $(\mathbf{\Sigma}_1, \mathbf{\Sigma}_2)$, i.e.

$$\mathbf{\Sigma}_1 \mathbf{A} = \mathbf{\Sigma}_2 \mathbf{A} \mathbf{\Lambda} \iff \mathbf{\Sigma}_2^\dagger \mathbf{\Sigma}_1 \mathbf{A} = \mathbf{A} \mathbf{\Lambda}$$

Since $\mathcal{R}(\mathbf{S}_b) = \mathcal{R}(\mathbf{X}^T \mathbf{H} \mathbf{E} \mathbf{\Pi}^{-1} \mathbf{E}^T \mathbf{H} \mathbf{X}) = \mathcal{R}(\mathbf{X}^T \mathbf{H} \mathbf{E} \mathbf{\Pi}^{-1/2})$ and $\mathcal{R}(\mathbf{S}_t) = \mathcal{R}(\mathbf{X}^T \mathbf{H} \mathbf{H} \mathbf{X}) = \mathcal{R}(\mathbf{X}^T \mathbf{H})$, we have $\mathcal{R}(\mathbf{S}_b) \subseteq \mathcal{R}(\mathbf{S}_t)$. By Theorem 1.1, equation (1) is equivalent to

$$\mathbf{S}_t^\dagger \mathbf{S}_b \mathbf{g} = \lambda \mathbf{g}$$

Now we move to prove Theorem 1.1.

Proof. Let $\Sigma_1 = \mathbf{U}_1 \Gamma_1 \mathbf{V}_1^T$ and $\Sigma_2 = \mathbf{U}_2 \Gamma_2 \mathbf{V}_2^T$ be the condensed SVD of Σ_1 and Σ_2 . Then we have $\mathcal{R}(\Sigma_1) = \mathcal{R}(\mathbf{U}_1)$ and $\mathcal{R}(\Sigma_2) = \mathcal{R}(\mathbf{U}_2)$. The psudo-inverse is $\Sigma_2^\dagger = \mathbf{V}_2 \Gamma_2^{-1} \mathbf{U}_2^T$, thus $\Sigma_2 \Sigma_2^\dagger = \mathbf{U}_2 \mathbf{U}_2^T$. Given $\mathcal{R}(\Sigma_1) \subseteq \mathcal{R}(\Sigma_2)$, we have $\mathcal{R}(\mathbf{U}_1) \subseteq \mathcal{R}(\mathbf{U}_2)$. Further, we can assume there is some \mathbf{Q} such that $\mathbf{U}_1 = \mathbf{U}_2 \mathbf{Q}$

$$\begin{aligned}\Sigma_2 \Sigma_2^\dagger \Sigma_1 &= \mathbf{U}_2 \mathbf{U}_2^T \mathbf{U}_1 \Gamma_1 \mathbf{V}_1^T \\ &= \mathbf{U}_2 \mathbf{U}_2^T \mathbf{U}_2 \mathbf{Q} \Gamma_1 \mathbf{V}_1^T \\ &= \mathbf{U}_1 \Gamma_1 \mathbf{V}_1^T \\ &= \Sigma_1\end{aligned}$$

Therefore,

$$\begin{aligned}\Sigma_2^\dagger \Sigma_1 \mathbf{A} &= \mathbf{A} \Lambda \\ \Sigma_2 \Sigma_2^\dagger \Sigma_1 \mathbf{A} &= \Sigma_2 \mathbf{A} \Lambda \\ \Sigma_1 \mathbf{A} &= \Sigma_2 \mathbf{A} \Lambda\end{aligned}$$

□

2 Method 2: Complete Orthogonal Decomposition

Definition 2.1 (Generalized Singular Value Decomposition). For $\mathbf{A} \in \mathbb{R}^{n \times m}$, $\mathbf{B} \in \mathbb{R}^{k \times m}$, their GSVD is given by:

$$\begin{aligned}\mathbf{U}^T \mathbf{A} \mathbf{X} = \mathbf{C} &= \text{diag}(\alpha_1, \dots, \alpha_m) = [\Sigma_{\mathbf{A}}, \mathbf{0}], \alpha_i \geq 0 \\ \mathbf{V}^T \mathbf{B} \mathbf{X} = \mathbf{S} &= \text{diag}(\beta_1, \dots, \beta_q) = [\Sigma_{\mathbf{B}}, \mathbf{0}], \beta_i \geq 0, q = \min(k, m),\end{aligned}$$

where $\mathbf{U} \in \mathbb{R}^{n \times n}$, $\mathbf{V} \in \mathbb{R}^{k \times k}$ and $\mathbf{U}^T \mathbf{U} = \mathbf{I}_n$, $\mathbf{V}^T \mathbf{V} = \mathbf{I}_k$, $\mathbf{X} \in \mathbb{R}^{m \times m}$ is nonsingular. It holds that $\Sigma_{\mathbf{A}}^T \Sigma_{\mathbf{A}} + \Sigma_{\mathbf{B}}^T \Sigma_{\mathbf{B}} = \mathbf{I}_m$.

Proposition 2.1. Application of GSVD. For $\mathbf{A} \in \mathbb{R}^{n \times m}$, $\mathbf{B} \in \mathbb{R}^{k \times m}$, their GSVD is

$$\begin{aligned}\mathbf{U}^T \mathbf{A} \mathbf{X} &= [\Sigma_{\mathbf{A}}, \mathbf{0}], \\ \mathbf{V}^T \mathbf{B} \mathbf{X} &= [\Sigma_{\mathbf{B}}, \mathbf{0}]\end{aligned}$$

$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_m)$, then first r vectors of \mathbf{X} are the generalized eigenvectors of $\mathbf{A}^T \mathbf{A}$ and $\mathbf{B}^T \mathbf{B}$, and the corresponding eigenvalue is $\frac{\alpha_i^2}{\beta_i^2}$.

Proof: From GSVD:

$$\begin{aligned}\mathbf{A} = \mathbf{U} [\Sigma_{\mathbf{A}}, \mathbf{0}] \mathbf{X}^{-1} &\implies \mathbf{A}^T \mathbf{A} = \mathbf{X}^{-T} [\Sigma_{\mathbf{A}}, \mathbf{0}]^T [\Sigma_{\mathbf{A}}, \mathbf{0}] \mathbf{X}^{-1} \\ \mathbf{B} = \mathbf{V} [\Sigma_{\mathbf{B}}, \mathbf{0}] \mathbf{X}^{-1} &\implies \mathbf{B}^T \mathbf{B} = \mathbf{X}^{-T} [\Sigma_{\mathbf{B}}, \mathbf{0}]^T [\Sigma_{\mathbf{B}}, \mathbf{0}] \mathbf{X}^{-1}\end{aligned}$$

Because \mathbf{X} is nonsingular,

$$\begin{aligned}\mathbf{A}^T \mathbf{A} \mathbf{X} &= \mathbf{X}^{-T} \begin{pmatrix} \Sigma_{\mathbf{A}}^T \Sigma_{\mathbf{A}} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \\ \mathbf{B}^T \mathbf{B} \mathbf{X} &= \mathbf{X}^{-T} \begin{pmatrix} \Sigma_{\mathbf{B}}^T \Sigma_{\mathbf{B}} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}\end{aligned}$$

Let $\mathbf{x}_i, i \leq r$, then

$$\begin{aligned}\mathbf{A}^T \mathbf{A} \mathbf{x}_i &= \mathbf{X}^{-T} \begin{pmatrix} \alpha_i^2 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \\ \mathbf{B}^T \mathbf{B} \mathbf{x}_i &= \mathbf{X}^{-T} \begin{pmatrix} \beta_i^2 \\ 0 \\ \vdots \\ 0 \end{pmatrix}\end{aligned}$$

We get $\mathbf{A}^T \mathbf{A} \mathbf{x}_i = \frac{\alpha_i^2}{\beta_i^2} \mathbf{B}^T \mathbf{B} \mathbf{x}_i$. This is the solution to the generalized eigenvalue problems.

Lemma 2.1. *The CS Decomposition: Consider matrix*

$$\mathbf{Q} = \begin{pmatrix} \mathbf{Q}_1 \\ \mathbf{Q}_2 \end{pmatrix}, \mathbf{Q}_1 \in \mathbb{R}^{m_1 \times n}, \mathbf{Q}_2 \in \mathbb{R}^{m_2 \times n},$$

where $m_1 \geq n, m_2 \geq n$, if the columns of \mathbf{Q} are orthogonal, then exist orthogonal matrices $\mathbf{U}_1 \in \mathbb{R}^{m_1 \times m_1}, \mathbf{U}_2 \in \mathbb{R}^{m_2 \times m_2}, \mathbf{V}_1 \in \mathbb{R}^{n \times n}$, such that

$$\mathbf{U}_1^T \mathbf{Q}_1 \mathbf{V}_1 = \mathbf{C}, \quad \mathbf{U}_2^T \mathbf{Q}_2 \mathbf{V}_1 = \mathbf{S}, \quad \mathbf{C}^T \mathbf{C} + \mathbf{S}^T \mathbf{S} = \mathbf{I}_n$$

Proposition 2.2. *Using QR decomposition to solve GSVD.*

Proof: Let $\mathbf{A} \in \mathbb{R}^{n \times m}, \mathbf{B} \in \mathbb{R}^{k \times m}$,

$$\mathbf{C} = \begin{pmatrix} \mathbf{A} \\ \mathbf{B} \end{pmatrix}$$

and $t = \text{rank}(\mathbf{C})$, perform QR to \mathbf{C} , gets

$$\mathbf{P}^T \mathbf{C} \mathbf{Q} = \begin{pmatrix} \mathbf{R} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}$$

Where \mathbf{Q} is a permutation matrix, \mathbf{P} is an orthogonal matrix. $\mathbf{R}_{t \times t}$ is nonsingular. Let

$$\mathbf{P} = \begin{pmatrix} \mathbf{P}_{11} & \mathbf{P}_{12} \\ \mathbf{P}_{21} & \mathbf{P}_{22} \end{pmatrix}$$

Where $\mathbf{P}_{11} \in \mathbb{R}^{n \times t}, \mathbf{P}_{21} \in \mathbb{R}^{k \times t}, \|\mathbf{P}\| \leq 1, \|\mathbf{P}_{11}\| \leq 1$.

From the lemma, first apply SVD on \mathbf{P}_{11} , we will have $\mathbf{U}^T \mathbf{P}_{11} \mathbf{W} = \mathbf{\Sigma}_A$, then apply QR to $\mathbf{P}_{21} \mathbf{W}$, we will have $\mathbf{P}_{21} \mathbf{W} = \mathbf{V} \mathbf{L}$. So

$$\begin{pmatrix} \mathbf{\Sigma}_A \\ \mathbf{L} \end{pmatrix} = \begin{pmatrix} \mathbf{U}^T \mathbf{P}_{11} \\ \mathbf{V}^T \mathbf{P}_{21} \end{pmatrix} \mathbf{W} = \begin{pmatrix} \mathbf{U}^T & \mathbf{0} \\ \mathbf{0} & \mathbf{V}^T \end{pmatrix} \begin{pmatrix} \mathbf{P}_{11} \\ \mathbf{P}_{21} \end{pmatrix} \mathbf{W}$$

From Lemma, we have

$$(\Sigma_{\mathbf{A}}^T \quad \mathbf{L}^T) \begin{pmatrix} \Sigma_{\mathbf{A}} \\ \mathbf{L} \end{pmatrix} = \Sigma_{\mathbf{A}}^T \Sigma_{\mathbf{A}} + \mathbf{L}^T \mathbf{L} = \mathbf{I}$$

Because $\Sigma_{\mathbf{B}}$ need not to be diagonal matrix, we use $\Sigma_{\mathbf{B}}$ to denote \mathbf{L} . To simplify:

$$\begin{pmatrix} \mathbf{A} \\ \mathbf{B} \end{pmatrix} \mathbf{Q} = \begin{pmatrix} \mathbf{P}_{11} & \mathbf{P}_{12} \\ \mathbf{P}_{21} & \mathbf{P}_{22} \end{pmatrix} \begin{pmatrix} \mathbf{R} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} = \begin{pmatrix} \mathbf{P}_{11}\mathbf{R} & \mathbf{0} \\ \mathbf{P}_{21}\mathbf{R} & \mathbf{0} \end{pmatrix} = \begin{pmatrix} \mathbf{U}\Sigma_{\mathbf{A}}\mathbf{W}^T\mathbf{R} & \mathbf{0} \\ \mathbf{V}\Sigma_{\mathbf{B}}\mathbf{W}^T\mathbf{R} & \mathbf{0} \end{pmatrix}$$

From above, we get

$$\begin{aligned} \mathbf{A}\mathbf{Q} &= [\mathbf{U}\Sigma_{\mathbf{A}}\mathbf{W}^T\mathbf{R}, \mathbf{0}] \\ \mathbf{B}\mathbf{Q} &= [\mathbf{V}\Sigma_{\mathbf{B}}\mathbf{W}^T\mathbf{R}, \mathbf{0}] \end{aligned}$$

Change form:

$$\begin{aligned} \mathbf{U}^T \mathbf{A}\mathbf{Q} &= [\Sigma_{\mathbf{A}}\mathbf{W}^T\mathbf{R}, \mathbf{0}] = [\Sigma_{\mathbf{A}}, \mathbf{0}] \begin{pmatrix} \mathbf{W}^T\mathbf{R} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{pmatrix} \\ \mathbf{V}^T \mathbf{B}\mathbf{Q} &= [\Sigma_{\mathbf{B}}\mathbf{W}^T\mathbf{R}, \mathbf{0}] = [\Sigma_{\mathbf{B}}, \mathbf{0}] \begin{pmatrix} \mathbf{W}^T\mathbf{R} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{pmatrix} \end{aligned}$$

Because $\mathbf{W}^T\mathbf{R}$ is invertible, we can set

$$\mathbf{X} = \mathbf{Q} \begin{pmatrix} \mathbf{R}^{-1}\mathbf{W} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{pmatrix}$$

Then come the solution of GSVD:

$$\begin{aligned} \mathbf{U}^T \mathbf{A}\mathbf{X} &= [\Sigma_{\mathbf{A}}, \mathbf{0}] \\ \mathbf{V}^T \mathbf{B}\mathbf{X} &= [\Sigma_{\mathbf{B}}, \mathbf{0}] \end{aligned}$$

Example 2.1. *The Step to Solve FDA.*

1. Compute $\mathbf{S}_t = \mathbf{X}^T \mathbf{H} \mathbf{H} \mathbf{X}$, $\mathbf{S}_b = \mathbf{X}^T \mathbf{H} \mathbf{E} \mathbf{H} \mathbf{E}^T \mathbf{H} \mathbf{X}$.
2. Let $\mathbf{A} = \mathbf{H} \mathbf{H} \mathbf{X}$, $\mathbf{B} = \mathbf{H} \mathbf{X}$.
3. Let

$$\mathbf{C} = \begin{pmatrix} \mathbf{A} \\ \mathbf{B} \end{pmatrix}$$

4. Apply QR to \mathbf{C} , get

$$\mathbf{P}^T \mathbf{C} \mathbf{Q} = \begin{pmatrix} \mathbf{R} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}$$

5. Perform the SVD of \mathbf{P}_{11} , get $\mathbf{U}^T \mathbf{P}_{11} \mathbf{W} = \Sigma_{\mathbf{A}}$.
6. The corresponding vectors is first $c - 1$ vectors of

$$\mathbf{X} = \mathbf{Q} \begin{pmatrix} \mathbf{R}^{-1}\mathbf{W} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{pmatrix}$$

.

3 Method 3: Regularized Discriminant Analysis

If \mathbf{S}_t is singular, we can add perturbation to \mathbf{S}_t , the problem then change to the following form.

$$\begin{aligned} (\mathbf{S}_t + \sigma^2 \mathbf{I}_p)^{-1} \mathbf{S}_b \mathbf{A} &= \mathbf{A} \mathbf{\Lambda} \\ (\mathbf{X}^T \mathbf{H} \mathbf{X} + \sigma^2 \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{H} \mathbf{E} \mathbf{\Pi}^{-\frac{1}{2}} \mathbf{\Pi}^{-\frac{1}{2}} \mathbf{E}^T \mathbf{H} \mathbf{X} \mathbf{A} &= \mathbf{A} \mathbf{\Lambda} \end{aligned}$$

For computation efficiency, If $n \gg p$, let $\mathbf{\Phi} = (\mathbf{X}^T \mathbf{H} \mathbf{X} + \sigma^2 \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{H} \mathbf{E} \mathbf{\Pi}^{-\frac{1}{2}}$.

If $p \gg n$, change it to $\mathbf{\Phi} = \mathbf{X}^T \mathbf{H} (\mathbf{H} \mathbf{X} \mathbf{X}^T \mathbf{H} + \sigma^2 \mathbf{I}_n)^{-1} \mathbf{E} \mathbf{\Pi}^{-\frac{1}{2}}$. The equation can be rewrite to $\mathbf{\Phi} \mathbf{\Pi}^{-\frac{1}{2}} \mathbf{E}^T \mathbf{H} \mathbf{X} \mathbf{A} = \mathbf{A} \mathbf{\Lambda}$. Let $\mathbf{\Psi} = \mathbf{\Pi}^{-\frac{1}{2}} \mathbf{E}^T \mathbf{H} \mathbf{X} \mathbf{\Phi}$, $\mathbf{B} = \mathbf{\Pi}^{-\frac{1}{2}} \mathbf{E}^T \mathbf{H} \mathbf{X}$.

Proposition 3.1. *The eigenvectors \mathbf{A} in $\mathbf{\Phi} \mathbf{B} \mathbf{A} = \mathbf{A} \mathbf{\Lambda}$ is $\mathbf{\Phi} \mathbf{V}_{\Psi} \mathbf{\Gamma}_{\Psi}^{-\frac{1}{2}}$, and the eigenvalues are $\mathbf{\Gamma}_{\Psi}$. Where $\mathbf{\Psi} = \mathbf{V}_{\Psi} \mathbf{\Gamma}_{\Psi} \mathbf{V}_{\Psi}^T$.*

Proof:

$$\begin{aligned} \mathbf{\Psi} &= \mathbf{B} \mathbf{\Phi} = \mathbf{V}_{\Psi} \mathbf{\Gamma}_{\Psi} \mathbf{V}_{\Psi}^T \\ \implies \mathbf{\Phi} \mathbf{B} \mathbf{\Phi} &= \mathbf{\Phi} \mathbf{V}_{\Psi} \mathbf{\Gamma}_{\Psi} \mathbf{V}_{\Psi}^T \\ \implies \mathbf{\Phi} \mathbf{B} \mathbf{\Phi} \mathbf{V}_{\Psi} \mathbf{\Gamma}_{\Psi}^{-\frac{1}{2}} &= \mathbf{\Phi} \mathbf{V}_{\Psi} \mathbf{\Gamma}_{\Psi}^{-\frac{1}{2}} \mathbf{\Gamma}_{\Psi} \\ \implies (\mathbf{\Phi} \mathbf{B}) (\mathbf{\Phi} \mathbf{V}_{\Psi} \mathbf{\Gamma}_{\Psi}^{-\frac{1}{2}}) &= (\mathbf{\Phi} \mathbf{V}_{\Psi} \mathbf{\Gamma}_{\Psi}^{-\frac{1}{2}}) \mathbf{\Gamma}_{\Psi} \\ \implies \mathbf{A} &= \mathbf{\Phi} \mathbf{V}_{\Psi} \mathbf{\Gamma}_{\Psi}^{-\frac{1}{2}}, \mathbf{\Lambda} = \mathbf{\Gamma}_{\Psi} \end{aligned}$$

Example 3.1. *The Step to Solve FDA.*

1. Compute $\mathbf{\Phi}, \mathbf{\Psi}$.
2. Perform the SVD of $\mathbf{\Psi}$. Get $\mathbf{\Psi} = \mathbf{V}_{\Psi} \mathbf{\Gamma}_{\Psi} \mathbf{V}_{\Psi}^T$.
3. $\mathbf{G} = \mathbf{\Phi} \mathbf{V}_{\Psi} \mathbf{\Gamma}_{\Psi}^{-\frac{1}{2}}$.

4 RFDA and Rigne Regression

Let $\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n]^T = \mathbf{E} \mathbf{\Pi}^{-1/2} \mathbf{H} \mathbf{\pi}$, where $\mathbf{H} \mathbf{\pi} = \mathbf{I}_c - \frac{1}{n} \sqrt{\mathbf{\pi}} \sqrt{\mathbf{\pi}}^T$ and $\sqrt{\mathbf{\pi}} = (\sqrt{n_1}, \sqrt{n_2}, \dots, \sqrt{n_c})^T$ be a vector associated with the square root of number of nodes in each cluster.

For each row in \mathbf{Y} , $\mathbf{y}_i = (y_{i1}, y_{i2}, \dots, y_{ic})$ where

$$y_{ij} = \begin{cases} \frac{n-n_j}{n\sqrt{n_j}} & \text{if } i \in V_j \\ \frac{\sqrt{n_j}}{n} & \text{if otherwise} \end{cases}$$

The goal is to minimize the following Lagrangian function:

$$\begin{aligned} \min_{\mathbf{w}_0, \mathbf{W}} L(\mathbf{w}_0, \mathbf{W}) &= \frac{1}{2} \|\mathbf{Y} - \mathbf{1}_n \mathbf{w}_0^T - \mathbf{X} \mathbf{W}\|_F^2 + \frac{\sigma^2}{2} \text{tr}(\mathbf{W}^T \mathbf{W}) \\ &= \frac{1}{2} \sum_{i=1}^n \|\mathbf{y}_i - \mathbf{w}_0 - \mathbf{W}^T \mathbf{x}_i\|^2 + \frac{\sigma^2}{2} \text{tr}(\mathbf{W}^T \mathbf{W}) \end{aligned}$$

By taking partial derivatives, we have

$$\frac{\partial L}{\mathbf{w}_0} = n\mathbf{w}_0 + \mathbf{W}^T \mathbf{X}^T \mathbf{1}_n - \mathbf{Y}^T \mathbf{1}_n = 0 \quad (2)$$

$$\frac{\partial L}{\mathbf{W}} = (\mathbf{X}^T \mathbf{X} + \sigma^2 \mathbf{I}_p) \mathbf{W} + \mathbf{X}^T \mathbf{1}_n \mathbf{w}_0^T - \mathbf{X}^T \mathbf{Y} = 0 \quad (3)$$

By solving equation (2), we have

$$\mathbf{w}_0 = -\mathbf{W}^T \mathbf{m}$$

where $\mathbf{m} = \frac{1}{n} \mathbf{X}^T \mathbf{1}_n$. Then in equation (3), we have

$$\begin{aligned} (\mathbf{X}^T \mathbf{X} + \sigma^2 \mathbf{I}_p) \mathbf{W} &= \mathbf{X}^T \mathbf{1}_n \mathbf{m}^T \mathbf{W} + \mathbf{X}^T \mathbf{Y} \\ \implies (\mathbf{X}^T \mathbf{X} - n \mathbf{m} \mathbf{m}^T + \sigma^2 \mathbf{I}_p) \mathbf{W} &= \mathbf{X}^T \mathbf{Y} \\ \implies (\mathbf{X}^T \mathbf{H} \mathbf{X} - \sigma^2 \mathbf{I}_p) \mathbf{W} &= \mathbf{X}^T \mathbf{Y} \\ \implies \mathbf{W} &= (\mathbf{X}^T \mathbf{H} \mathbf{X} - \sigma^2 \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{Y} \end{aligned}$$

Notice that $\mathbf{Y} = \mathbf{E} \mathbf{\Pi}^{-1/2} \mathbf{H} \boldsymbol{\pi} = \mathbf{E} \mathbf{\Pi}^{-1/2} (\mathbf{I}_c - \frac{1}{n} \sqrt{\boldsymbol{\pi}} \sqrt{\boldsymbol{\pi}}^T) = (\mathbf{I}_p - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T) \mathbf{E} \mathbf{\Pi}^{-1/2}$, we have

$$\mathbf{W} = (\mathbf{X}^T \mathbf{H} \mathbf{X} + \sigma^2 \mathbf{I}_p)^{-1} \mathbf{H} \mathbf{E} \mathbf{\Pi}^{1/2} \quad (4)$$

The result in equation (4) is exactly the same as $\boldsymbol{\Phi}$.