

Fisher's Linear Discriminant Analysis

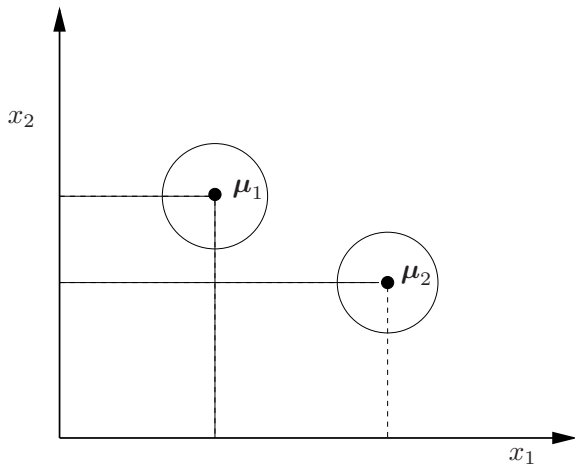
Seungjin Choi

Department of Computer Science and Engineering
Pohang University of Science and Technology
77 Cheongam-ro, Nam-gu, Pohang 37673, Korea
seungjin@postech.ac.kr

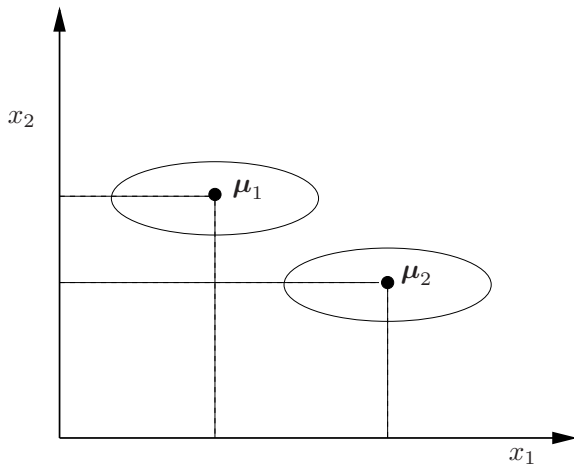
FLD or LDA

- ▶ Introduced by Fisher (1936)
- ▶ One of widely-used **linear discriminant analysis (LDA)** methods
- ▶ Curse of dimensionality
 - ▶ Linear dimensionality reduction: PCA, ICA, FLD, MDS
 - ▶ Nonlinear dimensionality reduction: Isomap, LLE, Laplacian eigenmap
- ▶ FLD aims at achieving an **optimal linear dimensionality reduction** for **classification**

An Example: Isotropic Case



FLD: A Graphical Illustration



Two Classes

Given a set of data points, $\{\mathbf{x} \in \mathbb{R}^D\}$, one wished to find a linear projection of the data onto a 1-dimensional space, $y = \mathbf{w}^\top \mathbf{x}$.

Sample means for \mathbf{x} :

$$\boldsymbol{\mu}_i = \frac{1}{N_i} \sum_{\mathbf{x} \in C_i} \mathbf{x}.$$

Sample means for the projected points:

$$\tilde{\mu}_i = \frac{1}{N_i} \sum_{y \in \mathcal{Y}_i} y = \frac{1}{N_i} \sum_{\mathbf{x} \in C_i} \mathbf{w}^\top \mathbf{x} = \mathbf{w}^\top \boldsymbol{\mu}_i.$$

We know that the difference between sample means is not always a good measure of the separation between projected points:

$$|\tilde{\mu}_1 - \tilde{\mu}_2| = |\mathbf{w}^\top (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)|.$$

Scale $\|\mathbf{w}\| \uparrow \Rightarrow |\tilde{\mu}_1 - \tilde{\mu}_2| \uparrow$ (not desirable!).

FLD: Two Classes

Define the **within-class scatter** for projected samples by $\tilde{s}_1^2 + \tilde{s}_2^2$, where

$$\tilde{s}_i^2 = \sum_{y \in \mathcal{Y}_i} (y - \tilde{\mu}_i)^2 = \mathbf{w}^\top \underbrace{\left[\sum_{\mathbf{x} \in C_i} (\mathbf{x} - \boldsymbol{\mu}_i)(\mathbf{x} - \boldsymbol{\mu}_i)^\top \right]}_{\mathbf{S}_i} \mathbf{w}.$$

FLD finds

$$\mathbf{w} = \arg \max_{\mathbf{w}} \frac{|\tilde{\mu}_1 - \tilde{\mu}_2|^2}{\tilde{s}_1^2 + \tilde{s}_2^2} = \arg \max_{\mathbf{w}} \frac{\mathbf{w}^\top \mathbf{S}_B \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_W \mathbf{w}},$$

where $\mathbf{S}_W = \mathbf{S}_1 + \mathbf{S}_2$ (**within-class scatter matrix**) and $\mathbf{S}_B = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top$ (**between-class scatter matrix**).

$$\boxed{\arg \max_{\mathbf{w}} \frac{\mathbf{w}^\top \mathbf{S}_B \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_W \mathbf{w}}} \Rightarrow \boxed{\mathbf{S}_B \mathbf{w} = \lambda \mathbf{S}_W \mathbf{w}}$$

(generalized eigenvalue problem).

Multiple Discriminant Functions

For the case of K classes, FLD involves $K - 1$ discriminant functions, i.e., the projection is from \mathbb{R}^D to \mathbb{R}^{K-1} .

Given a set of data $\{\mathbf{x} \in \mathbb{R}^D\}$, one wishes to find a linear lower-dimensional embedding $\mathbf{W}^\top \in \mathbb{R}^{(K-1) \times D}$ such that $\{\mathbf{y} = \mathbf{W}^\top \mathbf{x}\}$ are classified as well as possible in the lower-dimensional space.

$$\underbrace{\begin{bmatrix} y_1 \\ \vdots \\ y_{K-1} \end{bmatrix}}_{\mathbf{y}} = \underbrace{\begin{bmatrix} \mathbf{w}_1^\top \\ \vdots \\ \mathbf{w}_{K-1}^\top \end{bmatrix}}_{\mathbf{W}^\top} \underbrace{\begin{bmatrix} x_1 \\ \vdots \\ x_D \end{bmatrix}}_{\mathbf{x}}.$$

Scatter Matrices

Within-class scatter matrix

$$\mathbf{S}_W = \sum_{i=1}^K \sum_{\mathbf{x} \in C_i} (\mathbf{x} - \boldsymbol{\mu}_i)(\mathbf{x} - \boldsymbol{\mu}_i)^\top.$$

Between-class scatter matrix

$$\mathbf{S}_B = \sum_{i=1}^K \sum_{C_i} (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^\top = \sum_{i=1}^K N_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^\top.$$

Total scatter matrix: $\mathbf{S}_T = \mathbf{S}_W + \mathbf{S}_B$

$$\mathbf{S}_T = \sum_{\mathbf{x}} (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top.$$

$$\text{Rank}(\mathbf{S}_B) \leq K - 1, \quad \text{Rank}(\mathbf{S}_W) \leq N - K, \quad \text{Rank}(\mathbf{S}_T) \leq N - 1.$$

Total Scatter Matrix

Define $\mathbf{X} = [\mathbf{X}_1, \dots, \mathbf{X}_K]$ where \mathbf{X}_i is a matrix whose columns are associated with data vectors belonging to C_i .

Define

$$\begin{aligned}\mathbf{H}_W &= [\mathbf{X}_1 - \mu_1 \mathbf{e}_1^\top, \dots, \mathbf{X}_K - \mu_K \mathbf{e}_K^\top], \\ \mathbf{H}_B &= [(\mu_1 - \mu) \mathbf{e}_1^\top, \dots, (\mu_K - \mu) \mathbf{e}_K^\top], \\ \mathbf{H}_T &= [\mathbf{x}_1 - \mu, \dots, \mathbf{x}_N - \mu].\end{aligned}$$

One can easily see that $\mathbf{H}_T = \mathbf{X} - \mu \mathbf{e}^\top = \mathbf{H}_W + \mathbf{H}_B$.

We also have $\mathbf{S}_W = \mathbf{H}_W \mathbf{H}_W^\top$, $\mathbf{S}_B = \mathbf{H}_B \mathbf{H}_B^\top$, $\mathbf{S}_T = \mathbf{H}_T \mathbf{H}_T^\top$.

Since $\mathbf{H}_W \mathbf{H}_B^\top = 0$, we have

$$\mathbf{S}_T = (\mathbf{H}_W + \mathbf{H}_B)(\mathbf{H}_W + \mathbf{H}_B)^\top = \mathbf{S}_W + \mathbf{S}_B.$$

The column vectors of \mathbf{S}_W and \mathbf{S}_B are linear combinations of centered data samples.

FLD: Multiple Classes

Define

$$\tilde{\mathbf{S}}_W = \sum_{i=1}^K \sum_{\mathbf{y} \in \mathcal{Y}_i} (\mathbf{y} - \tilde{\boldsymbol{\mu}}_i) (\mathbf{y} - \tilde{\boldsymbol{\mu}}_i)^\top$$

$$\tilde{\mathbf{S}}_B = \sum_{i=1}^K N_i (\tilde{\boldsymbol{\mu}}_i - \tilde{\boldsymbol{\mu}}) (\tilde{\boldsymbol{\mu}}_i - \tilde{\boldsymbol{\mu}})^\top.$$

One can easily show that

$$\begin{aligned}\tilde{\mathbf{S}}_W &= \mathbf{W}^\top \mathbf{S}_W \mathbf{W}, \\ \tilde{\mathbf{S}}_B &= \mathbf{W}^\top \mathbf{S}_B \mathbf{W}.\end{aligned}$$

FLD seeks $K - 1$ discriminant functions \mathbf{W} such that $\mathbf{y} = \mathbf{W}^\top \mathbf{x}$:

$$\begin{aligned}\mathbf{W} &= \arg \max_{\mathbf{W}} \mathcal{J}_{FLD} \\ &= \arg \max_{\mathbf{W}} \operatorname{tr} \left\{ \tilde{\mathbf{S}}_W^{-1} \tilde{\mathbf{S}}_B \right\} \\ &= \arg \max_{\mathbf{W}} \operatorname{tr} \left\{ \left(\mathbf{W}^\top \mathbf{S}_W \mathbf{W} \right)^{-1} \left(\mathbf{W}^\top \mathbf{S}_B \mathbf{W} \right) \right\},\end{aligned}$$

leading to

$$\boxed{\arg \max_{\mathbf{W}} \mathcal{J}_{FLD}} \Rightarrow \boxed{\mathbf{S}_B \mathbf{w}_i = \lambda_i \mathbf{S}_W \mathbf{w}_i}.$$

generalized eigenvalue problem

Rayleigh Quotient

Definition

Let $\mathbf{A} \in \mathbb{R}^{m \times m}$ be symmetric. The Rayleigh quotient $R(\mathbf{x}, \mathbf{A})$ is defined by

$$R(\mathbf{x}, \mathbf{A}) = \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}.$$

Theorem

Let $\mathbf{A} \in \mathbb{R}^{m \times m}$ be symmetric with its eigenvalues being $\{\lambda_1 \geq \dots \geq \lambda_m\}$. For $\forall \mathbf{x} \neq \mathbf{0} \in \mathbb{R}^m$, we have

$$\lambda_m \leq \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} \leq \lambda_1,$$

and in particular,

$$\lambda_m = \min_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}, \quad \lambda_1 = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}.$$

An Extremal Property of Generalized Eigenvalues

Theorem

Let \mathbf{A} and \mathbf{B} be $m \times m$ matrices, with \mathbf{A} being nonnegative definite and \mathbf{B} positive definite. For $h = 1, \dots, m$, define

$$\mathbf{X}_h = [\mathbf{x}_1, \dots, \mathbf{x}_h], \quad \mathbf{Y}_h = [\mathbf{x}_h, \dots, \mathbf{x}_m],$$

where $\mathbf{x}_1, \dots, \mathbf{x}_m$ are linear independent eigenvectors of $\mathbf{B}^{-1}\mathbf{A}$ corresponding to the eigenvalues

$$\lambda_1(\mathbf{B}^{-1}\mathbf{A}) \geq \dots \geq \lambda_m(\mathbf{B}^{-1}\mathbf{A}).$$

Then

$$\lambda_m(\mathbf{B}^{-1}\mathbf{A}) = \min_{\mathbf{Y}_{h+1}^\top \mathbf{B} \mathbf{x} = 0} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{B} \mathbf{x}},$$
$$\lambda_1(\mathbf{B}^{-1}\mathbf{A}) = \max_{\mathbf{X}_{h-1}^\top \mathbf{B} \mathbf{x} = 0} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{B} \mathbf{x}},$$

where $\mathbf{x} = 0$ is excluded.

Relation to Least Squares Regression: Binary Class

Given a training set $\{\mathbf{x}_i, y_i\}_{i=1}^N$, where $\mathbf{x}_i \in \mathbb{R}^D$ and $y_i \in \{1, -1\}$, consider a linear discriminant function:

$$f(\mathbf{x}_i) = \mathbf{w}^\top \mathbf{x}_i + b.$$

Partition the data matrix into two groups, each group of which contains examples in class 1 or class 2, i.e., $\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2]$, where $\mathbf{X}_1 \in \mathbb{R}^{D \times N_1}$ and $\mathbf{X}_2 \in \mathbb{R}^{D \times N_2}$.

Define binary label vector $\mathbf{y} \in \mathbb{R}^N$, then LS regression is formulated as

$$\arg \min_{\mathbf{w}, b} \|\mathbf{y} - \mathbf{X}^\top \mathbf{w} - b \mathbf{1}_N\|^2,$$

where $\mathbf{1}_N$ is the N -dimensional vector of all ones, which can be re-written as

$$\arg \min_{\mathbf{w}, b} \left\| \begin{bmatrix} \mathbf{X}_1^\top & \mathbf{1}_{N_1} \\ \mathbf{X}_2^\top & \mathbf{1}_{N_2} \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix} - \begin{bmatrix} \mathbf{1}_{N_1} \\ -\mathbf{1}_{N_2} \end{bmatrix} \right\|^2.$$

The solution to this LS problem satisfies the **normal equation**:

$$\begin{bmatrix} \mathbf{X}_1 & \mathbf{X}_2 \\ \mathbf{1}_{N_1}^\top & \mathbf{1}_{N_2}^\top \end{bmatrix} \begin{bmatrix} \mathbf{X}_1^\top & \mathbf{1}_{N_1} \\ \mathbf{X}_2^\top & \mathbf{1}_{N_2} \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix} = \begin{bmatrix} \mathbf{X}_1 & \mathbf{X}_2 \\ \mathbf{1}_{N_1}^\top & \mathbf{1}_{N_2}^\top \end{bmatrix} \begin{bmatrix} \mathbf{1}_{N_1} \\ -\mathbf{1}_{N_2} \end{bmatrix},$$

which is written as

$$\begin{bmatrix} \mathbf{X}_1\mathbf{X}_1^\top + \mathbf{X}_2\mathbf{X}_2^\top & \mathbf{X}_1\mathbf{1}_{N_1} + \mathbf{X}_2\mathbf{1}_{N_2} \\ \mathbf{1}_{N_1}^\top\mathbf{X}_1^\top + \mathbf{1}_{N_2}^\top\mathbf{X}_2^\top & \mathbf{1}_{N_1}^\top\mathbf{1}_{N_1} + \mathbf{1}_{N_2}^\top\mathbf{1}_{N_2} \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix} = \begin{bmatrix} \mathbf{X}_1\mathbf{1}_{N_1} - \mathbf{X}_2\mathbf{1}_{N_2} \\ \mathbf{1}_{N_1}^\top\mathbf{1}_{N_1} - \mathbf{1}_{N_2}^\top\mathbf{1}_{N_2} \end{bmatrix}.$$

Recall

$$\begin{aligned} \mathbf{S}_B &= (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \\ \mathbf{S}_W &= (\mathbf{X}_1 - \boldsymbol{\mu}_1\mathbf{1}_{N_1}^\top)(\mathbf{X}_1 - \boldsymbol{\mu}_1\mathbf{1}_{N_1}^\top)^\top + (\mathbf{X}_2 - \boldsymbol{\mu}_2\mathbf{1}_{N_2}^\top)(\mathbf{X}_2 - \boldsymbol{\mu}_2\mathbf{1}_{N_2}^\top)^\top \\ &= \mathbf{X}_1\mathbf{X}_1^\top - N_1\boldsymbol{\mu}_1\boldsymbol{\mu}_1^\top + \mathbf{X}_2\mathbf{X}_2^\top - N_2\boldsymbol{\mu}_2\boldsymbol{\mu}_2^\top. \end{aligned}$$

With \mathbf{S}_B and \mathbf{S}_W , the normal equation is written as

$$\begin{bmatrix} \mathbf{S}_W + N_1\boldsymbol{\mu}_1\boldsymbol{\mu}_1^\top + N_2\boldsymbol{\mu}_2\boldsymbol{\mu}_2^\top & N_1\boldsymbol{\mu}_1 + N_2\boldsymbol{\mu}_2 \\ (N_1\boldsymbol{\mu}_1 + N_2\boldsymbol{\mu}_2)^\top & N_1 + N_2 \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix} = \begin{bmatrix} N_1\boldsymbol{\mu}_1 - N_2\boldsymbol{\mu}_2 \\ N_1 - N_2 \end{bmatrix}.$$

Solve the 2nd equation for b to obtain

$$b = \frac{(N_1 - N_2) - (N_1\boldsymbol{\mu}_1 + N_2\boldsymbol{\mu}_2)^\top \mathbf{w}}{N_1 + N_2}.$$

Substitute this into the 1st equation to obtain

$$\left[\mathbf{S}_W + \frac{N_1 N_2}{N_1 + N_2} \mathbf{S}_B \right] \mathbf{w} = 2N_1 N_2 (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2).$$

Note that the vector $\mathbf{S}_B \mathbf{w}$ is in the direction of $\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2$ for $\forall \mathbf{w}$, since $\mathbf{S}_B \mathbf{w} = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \mathbf{w}$.

Thus we write

$$\frac{N_1 N_2}{N_1 + N_2} \mathbf{S}_B \mathbf{w} = (2N_1 N_2 - \alpha)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2).$$

Then we have

$$\mathbf{w} = \alpha \mathbf{S}_W^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2),$$

which is identical to FLD solutions except for scaling factor.

Simultaneous Diagonalization

The goal: Given two symmetric matrices, Σ_1 and Σ_2 , find a linear transformation \mathbf{W} such that

$$\begin{aligned}\mathbf{W}^\top \Sigma_1 \mathbf{W} &= \mathbf{I}, \\ \mathbf{W}^\top \Sigma_2 \mathbf{W} &= \Lambda. \quad (\text{diagonal})\end{aligned}$$

Methods: It turns out that simultaneous diagonalization involves the **generalized eigen-decomposition**.

- ▶ Two-stage method
 1. whitening
 2. unitary transformation
- ▶ Single-stage method: generalized eigenvalue decomposition

Simultaneous Diagonalization: Algorithm Outline

1. First, whiten Σ_1 , i.e.,

$$\begin{aligned}D^{-\frac{1}{2}} U_1^\top \Sigma_1 U_1 D^{-\frac{1}{2}} &= I, \\D^{-\frac{1}{2}} U_1^\top \Sigma_2 U_1 D^{-\frac{1}{2}} &= K, \quad (\text{not diagonal}),\end{aligned}$$

where $\Sigma_1 = U_1 D U_1^\top$.

2. Second, apply an unitary transformation to diagonalize K , i.e.,

$$\begin{aligned}U_2^\top I U_2 &= I, \\U_2^\top K U_2 &= \Lambda,\end{aligned}$$

where $K = U_2 \Lambda U_2^\top$.

Then, the transformation W which simultaneously diagonalizes Σ_1 and Σ_2 , is given by, $W = U_1 D^{-\frac{1}{2}} U_2$, such that $W^\top \Sigma_1 W = I$ and $W^\top \Sigma_2 W = \Lambda$

Simultaneous Diagonalization: Generalized Eigen-Decomposition

Alternatively we can diagonalize two symmetric matrices Σ_1 and Σ_2 as

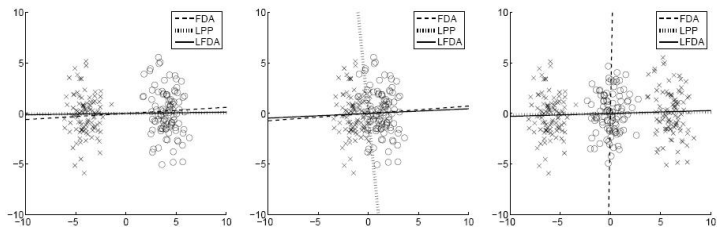
$$\begin{aligned}W^T \Sigma_1 W &= I, \\W^T \Sigma_2 W &= \Lambda, \quad (\text{diagonal})\end{aligned}$$

where Λ , W are eigenvalues and eigenvectors of $\Sigma_1^{-1}\Sigma_2$, i.e.,

$$\Sigma_1^{-1}\Sigma_2 W = W\Lambda.$$

Prove it!

Example: Multi-Modal Data



Alternative Expressions of \mathbf{S}_W and \mathbf{S}_B

Alternatively, \mathbf{S}_W and \mathbf{S}_B are expressed as

$$\mathbf{S}_W = \frac{1}{2} \sum_i \sum_j A_{ij}^W (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top,$$

$$\mathbf{S}_B = \frac{1}{2} \sum_i \sum_j A_{ij}^B (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top,$$

$$A_{ij}^W = \begin{cases} \frac{1}{N_k} & \text{if } \mathbf{x}_i \in C_k \text{ and } \mathbf{x}_j \in C_k, \\ 0 & \text{if } \mathbf{x}_i \text{ and } \mathbf{x}_j \text{ are in different classes,} \end{cases}$$

$$A_{ij}^B = \begin{cases} \frac{1}{N} - \frac{1}{N_k} & \text{if } \mathbf{x}_i \in C_k \text{ and } \mathbf{x}_j \in C_k, \\ \frac{1}{N} & \text{if } \mathbf{x}_i \text{ and } \mathbf{x}_j \text{ are in different classes.} \end{cases}$$

$$\begin{aligned}
\mathbf{S}_W &= \sum_{i=1}^K \sum_{\mathbf{x} \in C_i} (\mathbf{x} - \boldsymbol{\mu}_i)(\mathbf{x} - \boldsymbol{\mu}_i)^\top \\
&= \sum_{i=1}^K \sum_{\mathbf{x}_j \in C_i} \left(\mathbf{x}_j - \frac{1}{N_i} \sum_{\mathbf{x}_u \in C_i} \mathbf{x}_u \right) \left(\mathbf{x}_j - \frac{1}{N_i} \sum_{\mathbf{x}_v \in C_i} \mathbf{x}_v \right)^\top \\
&= \sum_{i=1}^K \sum_{\mathbf{x}_j \in C_i} \left\{ \mathbf{x}_j \mathbf{x}_j^\top - \frac{1}{N_i} \sum_{\mathbf{x}_v \in C_i} \mathbf{x}_j \mathbf{x}_v^\top - \frac{1}{N_i} \sum_{\mathbf{x}_u \in C_i} \mathbf{x}_u \mathbf{x}_j^\top + \frac{1}{N_i^2} \sum_{\mathbf{x}_u \in C_i} \sum_{\mathbf{x}_v \in C_i} \mathbf{x}_u \mathbf{x}_v^\top \right\} \\
&= \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^\top - \sum_{i=1}^K \frac{1}{N_i} \sum_{\mathbf{x}_u \in C_i} \sum_{\mathbf{x}_v \in C_i} \mathbf{x}_u \mathbf{x}_v^\top \\
&= \sum_{i=1}^N \left(\sum_{j=1}^N A_{ij}^W \right) \mathbf{x}_i \mathbf{x}_i^\top - \sum_{i=1}^N \sum_{j=1}^N A_{ij}^W \mathbf{x}_i \mathbf{x}_j^\top \\
&= \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N A_{ij}^W \left(\mathbf{x}_i \mathbf{x}_i^\top + \mathbf{x}_j \mathbf{x}_j^\top - \mathbf{x}_i \mathbf{x}_j^\top - \mathbf{x}_j \mathbf{x}_i^\top \right) \\
&= \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N A_{ij}^W (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top.
\end{aligned}$$

$$\begin{aligned}
\mathbf{S}_B &= \mathbf{S}_T - \mathbf{S}_W \\
&= \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top - \mathbf{S}_W \\
&= \sum_{i=1}^N \left(\sum_{j=1}^N \frac{1}{N} \right) \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N \mathbf{x}_i \mathbf{x}_j^\top \\
&\quad - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N A_{ij}^W \left(\mathbf{x}_i \mathbf{x}_i^\top + \mathbf{x}_j \mathbf{x}_j^\top - \mathbf{x}_i \mathbf{x}_j^\top - \mathbf{x}_j \mathbf{x}_i^\top \right) \\
&= \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \left(\frac{1}{N} - A_{ij}^W \right) \left(\mathbf{x}_i \mathbf{x}_i^\top + \mathbf{x}_j \mathbf{x}_j^\top - \mathbf{x}_i \mathbf{x}_j^\top - \mathbf{x}_j \mathbf{x}_i^\top \right) \\
&= \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \left(\frac{1}{N} - A_{ij}^W \right) (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top.
\end{aligned}$$

Local Within-Class and Between-Class Scatter

Given weighted adjacency matrix $[A_{ij}]$, introduce **local within-class scatter** and **local between-class scatter**:

$$\bar{\mathbf{S}}_W = \frac{1}{2} \sum_i \sum_j \bar{A}_{ij}^W (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top,$$

$$\bar{\mathbf{S}}_B = \frac{1}{2} \sum_i \sum_j \bar{A}_{ij}^B (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top,$$

$$\bar{A}_{ij}^W = \begin{cases} \frac{A_{ij}}{N_k} & \text{if } \mathbf{x}_i \in C_k \text{ and } \mathbf{x}_j \in C_k, \\ 0 & \text{if } \mathbf{x}_i \text{ and } \mathbf{x}_j \text{ are in different classes,} \end{cases}$$

$$\bar{A}_{ij}^B = \begin{cases} A_{ij} \left(\frac{1}{N} - \frac{1}{N_k} \right) & \text{if } \mathbf{x}_i \in C_k \text{ and } \mathbf{x}_j \in C_k, \\ \frac{1}{N} & \text{if } \mathbf{x}_i \text{ and } \mathbf{x}_j \text{ are in different classes.} \end{cases}$$

Local Fisher Discriminant Analysis (LFDA)

Proposed by M. Sugiyama (ICML-2006).

LFDA seeks $K - 1$ discriminant functions \mathbf{W} such that $\mathbf{y} = \mathbf{W}^\top \mathbf{x}$:

$$\arg \max_{\mathbf{W}} \text{tr} \left\{ \left(\mathbf{W}^\top \bar{\mathbf{S}}_W \mathbf{W} \right)^{-1} \left(\mathbf{W}^\top \bar{\mathbf{S}}_B \mathbf{W} \right) \right\},$$

Local within-class scatter matrix

$$\bar{\mathbf{S}}_W = \frac{1}{2} \sum_i \sum_j \bar{A}_{ij}^W (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top,$$

Local between-class scatter matrix

$$\bar{\mathbf{S}}_B = \frac{1}{2} \sum_i \sum_j \bar{A}_{ij}^B (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top.$$