

Fisher Linear Discriminant Analysis (2 Classes)

April 27, 2025

Contents

1	Introduction to Model	2
2	Model Fitting	4
2.1	Reasoning	4
2.2	Algorithm	5
2.3	Proof of Algorithm Feasibility	6
2.4	Proof of Algorithm Optimality	6
3	Model-based Prediction	18
3.1	Simplified Model	18
3.2	Proof of Simplified Model	18
4	Example	20
4.1	Example 1	20
4.2	Example 2	22
4.3	Example 3	23
4.4	Example 4	25
5	Some other Important Properties	27
6	Appendix: Calculus and Linear Algebra	32

Chapter 1

Introduction to Model

Population dataset \mathcal{P} exists. Given available dataset \mathcal{A} and its splits training dataset \mathcal{I} and test dataset \mathcal{S} , each dataset contains samples with features $\mathbf{x} = [[\mathbf{x}]_1, [\mathbf{x}]_2, \dots, [\mathbf{x}]_p]^\top$ and $y \in \{0, 1\}$. We want to establish a model:

$$M(\mathbf{x}) \tag{1.1}$$

which:

$$\text{maximize} \quad \sum_{i=1}^{N(\mathcal{P})} \text{I}_{\text{acc}}(M(\mathbf{x}|_{\mathcal{I}_i}), y|_{\mathcal{I}_i}) \quad (\text{Accuracy}) \tag{1.2}$$

$$\text{I}_{\text{acc}}(a, b) = \begin{cases} 1 & \text{if } a = b \\ 0 & \text{otherwise} \end{cases} \tag{1.3}$$

Fisher LDA is a **supervised** dimensionality reduction method that projects data onto a lower-dimensional space while:

- Maximizing **between-class scatter** (separation between class means)
- Minimizing **within-class scatter** (variance within each class)

which tends to fit a good classifier.

Mathematically, $\mathbf{w}^\top \mathbf{x}$ means projecting \mathbf{x} on a direction \mathbf{w} .

Define \mathcal{I}_0 as the training data set that contains only samples from class 0 ($y|_{\mathcal{I}_i} = 0$), and \mathcal{I}_1 as the training data set that contains only samples from class 1 ($y|_{\mathcal{I}_i} = 1$).

The objective is to find a projection matrix \mathbf{w} that maximizes the ratio between **between-class scatter** and **within-class scatter**:

$$J_1(\mathbf{w}) = \frac{\|\mathbf{w}^\top \boldsymbol{\mu}_0 - \mathbf{w}^\top \boldsymbol{\mu}_1\|_2^2}{\mathbf{w}^\top \boldsymbol{\Sigma}_0 \mathbf{w} + \mathbf{w}^\top \boldsymbol{\Sigma}_1 \mathbf{w}} \tag{1.4}$$

where:

$$\boldsymbol{\mu}_0 = \frac{1}{N(\mathcal{I}_0)} \sum_{i=1}^{N(\mathcal{I}_0)} \mathbf{x}|_{\mathcal{I}_i}^{\mathcal{I}_0} \tag{1.5}$$

$$\boldsymbol{\mu}_1 = \frac{1}{N(\mathcal{I}_1)} \sum_{i=1}^{N(\mathcal{I}_1)} \mathbf{x}|_{\mathcal{I}_i}^{\mathcal{I}_1} \tag{1.6}$$

$$\boldsymbol{\Sigma}_0 = \sum_{i=1}^{N(\mathcal{I}_0)} (\mathbf{x}|_{\mathcal{I}_i}^{\mathcal{I}_0} - \boldsymbol{\mu}_0) (\mathbf{x}|_{\mathcal{I}_i}^{\mathcal{I}_0} - \boldsymbol{\mu}_0)^\top \tag{1.7}$$

$$\boldsymbol{\Sigma}_1 = \sum_{i=1}^{N(\mathcal{I}_1)} (\mathbf{x}|_{\mathcal{I}_i}^{\mathcal{I}_1} - \boldsymbol{\mu}_1) (\mathbf{x}|_{\mathcal{I}_i}^{\mathcal{I}_1} - \boldsymbol{\mu}_1)^\top \tag{1.8}$$

- μ_0 is the class center of class 0 in training dataset \mathcal{I} .
- μ_1 is the class center of class 1 in training dataset \mathcal{I} .
- Σ_0 is the class variance of class 0 in training dataset \mathcal{I} .
- Σ_1 is the class variance of class 1 in training dataset \mathcal{I} .

Define the within-class scatter matrix:

$$\mathbf{S}_w = \Sigma_0 + \Sigma_1 \quad (1.9)$$

We require \mathbf{S}_w to be non-zero as this is not naturally guaranteed, so that the mathematical programming problem (P1) can be well-defined.

Define the between-class scatter matrix:

$$\mathbf{S}_b = (\mu_0 - \mu_1)(\mu_0 - \mu_1)^\top \quad (1.10)$$

Substitute the (1.9) and (1.10) into (1.11) yields:

$$J_1(\mathbf{w}) = \frac{\|\mathbf{w}^\top \mu_0 - \mathbf{w}^\top \mu_1\|_2^2}{\mathbf{w}^\top \Sigma_0 \mathbf{w} + \mathbf{w}^\top \Sigma_1 \mathbf{w}} = \frac{\mathbf{w}^\top (\mu_0 - \mu_1)(\mu_0 - \mu_1)^\top \mathbf{w}}{\mathbf{w}^\top (\Sigma_0 + \Sigma_1) \mathbf{w}} = \frac{\mathbf{w}^\top \mathbf{S}_b \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_w \mathbf{w}} \quad (1.11)$$

The J_1 is the objective we want to optimize. Define the problem:

$$\begin{array}{ll} \text{(P1) Maximize}_{\mathbf{w}} & J_1(\mathbf{w}) = \frac{\mathbf{w}^\top \mathbf{S}_b \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_w \mathbf{w}} \\ \text{Subject to} & \mathbf{w}^\top \mathbf{S}_w \mathbf{w} > 0 \end{array} \quad (1.12)$$

The problem (P1) is the target of LDA model fitting.

We need to admit that there can be more than one optimal solution of (P1) when \mathbf{S}_w is singular. In this case, we only need to find one of the optimal solutions for (P1). However, for interest, we also discuss about the optimal solution sets in the following several chapters.

In the next section we are going to talk about how to maximize J_1 .

Problem Solution Status

We define the following three problem solution status:

- Non-trivial Optimal Solution Status: (P1) has a non-empty optimal solution set, and the optimal solution set is different from the feasible solution set.
- Trivial Optimal Solution Status: (P1) has a non-empty optimal solution set, and the optimal solution set is equivalent to the feasible solution set.
- Infinite Optimal Solution Status: (P1) has an empty optimal solution set, and the objective value has no upper bound.

In details:

Status	Optimal Set	Optimal Objective Value	Relationship between Feasible and Optimal Sets
Non-trivial	Non-empty	Finite	Different
Trivial	Non-empty	Finite	Same
Infinite	Empty	Unbounded	Different

The trivial optimal solution status implies that (P1) is meaningless at all, as all feasible solutions can be optimal solutions. The infinite optimal solution status implies that (P1) has no upper bound.

You may find the examples in the Example Chapter.

Chapter 2

Model Fitting

2.1 Reasoning

The model fitting aims to find an optimal solution, and the solution status for (P1).

We discuss the following two cases:

$\exists \gamma \in \mathbb{R}$ such that $\mathbf{S}_b = \gamma \mathbf{S}_w$:

When there $\exists \gamma \in \mathbb{R}$ such that $\mathbf{S}_b = \gamma \mathbf{S}_w$, the objective function $J_1(\mathbf{w}) = \gamma$ for all feasible solutions, which means that the optimal solution status is trivial. Therefore, $\mathbf{J}_1(\mathbf{w}) = 0$ for any feasible \mathbf{w} , which implies that the problem has a trivial optimal solution.

In this case, the feasible solution satisfies $\mathbf{w}^\top \mathbf{S}_w \mathbf{w} > 0$, which is equivalent to $\mathbf{S}_w \mathbf{w} \neq \mathbf{0}$.

We know that since \mathbf{S}_w is non-zero, then there must exists at least one non-zero column. We just choose a non-zero column of \mathbf{S}_w to be a feasible as well as an optimal solution of (P1).

$\forall \gamma \in \mathbb{R}$ such that $\mathbf{S}_b \neq \gamma \mathbf{S}_w$:

$\forall \gamma$, $\mathbf{S}_b \neq \gamma \mathbf{S}_w$, implies that the $\mathbf{S}_b \neq \mathbf{0}$. According to the definition of \mathbf{S}_b : $\mathbf{S}_b = (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top$, $\boldsymbol{\mu}_0 \neq \boldsymbol{\mu}_1$.

We notice that the numerator and denominator are the quadratic form of \mathbf{w} , which means that the solution is irrelevant to the magnitude of \mathbf{w} but its direction. Without loss of generality, we let $\mathbf{w}^\top \mathbf{S}_w \mathbf{w} = 1$. Then we can remove the condition $\mathbf{w}^\top \mathbf{S}_w \mathbf{w} > 0$ because it is already included in $\mathbf{w}^\top \mathbf{S}_w \mathbf{w} = 1$. Then we can solve the (P2) instead of (P1):

$$\boxed{\begin{array}{ll} \text{(P2)} & \underset{\mathbf{w}}{\text{Minimize}} \quad J_2(\mathbf{w}) = -\mathbf{w}^\top \mathbf{S}_b \mathbf{w} \\ & \text{Subject to} \quad \mathbf{w}^\top \mathbf{S}_w \mathbf{w} = 1 \end{array}} \quad (2.1)$$

$$\text{Subject to} \quad \mathbf{w}^\top \mathbf{S}_w \mathbf{w} = 1 \quad (2.2)$$

Then we can apply the Generalized Fritz Johns Condition (**Theorem 4**):

$$\lambda \cdot \nabla_{\mathbf{w}}(-\mathbf{w}^\top \mathbf{S}_b \mathbf{w}) + v \cdot \nabla_{\mathbf{w}}(\mathbf{w}^\top \mathbf{S}_w \mathbf{w} - 1) = 0 \quad (2.3)$$

$$-\lambda \mathbf{S}_b \mathbf{w} + v \mathbf{S}_w \mathbf{w} = 0 \quad (2.4)$$

$$v \mathbf{S}_w \mathbf{w} = \lambda \mathbf{S}_b \mathbf{w} \quad (2.5)$$

and also the:

$$\lambda \geq 0 \quad (2.6)$$

$$[\lambda, v]^\top \neq \mathbf{0} \quad (2.7)$$

yields an equivalent problem (P3) of (P2):

$$\boxed{\begin{array}{ll} \text{(P3)} & \underset{\mathbf{w}, \lambda, v}{\text{Minimize}} \quad J_3(\mathbf{w}, \lambda, v) = -\mathbf{w}^\top \mathbf{S}_b \mathbf{w} \\ & \text{Subject to} \quad \lambda \mathbf{S}_b \mathbf{w} = v \mathbf{S}_w \mathbf{w} \end{array}} \quad (2.8)$$

$$\text{Subject to} \quad \lambda \mathbf{S}_b \mathbf{w} = v \mathbf{S}_w \mathbf{w} \quad (2.9)$$

$$\mathbf{w}^\top \mathbf{S}_w \mathbf{w} = 1 \quad (2.10)$$

$$\lambda \geq 0 \quad (2.11)$$

$$[\lambda, v]^\top \neq \mathbf{0} \quad (2.12)$$

$\mathbf{w}^\top \mathbf{S}_w \mathbf{w} = 1 > 0$ implies that $\mathbf{S}_w \mathbf{w} \neq \mathbf{0}$ (**Proposition 4**). We claim that $\lambda \neq 0$. If $\lambda = 0$, then $v \mathbf{S}_w \mathbf{w} = \lambda \mathbf{S}_b \mathbf{w} = \mathbf{0}$. Since $\mathbf{S}_w \mathbf{w} \neq \mathbf{0}$ then $v = 0$, which contradicts that $[\lambda, v]^\top \neq \mathbf{0}$. Therefore, $\lambda > 0$, $\lambda \mathbf{S}_b \mathbf{w} = v \mathbf{S}_w \mathbf{w}$ can be converted $\mathbf{S}_b \mathbf{w} = \frac{v}{\lambda} \mathbf{S}_w \mathbf{w}$.

Since $\lambda > 0$, then $[\lambda, v]^\top \neq \mathbf{0}$ holds for all $v \in \mathbb{R}$. We can remove this condition.

Now, since $\frac{v}{\lambda}$ ranges in \mathbb{R} , we can replace it with α , which yields (P4):

$$\boxed{\begin{array}{ll} \text{(P4)} & \text{Minimize}_{\mathbf{w}, \alpha} \quad J_4(\mathbf{w}, \alpha) = -\mathbf{w}^\top \mathbf{S}_b \mathbf{w} \end{array}} \quad (2.13)$$

$$\text{Subject to} \quad \mathbf{S}_b \mathbf{w} = \alpha \mathbf{S}_w \mathbf{w} \quad (2.14)$$

$$\mathbf{w}^\top \mathbf{S}_w \mathbf{w} = 1 \quad (2.15)$$

Then:

$$\mathbf{S}_b \mathbf{w} = \alpha \mathbf{S}_w \mathbf{w} \quad (2.16)$$

$$(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w} = \alpha \mathbf{S}_w \mathbf{w} \quad (2.17)$$

Suppose that $(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w} \neq 0$:

$$\alpha \mathbf{S}_w \mathbf{w} = (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w} \quad (2.18)$$

$$\frac{\alpha}{(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}} \mathbf{S}_w \mathbf{w} = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1 \quad (2.19)$$

Notice that $\frac{\alpha}{(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}}$ is a scalar and therefore we know that $\mathbf{S}_w \mathbf{w}$ is collinear with $\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1$. Since in (P1) the magnitude of \mathbf{w} does not affect the feasibility and optimality of the result, then we solve the (P5) instead:

$$\boxed{\begin{array}{ll} \text{(P5)} & \text{Solve}_{\mathbf{w}, \beta} \quad \beta \cdot (\mathbf{S}_w \mathbf{w}) = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1 \end{array}} \quad (2.20)$$

Since we just want to find a solution for (P5), then we can let $\beta = 1$ and solve the following (P6) instead:

$$\boxed{\begin{array}{ll} \text{(P6)} & \text{Solve}_{\mathbf{w}} \quad \mathbf{S}_w \mathbf{w} = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1 \end{array}} \quad (2.21)$$

Then, a solution of (P6) is the optimal solution of (P1).

However, the (P6) is not guaranteed to have a solution. When (P6) does not have a solution we can show that the objective value of (P1) is unbounded, and therefore have an infinite optimal solution status.

If (P6) has a solution, we still need to discuss whether the optimal solution status of (P6) is trivial or not.

By observation, we notice that if there $\exists \gamma \in \mathbb{R}$, such that $\mathbf{S}_b = \gamma \mathbf{S}_w$, the $J_1(\mathbf{w})$ is not relevant to the \mathbf{w} , which implies that this is a key condition for trivial optimal solution.

2.2 Algorithm

The optimal solution status and an optimal solution can be described as follows:

- $\exists \gamma \in \mathbb{R}$, such that $\mathbf{S}_b = \gamma \mathbf{S}_w$: (P1) has trivial optimal solution status, and an optimal solution is a non-zero column of \mathbf{S}_w .
- $\forall \gamma \in \mathbb{R}$, such that $\mathbf{S}_b \neq \gamma \mathbf{S}_w$, and (P6) has a solution: (P1) has non-trivial optimal solution status, and an optimal solution is a solution of (P6).
- $\forall \gamma \in \mathbb{R}$, such that $\mathbf{S}_b \neq \gamma \mathbf{S}_w$, and (P6) has no solution: (P1) has infinite optimal solution status, and the optimal solution set of (P1) is empty.

The proof can be found in **Claim 14**, **Claim 15**, **Claim 16**.

Based on this, the algorithm is described as follows:

- Step 1: Split the training dataset \mathcal{I} into \mathcal{I}_0 and \mathcal{I}_1
- Step 2: Compute $\boldsymbol{\mu}_0$, $\boldsymbol{\mu}_1$, \mathbf{S}_w . Check if $\mathbf{S}_w = \mathbf{0}$. If yes, return that (P1) is ill-defined. Else if no, turn to Step 3.

- Step 3: Compute \mathbf{S}_b . Check if there exists $\gamma \in \mathbb{R}$ such that $\mathbf{S}_b = \gamma \mathbf{S}_w$. If yes, return that (P1) has trivial optimal solution status, and a non-zero column of \mathbf{S}_w is an optimal solution of (P1). Else if no, turn to Step 4.
- Step 4: Solve (P6), if (P6) has solutions, then the (P1) has non-trivial optimal solution and the solution of (P6) is the optimal solution of (P1). Else if (P6) has no solution, then (P1) has infinite optimal solution status, and the optimal solution set is \emptyset .

2.3 Proof of Algorithm Feasibility

It is obvious that all steps in the algorithm are feasible.

2.4 Proof of Algorithm Optimality

Claim 1 \mathbf{S}_w is symmetric ($\mathbf{S}_w = \mathbf{S}_w^\top$)

Proof:

$$\mathbf{S}_w^\top = \left(\sum_{i=1}^{N(\mathcal{I}_0)} (\mathbf{x}_{\diamond i}^{\mathcal{I}_0} - \boldsymbol{\mu}_0) (\mathbf{x}_{\diamond i}^{\mathcal{I}_0} - \boldsymbol{\mu}_0)^\top + \sum_{i=1}^{N(\mathcal{I}_1)} (\mathbf{x}_{\diamond i}^{\mathcal{I}_1} - \boldsymbol{\mu}_1) (\mathbf{x}_{\diamond i}^{\mathcal{I}_1} - \boldsymbol{\mu}_1)^\top \right)^\top \quad (2.22)$$

$$\mathbf{S}_w^\top = \sum_{i=1}^{N(\mathcal{I}_0)} \left((\mathbf{x}_{\diamond i}^{\mathcal{I}_0} - \boldsymbol{\mu}_0) (\mathbf{x}_{\diamond i}^{\mathcal{I}_0} - \boldsymbol{\mu}_0)^\top \right)^\top + \sum_{i=1}^{N(\mathcal{I}_1)} \left((\mathbf{x}_{\diamond i}^{\mathcal{I}_1} - \boldsymbol{\mu}_1) (\mathbf{x}_{\diamond i}^{\mathcal{I}_1} - \boldsymbol{\mu}_1)^\top \right)^\top \quad (2.23)$$

$$\mathbf{S}_w^\top = \sum_{i=1}^{N(\mathcal{I}_0)} (\mathbf{x}_{\diamond i}^{\mathcal{I}_0} - \boldsymbol{\mu}_0) (\mathbf{x}_{\diamond i}^{\mathcal{I}_0} - \boldsymbol{\mu}_0)^\top + \sum_{i=1}^{N(\mathcal{I}_1)} (\mathbf{x}_{\diamond i}^{\mathcal{I}_1} - \boldsymbol{\mu}_1) (\mathbf{x}_{\diamond i}^{\mathcal{I}_1} - \boldsymbol{\mu}_1)^\top \quad (2.24)$$

$$\mathbf{S}_w^\top = \mathbf{S}_w \quad (2.25)$$

- (2.22): This comes directly from the definition: (1.9),(1.7), (1.8).
- (2.22) \Rightarrow (2.23): Transpose of summation equals to the summation of transpose.
- (2.23) \Rightarrow (2.24): Rearrangement of terms.
- (2.24) \Rightarrow (2.25): This also comes from the definition: (1.9),(1.7), (1.8).

which completes the proof.

Claim 2 \mathbf{S}_w is positive semi-definite.

Proof:

For any $\mathbf{v} \in \mathbb{R}^p$:

$$\mathbf{v}^\top \mathbf{S}_w \mathbf{v} = \mathbf{v}^\top \left(\sum_{i=1}^{N(\mathcal{I}_0)} \left((\mathbf{x}_{\diamond i}^{\mathcal{I}_0} - \boldsymbol{\mu}_0) (\mathbf{x}_{\diamond i}^{\mathcal{I}_0} - \boldsymbol{\mu}_0)^\top \right) + \sum_{i=1}^{N(\mathcal{I}_1)} \left((\mathbf{x}_{\diamond i}^{\mathcal{I}_1} - \boldsymbol{\mu}_1) (\mathbf{x}_{\diamond i}^{\mathcal{I}_1} - \boldsymbol{\mu}_1)^\top \right) \right) \mathbf{v} \quad (2.26)$$

$$\mathbf{v}^\top \mathbf{S}_w \mathbf{v} = \sum_{i=1}^{N(\mathcal{I}_0)} \mathbf{v}^\top (\mathbf{x}_{\diamond i}^{\mathcal{I}_0} - \boldsymbol{\mu}_0) (\mathbf{x}_{\diamond i}^{\mathcal{I}_0} - \boldsymbol{\mu}_0)^\top \mathbf{v} + \sum_{i=1}^{N(\mathcal{I}_1)} \mathbf{v}^\top (\mathbf{x}_{\diamond i}^{\mathcal{I}_1} - \boldsymbol{\mu}_1) (\mathbf{x}_{\diamond i}^{\mathcal{I}_1} - \boldsymbol{\mu}_1)^\top \mathbf{v} \quad (2.27)$$

$$\mathbf{v}^\top \mathbf{S}_w \mathbf{v} = \sum_{i=1}^{N(\mathcal{I}_0)} \left((\mathbf{x}_{\diamond i}^{\mathcal{I}_0} - \boldsymbol{\mu}_0)^\top \mathbf{v} \right)^\top \left((\mathbf{x}_{\diamond i}^{\mathcal{I}_0} - \boldsymbol{\mu}_0)^\top \mathbf{v} \right) + \sum_{i=1}^{N(\mathcal{I}_1)} \left((\mathbf{x}_{\diamond i}^{\mathcal{I}_1} - \boldsymbol{\mu}_1)^\top \mathbf{v} \right)^\top \left((\mathbf{x}_{\diamond i}^{\mathcal{I}_1} - \boldsymbol{\mu}_1)^\top \mathbf{v} \right) \quad (2.28)$$

$$\mathbf{v}^\top \mathbf{S}_w \mathbf{v} = \sum_{i=1}^{N(\mathcal{I}_0)} \left((\mathbf{x}_{|\mathcal{I}_0}^i - \boldsymbol{\mu}_0)^\top \mathbf{v} \right) \left((\mathbf{x}_{|\mathcal{I}_0}^i - \boldsymbol{\mu}_0)^\top \mathbf{v} \right) + \sum_{i=1}^{N(\mathcal{I}_1)} \left((\mathbf{x}_{|\mathcal{I}_1}^i - \boldsymbol{\mu}_1)^\top \mathbf{v} \right) \left((\mathbf{x}_{|\mathcal{I}_1}^i - \boldsymbol{\mu}_1)^\top \mathbf{v} \right) \quad (2.29)$$

$$\mathbf{v}^\top \mathbf{S}_w \mathbf{v} = \sum_{i=1}^{N(\mathcal{I}_0)} \left((\mathbf{x}_{|\mathcal{I}_0}^i - \boldsymbol{\mu}_0)^\top \mathbf{v} \right)^2 + \sum_{i=1}^{N(\mathcal{I}_1)} \left((\mathbf{x}_{|\mathcal{I}_1}^i - \boldsymbol{\mu}_1)^\top \mathbf{v} \right)^2 \quad (2.30)$$

$$\mathbf{v}^\top \mathbf{S}_w \mathbf{v} \geq 0 \quad (2.31)$$

- (2.26): This comes directly from the definition: (1.9), (1.7), (1.8).
- (2.26) \Rightarrow (2.27): By Distributive Law of Matrix (**Proposition 2**).
- (2.28) \Rightarrow (2.29): Since $(\mathbf{x}_{|\mathcal{I}_0}^i - \boldsymbol{\mu}_0)^\top \mathbf{v}$ is a scalar, its transpose equals to itself.

Combining $\mathbf{v}^\top \mathbf{S}_w \mathbf{v} \geq 0$ (2.31) and \mathbf{S}_w is symmetric (**Claim 1**), the \mathbf{S}_w is positive semi-definite by **Definition 2**.

Claim 3 \mathbf{S}_b is symmetric ($\mathbf{S}_b = \mathbf{S}_b^\top$)

Proof:

$$\mathbf{S}_b^\top = \left((\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \right)^\top \quad (2.32)$$

$$\mathbf{S}_b^\top = (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \quad (2.33)$$

$$\mathbf{S}_b^\top = \mathbf{S}_b \quad (2.34)$$

- (2.32): This comes directly (1.10).
- (2.33) \Rightarrow (2.34): This comes directly (1.10).

which completes the proof.

Claim 4 \mathbf{S}_b is positive semi-definite.

Proof: For any $\mathbf{v} \in \mathbb{R}^p$,

$$\mathbf{v}^\top \mathbf{S}_b \mathbf{v} = \mathbf{v}^\top (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{v} \quad (2.35)$$

$$\mathbf{v}^\top \mathbf{S}_b \mathbf{v} = \left((\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{v} \right)^\top \left((\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{v} \right) \quad (2.36)$$

$$\mathbf{v}^\top \mathbf{S}_b \mathbf{v} = \left((\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{v} \right) \left((\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{v} \right) \quad (2.37)$$

$$\mathbf{v}^\top \mathbf{S}_b \mathbf{v} = \left((\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{v} \right)^2 \quad (2.38)$$

$$\mathbf{v}^\top \mathbf{S}_b \mathbf{v} \geq 0 \quad (2.39)$$

- (2.35): This comes directly (1.10).
- (2.35) \Rightarrow (2.36): Rearrangement of terms.
- (2.36) \Rightarrow (2.37): Since $(\mathbf{x}_{|\mathcal{I}_0}^i - \boldsymbol{\mu}_0)^\top \mathbf{v}$ is a scalar, its transpose equals to itself.

Combining $\mathbf{v}^\top \mathbf{S}_b \mathbf{v} \geq 0$ (2.39) and \mathbf{S}_b is symmetric (**Claim 3**), the \mathbf{S}_b is positive semi-definite by **Definition 2**.

Claim 5 (Cauchy Schwarz Inequality) If S_w is positive semi-definite, then $(u^\top S_w v)^2 \leq (u^\top S_w u) \cdot (v^\top S_w v)$. The equality holds when $u = v$.

Proof: Let $t \in \mathbb{R}$, since S_w is positive semi-definite:

$$(u + tv)^\top S_w (u + tv) \geq 0 \quad (2.40)$$

$$u^\top S_w u + tu^\top S_w v + tv^\top S_w u + t^2 v^\top S_w v \geq 0 \quad (2.41)$$

Since $u^\top S_w v$ is a scalar, therefore its transpose is itself and S_w is symmetry, $u^\top S_w v = (u^\top S_w v)^\top = v^\top S_w u$. Then,

$$u^\top S_w u + tu^\top S_w v + tv^\top S_w u + t^2 v^\top S_w v \geq 0 \quad (2.42)$$

$$u^\top S_w u + tu^\top S_w v + tu^\top S_w v + t^2 v^\top S_w v \geq 0 \quad (2.43)$$

$$u^\top S_w u + 2tu^\top S_w v + t^2 v^\top S_w v \geq 0 \quad (2.44)$$

must hold for all $t \in \mathbb{R}$, as for all $t \in \mathbb{R}$, $(u + tv) \in \mathbb{R}^p$ (by positive semi-definite).

Since $u^\top S_w u + 2tu^\top S_w v + t^2 v^\top S_w v$ is a quadratic function. Therefore, $\Delta = (2u^\top S_w v)^2 - 4 \cdot (u^\top S_w u) \cdot (v^\top S_w v) \leq 0$:

$$(2u^\top S_w v)^2 - 4 \cdot (u^\top S_w u) \cdot (v^\top S_w v) \leq 0 \quad (2.45)$$

$$4(u^\top S_w v)^2 - 4 \cdot (u^\top S_w u) \cdot (v^\top S_w v) \leq 0 \quad (2.46)$$

$$(u^\top S_w v)^2 \leq (u^\top S_w u) \cdot (v^\top S_w v) \quad (2.47)$$

which completes the proof.

Claim 6 Let $\mu_0 \neq \mu_1$ and (P5) has solution, then:

- (P1) has optimal solutions
- Denote $(\hat{w}, \hat{\beta})$ to be the solution of (P5), then \hat{w} is an optimal solution of (P1).
- The optimal value of (P1) is strictly positive.

Proof:

$$\mu_0 - \mu_1 \neq 0 \quad (2.48)$$

$$\hat{\beta} \cdot S_w \hat{w} \neq 0 \quad (2.49)$$

$$S_w \hat{w} \neq 0 \quad (2.50)$$

$$\hat{w}^\top S_w \hat{w} > 0 \quad (2.51)$$

- (2.48): This comes directly from the assumption $\mu_0 \neq \mu_1$.
- (2.48) \Rightarrow (2.49): By assumption $(\hat{w}, \hat{\beta})$ is a solution of (P5), then $\hat{\beta} \cdot S_w \hat{w} = \mu_0 - \mu_1$.
- (2.49) \Rightarrow (2.50): If $S_w \hat{w} = 0$, then $\hat{\beta} \cdot S_w \hat{w} = 0$, which is a contradiction. Therefore, $S_w \hat{w} \neq 0$.
- (2.50) \Rightarrow (2.51): Since S_w is positive semi-definite by **Claim 2**, and $S_w \hat{w} \neq 0$, by **Proposition 4** we know that $\hat{w}^\top S_w \hat{w} > 0$.

(2.51) shows that \hat{w} is a feasible solution of (P1).

Denote any feasible solution of (P1) to be w then:

$$w^\top S_b w = w^\top (\mu_0 - \mu_1) (\mu_0 - \mu_1)^\top w \quad (2.52)$$

$$\mathbf{w}^\top \mathbf{S}_b \mathbf{w} = \mathbf{w}^\top (\hat{\beta} \mathbf{S}_w \hat{\mathbf{w}}) (\hat{\beta} \mathbf{S}_w \hat{\mathbf{w}})^\top \mathbf{w} \quad (2.53)$$

$$\mathbf{w}^\top \mathbf{S}_b \mathbf{w} = \hat{\beta}^2 (\mathbf{w}^\top \mathbf{S}_w \hat{\mathbf{w}}) \cdot (\mathbf{w}^\top \mathbf{S}_w \hat{\mathbf{w}})^\top \quad (2.54)$$

$$\mathbf{w}^\top \mathbf{S}_b \mathbf{w} = \hat{\beta}^2 (\mathbf{w}^\top \mathbf{S}_w \hat{\mathbf{w}}) \cdot (\mathbf{w}^\top \mathbf{S}_w \hat{\mathbf{w}}) \quad (2.55)$$

$$\mathbf{w}^\top \mathbf{S}_b \mathbf{w} = \hat{\beta}^2 (\mathbf{w}^\top \mathbf{S}_w \hat{\mathbf{w}})^2 \quad (2.56)$$

- (2.52): This comes directly from (1.10).
- (2.52) \Rightarrow (2.53): By assumption $(\hat{\mathbf{w}}, \hat{\beta})$ is the the feasible solution of (P5), then $\hat{\beta} \cdot (\mathbf{S}_w \hat{\mathbf{w}}) = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1$.
- (2.54) \Rightarrow (2.55): $\mathbf{w}^\top \mathbf{S}_w \hat{\mathbf{w}}$ is a scalar, and therefore its transpose $(\mathbf{w}^\top \mathbf{S}_w \hat{\mathbf{w}})^\top$ is still a scalar and equals to itself.

Then:

$$(\mathbf{w}^\top \mathbf{S}_w \hat{\mathbf{w}})^2 \leq (\mathbf{w}^\top \mathbf{S}_w \mathbf{w})(\hat{\mathbf{w}}^\top \mathbf{S}_w \hat{\mathbf{w}}) \quad (2.57)$$

$$\hat{\beta}^2 (\mathbf{w}^\top \mathbf{S}_w \hat{\mathbf{w}})^2 \leq \hat{\beta}^2 (\mathbf{w}^\top \mathbf{S}_w \mathbf{w})(\hat{\mathbf{w}}^\top \mathbf{S}_w \hat{\mathbf{w}}) \quad (2.58)$$

$$\mathbf{w}^\top \mathbf{S}_b \mathbf{w} \leq \hat{\beta}^2 (\mathbf{w}^\top \mathbf{S}_w \mathbf{w})(\hat{\mathbf{w}}^\top \mathbf{S}_w \hat{\mathbf{w}}) \quad (2.59)$$

$$\frac{\mathbf{w}^\top \mathbf{S}_b \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_w \mathbf{w}} \leq \hat{\beta}^2 (\hat{\mathbf{w}}^\top \mathbf{S}_w \hat{\mathbf{w}}) \quad (2.60)$$

$$J_1(\mathbf{w}) \leq \hat{\beta}^2 (\hat{\mathbf{w}}^\top \mathbf{S}_w \hat{\mathbf{w}}) \quad (2.61)$$

- (2.57): By Cauchy Schwarz Inequality (**Claim 5**).
- (2.57) \Rightarrow (2.58): Multiply both sides with $\hat{\beta}^2$.
- (2.58) \Rightarrow (2.59): $\hat{\beta}^2 (\mathbf{w}^\top \mathbf{S}_w \hat{\mathbf{w}})^2 = \mathbf{w}^\top \mathbf{S}_b \mathbf{w}$ comes from (2.56).
- (2.59) \Rightarrow (2.60): Since \mathbf{w} is a feasible solution of (P1), then by (1.12), $\mathbf{w}^\top \mathbf{S}_w \mathbf{w} > 0$. Then divide both sides by $\mathbf{w}^\top \mathbf{S}_w \mathbf{w}$.
- (2.60) \Rightarrow (2.61): This comes directly from (1.11).

Similarly:

$$\hat{\mathbf{w}}^\top \mathbf{S}_b \hat{\mathbf{w}} = \hat{\beta}^2 (\hat{\mathbf{w}}^\top \mathbf{S}_w \hat{\mathbf{w}})(\hat{\mathbf{w}}^\top \mathbf{S}_w \hat{\mathbf{w}}) \quad (2.62)$$

$$\frac{\hat{\mathbf{w}}^\top \mathbf{S}_b \hat{\mathbf{w}}}{\hat{\mathbf{w}}^\top \mathbf{S}_w \hat{\mathbf{w}}} = \hat{\beta}^2 (\hat{\mathbf{w}}^\top \mathbf{S}_w \hat{\mathbf{w}}) \quad (2.63)$$

$$J_1(\hat{\mathbf{w}}) = \hat{\beta}^2 (\hat{\mathbf{w}}^\top \mathbf{S}_w \hat{\mathbf{w}}) \quad (2.64)$$

- (2.57): (2.51) shows that $\hat{\mathbf{w}}$ is a feasible solution of (P1), then we can apply (2.56), by substitute \mathbf{w} with $\hat{\mathbf{w}}$.
- (2.62) \Rightarrow (2.63): Since $\hat{\mathbf{w}}$ is a feasible solution of (P1), then by (1.12), $\hat{\mathbf{w}}^\top \mathbf{S}_w \hat{\mathbf{w}} > 0$. Then divide both sides by $\hat{\mathbf{w}}^\top \mathbf{S}_w \hat{\mathbf{w}}$.
- (2.63) \Rightarrow (2.64): This comes directly from (1.11).

Since (2.51) shows that $\hat{\mathbf{w}}$ is a feasible solution of (P1), (2.61) shows that $J_1(\mathbf{w})$ has a upper bound and (2.64) shows that the upper bound is achieved at solutions $\hat{\mathbf{w}}$, which is exactly from the solution of (P5).

Therefore, we proved that if $(\hat{\mathbf{w}}, \hat{\beta})$ is a solution of (P5), then $\hat{\mathbf{w}}$ is an optimal solution of (P1).

Since $(\hat{\mathbf{w}}, \hat{\beta})$ is a solution of (P5), it must satisfy (2.20): $\hat{\beta} \mathbf{S}_w \hat{\mathbf{w}} = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1$. Since $\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1 \neq \mathbf{0}$, then $\hat{\beta} \neq 0$.

The optimal value $\hat{\beta}^2 (\hat{\mathbf{w}}^\top \mathbf{S}_w \hat{\mathbf{w}}) > 0$ because $\hat{\beta} \neq 0$ and we have already shown that $\hat{\mathbf{w}}^\top \mathbf{S}_w \hat{\mathbf{w}} > 0$, which completes the proof.

Claim 7 If \mathbf{w}^* is an optimal solution of (P1), then there $\exists \lambda^*, v^* \in \mathbb{R}$, such that $\left(\frac{\mathbf{w}^*}{\sqrt{\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^*}}, \lambda^*, v^* \right)$ is an optimal solution of (P3).

Proof:

For all feasible solutions \mathbf{w} in (P2), since it satisfies $\mathbf{w}^\top \mathbf{S}_w \mathbf{w} = 1$ (2.2), it must satisfy $\mathbf{w}^\top \mathbf{S}_w \mathbf{w} > 0$ (1.12), which means that \mathbf{w} is also the feasible solution of (P1).

$$J_1(\mathbf{w}^*) \geq J_1(\mathbf{w}) \quad (2.65)$$

$$\frac{\mathbf{w}^{*\top} \mathbf{S}_b \mathbf{w}^*}{\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^*} \geq \frac{\mathbf{w}^\top \mathbf{S}_b \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_w \mathbf{w}} \quad (2.66)$$

$$\left(\frac{\mathbf{w}^*}{\sqrt{\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^*}} \right)^\top \mathbf{S}_b \left(\frac{\mathbf{w}^*}{\sqrt{\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^*}} \right) \geq \frac{\mathbf{w}^\top \mathbf{S}_b \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_w \mathbf{w}} \quad (2.67)$$

$$\left(\frac{\mathbf{w}^*}{\sqrt{\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^*}} \right)^\top \mathbf{S}_b \left(\frac{\mathbf{w}^*}{\sqrt{\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^*}} \right) \geq \frac{\mathbf{w}^\top \mathbf{S}_b \mathbf{w}}{1} \quad (2.68)$$

$$\left(\frac{\mathbf{w}^*}{\sqrt{\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^*}} \right)^\top \mathbf{S}_b \left(\frac{\mathbf{w}^*}{\sqrt{\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^*}} \right) \geq \mathbf{w}^\top \mathbf{S}_b \mathbf{w} \quad (2.69)$$

$$-\left(\frac{\mathbf{w}^*}{\sqrt{\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^*}} \right)^\top \mathbf{S}_b \left(\frac{\mathbf{w}^*}{\sqrt{\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^*}} \right) \leq -\mathbf{w}^\top \mathbf{S}_b \mathbf{w} \quad (2.70)$$

$$J_2 \left(\frac{\mathbf{w}^*}{\sqrt{\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^*}} \right) \leq J_2(\mathbf{w}) \quad (2.71)$$

- (2.65): Since \mathbf{w}^* is the optimal solution of (P1) and \mathbf{w} is the feasible solution of (P1), $J_1(\mathbf{w}^*) \geq J_1(\mathbf{w})$
- (2.65) \Rightarrow (2.66): This comes directly from the definition of $J_1(\mathbf{w})$ (1.11).
- (2.67) \Rightarrow (2.68): Since \mathbf{w} is defined to be a feasible solution of (P2), it must satisfy $\mathbf{w}^\top \mathbf{S}_w \mathbf{w} = 1$ (2.2)
- (2.70) \Rightarrow (2.71): This comes directly from the definition of $J_2(\mathbf{w}^*)$ (2.1).

Also:

$$\frac{\mathbf{w}^{*\top}}{\sqrt{\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^*}} \cdot \mathbf{S}_w \cdot \frac{\mathbf{w}^*}{\sqrt{\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^*}} = 1 \quad (2.72)$$

(2.72) implies that $\frac{\mathbf{w}^*}{\sqrt{\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^*}}$ satisfies the (P2) constraint (2.2), which means that $\frac{\mathbf{w}^*}{\sqrt{\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^*}}$ is a feasible solution of

(P2). Also, we have shown that for all \mathbf{w} which are feasible solutions of (P2), there is $J_2 \left(\frac{\mathbf{w}^*}{\sqrt{\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^*}} \right) \leq J_2(\mathbf{w})$ (2.71),

which means that $\frac{\mathbf{w}^*}{\sqrt{\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^*}}$ has minimality among all feasible solutions. Combine the feasibility and minimality,

$\frac{\mathbf{w}^*}{\sqrt{\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^*}}$ is an optimal solution of (P2).

By Generalized Fritz John condition (**Proposition 1**), since $\frac{\mathbf{w}^*}{\sqrt{\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^*}}$ is an optimal solution of (P2), we conclude

that there exists $\lambda^*, v^* \in \mathbb{R}$ such that $\left(\frac{\mathbf{w}^*}{\sqrt{\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^*}}, \lambda^*, v^* \right)$ is an optimal solution of (P3).

This completes the proof.

Claim 8 Let $\mu_0 \neq \mu_1$ and (P1) has optimal solutions: If w^* is an optimal solution of (P1), then there $\exists \beta^* \in \mathbb{R}$, such that (w^*, β^*) is a solution of (P5).

Proof:

Since w^* is an optimal solution of (P1), by **Claim 7** we conclude that there exists $\lambda^*, v^* \in \mathbb{R}$ such that $\left(\frac{w^*}{\sqrt{w^{*\top} S_w w^*}}, \lambda^*, v^*\right)$ is an optimal solution of (P3). Therefore, it should be a feasible solution of (P3) as well, which means that it should satisfy the constraints (2.9), (2.12) ((2.10), (2.11) are ignored because they are useless here):

$$\lambda^* S_b \cdot \frac{w^*}{\sqrt{w^{*\top} S_w w^*}} = v^* S_w \cdot \frac{w^*}{\sqrt{w^{*\top} S_w w^*}} \quad (2.73)$$

$$[\lambda^*, v^*]^\top \neq 0 \quad (2.74)$$

We also have:

$$J_1(w^*) > 0 \quad (2.75)$$

$$\frac{w^{*\top} S_b w^*}{w^{*\top} S_w w^*} > 0 \quad (2.76)$$

$$w^{*\top} S_b w^* > 0 \quad (2.77)$$

$$w^{*\top} (\mu_0 - \mu_1) (\mu_0 - \mu_1)^\top w^* > 0 \quad (2.78)$$

$$((\mu_0 - \mu_1)^\top w^*)^\top ((\mu_0 - \mu_1)^\top w^*) > 0 \quad (2.79)$$

$$((\mu_0 - \mu_1)^\top w^*) ((\mu_0 - \mu_1)^\top w^*) > 0 \quad (2.80)$$

$$((\mu_0 - \mu_1)^\top w^*)^2 > 0 \quad (2.81)$$

$$(\mu_0 - \mu_1)^\top w^* \neq 0 \quad (2.82)$$

- (2.75): In **Claim 6** we have shown that the optimal value of (P1) is strictly positive. Since w^* is an optimal solution of (P1), $J_1(w^*) > 0$.
- (2.75) \Rightarrow (2.76): This comes directly from the definition of $J_1(w)$ (1.11).
- (2.76) \Rightarrow (2.77): Since w^* is the optimal solution of (P1), it must be feasible solution of (P1) and therefore satisfy $w^{*\top} S_w w^* > 0$ (1.12).
- (2.77) \Rightarrow (2.78): This comes directly from $S_b = (\mu_0 - \mu_1) (\mu_0 - \mu_1)^\top$ (1.10).
- (2.78) \Rightarrow (2.79): Rearrangement of terms.
- (2.79) \Rightarrow (2.80): Since $(\mu_0 - \mu_1)^\top w^*$ is a scalar, its transpose equals to itself.
- (2.81) \Rightarrow (2.82): If the square of a scalar is strictly greater than 0, then this scalar cannot be 0.

We claim that $\lambda^* \neq 0$ and prove it by contradiction:

If $\lambda^* = 0$, then:

$$v^* S_w \cdot \frac{w^*}{\sqrt{w^{*\top} S_w w^*}} = \lambda^* S_b \cdot \frac{w^*}{\sqrt{w^{*\top} S_w w^*}} \quad (2.83)$$

$$v^* S_w w^* = \lambda^* S_b w^* \quad (2.84)$$

$$v^* S_w w^* = 0 \cdot S_b w^* \quad (2.85)$$

$$v^* S_w w^* = 0 \quad (2.86)$$

$$v^* = 0 \quad (2.87)$$

- (2.83): This comes directly from (2.73).

- (2.83) \Rightarrow (2.84): Multiply both sides with $\sqrt{\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^*}$.
- (2.84) \Rightarrow (2.85): $\lambda^* = 0$ as the contradiction assumption.
- (2.86) \Rightarrow (2.87): Since \mathbf{w}^* is an optimal solution of (P1) then it must be feasible satisfying (1.12) $\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^* > 0$. By **Proposition 4** we know that $\mathbf{S}_w \mathbf{w}^* \neq \mathbf{0}$. Therefore, $v^* = 0$ must hold so that $v^* \mathbf{S}_w \mathbf{w}^* = \mathbf{0}$.

However, now $[\lambda^*, v^*]^\top = [0, 0]^\top = \mathbf{0}$ contradicts (2.74). Therefore:

$$\lambda^* \neq 0 \quad (2.88)$$

Then:

$$\lambda^* \mathbf{S}_b \frac{\mathbf{w}^*}{\sqrt{\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^*}} = v^* \mathbf{S}_w \frac{\mathbf{w}^*}{\sqrt{\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^*}} \quad (2.89)$$

$$\lambda^* \mathbf{S}_b \mathbf{w}^* = v^* \mathbf{S}_w \mathbf{w}^* \quad (2.90)$$

$$\lambda^* (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}^* = v^* \mathbf{S}_w \mathbf{w}^* \quad (2.91)$$

$$\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1 = \frac{v^*}{\lambda^* (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}^*} \mathbf{S}_w \mathbf{w}^* \quad (2.92)$$

- (2.89): This comes directly from (2.73).
- (2.89) \Rightarrow (2.90): Multiply both sides with $\sqrt{\mathbf{w}^{*\top} \mathbf{S}_w \mathbf{w}^*}$.
- (2.90) \Rightarrow (2.91): This comes directly from (1.10).
- (2.91) \Rightarrow (2.92): Since $\lambda^* \neq 0$ (2.88) and $(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}^* \neq 0$ (2.82), we can divide both sides by $\lambda^* (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}^*$

By observation we see that $(\mathbf{w}, \beta) = \left(\mathbf{w}^*, \frac{v^*}{\lambda^* (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}^*} \right)$ is a solution of (P5).

This shows that if \mathbf{w}^* is the optimal solution of (P1), then there $\exists \beta^* \in \mathbb{R}$, such that (\mathbf{w}^*, β^*) is a solution of (P5).

Claim 9 Let $\boldsymbol{\mu}_0 \neq \boldsymbol{\mu}_1$ and (P5) has solutions, then:

- (P1) has optimal solutions.
- Denote \mathbf{w}^* to be the optimal solution of (P1), then there $\exists \beta^* \in \mathbb{R}$, such that (\mathbf{w}^*, β^*) is a solution of (P5).

Proof:

In **Claim 6** we have shown that if (P5) has solutions, then (P1) must have optimal solutions.

Since (P1) has optimal solutions, we can define any optimal solution to be \mathbf{w}^* .

Since \mathbf{w}^* is an optimal solution of (P1), by **Claim 8** we conclude that there exists $\beta^* \in \mathbb{R}$ such that (\mathbf{w}^*, β^*) is an optimal solution of (P5), which completes the proof.

Claim 10 If (P5) does not have a solution, then (P1) has infinite maximum.

Proof: Since (P5) does not have solution, then a special case when $\beta = 1$ of (P5) does not have solution as well: $\mathbf{S}_w \mathbf{w} = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1$ does not have a solution, this means that $\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1 \notin \text{Col}(\mathbf{S}_w) = \text{Row}(\mathbf{S}_w)$, as \mathbf{S}_w is symmetric.

By **Theorem 1** we have $\text{Row}(\mathbf{S}_w)^\perp = \text{Ker}(\mathbf{S}_w)$, then by **Theorem 2** we have $\text{Row}(\mathbf{S}_w) = (\text{Row}(\mathbf{S}_w)^\perp)^\perp = \text{Ker}(\mathbf{S}_w)^\perp$. Therefore, $\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1 \notin \text{Col}(\mathbf{S}_w) = \text{Row}(\mathbf{S}_w) = \text{Ker}(\mathbf{S}_w)^\perp$.

Recall the definition of \perp , $\text{Ker}(\mathbf{S}_w)^\perp = \{\mathbf{y} \in \mathbb{R}^p \mid \forall \mathbf{v} \in \text{Ker}(\mathbf{S}_w), \mathbf{y}^\top \mathbf{v} = 0\}$. Since $\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1 \notin \text{Ker}(\mathbf{S}_w)^\perp$, by definition, there $\exists \mathbf{w}_0 \in \text{Ker}(\mathbf{S}_w)$ such that:

$$(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}_0 \neq 0 \quad (2.93)$$

Since $\mathbf{w}_0 \in \text{Ker}(\mathbf{S}_w)$, by definition:

$$\mathbf{S}_w \mathbf{w}_0 = \mathbf{0} \quad (2.94)$$

Since $\text{rank}(\mathbf{S}_w) \geq 1$ (as \mathbf{S}_w is nonzero required by algorithm). Then there must $\exists \mathbf{v}_0$ such that $\mathbf{S}_w \mathbf{v}_0 \neq \mathbf{0}$. Since \mathbf{S}_w is positive semi-definite, by **Proposition 4**, we have $\mathbf{v}_0^\top \mathbf{S}_w \mathbf{v}_0 > 0$.

Then we can define $\mathbf{w}_\epsilon := \mathbf{w}_0 + \epsilon \mathbf{v}_0$. Then:

$$\lim_{\epsilon \rightarrow 0} \mathbf{w}_\epsilon^\top \mathbf{S}_b \mathbf{w}_\epsilon = \lim_{\epsilon \rightarrow 0} ((\mathbf{w}_0 + \epsilon \mathbf{v}_0)^\top \mathbf{S}_b (\mathbf{w}_0 + \epsilon \mathbf{v}_0)) \quad (2.95)$$

$$\lim_{\epsilon \rightarrow 0} \mathbf{w}_\epsilon^\top \mathbf{S}_b \mathbf{w}_\epsilon = \lim_{\epsilon \rightarrow 0} (\mathbf{w}_0^\top \mathbf{S}_b \mathbf{w}_0 + \epsilon \mathbf{w}_0^\top \mathbf{S}_b \mathbf{v}_0 + \epsilon \mathbf{v}_0^\top \mathbf{S}_b \mathbf{w}_0 + \epsilon^2 \mathbf{v}_0^\top \mathbf{S}_b \mathbf{v}_0) \quad (2.96)$$

$$\lim_{\epsilon \rightarrow 0} \mathbf{w}_\epsilon^\top \mathbf{S}_b \mathbf{w}_\epsilon = \lim_{\epsilon \rightarrow 0} \mathbf{w}_0^\top \mathbf{S}_b \mathbf{w}_0 + \lim_{\epsilon \rightarrow 0} \epsilon \mathbf{w}_0^\top \mathbf{S}_b \mathbf{v}_0 + \lim_{\epsilon \rightarrow 0} \epsilon \mathbf{v}_0^\top \mathbf{S}_b \mathbf{w}_0 + \lim_{\epsilon \rightarrow 0} \epsilon^2 \mathbf{v}_0^\top \mathbf{S}_b \mathbf{v}_0 \quad (2.97)$$

$$\lim_{\epsilon \rightarrow 0} \mathbf{w}_\epsilon^\top \mathbf{S}_b \mathbf{w}_\epsilon = \mathbf{w}_0^\top \mathbf{S}_b \mathbf{w}_0 + 0 + 0 + 0 \quad (2.98)$$

$$\lim_{\epsilon \rightarrow 0} \mathbf{w}_\epsilon^\top \mathbf{S}_b \mathbf{w}_\epsilon = \mathbf{w}_0^\top (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}_0 \quad (2.99)$$

$$\lim_{\epsilon \rightarrow 0} \mathbf{w}_\epsilon^\top \mathbf{S}_b \mathbf{w}_\epsilon = ((\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}_0)^\top ((\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}_0) \quad (2.100)$$

$$\lim_{\epsilon \rightarrow 0} \mathbf{w}_\epsilon^\top \mathbf{S}_b \mathbf{w}_\epsilon = ((\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}_0) ((\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}_0) \quad (2.101)$$

$$\lim_{\epsilon \rightarrow 0} \mathbf{w}_\epsilon^\top \mathbf{S}_b \mathbf{w}_\epsilon = ((\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}_0)^2 \quad (2.102)$$

$$\lim_{\epsilon \rightarrow 0} \mathbf{w}_\epsilon^\top \mathbf{S}_b \mathbf{w}_\epsilon > 0 \quad (2.103)$$

$$(2.104)$$

- (2.95): This comes directly from the definition $\mathbf{w}_\epsilon := \mathbf{w}_0 + \epsilon \mathbf{v}_0$.
- (2.96) \Rightarrow (2.97): Distributive Law of Limit Addition.
- (2.98) \Rightarrow (2.99): This comes directly from the (1.10).
- (2.99) \Rightarrow (2.100): Rearrangement of terms.
- (2.100) \Rightarrow (2.101): Since $(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}_0$ is a scalar, the transpose of a scalar equals to itself.
- (2.102) \Rightarrow (2.103): Since (2.93), then the square of a non-zero scalar is greater than 0.

There is an important property to show:

$$\mathbf{w}_0^\top \mathbf{S}_w \mathbf{v}_0 = (\mathbf{w}_0^\top \mathbf{S}_w \mathbf{v}_0)^\top \quad (2.105)$$

$$\mathbf{w}_0^\top \mathbf{S}_w \mathbf{v}_0 = \mathbf{v}_0^\top \mathbf{S}_w^\top \mathbf{w}_0 \quad (2.106)$$

$$\mathbf{w}_0^\top \mathbf{S}_w \mathbf{v}_0 = \mathbf{v}_0^\top \mathbf{S}_w \mathbf{w}_0 \quad (2.107)$$

- (2.105): Since $\mathbf{w}_0^\top \mathbf{S}_w \mathbf{v}_0$ is a scalar, the transpose of scalar equals to itself.
- (2.106) \Rightarrow (2.107): This comes from **Claim 1**.

$$\lim_{\epsilon \rightarrow 0} \mathbf{w}_\epsilon^\top \mathbf{S}_w \mathbf{w}_\epsilon = \lim_{\epsilon \rightarrow 0} ((\mathbf{w}_0 + \epsilon \mathbf{v}_0)^\top \mathbf{S}_w (\mathbf{w}_0 + \epsilon \mathbf{v}_0)) \quad (2.108)$$

$$\lim_{\epsilon \rightarrow 0} \mathbf{w}_\epsilon^\top \mathbf{S}_w \mathbf{w}_\epsilon = \lim_{\epsilon \rightarrow 0} (\mathbf{w}_0^\top \mathbf{S}_w \mathbf{w}_0 + \epsilon \mathbf{w}_0^\top \mathbf{S}_w \mathbf{v}_0 + \epsilon \mathbf{v}_0^\top \mathbf{S}_w \mathbf{w}_0 + \epsilon^2 \mathbf{v}_0^\top \mathbf{S}_w \mathbf{v}_0) \quad (2.109)$$

$$\lim_{\epsilon \rightarrow 0} \mathbf{w}_\epsilon^\top \mathbf{S}_w \mathbf{w}_\epsilon = \lim_{\epsilon \rightarrow 0} (\mathbf{w}_0^\top \mathbf{S}_w \mathbf{w}_0 + \epsilon \mathbf{v}_0^\top \mathbf{S}_w \mathbf{w}_0 + \epsilon \mathbf{v}_0^\top \mathbf{S}_w \mathbf{w}_0 + \epsilon^2 \mathbf{v}_0^\top \mathbf{S}_w \mathbf{v}_0) \quad (2.110)$$

$$\lim_{\epsilon \rightarrow 0} \mathbf{w}_\epsilon^\top \mathbf{S}_w \mathbf{w}_\epsilon = \lim_{\epsilon \rightarrow 0} (\mathbf{w}_0^\top \mathbf{0} + \epsilon \mathbf{v}_0^\top \mathbf{0} + \epsilon \mathbf{v}_0^\top \mathbf{0} + \epsilon^2 \mathbf{v}_0^\top \mathbf{S}_w \mathbf{v}_0) \quad (2.111)$$

$$\lim_{\epsilon \rightarrow 0} \mathbf{w}_\epsilon^\top \mathbf{S}_w \mathbf{w}_\epsilon = \lim_{\epsilon \rightarrow 0} \epsilon^2 \mathbf{v}_0^\top \mathbf{S}_w \mathbf{v}_0 \quad (2.112)$$

$$\lim_{\epsilon \rightarrow 0} \mathbf{w}_\epsilon^\top \mathbf{S}_w \mathbf{w}_\epsilon = 0 \quad (2.113)$$

- (2.108): This comes directly from the definition $\mathbf{w}_\epsilon := \mathbf{w}_0 + \epsilon \mathbf{v}_0$.
- (2.109) \Rightarrow (2.110): This comes from $\mathbf{w}_0^\top \mathbf{S}_w \mathbf{v}_0 = \mathbf{v}_0^\top \mathbf{S}_w \mathbf{w}_0$ (2.107).
- (2.110) \Rightarrow (2.111): This comes from (2.94)

Then:

$$\lim_{\epsilon \rightarrow 0} \frac{\mathbf{w}_\epsilon^\top \mathbf{S}_b \mathbf{w}_\epsilon}{\mathbf{w}_\epsilon^\top \mathbf{S}_w \mathbf{w}_\epsilon} = +\infty \quad (2.114)$$

$$\lim_{\epsilon \rightarrow 0} J_1(\mathbf{w}_\epsilon) = +\infty \quad (2.115)$$

- (2.114): Since \mathbf{S}_w and \mathbf{S}_b are positive semi-definite, $\mathbf{w}_\epsilon^\top \mathbf{S}_b \mathbf{w}_\epsilon \geq 0$ and $\mathbf{w}_\epsilon^\top \mathbf{S}_w \mathbf{w}_\epsilon \geq 0$ for $\forall \epsilon \in \mathbb{R}$, $\lim_{\epsilon \rightarrow 0} \mathbf{w}_\epsilon^\top \mathbf{S}_b \mathbf{w}_\epsilon > 0$, and $\lim_{\epsilon \rightarrow 0} \mathbf{w}_\epsilon^\top \mathbf{S}_w \mathbf{w}_\epsilon = 0$.
- (2.114) \Rightarrow (2.115): This comes directly from the definition of $J_1(\mathbf{w})$ (1.11).

which completes the proof.

Claim 11 *There exists $\gamma \in \mathbb{R}$, such that $\mathbf{S}_b = \gamma \mathbf{S}_w \Rightarrow$ (P1) has trivial solution status.*

Proof:

Since $\mathbf{w}^\top \mathbf{S}_w \mathbf{w} > 0$ by constraint of (P1) we know that the denominator of $J_1(\mathbf{w})$ should be non-zero and therefore well-defined.

$$J_1(\mathbf{w}) = \frac{\mathbf{w}^\top \mathbf{S}_b \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_w \mathbf{w}} \quad (2.116)$$

$$J_1(\mathbf{w}) = \frac{\mathbf{w}^\top (\gamma \mathbf{S}_w) \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_w \mathbf{w}} \quad (2.117)$$

$$J_1(\mathbf{w}) = \gamma \cdot \frac{\mathbf{w}^\top \mathbf{S}_w \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_w \mathbf{w}} \quad (2.118)$$

$$J_1(\mathbf{w}) = \gamma \quad (2.119)$$

- (2.116): This comes directly from the definition of J_1 (1.11).
- (2.116) \Rightarrow (2.117): This comes from the assumption that $\mathbf{S}_b = \gamma \mathbf{S}_w$.

which shows that $J_1(\mathbf{w})$ for all feasible \mathbf{w} has the same value γ .

Since all feasible solutions have the same value, all feasible solutions are optimal solutions. Therefore, the feasible solution set is the same as the optimal solution set, which means that (P1) has trivial solution status.

Claim 12 *(P1) has trivial solution status \Rightarrow There exists $\gamma \in \mathbb{R}$, such that $\mathbf{S}_b = \gamma \mathbf{S}_w$.*

Proof:

Case $\mu_0 = \mu_1$:

If $\mu_0 = \mu_1$, then $\mathbf{S}_b = (\mu_0 - \mu_1)(\mu_0 - \mu_1)^\top = \mathbf{0}$.

This means that there exists $\gamma = 0$ such that $\mathbf{S}_b = \gamma \mathbf{S}_w$

Case $\mu_0 \neq \mu_1$:

If $\mu_0 \neq \mu_1$: The (P1) has trivial solution status implies that $J_1(\mathbf{w})$ is the same for all feasible solutions (In other words, \mathbf{w} satisfying $\mathbf{w}^\top \mathbf{S}_w \mathbf{w} > 0$). Then this means that (P1) has optimal solutions and all \mathbf{w} satisfying $\mathbf{w}^\top \mathbf{S}_w \mathbf{w} > 0$ are optimal solutions. Let us denote the same objective value to be γ_0 .

Since $\mu_0 \neq \mu_1$ all \mathbf{w} satisfying $\mathbf{w}^\top \mathbf{S}_w \mathbf{w} > 0$ are optimal solutions of (P1), by **Claim 9** we conclude that for each \mathbf{w} , there exists a corresponding $\beta \in \mathbb{R}$ such that (\mathbf{w}, β) is a solution of (P5) (2.20), which means that:

$$\beta \cdot (\mathbf{S}_w \mathbf{w}) = \mu_0 - \mu_1 \quad (2.120)$$

Then:

$$\beta \mathbf{S}_w \mathbf{w} = \mu_0 - \mu_1 \quad (2.121)$$

$$\beta \mathbf{S}_w \mathbf{w} \neq \mathbf{0} \quad (2.122)$$

$$\beta \neq 0 \quad (2.123)$$

- (2.121): This comes directly from (2.120).
- (2.121) \Rightarrow (2.122): Now we are in the case that $\mu_0 \neq \mu_1$.
- (2.122) \Rightarrow (2.123): $\beta \mathbf{S}_w \mathbf{w} \neq \mathbf{0}$ implies that $\beta \neq 0$ and $\mathbf{S}_w \mathbf{w} \neq \mathbf{0}$.

Now, we show by contradiction that $\text{rank}(\mathbf{S}_w) \geq 2$ does not hold:

Suppose that $\text{rank}(\mathbf{S}_w) \geq 2$, then there exists two non-zero linearly independent columns \mathbf{s}_i and \mathbf{s}_j of \mathbf{S}_w :

$$\mathbf{s}_i = \mathbf{S}_w \mathbf{e}_i \quad (2.124)$$

$$\mathbf{s}_j = \mathbf{S}_w \mathbf{e}_j \quad (2.125)$$

where \mathbf{e}_i is a vector that the i -th element is 1 and the others are 0s, and \mathbf{e}_j is a vector that the j -th element is 1 and the others are 0s.

Continue:

$$\mathbf{S}_w \mathbf{e}_i = \mathbf{s}_i \quad (2.126)$$

$$\mathbf{S}_w \mathbf{e}_i \neq \mathbf{0} \quad (2.127)$$

$$\mathbf{e}_i^\top \mathbf{S}_w \mathbf{e}_i > 0 \quad (2.128)$$

- (2.126): This comes directly from (2.124).
- (2.126) \Rightarrow (2.127): This is because \mathbf{s}_i is non-zero.
- (2.127) \Rightarrow (2.128): $\mathbf{S}_w \mathbf{e}_i \neq \mathbf{0}$, then $\mathbf{e}_i^\top \mathbf{S}_w \mathbf{e}_i > 0$ by **Proposition 4**.

$$\mathbf{S}_w \mathbf{e}_j = \mathbf{s}_j \quad (2.129)$$

$$\mathbf{S}_w \mathbf{e}_j \neq \mathbf{0} \quad (2.130)$$

$$\mathbf{e}_j^\top \mathbf{S}_w \mathbf{e}_j > 0 \quad (2.131)$$

- (2.129): This comes directly from (2.125).
- (2.129) \Rightarrow (2.130): This is because \mathbf{s}_j is non-zero.
- (2.130) \Rightarrow (2.131): $\mathbf{S}_w \mathbf{e}_j \neq \mathbf{0}$, then $\mathbf{e}_j^\top \mathbf{S}_w \mathbf{e}_j > 0$ by **Proposition 4**.

Since we have already concluded that for each optimal solution \mathbf{w} of (P1) there exists $\beta \in \mathbb{R}$ such that (2.120) holds, and \mathbf{e}_i and \mathbf{e}_j are optimal solutions, the following two equations hold:

$$\beta_i \mathbf{S}_w \mathbf{e}_i = \mu_0 - \mu_1 \quad (2.132)$$

$$\beta_j \mathbf{S}_w \mathbf{e}_j = \mu_0 - \mu_1 \quad (2.133)$$

where (\mathbf{e}_i, β_i) and (\mathbf{e}_j, β_j) are the solution of (P5) and $\beta_i, \beta_j \neq 0$.

Then:

$$\beta_i \mathbf{S}_w \mathbf{e}_i = \beta_j \mathbf{S}_w \mathbf{e}_j \quad (2.134)$$

$$\mathbf{S}_w \mathbf{e}_i = \frac{\beta_j}{\beta_i} \mathbf{S}_w \mathbf{e}_j \quad (2.135)$$

$$\mathbf{s}_i = \frac{\beta_j}{\beta_i} \mathbf{s}_j \quad (2.136)$$

- (2.134): Combine (2.132) and (2.133).
- (2.134) \Rightarrow (2.135): In (2.123) we have proved that $\beta \neq 0$, and therefore β_i and β_j are special case of it, and therefore we can divide β_i and β_j on both sides.
- (2.135) \Rightarrow (2.136): This comes directly from (2.124) and (2.125).

In (2.136) we proved that $\mathbf{s}_i = \frac{\beta_j}{\beta_i} \mathbf{s}_j$, which means that \mathbf{s}_i and \mathbf{s}_j are collinear. However, \mathbf{s}_i and \mathbf{s}_j are linearly independent (because \mathbf{s}_i and \mathbf{s}_j are two non-zero linearly independent columns of \mathbf{S}_w), which results in a contradiction. Therefore, $\text{rank}(\mathbf{S}_w) \geq 2$ does not hold.

Therefore, $\text{rank}(\mathbf{S}_w) = 1$: coupled with the fact that \mathbf{S}_w is positive semi-definite (**Claim 2**) then \mathbf{S}_w can be represented as $\mathbf{S}_w = \mathbf{u}\mathbf{u}^\top$ (**Proposition 3**).

Then:

$$\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1 = \beta \mathbf{S}_w \mathbf{w} \quad (2.137)$$

$$\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1 = \beta \mathbf{u}\mathbf{u}^\top \mathbf{w} \quad (2.138)$$

$$\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1 = \beta \mathbf{u}^\top \mathbf{w} \mathbf{u} \quad (2.139)$$

$$(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top = (\beta \mathbf{u}^\top \mathbf{w})^2 \mathbf{u}\mathbf{u}^\top \quad (2.140)$$

$$\mathbf{S}_b = (\beta \mathbf{u}^\top \mathbf{w})^2 \mathbf{S}_w \quad (2.141)$$

- (2.137): This comes directly from (2.120).
- (2.137) \Rightarrow (2.138): Since the \mathbf{S}_w has been represented as: $\mathbf{S}_w = \mathbf{u}\mathbf{u}^\top$.
- (2.140) \Rightarrow (2.141): This comes from (1.10) and $\mathbf{S}_w = \mathbf{u}\mathbf{u}^\top$.

The $(\beta \mathbf{u}^\top \mathbf{w})^2$ is a scalar and therefore we prove that there $\exists \gamma = (\beta \mathbf{u}^\top \mathbf{w})^2$ such that $\mathbf{S}_b = \gamma \mathbf{S}_w$.

There are some interesting and special cases which satisfying that: there exists $\gamma \in \mathbb{R}$, such that $\mathbf{S}_b = \gamma \mathbf{S}_w$.

An interesting case is $\boldsymbol{\mu}_0 = \boldsymbol{\mu}_1$, which means that the centers of two classes are the same. Since $\mathbf{S}_b = (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top = \mathbf{0} \cdot \mathbf{0}^\top = \mathbf{0}$, and the corresponding $\gamma = 0$, $\mathbf{S}_b = \gamma \mathbf{S}_w$. Therefore, the (P1) is trivial when $\boldsymbol{\mu}_0 = \boldsymbol{\mu}_1$.

Another interesting case is that $p = 1$, which means that there is only one feature. In this case, \mathbf{S}_b and $\mathbf{S}_w \in \mathbb{R}^{1 \times 1}$. Then, there exists $\gamma = \frac{B_{1,1}}{S_{1,1}}$ such that $\mathbf{S}_b = [B_{1,1}] = [\gamma S_{1,1}] = \gamma [S_{1,1}] = \gamma \mathbf{S}_w$. Therefore, the (P1) is trivial when $p = 1$.

Claim 13 A non-zero column of \mathbf{S}_w satisfies $\mathbf{w}^\top \mathbf{S}_w \mathbf{w} > 0$.

Proof:

Let us denote the i -th column of \mathbf{S}_w to be $\mathbf{s}_i \in \mathbb{R}^p$. Since \mathbf{S}_w is symmetric **Claim 1**, then the i -th row of \mathbf{S}_w is $\mathbf{s}_i^\top \in \mathbb{R}^{1 \times p}$. By assumption, suppose that the k -th column of \mathbf{S}_w is non-zero, which implies that: $\mathbf{s}_k \neq \mathbf{0}$.

Therefore:

$$\mathbf{S}_w \mathbf{s}_k = \begin{bmatrix} \mathbf{s}_1^\top \\ \mathbf{s}_2^\top \\ \vdots \\ \mathbf{s}_k^\top \\ \vdots \\ \mathbf{s}_p^\top \end{bmatrix} \mathbf{s}_k = \begin{bmatrix} \star \\ \star \\ \vdots \\ \|\mathbf{s}_k\|_2^2 \\ \vdots \\ \star \end{bmatrix} \quad (2.142)$$

Since $\mathbf{s}_k \neq \mathbf{0}$, $\|\mathbf{s}_k\|_2^2 > 0$, which shows that $\mathbf{S}_w \mathbf{s}_k \neq \mathbf{0}$.

Since $\mathbf{S}_w \mathbf{s}_k \neq \mathbf{0}$ and \mathbf{S}_w is positive semi-definite (**Claim 2**), by **Proposition 4** we know that $\mathbf{w}^\top \mathbf{S}_w \mathbf{w} > \mathbf{0}$.

Claim 14 If $\exists \gamma \in \mathbb{R}$, such that $\mathbf{S}_b = \gamma \mathbf{S}_w$, then:

- (P1) has trivial optimal solution status.
- A non-zero column of \mathbf{S}_w is an optimal solution of (P1).

Proof:

By **Claim 11** we know that • (P1) has trivial optimal solution status holds.

By **Claim 13** a non-zero column of \mathbf{S}_w satisfies $\mathbf{w}^\top \mathbf{S}_w \mathbf{w} > \mathbf{0}$. Since $\mathbf{w}^\top \mathbf{S}_w \mathbf{w} > \mathbf{0}$ is the only constraint of (P1), then \mathbf{w} is a feasible solution of (P1).

Since (P1) has trivial optimal solution status, by definition all feasible solutions of (P1) are optimal solutions of (P1).

Therefore, \mathbf{w} is an optimal solution of (P1).

Claim 15 If $\forall \gamma \in \mathbb{R}$, $\mathbf{S}_b \neq \gamma \mathbf{S}_w$, and (P6) has a solution, then:

- (P1) has non-trivial optimal solution status.
- Denote $\hat{\mathbf{w}}$ to be the solution of (P6), then $\hat{\mathbf{w}}$ is an optimal solution of (P1).

Proof:

The contrapositive of **Claim 12** tells us that if $\forall \gamma \in \mathbb{R}$, $\mathbf{S}_b \neq \gamma \mathbf{S}_w$, then (P1) has non-trivial optimal solution status.

Since $\hat{\mathbf{w}}$ is a solution of (P6), let $\hat{\beta} = 1$ and it is obvious that $(\hat{\mathbf{w}}, \hat{\beta})$ is a solution of (P5). By **Claim 6** we know that $\hat{\mathbf{w}}$ is an optimal solution of (P1).

Claim 16 If $\forall \gamma \in \mathbb{R}$, $\mathbf{S}_b \neq \gamma \mathbf{S}_w$, and (P6) has no solution, then:

- (P1) has infinite optimal solution status.
- The optimal solution set of (P1) is empty.

Proof:

Suppose that (P5) has a solution $(\hat{\mathbf{w}}, \hat{\beta})$, then it satisfies: $\hat{\beta} \mathbf{S}_w \hat{\mathbf{w}} = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1$. Rearrange the terms we have: $\mathbf{S}_w(\hat{\beta} \hat{\mathbf{w}}) = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1$, which implies that $\hat{\beta} \hat{\mathbf{w}}$ is a solution of (P6).

In a nutshell, if (P5) has a solution, then (P6) has a solution. By contrapositive we know that if (P6) has no solution, then (P5) has no solution.

Since the assumption has said that (P5) has no solution, **Claim 10** tells us that (P1) has infinite optimal solution status, and it is natural that the optimal solution set is empty.

Chapter 3

Model-based Prediction

3.1 Simplified Model

Suppose that the minimizer of $J_1(\mathbf{w})$ is \mathbf{w}^* , then the model is formulated as:

$$M(\mathbf{x}) = \begin{cases} 0 & \text{if } |\mathbf{w}^{*\top} \mathbf{x} - \mathbf{w}^{*\top} \boldsymbol{\mu}_0| \leq |\mathbf{w}^{*\top} \mathbf{x} - \mathbf{w}^{*\top} \boldsymbol{\mu}_1| \\ 1 & \text{if } |\mathbf{w}^{*\top} \mathbf{x} - \mathbf{w}^{*\top} \boldsymbol{\mu}_0| > |\mathbf{w}^{*\top} \mathbf{x} - \mathbf{w}^{*\top} \boldsymbol{\mu}_1| \end{cases} \quad (3.1)$$

which classifies the samples according to the distance to each projected cluster center.

We can simplify the model.

Since $|\mathbf{w}^{*\top} \mathbf{x} - \mathbf{w}^{*\top} \boldsymbol{\mu}_0|$ measures the distance between $\mathbf{w}^{*\top} \mathbf{x}$ and $\mathbf{w}^{*\top} \boldsymbol{\mu}_0$, and $|\mathbf{w}^{*\top} \mathbf{x} - \mathbf{w}^{*\top} \boldsymbol{\mu}_1|$ measures the distance between $\mathbf{w}^{*\top} \mathbf{x}$ and $\mathbf{w}^{*\top} \boldsymbol{\mu}_1$. Then $|\mathbf{w}^{*\top} \mathbf{x} - \mathbf{w}^{*\top} \boldsymbol{\mu}_0| \leq |\mathbf{w}^{*\top} \mathbf{x} - \mathbf{w}^{*\top} \boldsymbol{\mu}_1|$ implies that $\mathbf{w}^{*\top} \mathbf{x}$ is closer to $\mathbf{w}^{*\top} \boldsymbol{\mu}_0$. Therefore, it is obvious that the decision boundary is the middle point $\frac{1}{2} \mathbf{w}^{*\top} (\boldsymbol{\mu}_0 + \boldsymbol{\mu}_1)$.

Therefore, the model is simplified as:

$$M(\mathbf{x}) = \begin{cases} 0 & \text{if } (2 \cdot \mathbf{w}^{*\top} \mathbf{x} - \mathbf{w}^{*\top} (\boldsymbol{\mu}_0 + \boldsymbol{\mu}_1)) \cdot \mathbf{w}^{*\top} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) \leq 0 \\ 1 & \text{if } (2 \cdot \mathbf{w}^{*\top} \mathbf{x} - \mathbf{w}^{*\top} (\boldsymbol{\mu}_0 + \boldsymbol{\mu}_1)) \cdot \mathbf{w}^{*\top} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) > 0 \end{cases} \quad (3.2)$$

3.2 Proof of Simplified Model

Claim 17 *Let*

$$M(\mathbf{x}) = \begin{cases} 0 & \text{if } |\mathbf{w}^{*\top} \mathbf{x} - \mathbf{w}^{*\top} \boldsymbol{\mu}_0| \leq |\mathbf{w}^{*\top} \mathbf{x} - \mathbf{w}^{*\top} \boldsymbol{\mu}_1| \\ 1 & \text{if } |\mathbf{w}^{*\top} \mathbf{x} - \mathbf{w}^{*\top} \boldsymbol{\mu}_0| > |\mathbf{w}^{*\top} \mathbf{x} - \mathbf{w}^{*\top} \boldsymbol{\mu}_1| \end{cases} \quad (3.1)$$

Then it can be rewritten as:

$$M(\mathbf{x}) = \begin{cases} 0 & \text{if } (2 \cdot \mathbf{w}^{*\top} \mathbf{x} - \mathbf{w}^{*\top} (\boldsymbol{\mu}_0 + \boldsymbol{\mu}_1)) \cdot \mathbf{w}^{*\top} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) \leq 0 \\ 1 & \text{if } (2 \cdot \mathbf{w}^{*\top} \mathbf{x} - \mathbf{w}^{*\top} (\boldsymbol{\mu}_0 + \boldsymbol{\mu}_1)) \cdot \mathbf{w}^{*\top} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) > 0 \end{cases} \quad (3.2)$$

Proof:

Let

$$z := \mathbf{w}^{*\top} \mathbf{x}, \quad (3.3)$$

$$m_0 := \mathbf{w}^{*\top} \boldsymbol{\mu}_0, \quad (3.4)$$

$$m_1 := \mathbf{w}^{*\top} \boldsymbol{\mu}_1. \quad (3.5)$$

We begin with the original definition of the classifier:

$$M(\mathbf{x}) = \begin{cases} 0 & \text{if } |z - m_0| \leq |z - m_1| \\ 1 & \text{if } |z - m_0| > |z - m_1| \end{cases} \quad (3.6)$$

To proceed, we first square both sides:

$$|z - m_0| \leq |z - m_1| \quad (3.7)$$

$$(z - m_0)^2 \leq (z - m_1)^2 \quad (3.8)$$

$$z^2 - 2zm_0 + m_0^2 \leq z^2 - 2zm_1 + m_1^2 \quad (3.9)$$

$$-2z(m_0 - m_1) + (m_0^2 - m_1^2) \leq 0 \quad (3.10)$$

$$-2z(m_0 - m_1) + (m_0 - m_1)(m_0 + m_1) \leq 0 \quad (3.11)$$

$$(2z - (m_0 + m_1))(m_1 - m_0) \leq 0 \quad (3.12)$$

Substituting back the original terms:

$$(2 \cdot \mathbf{w}^{*\top} \mathbf{x} - \mathbf{w}^{*\top} (\boldsymbol{\mu}_0 + \boldsymbol{\mu}_1)) \cdot \mathbf{w}^{*\top} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) \leq 0 \quad (3.13)$$

This completes the proof.

Chapter 4

Example

4.1 Example 1

Question

Given a training dataset \mathcal{I} :

- $\mathbf{x}|_{\diamond 1}^{\mathcal{I}} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \mathbf{x}|_{\diamond 2}^{\mathcal{I}} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \mathbf{x}|_{\diamond 3}^{\mathcal{I}} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}, \mathbf{x}|_{\diamond 4}^{\mathcal{I}} = \begin{bmatrix} 3 \\ 7 \end{bmatrix}$
- $\mathbf{y}|_{\diamond 1}^{\mathcal{I}} = 0, \mathbf{y}|_{\diamond 2}^{\mathcal{I}} = 0, \mathbf{y}|_{\diamond 3}^{\mathcal{I}} = 1, \mathbf{y}|_{\diamond 4}^{\mathcal{I}} = 1$

Model Fitting

Step 1: Split Dataset \mathcal{I} into \mathcal{I}_0 and \mathcal{I}_1

$$\mathcal{I}_0 = \{\diamond 1, \diamond 2\}, \quad \mathcal{I}_1 = \{\diamond 3, \diamond 4\}$$

Step 2: Compute $\boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \mathbf{S}_w$

$$\begin{aligned}\boldsymbol{\mu}_0 &= \frac{1}{2} (\mathbf{x}|_{\diamond 1}^{\mathcal{I}_0} + \mathbf{x}|_{\diamond 2}^{\mathcal{I}_0}) = \frac{1}{2} \cdot \begin{bmatrix} 1 \\ 3 \end{bmatrix} = \begin{bmatrix} 0.5 \\ 1.5 \end{bmatrix} \\ \boldsymbol{\mu}_1 &= \frac{1}{2} (\mathbf{x}|_{\diamond 3}^{\mathcal{I}_1} + \mathbf{x}|_{\diamond 4}^{\mathcal{I}_1}) = \frac{1}{2} \cdot \begin{bmatrix} 6 \\ 12 \end{bmatrix} = \begin{bmatrix} 3 \\ 6 \end{bmatrix} \\ \boldsymbol{\Sigma}_0 &= \sum_{i=1}^2 (\mathbf{x}|_{\diamond i}^{\mathcal{I}_0} - \boldsymbol{\mu}_0) (\mathbf{x}|_{\diamond i}^{\mathcal{I}_0} - \boldsymbol{\mu}_0)^\top \\ &= \begin{bmatrix} 0.25 & 0.25 \\ 0.25 & 0.25 \end{bmatrix} + \begin{bmatrix} 0.25 & 0.25 \\ 0.25 & 0.25 \end{bmatrix} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix} \\ \boldsymbol{\Sigma}_1 &= \sum_{i=1}^2 (\mathbf{x}|_{\diamond i}^{\mathcal{I}_1} - \boldsymbol{\mu}_1) (\mathbf{x}|_{\diamond i}^{\mathcal{I}_1} - \boldsymbol{\mu}_1)^\top \\ &= \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0 & 2 \end{bmatrix} \\ \mathbf{S}_w &= \boldsymbol{\Sigma}_0 + \boldsymbol{\Sigma}_1 = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 2.5 \end{bmatrix}\end{aligned}$$

Step 3: Compute \mathbf{S}_b , check if there exists $\gamma \in \mathbb{R}$ such that $\mathbf{S}_b = \gamma \mathbf{S}_w$

$$\mathbf{S}_b = (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top = \begin{bmatrix} -2.5 \\ -4.5 \end{bmatrix} \begin{bmatrix} -2.5 & -4.5 \end{bmatrix} = \begin{bmatrix} 6.25 & 11.25 \\ 11.25 & 20.25 \end{bmatrix}$$

Does not exist.

Step 4: Solve $S_w w = \mu_0 - \mu_1$

$$\begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 2.5 \end{bmatrix} \cdot w^* = \begin{bmatrix} -2.5 \\ -4.5 \end{bmatrix}$$

$$w^* = \begin{bmatrix} -4 \\ -1 \end{bmatrix}$$

(P1) has non-trivial optimal solution status, and w^* is an optimal solution of (P1).

Prediction

We consider the inequality:

$$(2 \cdot w^{*\top} x - w^{*\top}(\mu_0 + \mu_1)) \cdot w^{*\top}(\mu_1 - \mu_0) \leq 0$$

Given:

$$w^* = \begin{bmatrix} -4 \\ -1 \end{bmatrix}, \quad \mu_0 = \begin{bmatrix} 0.5 \\ 1.5 \end{bmatrix}, \quad \mu_1 = \begin{bmatrix} 3 \\ 6 \end{bmatrix}$$

Step 1:

Compute $w^{*\top} x$

$$w^{*\top} x = [-4, -1] \cdot \begin{bmatrix} [x]_1 \\ [x]_2 \end{bmatrix} = -4[x]_1 - [x]_2$$

Step 2:

Compute $w^{*\top}(\mu_0 + \mu_1)$

$$\mu_0 + \mu_1 = \begin{bmatrix} 0.5 + 3 \\ 1.5 + 6 \end{bmatrix} = \begin{bmatrix} 3.5 \\ 7.5 \end{bmatrix}$$

$$w^{*\top}(\mu_0 + \mu_1) = [-4, -1] \cdot \begin{bmatrix} 3.5 \\ 7.5 \end{bmatrix} = -4 \cdot 3.5 - 1 \cdot 7.5 = -14 - 7.5 = -21.5$$

Step 3:

Compute $w^{*\top}(\mu_1 - \mu_0)$

$$\mu_1 - \mu_0 = \begin{bmatrix} 3 - 0.5 \\ 6 - 1.5 \end{bmatrix} = \begin{bmatrix} 2.5 \\ 4.5 \end{bmatrix}$$

$$w^{*\top}(\mu_1 - \mu_0) = [-4, -1] \cdot \begin{bmatrix} 2.5 \\ 4.5 \end{bmatrix} = -4 \cdot 2.5 - 1 \cdot 4.5 = -10 - 4.5 = -14.5$$

Step 4:

Substitute all values into the inequality

$$(2(-4[x]_1 - [x]_2) - (-21.5)) \cdot (-14.5) \leq 0$$

Simplify:

$$16[x]_1 + 4[x]_2 - 43 \leq 0$$

Final result:

$$M(x) = \begin{cases} 0 & \text{if } 16[x]_1 + 4[x]_2 - 43 \leq 0 \\ 1 & \text{if } 16[x]_1 + 4[x]_2 - 43 > 0 \end{cases}$$

4.2 Example 2

Question

Given a training dataset \mathcal{I} :

- $\mathbf{x}|_{\diamond 1}^{\mathcal{I}} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \mathbf{x}|_{\diamond 2}^{\mathcal{I}} = \begin{bmatrix} 2 \\ 4 \end{bmatrix}, \mathbf{x}|_{\diamond 3}^{\mathcal{I}} = \begin{bmatrix} 3 \\ 6 \end{bmatrix}, \mathbf{x}|_{\diamond 4}^{\mathcal{I}} = \begin{bmatrix} 4 \\ 8 \end{bmatrix}$
- $\mathbf{y}|_{\diamond 1}^{\mathcal{I}} = 0, \mathbf{y}|_{\diamond 2}^{\mathcal{I}} = 0, \mathbf{y}|_{\diamond 3}^{\mathcal{I}} = 1, \mathbf{y}|_{\diamond 4}^{\mathcal{I}} = 1$

Model Fitting

Step 1: Split Dataset \mathcal{I} into \mathcal{I}_0 and \mathcal{I}_1

$$\mathcal{I}_0 = \mathcal{I}_{\diamond 1}, \quad \mathcal{I}_0 = \mathcal{I}_{\diamond 2}, \quad \mathcal{I}_1 = \mathcal{I}_{\diamond 3}, \quad \mathcal{I}_1 = \mathcal{I}_{\diamond 4}$$

Step 2: Compute $\boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \mathbf{S}_w$

$$\begin{aligned} \boldsymbol{\mu}_0 &= \frac{1}{2} (\mathbf{x}|_{\diamond 1}^{\mathcal{I}_0} + \mathbf{x}|_{\diamond 2}^{\mathcal{I}_0}) = \frac{1}{2} \left(\begin{bmatrix} 1 \\ 2 \end{bmatrix} + \begin{bmatrix} 2 \\ 4 \end{bmatrix} \right) = \frac{1}{2} \begin{bmatrix} 3 \\ 6 \end{bmatrix} = \begin{bmatrix} 1.5 \\ 3 \end{bmatrix} \\ \boldsymbol{\mu}_1 &= \frac{1}{2} (\mathbf{x}|_{\diamond 3}^{\mathcal{I}_1} + \mathbf{x}|_{\diamond 4}^{\mathcal{I}_1}) = \frac{1}{2} \left(\begin{bmatrix} 3 \\ 6 \end{bmatrix} + \begin{bmatrix} 4 \\ 8 \end{bmatrix} \right) = \frac{1}{2} \begin{bmatrix} 7 \\ 14 \end{bmatrix} = \begin{bmatrix} 3.5 \\ 7 \end{bmatrix} \\ \boldsymbol{\Sigma}_0 &= \sum_{i=1}^2 (\mathbf{x}|_{\diamond i}^{\mathcal{I}_0} - \boldsymbol{\mu}_0) (\mathbf{x}|_{\diamond i}^{\mathcal{I}_0} - \boldsymbol{\mu}_0)^\top \\ &= \begin{bmatrix} 0.25 & 0.5 \\ 0.5 & 1 \end{bmatrix} + \begin{bmatrix} 0.25 & 0.5 \\ 0.5 & 1 \end{bmatrix} = \begin{bmatrix} 0.5 & 1 \\ 1 & 2 \end{bmatrix} \\ \boldsymbol{\Sigma}_1 &= \sum_{i=1}^2 (\mathbf{x}|_{\diamond i}^{\mathcal{I}_1} - \boldsymbol{\mu}_1) (\mathbf{x}|_{\diamond i}^{\mathcal{I}_1} - \boldsymbol{\mu}_1)^\top \\ &= \begin{bmatrix} 0.25 & 0.5 \\ 0.5 & 1 \end{bmatrix} + \begin{bmatrix} 0.25 & 0.5 \\ 0.5 & 1 \end{bmatrix} = \begin{bmatrix} 0.5 & 1 \\ 1 & 2 \end{bmatrix} \\ \mathbf{S}_w &= \boldsymbol{\Sigma}_0 + \boldsymbol{\Sigma}_1 = \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix} \end{aligned}$$

Step 3: Compute \mathbf{S}_b , check if there exists $\gamma \in \mathbb{R}$ such that $\mathbf{S}_b = \gamma \mathbf{S}_w$

$$\begin{aligned} \mathbf{S}_b &= (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \\ &= \begin{bmatrix} 4 & 8 \\ 8 & 16 \end{bmatrix} \end{aligned}$$

Thus, there exists $\gamma = 4$ such that $\mathbf{S}_b = \gamma \mathbf{S}_w$.

(P1) has trivial optimal solution status, and a non-zero column of \mathbf{S}_w is an optimal solution of (P1).

$$\mathbf{w}^* = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

Prediction

We consider the inequality:

$$(2 \cdot \mathbf{w}^{*\top} \mathbf{x} - \mathbf{w}^{*\top} (\boldsymbol{\mu}_0 + \boldsymbol{\mu}_1)) \cdot \mathbf{w}^{*\top} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) \leq 0$$

Given:

$$\mathbf{w}^* = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \quad \boldsymbol{\mu}_0 = \begin{bmatrix} 1.5 \\ 3.5 \end{bmatrix}, \quad \boldsymbol{\mu}_1 = \begin{bmatrix} 3.5 \\ 7 \end{bmatrix}$$

Step 1:

Compute $\mathbf{w}^{*\top} \mathbf{x}$

$$\mathbf{w}^{*\top} \mathbf{x} = [1, 2] \cdot \begin{bmatrix} [\mathbf{x}]_1 \\ [\mathbf{x}]_2 \end{bmatrix} = 1[\mathbf{x}]_1 + 2[\mathbf{x}]_2$$

Step 2:

Compute $\mathbf{w}^{*\top}(\boldsymbol{\mu}_0 + \boldsymbol{\mu}_1)$

$$\begin{aligned} \boldsymbol{\mu}_0 + \boldsymbol{\mu}_1 &= \begin{bmatrix} 1.5 + 3.5 \\ 3.5 + 7 \end{bmatrix} = \begin{bmatrix} 5 \\ 10.5 \end{bmatrix} \\ \mathbf{w}^{*\top}(\boldsymbol{\mu}_0 + \boldsymbol{\mu}_1) &= [1, 2] \cdot \begin{bmatrix} 5 \\ 10.5 \end{bmatrix} = 1 \times 5 + 2 \times 10.5 = 5 + 21 = 26 \end{aligned}$$

Step 3:

Compute $\mathbf{w}^{*\top}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)$

$$\begin{aligned} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_0 &= \begin{bmatrix} 3.5 - 1.5 \\ 7 - 3.5 \end{bmatrix} = \begin{bmatrix} 2 \\ 3.5 \end{bmatrix} \\ \mathbf{w}^{*\top}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) &= [1, 2] \cdot \begin{bmatrix} 2 \\ 3.5 \end{bmatrix} = 1 \times 2 + 2 \times 3.5 = 2 + 7 = 9 \end{aligned}$$

Step 4:

Substitute all values into the inequality

$$(2([\mathbf{x}]_1 + 2[\mathbf{x}]_2) - 26) \cdot 9 \leq 0$$

Simplify:

$$(2[\mathbf{x}]_1 + 4[\mathbf{x}]_2 - 26) \times 9 \leq 0$$

which is equivalent to

$$[\mathbf{x}]_1 + 2[\mathbf{x}]_2 - 13 \leq 0$$

Final result:

$$M(\mathbf{x}) = \begin{cases} 0 & \text{if } [\mathbf{x}]_1 + 2[\mathbf{x}]_2 - 13 \leq 0 \\ 1 & \text{if } [\mathbf{x}]_1 + 2[\mathbf{x}]_2 - 13 > 0 \end{cases}$$

4.3 Example 3

Question

Given a training dataset \mathcal{I} :

- $\mathbf{x}|_{\diamond 1}^{\mathcal{I}} = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}$, $\mathbf{x}|_{\diamond 2}^{\mathcal{I}} = \begin{bmatrix} 0 \\ -1 \\ 0 \end{bmatrix}$, $\mathbf{x}|_{\diamond 3}^{\mathcal{I}} = \begin{bmatrix} -1 \\ 0 \\ 0 \end{bmatrix}$, $\mathbf{x}|_{\diamond 4}^{\mathcal{I}} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$
- $\mathbf{y}|_{\diamond 1}^{\mathcal{I}} = 0$, $\mathbf{y}|_{\diamond 2}^{\mathcal{I}} = 0$, $\mathbf{y}|_{\diamond 3}^{\mathcal{I}} = 1$, $\mathbf{y}|_{\diamond 4}^{\mathcal{I}} = 1$

Model Fitting

Step 1: Split Dataset \mathcal{I} into \mathcal{I}_0 and \mathcal{I}_1

$$\mathcal{I}_0 = \mathcal{I}_{\diamond 1}, \quad \mathcal{I}_0 = \mathcal{I}_{\diamond 2}, \quad \mathcal{I}_1 = \mathcal{I}_{\diamond 3}, \quad \mathcal{I}_1 = \mathcal{I}_{\diamond 4}$$

Step 2: Compute μ_0, μ_1, S_w

$$\begin{aligned}\mu_0 &= \frac{1}{2} (\mathbf{x}|_{\diamond_1}^{\mathcal{I}_0} + \mathbf{x}|_{\diamond_2}^{\mathcal{I}_0}) = \frac{1}{2} \cdot \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0 \\ 0 \end{bmatrix} \\ \mu_1 &= \frac{1}{2} (\mathbf{x}|_{\diamond_1}^{\mathcal{I}_1} + \mathbf{x}|_{\diamond_2}^{\mathcal{I}_1}) = \frac{1}{2} \cdot \begin{bmatrix} -1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} -0.5 \\ 0 \\ 0 \end{bmatrix} \\ \Sigma_0 &= \sum_{i=1}^2 (\mathbf{x}|_{\diamond_i}^{\mathcal{I}_0} - \mu_0) (\mathbf{x}|_{\diamond_i}^{\mathcal{I}_0} - \mu_0)^\top = \begin{bmatrix} 0.5 & 1 & 0 \\ 1 & 2 & 0 \\ 0 & 0 & 0 \end{bmatrix} \\ \Sigma_1 &= \sum_{i=1}^2 (\mathbf{x}|_{\diamond_i}^{\mathcal{I}_1} - \mu_1) (\mathbf{x}|_{\diamond_i}^{\mathcal{I}_1} - \mu_1)^\top = \begin{bmatrix} 0.5 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \\ S_w &= \Sigma_0 + \Sigma_1 = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 2 & 0 \\ 0 & 0 & 0 \end{bmatrix}\end{aligned}$$

Step 3: Compute S_b , check if there exists $\gamma \in \mathbb{R}$ such that $S_b = \gamma S_w$

$$S_b = (\mu_0 - \mu_1)(\mu_0 - \mu_1)^\top = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

Does not exist.

Step 4: Solve $S_w w = \mu_0 - \mu_1$

We solve:

$$\begin{bmatrix} 1 & 1 & 0 \\ 1 & 2 & 0 \\ 0 & 0 & 0 \end{bmatrix} \mathbf{w}^* = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$\mathbf{w}^* = \begin{bmatrix} 2 \\ -1 \\ 0 \end{bmatrix} + \xi \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}, \quad \xi \in \mathbb{R}.$$

We choose:

$$\mathbf{w}^* = \begin{bmatrix} 2 \\ -1 \\ 0 \end{bmatrix}$$

(P1) has non-trivial optimal solution status, and \mathbf{w}^* is an optimal solution of (P1).

Prediction

We consider the inequality:

$$(2 \cdot \mathbf{w}^{*\top} \mathbf{x} - \mathbf{w}^{*\top} (\mu_0 + \mu_1)) \cdot \mathbf{w}^{*\top} (\mu_1 - \mu_0) \leq 0$$

Given:

$$\mathbf{w}^* = \begin{bmatrix} 2 \\ -1 \\ 0 \end{bmatrix}, \quad \mu_0 = \begin{bmatrix} 0.5 \\ 0 \\ 0 \end{bmatrix}, \quad \mu_1 = \begin{bmatrix} -0.5 \\ 0 \\ 0 \end{bmatrix}$$

Step 1:

Compute $\mathbf{w}^{*\top} \mathbf{x}$

$$\mathbf{w}^{*\top} \mathbf{x} = [2, -1, 0] \cdot \begin{bmatrix} [\mathbf{x}]_1 \\ [\mathbf{x}]_2 \\ [\mathbf{x}]_3 \end{bmatrix} = 2[\mathbf{x}]_1 - [\mathbf{x}]_2$$

Step 2:

Compute $\mathbf{w}^{*\top}(\boldsymbol{\mu}_0 + \boldsymbol{\mu}_1)$

$$\boldsymbol{\mu}_0 + \boldsymbol{\mu}_1 = \begin{bmatrix} 0.5 + (-0.5) \\ 0 + 0 \\ 0 + 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\mathbf{w}^{*\top}(\boldsymbol{\mu}_0 + \boldsymbol{\mu}_1) = [2, -1, 0] \cdot \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} = 0$$

Step 3:

Compute $\mathbf{w}^{*\top}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)$

$$\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0 = \begin{bmatrix} -0.5 - 0.5 \\ 0 - 0 \\ 0 - 0 \end{bmatrix} = \begin{bmatrix} -1 \\ 0 \\ 0 \end{bmatrix}$$

$$\mathbf{w}^{*\top}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) = [2, -1, 0] \cdot \begin{bmatrix} -1 \\ 0 \\ 0 \end{bmatrix} = -2$$

Step 4:

Substitute all values into the inequality

$$(2(2[\mathbf{x}]_1 - [\mathbf{x}]_2) - 0) \times (-2) \leq 0$$

Simplify inside:

$$2[\mathbf{x}]_1 - [\mathbf{x}]_2 \geq 0$$

Final result:

$$M(\mathbf{x}) = \begin{cases} 0 & \text{if } 2[\mathbf{x}]_1 - [\mathbf{x}]_2 \geq 0 \\ 1 & \text{if } 2[\mathbf{x}]_1 - [\mathbf{x}]_2 < 0 \end{cases}$$

4.4 Example 4

Question

Given a training dataset \mathcal{I} :

- $\mathbf{x}|_{\diamond_1}^{\mathcal{I}} = \begin{bmatrix} 8 \\ 6 \end{bmatrix}$, $\mathbf{x}|_{\diamond_2}^{\mathcal{I}} = \begin{bmatrix} 6 \\ 2 \end{bmatrix}$, $\mathbf{x}|_{\diamond_3}^{\mathcal{I}} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$, $\mathbf{x}|_{\diamond_4}^{\mathcal{I}} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$
- $y|_{\diamond_1}^{\mathcal{I}} = 0$, $y|_{\diamond_2}^{\mathcal{I}} = 0$, $y|_{\diamond_3}^{\mathcal{I}} = 1$, $y|_{\diamond_4}^{\mathcal{I}} = 1$

Model Fitting

Step 1: Split Dataset \mathcal{I} into \mathcal{I}_0 and \mathcal{I}_1

$$\mathcal{I}_0 = \mathcal{I}_{\diamond_1}, \quad \mathcal{I}_0 = \mathcal{I}_{\diamond_2}, \quad \mathcal{I}_1 = \mathcal{I}_{\diamond_3}, \quad \mathcal{I}_1 = \mathcal{I}_{\diamond_4}$$

Step 2: Compute μ_0, μ_1, S_w

$$\begin{aligned}\mu_0 &= \frac{1}{2} (\mathbf{x}|_{\diamond_1}^{\mathcal{I}_0} + \mathbf{x}|_{\diamond_2}^{\mathcal{I}_0}) = \frac{1}{2} \left(\begin{bmatrix} 8 \\ 6 \end{bmatrix} + \begin{bmatrix} 6 \\ 2 \end{bmatrix} \right) = \frac{1}{2} \begin{bmatrix} 14 \\ 8 \end{bmatrix} = \begin{bmatrix} 7 \\ 4 \end{bmatrix} \\ \mu_1 &= \frac{1}{2} (\mathbf{x}|_{\diamond_1}^{\mathcal{I}_1} + \mathbf{x}|_{\diamond_2}^{\mathcal{I}_1}) = \frac{1}{2} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix} \right) = \frac{1}{2} \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \\ \Sigma_0 &= \sum_{i=1}^2 (\mathbf{x}|_{\diamond_i}^{\mathcal{I}_0} - \mu_0) (\mathbf{x}|_{\diamond_i}^{\mathcal{I}_0} - \mu_0)^\top = \begin{bmatrix} 2 & 4 \\ 4 & 8 \end{bmatrix} \\ \Sigma_1 &= \sum_{i=1}^2 (\mathbf{x}|_{\diamond_i}^{\mathcal{I}_1} - \mu_1) (\mathbf{x}|_{\diamond_i}^{\mathcal{I}_1} - \mu_1)^\top = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \\ S_w &= \Sigma_0 + \Sigma_1 = \begin{bmatrix} 2 & 4 \\ 4 & 8 \end{bmatrix}\end{aligned}$$

Step 3: Compute S_b , check if there exists $\gamma \in \mathbb{R}$ such that $S_b = \gamma S_w$

$$S_b = (\mu_0 - \mu_1)(\mu_0 - \mu_1)^\top = \begin{bmatrix} 49 & 28 \\ 28 & 16 \end{bmatrix}$$

Does not exist.

Step 4: Solve $S_w w = \mu_0 - \mu_1$

We solve:

$$\begin{bmatrix} 2 & 4 \\ 4 & 8 \end{bmatrix} w^* = \begin{bmatrix} 7 \\ 4 \end{bmatrix}$$

This problem has no solution.

Since (P5) $S_w w = \mu_0 - \mu_1$ has no solution, (P1) has infinite optimal solution status, and the optimal set is empty.

Chapter 5

Some other Important Properties

Solution Set of (P1)

Case $\mu_0 \neq \mu_1$ and (P6) has a solution:

In this case, there exists $\mathbf{v}_0, \mathbf{v}_1, \dots, \mathbf{v}_{p-\text{rank}(\mathbf{S}_w)} \in \mathbb{R}^p$, such that \mathbb{D}^p is the solution set of (P6) $\mathbf{S}_w \mathbf{w} = \mu_0 - \mu_1$, where

$$\mathbb{D}^p = \left\{ \mathbf{v}_0 + \sum_{k=1}^{p-\text{rank}(\mathbf{S}_w)} \xi_k \mathbf{v}_k \mid \xi_1, \dots, \xi_{p-\text{rank}(\mathbf{S}_w)} \in \mathbb{R} \right\} \quad (5.1)$$

The optimal solution set of (P1) is the following equivalent sets:

$$\mathbb{A}^p = \{ \mathbf{w} \in \mathbb{R}^p \mid \exists \beta \in \mathbb{R}, \text{ such that } \beta \mathbf{S}_w \mathbf{w} = \mathbf{S}_b \mathbf{w} \} \quad (5.2)$$

$$\mathbb{B}^p = \left\{ \xi_0 \mathbf{v}_0 + \sum_{k=1}^{p-\text{rank}(\mathbf{S}_w)} \xi_k \mathbf{v}_k \mid \text{where } \xi_0 \in \mathbb{R} \setminus \{0\} \text{ and } \xi_1, \dots, \xi_{p-\text{rank}(\mathbf{S}_w)} \in \mathbb{R} \right\} \quad (5.3)$$

Notice that \mathbb{A}^p is relevant to (P5).

Case $\mu_0 = \mu_1$:

In this case, there exists $\mathbf{v}_1, \dots, \mathbf{v}_{p-\text{rank}(\mathbf{S}_w)} \in \mathbb{R}^p$, such that \mathbb{W}^p is the solution set of (P6) $\mathbf{S}_w \mathbf{w} = \mu_0 - \mu_1$, where

$$\mathbb{W}^p = \left\{ \sum_{k=1}^{p-\text{rank}(\mathbf{S}_w)} \xi_k \mathbf{v}_k \mid \xi_1, \dots, \xi_{p-\text{rank}(\mathbf{S}_w)} \in \mathbb{R} \right\} \quad (5.4)$$

Then there exists $\mathbf{v}_{p-\text{rank}(\mathbf{S}_w)+1}, \dots, \mathbf{v}_p \in \{e_1, e_2, \dots, e_p\} \subseteq \mathbb{R}^p$ such that $\mathbf{v}_1, \dots, \mathbf{v}_p \in \mathbb{R}^p$ are linearly independent. The readers can select $\mathbf{v}_{p-\text{rank}(\mathbf{S}_w)+1}, \dots, \mathbf{v}_p$ from e_1, \dots, e_p via a brute-force combinatorial search as follows:

- Consider all subsets $\mathbb{I} \subset \{1, 2, \dots, p\}$ such that $|\mathbb{I}| = \text{rank}(\mathbf{S}_w)$.
- For each candidate subset $\mathbb{I} = \{i_1, i_2, \dots, i_{\text{rank}(\mathbf{S}_w)}\}$, form the set of vectors $\{e_{i_1}, e_{i_2}, \dots, e_{i_{\text{rank}(\mathbf{S}_w)}}\}$.
- Check whether

$$\{\mathbf{v}_1, \dots, \mathbf{v}_{p-\text{rank}(\mathbf{S}_w)}\} \cup \{e_i : i \in \mathbb{I}\}$$

is a linearly independent set.

- If such a subset \mathbb{I} is found, assign

$$(\mathbf{v}_{p-\text{rank}(\mathbf{S}_w)+1}, \dots, \mathbf{v}_p) = (e_{i_1}, \dots, e_{i_{\text{rank}(\mathbf{S}_w)}}),$$

and terminate the search.

The optimal solution set of (P1) is the following equivalent set:

$$\mathbb{U}^p = \{ \mathbf{w} \in \mathbb{R}^p \mid \mathbf{S}_w \mathbf{w} \neq \mathbf{0} \} \quad (5.5)$$

$$\mathbb{V}^p = \left\{ \sum_{k=1}^p \xi_k \mathbf{v}_k \mid \text{where } \xi_1, \dots, \xi_{p-\text{rank}(\mathbf{S}_w)} \in \mathbb{R} \setminus \{0\} \text{ and } \xi_{p-\text{rank}(\mathbf{S}_w)+1}, \dots, \xi_p \in \mathbb{R} \right\} \quad (5.6)$$

Claim 18 Let $\mu_0 \neq \mu_1$ and (P5) has solutions, then the optimal solution set of (P1) is: $\{\mathbf{w} \in \mathbb{R}^p : \exists \beta \in \mathbb{R}, \text{ such that } \beta \mathbf{S}_w \mathbf{w} = \mathbf{S}_b \mathbf{w}\}$

If (\mathbf{w}, β) is the solution of (P5). By **Claim 6**, we have \mathbf{w} is an optimal solution of (P1).

If for a specific \mathbf{w} , there does not exist $\beta \in \mathbb{R}$ such that (\mathbf{w}, β) is a solution of (P5), then by contrapositive of **Claim 9** we have that \mathbf{w} is not an optimal solution of (P1).

Therefore, the solution set is: $\{\mathbf{w} \in \mathbb{R}^p : \exists \beta \in \mathbb{R}, \text{ such that } \beta \mathbf{S}_w \mathbf{w} = \mathbf{S}_b \mathbf{w}\}$.

We have shown that if $\mu_0 \neq \mu_1$ and (P5) have solutions. Then (P1) must have solutions (**Claim 6**).

Combine **Claim 6** and **Claim 9** we conclude that $\hat{\mathbf{w}}$ is an optimal solution of (P1) if and only if there $\exists \hat{\beta} \in \mathbb{R}$ such that $(\hat{\mathbf{w}}, \hat{\beta})$ is a solution of (P5).

Claim 19 If there exists $\gamma \in \mathbb{R}$ such that $\mathbf{S}_b = \gamma \mathbf{S}_w$, then (P5) has solutions.

Proof:

Case $\mu_0 = \mu_1$:

We can have: $\mathbf{S}_b = \mathbf{0}$. Then there exists $\gamma = 0$ such that $\mathbf{S}_b = \gamma \mathbf{S}_w$.

$\mu_0 - \mu_1 = \mathbf{0}$. Therefore, $(\mathbf{0}, 0)$ can be a solution of (P5).

This shows that (P5) has solutions when $\mu_0 = \mu_1$.

Case $\mu_0 \neq \mu_1$:

$$\gamma \mathbf{S}_w = \mathbf{S}_b \tag{5.7}$$

$$\gamma \mathbf{S}_w = (\mu_0 - \mu_1)(\mu_0 - \mu_1)^\top \tag{5.8}$$

$$\gamma \mathbf{S}_w \neq \mathbf{0} \tag{5.9}$$

$$\gamma \neq 0 \tag{5.10}$$

- (5.7): This is by assumption.
- (5.7) \Rightarrow (5.8): This comes directly from (1.10).
- (5.8) \Rightarrow (5.9): Since $\mu_0 \neq \mu_1$, $(\mu_0 - \mu_1)(\mu_0 - \mu_1)^\top \neq \mathbf{0}$.
- (5.9) \Rightarrow (5.10): Since $\gamma \mathbf{S}_w \neq \mathbf{0}$, $\gamma \neq 0$.

Since \mathbf{S}_w is non-zero then there exists $\mathbf{w}_0 \in \mathbb{R}^p$ such that:

$$\mathbf{S}_w \mathbf{w}_0 \neq \mathbf{0} \tag{5.11}$$

Then:

$$\mathbf{S}_w \mathbf{w}_0 \neq \mathbf{0} \tag{5.12}$$

$$\gamma \mathbf{S}_w \mathbf{w}_0 \neq \mathbf{0} \tag{5.13}$$

$$\mathbf{S}_b \mathbf{w}_0 \neq \mathbf{0} \tag{5.14}$$

$$(\mu_0 - \mu_1)(\mu_0 - \mu_1)^\top \mathbf{w}_0 \neq \mathbf{0} \tag{5.15}$$

$$(\mu_0 - \mu_1)^\top \mathbf{w}_0 \neq 0 \tag{5.16}$$

- (5.12): This comes directly from (5.11).
- (5.12) \Rightarrow (5.13): Multiply both sides with γ and the equality still holds because $\gamma \neq 0$ by (5.10).
- (5.13) \Rightarrow (5.14): This comes from the assumption that $\mathbf{S}_b = \gamma \mathbf{S}_w$.
- (5.14) \Rightarrow (5.15): This comes directly from the definition of \mathbf{S}_b (1.10).
- (5.15) \Rightarrow (5.16): Since $(\mu_0 - \mu_1)(\mu_0 - \mu_1)^\top \mathbf{w}_0 \neq \mathbf{0}$, $(\mu_0 - \mu_1)^\top \mathbf{w}_0 \neq 0$.

Then:

$$\gamma \mathbf{S}_w \mathbf{w}_0 = \mathbf{S}_b \mathbf{w}_0 \quad (5.17)$$

$$\frac{\gamma}{(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}_0} \mathbf{S}_w \mathbf{w}_0 = \frac{1}{(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}_0} \mathbf{S}_b \mathbf{w}_0 \quad (5.18)$$

$$\frac{\gamma}{(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}_0} \mathbf{S}_w \mathbf{w}_0 = \frac{1}{(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}_0} (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}_0 \quad (5.19)$$

$$\frac{\gamma}{(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}_0} \mathbf{S}_w \mathbf{w}_0 = \frac{1}{(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}_0} (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}_0 (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) \quad (5.20)$$

$$\frac{\gamma}{(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}_0} \mathbf{S}_w \mathbf{w}_0 = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1 \quad (5.21)$$

- (5.17): This comes directly from the assumption that $\mathbf{S}_b = \gamma \mathbf{S}_w$
- (5.17) \Rightarrow (5.18): We can divide both sides with $(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}_0$ because $(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}_0 \neq 0$ (5.16).
- (5.18) \Rightarrow (5.19): This comes directly from the definition of \mathbf{S}_b (1.10).
- (5.19) \Rightarrow (5.20): Since $(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}_0$ is a scalar, we can rearrange it.

From (5.21) we can see that $\left(\frac{\gamma}{(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w}_0}, \mathbf{w}_0 \right)$ is a solution of (P5).

Claim 20 Let $\boldsymbol{\mu}_0 \neq \boldsymbol{\mu}_1$ and (P6) has a solution, then:

- There exists $\mathbf{v}_0, \mathbf{v}_1, \dots, \mathbf{v}_{p-\text{rank}(k)} \in \mathbb{R}^p$, such that \mathbb{D}^p is the solution set of (P6) $\mathbf{S}_w \mathbf{w} = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1$

$$\mathbb{D}^p = \left\{ \mathbf{v}_0 + \sum_{k=1}^{p-\text{rank}(\mathbf{S}_w)} \xi_k \mathbf{v}_k \mid \xi_1, \dots, \xi_{p-\text{rank}(\mathbf{S}_w)} \in \mathbb{R} \right\} \quad (5.1)$$

- The following two sets \mathbb{A}^p and \mathbb{B}^p are equivalent:

$$\mathbb{A}^p = \{ \mathbf{w} \in \mathbb{R}^p \mid \exists \beta, \text{ such that } \beta \mathbf{S}_w \mathbf{w} = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1 \} \quad (5.2)$$

$$\mathbb{B}^p = \left\{ \xi_0 \mathbf{v}_0 + \sum_{k=1}^{p-\text{rank}(\mathbf{S}_w)} \xi_k \mathbf{v}_k \mid \text{where } \xi_0 \in \mathbb{R} \setminus \{0\} \text{ and } \xi_1, \dots, \xi_{p-\text{rank}(\mathbf{S}_w)} \in \mathbb{R} \right\} \quad (5.3)$$

Proof:

The proof of the first bullet point is ignored, but you can refer to the linear algebra textbook.

$(\mathbf{w}_0 \in \mathbb{B}^p \Rightarrow \mathbf{w}_0 \in \mathbb{A}^p)$:

Suppose that $\mathbf{w}_0 \in \mathbb{B}^p$. Then there exists $\bar{\xi}_0 \in \mathbb{R} \setminus \{0\}$ and $\bar{\xi}_k \in \mathbb{R}$ for $k = 1, 2, \dots, p - \text{rank}(\mathbf{S}_w)$, such that:

$$\mathbf{w}_0 = \bar{\xi}_0 \mathbf{v}_0 + \sum_{k=1}^{p-\text{rank}(\mathbf{S}_w)} \bar{\xi}_k \mathbf{v}_k \quad (5.22)$$

Since $\xi_0 \neq 0$:

$$\frac{1}{\bar{\xi}_0} \mathbf{w}_0 = \mathbf{v}_0 + \sum_{k=1}^{p-\text{rank}(\mathbf{S}_w)} \frac{\bar{\xi}_k}{\bar{\xi}_0} \mathbf{v}_k \quad (5.23)$$

Since $\bar{\xi}_k \in \mathbb{R}$ and $\bar{\xi}_0 \in \mathbb{R} \setminus \{0\}$, $\frac{\bar{\xi}_k}{\bar{\xi}_0} \in \mathbb{R}$ as well. This implies that $\frac{1}{\bar{\xi}_0} \mathbf{w}_0 \in \mathbb{D}^p$, which is exactly the solution set of $\mathbf{S}_w \mathbf{w} = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1$. Therefore, $\frac{1}{\bar{\xi}_0} \mathbf{w}_0$ is the solution of $\mathbf{S}_w \mathbf{w} = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1$ the following must hold:

$$\mathbf{S}_w \left(\frac{1}{\bar{\xi}_0} \mathbf{w}_0 \right) = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1 \quad (5.24)$$

Then there $\exists \beta = \frac{1}{\bar{\xi}_0}$, such that:

$$\beta \mathbf{S}_w \mathbf{w}_0 = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1 \quad (5.25)$$

which shows that $\mathbf{w}_0 \in \mathbb{A}^p$.

($\mathbf{w}_0 \in \mathbb{A}^p \Rightarrow \mathbf{w}_0 \in \mathbb{B}^p$):

Suppose that $\mathbf{w}_0 \in \mathbb{A}^p$, then there exists β_0 such that:

$$\beta_0 \mathbf{S}_w \mathbf{w}_0 = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1 \quad (5.26)$$

$$\mathbf{S}_w (\beta_0 \mathbf{w}_0) = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1 \quad (5.27)$$

which shows that $\beta_0 \mathbf{w}_0$ is a solution of $\mathbf{S}_w \mathbf{w} = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1$, then it should have the form of (5.1), which means that there exists $\bar{\xi}_k \in \mathbb{R}$ for $k = 1, 2, \dots, p - \text{rank}(\mathbf{S}_w)$.

$$\beta_0 \mathbf{w}_0 = \mathbf{v}_0 + \sum_{k=1}^{p-\text{rank}(\mathbf{S}_w)} \bar{\xi}_k \mathbf{v}_k \quad (5.28)$$

By assumption $\boldsymbol{\mu}_0 \neq \boldsymbol{\mu}_1$, then $\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1 \neq \mathbf{0}$. Combine this with (5.26) we can conclude that $\beta_0 \neq 0$. Then:

$$\beta_0 \mathbf{w}_0 = \mathbf{v}_0 + \sum_{k=1}^{p-\text{rank}(\mathbf{S}_w)} \bar{\xi}_k \mathbf{v}_k \quad (5.29)$$

$$\mathbf{w}_0 = \frac{1}{\beta_0} \mathbf{v}_0 + \sum_{k=1}^{p-\text{rank}(\mathbf{S}_w)} \frac{\bar{\xi}_k}{\beta_0} \mathbf{v}_k \quad (5.30)$$

Since $\frac{1}{\beta_0} \in \mathbb{R} \setminus \{0\}$ and $\frac{\bar{\xi}_k}{\beta_0} \in \mathbb{R}$, we can conclude that $\mathbf{w}_0 \in \mathbb{B}^p$.

Since ($\mathbf{w}_0 \in \mathbb{B}^p \Rightarrow \mathbf{w}_0 \in \mathbb{A}^p$) and ($\mathbf{w}_0 \in \mathbb{A}^p \Rightarrow \mathbf{w}_0 \in \mathbb{B}^p$) we can conclude that \mathbb{A}^p and \mathbb{B}^p are equivalent.

Claim 21 Let $\boldsymbol{\mu}_0 \neq \boldsymbol{\mu}_1$, and $\mathbf{S}_w \mathbf{w} = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1$ has solutions, and the solution set of $\mathbf{S}_w \mathbf{w} = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1$ to be:

$$\mathbb{D}^p = \left\{ \mathbf{w} \in \mathbb{R}^p \left| \mathbf{w} = \mathbf{v}_0 + \sum_{k=1}^{p-\text{rank}(\mathbf{S}_w)} \xi_k \mathbf{v}_k, \text{ where } \xi_1, \dots, \xi_{p-\text{rank}(\mathbf{S}_w)} \in \mathbb{R} \right. \right\} \quad (5.1)$$

where:

$$\mathbf{S}_w \mathbf{v}_0 = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1 \quad (5.31)$$

$$\mathbf{S}_w \mathbf{v}_k = \mathbf{0} \quad (5.32)$$

then the optimal solution set of (P1) is the following two equivalent sets:

- $\mathbb{A} = \{ \mathbf{w} \in \mathbb{R}^p : \exists \beta \in \mathbb{R}, \text{ such that } \beta \mathbf{S}_w \mathbf{w} = \mathbf{S}_b \mathbf{w} \}$
- $\mathbb{B} = \left\{ \mathbf{w} \in \mathbb{R}^p \left| \mathbf{w} = \xi_0 \mathbf{v}_0 + \sum_{k=1}^{p-\text{rank}(\mathbf{S}_w)} \xi_k \mathbf{v}_k, \text{ where } \xi_0 \in \mathbb{R} \setminus \{0\} \text{ and } \xi_1, \dots, \xi_{p-\text{rank}(\mathbf{S}_w)} \in \mathbb{R} \right. \right\}$

Proof:

In the assumption $\mathbf{S}_w \mathbf{w} = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1$ has solutions, which also implies that (P5) $\beta \mathbf{S}_w \mathbf{w} = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1$ has solutions.

Since $\boldsymbol{\mu}_0 \neq \boldsymbol{\mu}_1$ and (P5) has solutions, we can apply **Claim 18** which claims that the solution set of (P1) is $\{\mathbf{w} \in \mathbb{R}^p : \exists \beta \in \mathbb{R}, \text{ such that } \beta \mathbf{S}_w \mathbf{w} = \mathbf{S}_b \mathbf{w}\}$.

Then we can apply **Claim 20** which shows that $\{\mathbf{w} \in \mathbb{R}^p : \exists \beta \in \mathbb{R}, \text{ such that } \beta \mathbf{S}_w \mathbf{w} = \mathbf{S}_b \mathbf{w}\}$ is equivalent to $\left\{ \mathbf{w} \in \mathbb{R}^p \mid \mathbf{w} = \xi_0 \mathbf{v}_0 + \sum_{k=1}^{p-\text{rank}(\mathbf{S}_w)} \xi_k \mathbf{v}_k, \text{ where } \xi_0 \in \mathbb{R} \setminus \{0\} \text{ and } \xi_1, \dots, \xi_{p-\text{rank}(\mathbf{S}_w)} \in \mathbb{R} \right\}$.

This completes the proof.

Chapter 6

Appendix: Calculus and Linear Algebra

Definition 1 Let A be an $n \times n$ square matrix over a field \mathbb{F} . A matrix S_b is called the **inverse** of A , denoted A^{-1} , if it satisfies:

$$AB = BA = I_n$$

where I_n is the $n \times n$ identity matrix.

A matrix possessing an inverse is said to be **invertible** or **non-singular**.

If no such matrix S_b exists for A , then A is called **non-invertible** or **singular**.

Definition 2 A real matrix $A \in \mathbb{R}^{n \times n}$ is called:

- **positive definite**, if $A = A^\top$ (symmetric) and $\mathbf{v}^\top A \mathbf{v} > 0$ for all non-zero $\mathbf{v} \in \mathbb{R}^n$.
- **positive semi-definite**, if $A = A^\top$ (symmetric) and $\mathbf{v}^\top A \mathbf{v} \geq 0$ for all non-zero $\mathbf{v} \in \mathbb{R}^n$;

Theorem 1 (Fundamental Subspace Theorem)

$$\text{Row}(A)^\perp = \text{Ker}(A) \tag{6.1}$$

Proof:

1. $\text{Ker}(A) \subseteq \text{Row}(A)^\perp$
2. $\text{Row}(A)^\perp \subseteq \text{Ker}(A)$

1. Proof of $\text{Ker}(A) \subseteq \text{Row}(A)^\perp$:

Let $\mathbf{x} \in \text{Ker}(A)$, which means $A\mathbf{x} = \mathbf{0}$.

For any $\mathbf{v} \in \text{Row}(A)$, we can write \mathbf{v} as a linear combination of the rows of A :

$$\mathbf{v} = c_1 \mathbf{a}_1 + c_2 \mathbf{a}_2 + \cdots + c_m \mathbf{a}_m$$

where \mathbf{a}_i are the row vectors of A .

Compute the inner product:

$$\mathbf{v} \cdot \mathbf{x} = (c_1 \mathbf{a}_1 + c_2 \mathbf{a}_2 + \cdots + c_m \mathbf{a}_m) \cdot \mathbf{x} = c_1 (\mathbf{a}_1 \cdot \mathbf{x}) + \cdots + c_m (\mathbf{a}_m \cdot \mathbf{x})$$

Since $A\mathbf{x} = \mathbf{0}$, we have $\mathbf{a}_i \cdot \mathbf{x} = 0$ for all i , thus:

$$\mathbf{v} \cdot \mathbf{x} = 0$$

This shows \mathbf{x} is orthogonal to every vector in $\text{Row}(A)$, so:

$$\text{Ker}(A) \subseteq \text{Row}(A)^\perp$$

2. Proof of $\text{Row}(A)^\perp \subseteq \text{Ker}(A)$:

Let $\mathbf{x} \in \text{Row}(A)^\perp$, meaning \mathbf{x} is orthogonal to all vectors in $\text{Row}(A)$. Since $\text{Row}(A)$ contains all row vectors of A , we have:

$$\mathbf{a}_i \cdot \mathbf{x} = 0 \quad (\text{for all rows } \mathbf{a}_i \text{ of } A)$$

This implies:

$$A\mathbf{x} = \mathbf{0}$$

Therefore $\mathbf{x} \in \text{Ker}(A)$, and we conclude:

$$\text{Row}(A)^\perp \subseteq \text{Ker}(A)$$

From both inclusions, we obtain:

$$\text{Row}(A)^\perp = \text{Ker}(A)$$

Theorem 2 (Double Orthogonal Complement Theorem) *For any subspace $V \subseteq \mathbb{R}^n$,*

$$(V^\perp)^\perp = V \tag{6.2}$$

Proof:

$$1. V \subseteq (V^\perp)^\perp$$

$$2. (V^\perp)^\perp \subseteq V$$

1. Proof of $V \subseteq (V^\perp)^\perp$:

Let $\mathbf{v} \in V$. By definition of orthogonal complement, for any $\mathbf{w} \in V^\perp$ we have $\mathbf{v} \cdot \mathbf{w} = 0$.

This means \mathbf{v} is orthogonal to every vector in V^\perp , which by definition implies $\mathbf{v} \in (V^\perp)^\perp$. Therefore:

$$V \subseteq (V^\perp)^\perp$$

2. Proof of $(V^\perp)^\perp \subseteq V$:

Let $\mathbf{x} \in (V^\perp)^\perp$. We can decompose \mathbb{R}^n as:

$$\mathbb{R}^n = V \oplus V^\perp$$

Thus, \mathbf{x} can be uniquely written as:

$$\mathbf{x} = \mathbf{v} + \mathbf{v}^\perp \quad \text{where } \mathbf{v} \in V, \mathbf{v}^\perp \in V^\perp$$

Since $\mathbf{x} \in (V^\perp)^\perp$, it must be orthogonal to all vectors in V^\perp . In particular:

$$\mathbf{x} \cdot \mathbf{v}^\perp = (\mathbf{v} + \mathbf{v}^\perp) \cdot \mathbf{v}^\perp = \mathbf{v} \cdot \mathbf{v}^\perp + \mathbf{v}^\perp \cdot \mathbf{v}^\perp = 0 + \|\mathbf{v}^\perp\|^2 = 0$$

This implies $\|\mathbf{v}^\perp\|^2 = 0$, so $\mathbf{v}^\perp = \mathbf{0}$. Therefore $\mathbf{x} = \mathbf{v} \in V$, proving:

$$(V^\perp)^\perp \subseteq V$$

Combining both inclusions, we conclude:

$$(V^\perp)^\perp = V$$

Theorem 3 (Cholesky Decomposition) *If a matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is positive semi-definite, then there exists a lower triangular matrix $\mathbf{L} \in \mathbb{R}^{n \times n}$ such that $\mathbf{A} = \mathbf{L}\mathbf{L}^\top$.*

Proof: Ignored.

Theorem 4 (Generalized Fritz John Necessary Conditions) *Consider the following mathematical programming problem:*

$$(F) \quad \underset{\mathbf{x}}{\text{Minimize}} \quad f(\mathbf{x}) \quad (6.3)$$

$$\text{Subject to} \quad g_i(\mathbf{x}) \leq 0, \quad i \in \{1, 2, \dots, m\} \quad (6.4)$$

$$h_j(\mathbf{x}) = 0, \quad j \in \{1, 2, \dots, n\} \quad (6.5)$$

where $f(\mathbf{x})$, $g_i(\mathbf{x})$, and $h_j(\mathbf{x})$ are functions defined on the n -dimensional Euclidean space \mathbb{R}^n and have continuous first partial derivatives on \mathbb{R}^n .

If this problem has minimum, and $\bar{\mathbf{x}}$ is an optimal solution, then there exist vectors

$$\bar{\beta} \in \mathbb{R}, \quad \bar{\mathbf{u}} = (\bar{u}_1, \dots, \bar{u}_m)^\top \in \mathbb{R}^m, \quad \bar{\mathbf{v}} = (\bar{v}_1, \bar{v}_2, \dots, \bar{v}_k)^\top \in \mathbb{R}^k$$

such that

$$\bar{\beta} \nabla \theta(\bar{\mathbf{x}}) + \sum_{i=1}^m \bar{u}_i \nabla g_i(\bar{\mathbf{x}}) + \sum_{j=1}^k \bar{v}_j \nabla h_j(\bar{\mathbf{x}}) = 0 \quad (6.6)$$

$$\sum_{i=1}^m \bar{u}_i g_i(\bar{\mathbf{x}}) = 0 \quad (6.7)$$

$$\bar{\beta} \geq 0 \quad (6.8)$$

$$\bar{u}_i \geq 0 \quad \forall i \in \{1, \dots, m\} \quad (6.9)$$

$$\begin{bmatrix} \bar{\beta} \\ \bar{\mathbf{u}} \\ \bar{\mathbf{v}} \end{bmatrix} \neq 0 \quad (6.10)$$

Proof: <https://pages.cs.wisc.edu/~olvi/oldpapers/MFCQ.pdf>

Proposition 1 Consider the following two mathematical programming problem:

$$(P) \quad \underset{\mathbf{x}}{\text{Minimize}} \quad f(\mathbf{x}) \quad (6.11)$$

$$\text{Subject to} \quad g_i(\mathbf{x}) \leq 0, \quad i \in \{1, 2, \dots, m\} \quad (6.12)$$

$$h_j(\mathbf{x}) = 0, \quad j \in \{1, 2, \dots, n\} \quad (6.13)$$

$$(D) \quad \underset{\mathbf{x}, \beta, \mathbf{u}, \mathbf{v}}{\text{Minimize}} \quad f(\mathbf{x}) \quad (6.14)$$

$$\text{Subject to} \quad \beta \cdot \nabla \theta(\mathbf{x}) + \sum_{i=1}^m [\mathbf{u}]_i \cdot \nabla g_i(\mathbf{x}) + \sum_{j=1}^k [\mathbf{v}]_j \cdot \nabla h_j(\mathbf{x}) = 0 \quad (6.15)$$

$$\sum_{i=1}^m [\mathbf{u}]_i \cdot g_i(\mathbf{x}) = 0 \quad (6.16)$$

$$g_i(\mathbf{x}) \leq 0, \quad i \in \{1, 2, \dots, m\} \quad (6.17)$$

$$h_j(\mathbf{x}) = 0, \quad j \in \{1, 2, \dots, n\} \quad (6.18)$$

$$\beta \geq 0 \quad (6.19)$$

$$[\mathbf{u}]_i \geq 0, \quad i \in \{1, 2, \dots, m\} \quad (6.20)$$

$$\begin{bmatrix} \beta \\ \mathbf{u} \\ \mathbf{v} \end{bmatrix} \neq 0 \quad (6.21)$$

where $f(\mathbf{x})$, $g_i(\mathbf{x})$, and $h_j(\mathbf{x})$ are functions defined on the n -dimensional Euclidean space \mathbb{R}^n and have continuous first partial derivatives on \mathbb{R}^n .

Suppose that (P) has optimal solutions:

- \Rightarrow : If \mathbf{x}^* is an optimal solution of (P), then there $\exists \beta^*, \mathbf{u}^*, \mathbf{v}^*$ such that $(\mathbf{x}^*, \beta^*, \mathbf{u}^*, \mathbf{v}^*)$ is an optimal solution of (D).
- \Leftarrow : If $(\mathbf{x}^*, \beta^*, \mathbf{u}^*, \mathbf{v}^*)$ is an optimal solution of (D), then \mathbf{x}^* is an optimal solution of (P).

Proof (\Rightarrow):

If \mathbf{x}^* is an optimal solution of (P), then by **Theorem 4** there exists $\beta^*, \mathbf{u}^*, \mathbf{v}^*$ such that $(\mathbf{x}^*, \beta^*, \mathbf{u}^*, \mathbf{v}^*)$ satisfying (6.6), (6.7), (6.8), (6.9), (6.10).

Since \mathbf{x}^* is an optimal solution of (P), it must be feasible as well, which means that \mathbf{x}^* should satisfy (6.12), (6.13).

Combine these we know that: If \mathbf{x}^* is an optimal solution of (P), then there exists $\beta^*, \mathbf{u}^*, \mathbf{v}^*$ such that $(\mathbf{x}^*, \beta^*, \mathbf{u}^*, \mathbf{v}^*)$ is the solution of (6.6), (6.7), (6.8), (6.9), (6.10), (6.12), (6.13) (which is equivalent to (6.15), (6.16), (6.19), (6.20), (6.21), (6.17), (6.18)). Therefore, we know that there exists $\beta^*, \mathbf{u}^*, \mathbf{v}^*$ such that $(\mathbf{x}^*, \beta^*, \mathbf{u}^*, \mathbf{v}^*)$ is a feasible solution of (D).

Denote all feasible solutions of (D) to be $(\mathbf{x}, \beta, \mathbf{u}, \mathbf{v})$, since the (6.17), and (6.18) hold (which are equivalent to (6.12), (6.13)), \mathbf{x} is a feasible solution of (P) as well.

Now, \mathbf{x} is a feasible solution of (P), \mathbf{x}^* is an optimal solution of (P), which implies that $f(\mathbf{x}^*) \leq f(\mathbf{x})$. This means that $f(\mathbf{x}^*) \leq f(\mathbf{x})$ holds for all feasible solutions in (D).

Therefore, $(\mathbf{x}^*, \beta^*, \mathbf{u}^*, \mathbf{v}^*)$ is an optimal solution of (D) because $f(\mathbf{x}^*) \leq f(\mathbf{x})$ holds for all feasible solutions in (D) and $(\mathbf{x}^*, \beta^*, \mathbf{u}^*, \mathbf{v}^*)$ is a feasible solution of (D).

Proof (\Leftarrow):

By **Theorem 4** and feasibility, we know that if a solution of (P) is optimal, it must satisfy (6.6), (6.7), (6.8), (6.9), (6.10), (6.12), (6.13). Then the optimal solutions are the solutions with minimal objective values and satisfying (6.6), (6.7), (6.8), (6.9), (6.10), (6.12), (6.13).

Therefore, the optimal solution of (D) has minimal objective values among those who exactly satisfying (6.6), (6.7), (6.8), (6.9), (6.10), (6.12), (6.13) since these are equivalent to (6.15) \sim (6.21).

This completes the proof.

To help the readers understand the relationships between these solutions, we draw a Venn Chart in Figure 6.1, where:

- All: All points in \mathbb{R}^n
- Feasible: All feasible solutions of (P).
- Stationary: All feasible solutions of (P) which satisfying (6.6), (6.7), (6.8), (6.9), (6.10).
- Optimal: The optimal solutions of (P).

Proposition 2 Let $\mathbf{A} \in \mathbb{F}^{m \times p}$, $\mathbf{C} \in \mathbb{F}^{q \times r}$, and $\mathbf{B}_1, \dots, \mathbf{B}_n \in \mathbb{F}^{p \times q}$. Then the following identity holds:

$$\mathbf{A} \left(\sum_{i=1}^n \mathbf{B}_i \right) \mathbf{C} = \sum_{i=1}^n \mathbf{A} \mathbf{B}_i \mathbf{C}$$

Proof:

We begin by expanding the left-hand side of the equation:

$$\mathbf{A} \left(\sum_{i=1}^n \mathbf{B}_i \right) \mathbf{C} = \mathbf{A} (\mathbf{B}_1 + \mathbf{B}_2 + \dots + \mathbf{B}_n) \mathbf{C}$$

Using the distributive property of matrix multiplication over addition, we can distribute \mathbf{A} and \mathbf{C} over the sum:

$$= \mathbf{A} \mathbf{B}_1 \mathbf{C} + \mathbf{A} \mathbf{B}_2 \mathbf{C} + \dots + \mathbf{A} \mathbf{B}_n \mathbf{C}$$

Thus, we obtain the right-hand side of the equation:

$$= \sum_{i=1}^n \mathbf{A} \mathbf{B}_i \mathbf{C}$$

This completes the proof.

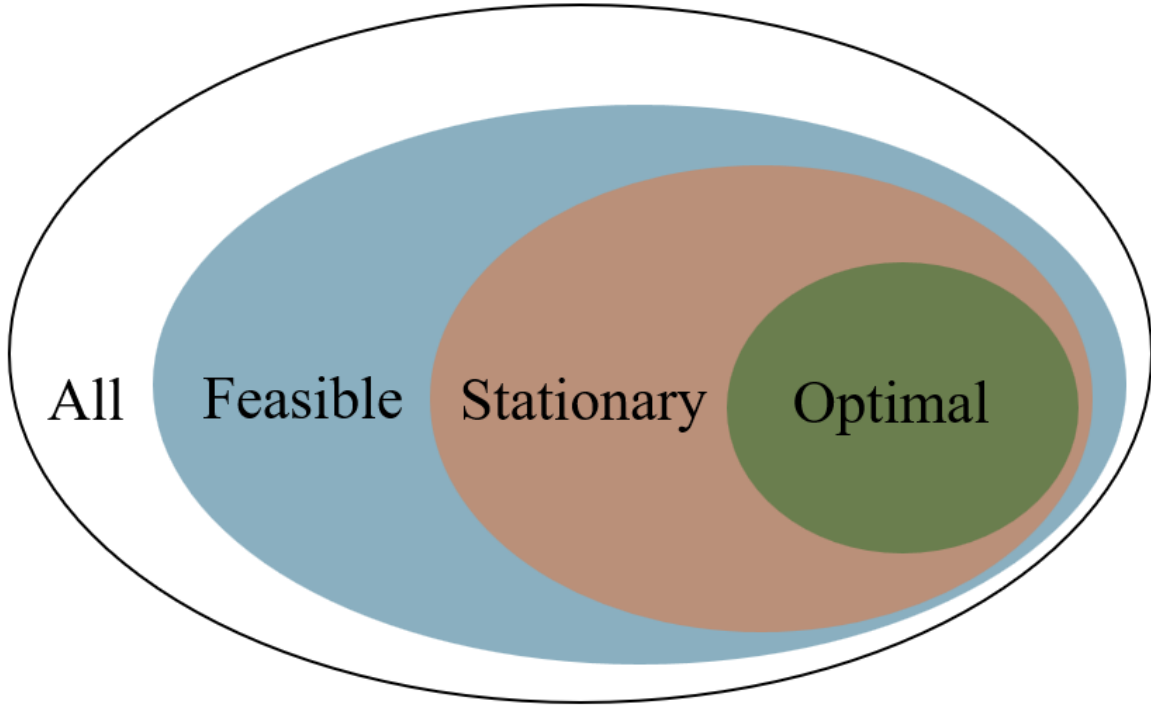


Figure 6.1: Venn Chart for Generalized Fritz John Condition

Proposition 3 *If a matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is symmetric, positive semi-definite, and of rank one, then there exists a vector $\mathbf{u} \in \mathbb{R}^n$ such that $\mathbf{A} = \mathbf{u}\mathbf{u}^\top$.*

Proof:

Since \mathbf{A} is symmetric and positive semi-definite, there exists a matrix $\mathbf{L} \in \mathbb{R}^{n \times n}$ such that

$$\mathbf{A} = \mathbf{L}\mathbf{L}^\top, \quad (6.22)$$

where \mathbf{L} is lower triangular with non-negative diagonal entries (by the Cholesky decomposition **Theorem 3**).

Given that $\text{rank}(\mathbf{A}) = 1$, it follows that \mathbf{A} has rank one, so \mathbf{L} must also have rank one. Thus, \mathbf{L} has exactly one nonzero column, say the first column, and all other columns are zero vectors. Denote this nonzero column by $\mathbf{u} \in \mathbb{R}^n$. Then \mathbf{L} can be written as

$$\mathbf{L} = \mathbf{u}\mathbf{e}_1^\top, \quad (6.23)$$

where $\mathbf{e}_1 = (1, 0, \dots, 0)^\top \in \mathbb{R}^n$ is the first standard basis vector.

Substituting back, we have

$$\mathbf{A} = \mathbf{L}\mathbf{L}^\top = (\mathbf{u}\mathbf{e}_1^\top)(\mathbf{e}_1\mathbf{u}^\top) = \mathbf{u}(\mathbf{e}_1^\top\mathbf{e}_1)\mathbf{u}^\top = \mathbf{u}\mathbf{u}^\top, \quad (6.24)$$

as desired.

Proposition 4 *Let \mathbf{A} to be a positive semi-definite matrix:*

- $\mathbf{A}\mathbf{v} = \mathbf{0}$ if and only if $\mathbf{v}^\top\mathbf{A}\mathbf{v} = 0$
- $\mathbf{A}\mathbf{v} \neq \mathbf{0}$ if and only if $\mathbf{v}^\top\mathbf{A}\mathbf{v} > 0$

Proof of $(\mathbf{A}\mathbf{v} = \mathbf{0} \Rightarrow \mathbf{v}^\top\mathbf{A}\mathbf{v} = 0)$:

If $\mathbf{A}\mathbf{v} = \mathbf{0}$, then $\mathbf{v}^\top\mathbf{A}\mathbf{v} = \mathbf{v}^\top\mathbf{0} = 0$, which completes the proof.

Proof of $(\mathbf{v}^\top \mathbf{A} \mathbf{v} = 0 \Rightarrow \mathbf{A} \mathbf{v} = \mathbf{0})$:

$\mathbf{v}^\top \mathbf{A} \mathbf{v} = 0$. Since \mathbf{A} is positive semi-definite, there exists \mathbf{L} such that $\mathbf{A} = \mathbf{L} \mathbf{L}^\top$ (**Theorem 3**):

$$\mathbf{v}^\top \mathbf{A} \mathbf{v} = 0 \quad (6.25)$$

$$\mathbf{v}^\top \mathbf{L} \mathbf{L}^\top \mathbf{v} = 0 \quad (6.26)$$

$$(\mathbf{L}^\top \mathbf{v})^\top (\mathbf{L}^\top \mathbf{v}) = 0 \quad (6.27)$$

$$(\mathbf{L}^\top \mathbf{v})^2 = 0 \quad (6.28)$$

$$\mathbf{L}^\top \mathbf{v} = \mathbf{0} \quad (6.29)$$

Continue:

$$\mathbf{A} \mathbf{w} = \mathbf{L} \mathbf{L}^\top \mathbf{w} = \mathbf{L} (\mathbf{L}^\top \mathbf{w}) = \mathbf{L} \cdot \mathbf{0} = \mathbf{0} \quad (6.30)$$

which completes the proof.

Proof of $(\mathbf{A} \mathbf{v} \neq \mathbf{0} \Rightarrow \mathbf{v}^\top \mathbf{A} \mathbf{v} > 0)$:

Since we have proved that if $\mathbf{v}^\top \mathbf{A} \mathbf{v} = 0$, then $\mathbf{A} \mathbf{v} = \mathbf{0}$. Therefore, when $\mathbf{A} \mathbf{v} \neq \mathbf{0}$, $\mathbf{v}^\top \mathbf{A} \mathbf{v} \neq 0$.

Since \mathbf{A} is positive semi-definite, by **Definition 2'** we have $\mathbf{v}^\top \mathbf{A} \mathbf{v} \geq 0$. Combine $\mathbf{v}^\top \mathbf{A} \mathbf{v} \neq 0$ and $\mathbf{v}^\top \mathbf{A} \mathbf{v} \geq 0$ we have $\mathbf{v}^\top \mathbf{A} \mathbf{v} > 0$, which completes the proof.

Proof of $(\mathbf{v}^\top \mathbf{A} \mathbf{v} > 0 \Rightarrow \mathbf{A} \mathbf{v} \neq \mathbf{0})$:

Since we have proved that $\mathbf{A} \mathbf{v} = \mathbf{0}$, $\mathbf{v}^\top \mathbf{A} \mathbf{v} = 0$. If $\mathbf{v}^\top \mathbf{A} \mathbf{v} > 0$, we know that $\mathbf{A} \mathbf{v} \neq \mathbf{0}$, which completes the proof.