

2) Binary classification problem;  $\theta \in \{0, 1\}$  denotes class label

$$X | (\theta=0) \sim N_p(\mu_0, \Sigma) \text{ and } X | (\theta=1) \sim N_p(\mu_1, \Sigma)$$

$\Sigma$  is +DEF, 0-1 Loss, prior dist of  $\theta$  is  $P(\theta=0)=1/2$  and  $P(\theta=1)=1/2$ , so  $\theta \sim \text{Ber}(1/2)$

(a) Derive the Bayes rule for classifying a new obs.  $X \in \mathbb{R}^p$ .

→ Bayes rules minimize post. exp. loss

Under 0-1 Loss, Bayes Rule classifies observation to the posterior mode → whichever occurs more,  $\theta=1$  or  $\theta=0$

$$\text{So, when } \frac{P(\theta=1|X)}{P(\theta=0|X)} > 1, \text{ then } X \text{ gets assigned to class 1}$$

$$\text{We know } P(\theta_1|X) \propto P(X|\theta_1)P(\theta_1) = \{(2\pi)^p |\Sigma|\}^{-1/2} \exp\left\{-\frac{(X-\mu_1)^T \Sigma^{-1} (X-\mu_1)}{2}\right\} \cdot \frac{1}{2}$$

$$\text{and } P(\theta_2|X) \propto P(X|\theta_2)P(\theta_2) = \{(2\pi)^p |\Sigma|\}^{-1/2} \exp\left\{-\frac{(X-\mu_0)^T \Sigma^{-1} (X-\mu_0)}{2}\right\} \cdot \frac{1}{2}$$

So, Bayes rule assigns  $X$  to  $\theta=1$  if:

$$\{(2\pi)^p |\Sigma|\}^{-1/2} \exp\left\{-\frac{(X-\mu_1)^T \Sigma^{-1} (X-\mu_1)}{2}\right\} \cdot \frac{1}{2} > \{(2\pi)^p |\Sigma|\}^{-1/2} \exp\left\{-\frac{(X-\mu_0)^T \Sigma^{-1} (X-\mu_0)}{2}\right\} \cdot \frac{1}{2}$$

$$\Leftrightarrow \exp\left\{-\frac{(X-\mu_1)^T \Sigma^{-1} (X-\mu_1)}{2}\right\} > \exp\left\{-\frac{(X-\mu_0)^T \Sigma^{-1} (X-\mu_0)}{2}\right\}$$

$$\Leftrightarrow \exp\left\{-\frac{(X-\mu_1)^T \Sigma^{-1} (X-\mu_1)}{2} + \frac{(X-\mu_0)^T \Sigma^{-1} (X-\mu_0)}{2}\right\} > 1$$

$$\Leftrightarrow -\frac{(X-\mu_1)^T \Sigma^{-1} (X-\mu_1)}{2} + \frac{(X-\mu_0)^T \Sigma^{-1} (X-\mu_0)}{2} > 0$$

$$\Leftrightarrow -\cancel{X^T \Sigma^{-1} X} + \underline{\mu_1^T \Sigma^{-1} X} + \underline{X^T \Sigma^{-1} \mu_1} - \mu_1^T \Sigma^{-1} \mu_1 + \cancel{X^T \Sigma^{-1} X} - \underline{X^T \Sigma^{-1} \mu_0} - \underline{\mu_0^T \Sigma^{-1} X} + \mu_0^T \Sigma^{-1} \mu_0 > 0$$

$$\Leftrightarrow X^T \Sigma^{-1} (\mu_1 - \mu_0) + \underbrace{(\mu_1^T - \mu_0^T) \Sigma^{-1} X}_{(\mu_1 - \mu_0)^T} + \underbrace{(-\mu_1^T + \mu_0^T) \Sigma^{-1} (\mu_1 + \mu_0)}_{(\mu_1 + \mu_0)^T} > 0$$

$$\text{Let } \delta = (\mu_1 - \mu_0) \Leftrightarrow X^T \Sigma^{-1} \delta - \delta^T \Sigma^{-1} X - \delta^T \Sigma^{-1} 2\bar{\mu} > 0$$

$$\text{and } \bar{\mu} = \frac{\mu_0 + \mu_1}{2} \quad (\delta^T \Sigma^{-1} X)^T - \delta^T \Sigma^{-1} X - 2\delta^T \Sigma^{-1} \bar{\mu} > 0$$

$$2\delta^T \Sigma^{-1} (X - \bar{\mu}) > 0$$

$$\Leftrightarrow \boxed{\delta^T \Sigma^{-1} (X - \bar{\mu}) > 0} \checkmark$$

probably would have stepped here...

↓ Had to look ahead to part (d) to see  $\delta$  and  $\bar{\mu}$  defined

(b) Derive the misclassification rate  $R^*$  of the Bayes rule.

$$R^* = R(0, d) + R(1, d) \leftarrow \text{BINARY! So, 2 risks to consider!}$$

$$R^* = P(\text{choose } X \text{ in } \theta=1 \mid \theta=0) \cdot P(\theta=0) + P(\text{choose } X \text{ in } \theta=0 \mid \theta=1) \cdot P(\theta=1)$$

$$R^* = \frac{1}{2} P(-\delta^T \Sigma^{-1}(X - \mu) > 0 \mid \theta=0) + \frac{1}{2} P(-\delta^T \Sigma^{-1}(X - \mu) \leq 0 \mid \theta=1)$$

$\swarrow$  since their  $\delta$  in (d) is  $\hat{\mu}_0 - \hat{\mu}_1$  and their  $\mu = \bar{\mu}$  from (a)

$$\frac{1}{2} P(-\delta^T \Sigma^{-1}(X - \mu) > 0 \mid \theta=0)$$

$$= \frac{1}{2} P(\delta^T \Sigma^{-1}(X - \mu) \leq 0 \mid \theta=0)$$

$$= \frac{1}{2} P(\delta^T \Sigma^{-1}(X - \mu_0 + \mu_0 - \mu) \leq 0 \mid \theta=0) ; \quad X \mid \theta=0 \sim N_p(\mu_0, \Sigma)$$

$$= \frac{1}{2} P(\delta^T \Sigma^{-1}(X - \mu_0) \leq \delta^T \Sigma^{-1}(\mu - \mu_0) \mid \theta=0)$$

$$= \frac{1}{2} P\left(\frac{\delta^T \Sigma^{-1}(X - \mu_0)}{\sqrt{\delta^T \Sigma^{-1} \delta}} \leq \frac{\delta^T \Sigma^{-1}(\mu - \mu_0)}{\sqrt{\delta^T \Sigma^{-1} \delta}}\right)$$

$$\left. \begin{aligned} &\Rightarrow X - \mu_0 \mid \theta=0 \sim N_p(0, \Sigma) \\ &\Rightarrow \delta^T \Sigma^{-1}(X - \mu_0) \mid \theta=0 \sim N_p(0, \delta^T \Sigma^{-1} \delta) \\ &\Rightarrow \frac{\delta^T \Sigma^{-1}(X - \mu_0)}{\sqrt{\delta^T \Sigma^{-1} \delta}} \mid \theta=0 \sim N_p(0, 1) \end{aligned} \right\}$$

$$= \frac{1}{2} \Phi\left(\frac{\delta^T \Sigma^{-1}(\mu - \mu_0)}{\sqrt{\delta^T \Sigma^{-1} \delta}}\right), \text{ where } \delta = \mu_0 - \mu_1 \text{ and } \mu = \frac{\mu_0 + \mu_1}{2} ; \quad \mu - \mu_0 = \frac{\mu_1 - \mu_0}{2} = -\frac{1}{2} \delta$$

$$\text{Similarly, } \frac{1}{2} P(-\delta^T \Sigma^{-1}(X - \mu) \leq 0 \mid \theta=1) = \frac{1}{2} \Phi\left(-\frac{\delta^T \Sigma^{-1}(\mu - \mu_1)}{\sqrt{\delta^T \Sigma^{-1} \delta}}\right)$$

$$= \frac{1}{2} \Phi\left(-\frac{\delta^T \Sigma^{-1}(-\frac{1}{2} \delta)}{\sqrt{\delta^T \Sigma^{-1} \delta}}\right) = \frac{1}{2} \Phi\left(\frac{1}{2} (\delta^T \Sigma^{-1} \delta)^{1/2}\right)$$

$$\text{So, } R^* = \frac{1}{2} \Phi\left(\frac{1}{2} (\delta^T \Sigma^{-1} \delta)^{1/2}\right) + \frac{1}{2} \Phi\left(\frac{1}{2} (\delta^T \Sigma^{-1} \delta)^{1/2}\right)$$

$$= \Phi\left(\frac{1}{2} (\delta^T \Sigma^{-1} \delta)^{1/2}\right)$$

(C) Let  $X_{0i}$  ( $i=1, \dots, n_0$ ) be iid samples class  $\theta=0$  and  $X_{1i}$  ( $i=1, \dots, n_1$ ) be iid samples class  $\theta=1$ .

$X_{0i} \perp\!\!\!\perp X_{1i}$ . Derive MLEs of  $(\mu_0, \mu_1, \Sigma)$ .

The likelihood of  $X_{0i} = \prod_{i=1}^{n_0} \{ (2\pi)^p |\Sigma| \}^{-1/2} \exp \{ - (X_{0i} - \mu_0)^T \Sigma^{-1} (X_{0i} - \mu_0) / 2 \}$

$$\Rightarrow \text{Log-likelihood: } \ln(\mu_0, \Sigma) = \sum_{i=1}^{n_0} \left\{ -\frac{1}{2} \log[(2\pi)^p |\Sigma|] - \frac{(X_{0i} - \mu_0)^T \Sigma^{-1} (X_{0i} - \mu_0)}{2} \right\}$$

$$\text{Then, } \frac{\partial \ln(\mu_0, \Sigma)}{\partial \mu_0} = \sum_{i=1}^{n_0} \frac{2(X_{0i} - \mu_0) \Sigma^{-1} (-1)}{2} \stackrel{\text{set}}{=} 0 \rightarrow \sum_{i=1}^{n_0} \hat{\mu}_0 - X_{0i} = 0 \rightarrow n_0 \hat{\mu}_0 = \sum_{i=1}^{n_0} X_{0i}$$

$$\Rightarrow \boxed{\hat{\mu}_0 = \frac{1}{n_0} \sum_{i=1}^{n_0} X_{0i}} \quad \text{and similarly,} \quad \boxed{\hat{\mu}_1 = \frac{1}{n_1} \sum_{i=1}^{n_1} X_{1i}}$$

$= \bar{X}_0 \qquad \qquad \qquad = \bar{X}_1$

Now, to find the MLE of  $\Sigma$ , we rewrite the Log-likelihood in terms of trace:

$$\ln(\mu_0, \Sigma) = -\frac{p n_0}{2} \log(2\pi) - \frac{n_0}{2} \log(|\Sigma|) - \sum_{i=1}^{n_0} (X_{0i} - \mu_0)^T \Sigma^{-1} (X_{0i} - \mu_0)$$

$$\propto -\frac{n_0}{2} \log(|\Sigma|) - \frac{1}{2} \text{tr} \left( \sum_{i=1}^{n_0} (X_{0i} - \mu_0)^T \Sigma^{-1} (X_{0i} - \mu_0) \right)$$

← b/c scalar

$$= -\frac{n_0}{2} \log\left(\frac{1}{|\Sigma^{-1}|}\right) - \frac{1}{2} \text{tr} \left( \Sigma^{-1} \underbrace{\sum_{i=1}^{n_0} (X_{0i} - \mu_0)(X_{0i} - \mu_0)^T}_{\text{Symmetric}} \right)$$

Note, as given:  $\frac{\partial \log(|\Sigma|)}{\partial \Sigma} = \Sigma^{-1}$  and  $\frac{\partial \text{tr} \left( \Sigma^{-1} \sum_{i=1}^{n_0} (X_{0i} - \mu_0)(X_{0i} - \mu_0)^T \right)}{\partial \Sigma^{-1}} = \sum_{i=1}^{n_0} (X_{0i} - \hat{\mu}_0)(X_{0i} - \hat{\mu}_0)^T$

← plug in  $\hat{\mu}_0$

This makes this part much easier, so we take derivative wrt  $\Sigma^{-1}$

Note  $\frac{\partial \log(\frac{1}{|\Sigma^{-1}|})}{\partial \Sigma^{-1}} = -\frac{\partial \log(|\Sigma^{-1}|)}{\partial \Sigma^{-1}} = -(\Sigma^{-1})^{-1} = -\Sigma$

Since  $X_{0i} \perp\!\!\!\perp X_{1i}$ , then  $\ln(\mu_0, \mu_1, \Sigma) = \ln(\mu_0, \Sigma) + \ln(\mu_1, \Sigma)$  \*

$$= -\frac{n_0}{2} \log\left(\frac{1}{|\Sigma^{-1}|}\right) - \frac{1}{2} \text{tr} \left( \Sigma^{-1} \sum_{i=1}^{n_0} (X_{0i} - \mu_0)(X_{0i} - \mu_0)^T \right) - \frac{n_1}{2} \log\left(\frac{1}{|\Sigma^{-1}|}\right) - \frac{1}{2} \text{tr} \left( \Sigma^{-1} \sum_{i=1}^{n_1} (X_{1i} - \mu_1)(X_{1i} - \mu_1)^T \right)$$

$$\Rightarrow \frac{\partial \ln(\mu_0, \mu_1, \Sigma)}{\partial \Sigma^{-1}} = -\frac{n_0}{2} (-\Sigma) - \frac{n_1}{2} (-\Sigma) - \frac{1}{2} \sum_{i=1}^{n_0} (X_{0i} - \bar{X}_0)(X_{0i} - \bar{X}_0)^T - \frac{1}{2} \sum_{i=1}^{n_1} (X_{1i} - \bar{X}_1)(X_{1i} - \bar{X}_1)^T \stackrel{\text{set}}{=} 0$$

$$\hat{\Sigma} = \left( \frac{1}{n_0 + n_1} \right) \left( \sum_{i=1}^{n_0} (X_{0i} - \bar{X}_0)(X_{0i} - \bar{X}_0)^T + \sum_{i=1}^{n_1} (X_{1i} - \bar{X}_1)(X_{1i} - \bar{X}_1)^T \right)$$



(d) If we replace  $(\mu_0, \mu_1, \Sigma)$  in Bayes rule with  $(\hat{\mu}_0, \hat{\mu}_1, \hat{\Sigma})$ , prove that the misclassification rate of the resulting rule, i.e. the probability of classifying  $x$  to a wrong class given the training data  $\{x_{0i}\}_{i=1}^{n_0}$  and  $\{x_{1i}\}_{i=1}^{n_1}$ , is given by:  $\frac{1}{2} \Phi\left(\frac{\hat{\delta}^T \hat{\Sigma}^{-1}(\mu_1 - \hat{\mu})}{\sqrt{\hat{\delta}^T \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} \hat{\delta}}}\right) + \frac{1}{2} \Phi\left(-\frac{\hat{\delta}^T \hat{\Sigma}^{-1}(\mu_0 - \hat{\mu})}{\sqrt{\hat{\delta}^T \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} \hat{\delta}}}\right)$

From (b), we have  $R^* = \frac{1}{2} P(\delta^T \Sigma^{-1}(x - \bar{\mu}) > 0 | \theta = 0) + \frac{1}{2} P(\delta^T \Sigma^{-1}(x - \bar{\mu}) \leq 0 | \theta = 1)$

→ replacing  $(\mu_0, \mu_1, \Sigma)$  with  $(\hat{\mu}_0, \hat{\mu}_1, \hat{\Sigma})$ , we have:

$$\begin{aligned} R^* &= \frac{1}{2} P\left((\hat{\mu}_1 - \hat{\mu}_0)^T \hat{\Sigma}^{-1} \left(x - \frac{\hat{\mu}_1 + \hat{\mu}_0}{2}\right) > 0 \mid \theta = 0\right) + \frac{1}{2} P\left((\hat{\mu}_1 - \hat{\mu}_0)^T \hat{\Sigma}^{-1} \left(x - \frac{\hat{\mu}_1 + \hat{\mu}_0}{2}\right) \leq 0 \mid \theta = 1\right) \\ &= \frac{1}{2} P\left(\hat{\delta}^T \hat{\Sigma}^{-1} (x - \hat{\mu}) > 0 \mid \theta = 0\right) \\ &= \frac{1}{2} P\left(\hat{\delta}^T \hat{\Sigma}^{-1} (x - \mu_0 + \mu_0 - \hat{\mu}) > 0 \mid \theta = 0\right) ; \quad x \mid \theta = 0 \sim N_p(\mu_0, \Sigma) \Rightarrow x - \mu_0 \mid \theta = 0 \sim N(0, \Sigma) \\ &\quad \Rightarrow \hat{\Sigma}^{-1} (x - \mu_0) \mid \theta = 0 \sim N(0, \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1}) \\ &\quad \Rightarrow \hat{\delta}^T \hat{\Sigma}^{-1} (x - \mu_0) \mid \theta = 0 \sim N(0, \hat{\delta}^T \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} \hat{\delta}) \\ &= \frac{1}{2} P\left(\hat{\delta}^T \hat{\Sigma}^{-1} (x - \mu_0) > \hat{\delta}^T \hat{\Sigma}^{-1} (\hat{\mu} - \mu_0) \mid \theta = 0\right) \\ &= \frac{1}{2} P\left(\frac{\hat{\delta}^T \hat{\Sigma}^{-1} (x - \mu_0)}{\sqrt{\hat{\delta}^T \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} \hat{\delta}}} > \frac{\hat{\delta}^T \hat{\Sigma}^{-1} (\hat{\mu} - \mu_0)}{\sqrt{\hat{\delta}^T \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} \hat{\delta}}} \mid \theta = 0\right) \\ &\quad \Rightarrow \frac{\hat{\delta}^T \hat{\Sigma}^{-1} (x - \mu_0)}{\sqrt{\hat{\delta}^T \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} \hat{\delta}}} \mid \theta = 0 \sim N(0, 1) \\ &= \frac{1}{2} \Phi\left(-\frac{\hat{\delta}^T \hat{\Sigma}^{-1} (\hat{\mu} - \mu_0)}{\sqrt{\hat{\delta}^T \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} \hat{\delta}}}\right) \end{aligned}$$

\* Note: My  $\hat{\delta}$  is  $-\hat{\delta}$  of what is given

Similarly,  $\frac{1}{2} P((\hat{\mu}_1 - \hat{\mu}_0)^T \hat{\Sigma}^{-1} (x - \hat{\mu}) \leq 0 \mid \theta = 1) = \frac{1}{2} \Phi\left(\frac{-\hat{\delta}^T \hat{\Sigma}^{-1} (\hat{\mu} - \mu_1)}{\sqrt{\hat{\delta}^T \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} \hat{\delta}}}\right)$

Thus, the misclassification rate is:

$$R^* = \frac{1}{2} \Phi\left(\frac{\hat{\delta}^T \hat{\Sigma}^{-1} (\mu_1 - \hat{\mu})}{\sqrt{\hat{\delta}^T \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} \hat{\delta}}}\right) + \frac{1}{2} \Phi\left(-\frac{\hat{\delta}^T \hat{\Sigma}^{-1} (\mu_0 - \hat{\mu})}{\sqrt{\hat{\delta}^T \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} \hat{\delta}}}\right)$$

(e) We propose another classification rule that assigns  $x$  to class of  $\theta = 0$  iff  $\hat{\beta}^T(x - \hat{\mu}) \geq 0$   
 where  $\hat{\mu} = (\hat{\mu}_0 + \hat{\mu}_1)/2$  and  $\hat{\beta}$  solves  $\hat{\beta} = \arg\min_{\beta \in \mathbb{R}^p} \frac{1}{2} \beta^T \hat{\Sigma} \beta - (\hat{\mu}_0 - \hat{\mu}_1)^T \beta + \lambda \sum_{j=1}^p |\beta_j|$

Derive the Majorization-Minimization algorithm for solving  $\hat{\beta}$ .

Give an explicit choice of step size and closed-form expressions on how iterations need to be done.

[MAJORIZATION]

$$\text{Let } l(\beta) = \frac{1}{2} \beta^T \hat{\Sigma} \beta - (\hat{\mu}_0 - \hat{\mu}_1)^T \beta$$

$$\text{Then, } \nabla l(\beta) = \hat{\Sigma} \beta - (\hat{\mu}_0 - \hat{\mu}_1) \quad ; \quad \nabla^2 l(\beta) = \hat{\Sigma}$$

By 2nd-order Taylor Expansion of  $l(\beta)$  around  $\tilde{\beta}$ :

$$\begin{aligned} l(\beta) &= l(\tilde{\beta}) + \nabla l(\tilde{\beta})^T (\beta - \tilde{\beta}) + \frac{1}{2} (\beta - \tilde{\beta})^T \nabla^2 l(\tilde{\beta}) (\beta - \tilde{\beta}) \\ &\leq \underbrace{l(\tilde{\beta}) + \nabla l(\tilde{\beta})^T (\beta - \tilde{\beta}) + c (\beta - \tilde{\beta})^T (\beta - \tilde{\beta})}_{\text{Call this } l_q(\beta)}, \text{ where } c \geq \lambda_{\max}(\hat{\Sigma}) \end{aligned}$$

[MINIMIZATION]

Now, let  $p(\beta) = l_q(\beta) + \lambda \|\beta\|_1$ ; we want to minimize  $p(\beta)$  w.r.t  $\beta$  to find  $\hat{\beta}^{(new)}$

$$p(\beta) = l_q(\beta) + \lambda \|\beta\|_1 = l(\tilde{\beta}) + \nabla l(\tilde{\beta})^T (\beta - \tilde{\beta}) + c (\beta - \tilde{\beta})^T (\beta - \tilde{\beta}) + \lambda \|\beta\|_1$$

$$\frac{\partial p(\beta)}{\partial \beta} = \nabla l(\tilde{\beta}) + 2c (\beta - \tilde{\beta}) + \lambda \frac{\partial \|\beta\|_1}{\partial \beta}$$

$$\cdot \text{ If } \beta \neq 0, \text{ then } \frac{\partial \|\beta\|_1}{\partial \beta} = \frac{\partial}{\partial \beta} (\sqrt{\beta^2}) = \frac{1}{2} (\beta^2)^{-1/2} \cdot 2\beta = \frac{\beta}{\|\beta\|_1}$$

$$\cdot \text{ If } \beta = 0, \text{ then we have previously shown } \frac{\partial \|\beta\|_1}{\partial \beta} = [-1, 1]$$

$$\text{So, } \frac{\partial \|\beta\|_1}{\partial \beta} = \begin{cases} 1 & \text{if } \beta > 0 \\ [-1, 1] & \text{if } \beta = 0 \\ -1 & \text{if } \beta < 0 \end{cases}$$

(and we know  $0 \in \frac{\partial p(\beta)}{\partial \beta} \Big|_{\hat{\beta}^{new}}$  by KKT)

→ We will look @ all 3 cases:

• When  $\beta > 0 \Rightarrow \tilde{\beta}^{(new)} > 0$

$$\frac{d\rho(\beta)}{d\beta} = \hat{\Sigma} \tilde{\beta} - (\hat{\mu}_0 - \hat{\mu}_1) + 2c(\beta - \tilde{\beta}) + \lambda \stackrel{\text{set}}{=} 0$$

$$-\hat{\Sigma} \tilde{\beta} + (\hat{\mu}_0 - \hat{\mu}_1) + 2c\tilde{\beta} - \lambda = 2c\tilde{\beta}^{(new)}$$

$$\Rightarrow \tilde{\beta}^{(new)} = \tilde{\beta} + \frac{1}{2c} \left( \hat{\Sigma} \tilde{\beta} - (\hat{\mu}_0 - \hat{\mu}_1) + \lambda \right)$$

as long as this is  $> 0$

• When  $\beta < 0 \Rightarrow \tilde{\beta}^{(new)} < 0$

$$\text{Similarly, } \tilde{\beta}^{(new)} = \tilde{\beta} - \frac{1}{2c} \left( \hat{\Sigma} \tilde{\beta} - (\hat{\mu}_0 - \hat{\mu}_1) + \lambda \right)$$

as long as this is  $< 0$

• When  $\beta = 0 \Rightarrow \tilde{\beta}^{(new)} = 0$

$$\frac{d\rho(\beta)}{d\beta} = \hat{\Sigma} \tilde{\beta} - (\hat{\mu}_0 - \hat{\mu}_1) + 2c(\beta - \tilde{\beta}) + \lambda [-1, 1] \Big|_{\beta=0}$$

$$= \left[ \underbrace{\hat{\Sigma} \tilde{\beta} - (\hat{\mu}_0 - \hat{\mu}_1) - 2c\tilde{\beta} - \lambda}_{\text{must be } < 0}, \underbrace{\hat{\Sigma} \tilde{\beta} - (\hat{\mu}_0 - \hat{\mu}_1) - 2c\tilde{\beta} + \lambda}_{\text{must be } > 0} \right]$$

as long as  $(\hat{\Sigma} \tilde{\beta} - (\hat{\mu}_0 - \hat{\mu}_1) - 2c\tilde{\beta}) \in [-\lambda, \lambda]$

$$\tilde{\beta} - \frac{1}{2c} (\hat{\Sigma} \tilde{\beta} - (\hat{\mu}_0 - \hat{\mu}_1)) \in \left[-\frac{1}{2c}\lambda, \frac{1}{2c}\lambda\right]$$

$$\text{So, } \tilde{\beta}^{(new)} = S\left(\underbrace{\tilde{\beta} - \frac{1}{2c} (\hat{\Sigma} \tilde{\beta} - (\hat{\mu}_0 - \hat{\mu}_1))}_A, \frac{1}{2c} \lambda\right)$$

$$= \begin{cases} A - \lambda & \text{if } |A| > \lambda \text{ and } A > 0 \\ 0 & \text{if } |A| \leq \lambda \\ A + \lambda & \text{if } |A| > \lambda \text{ and } A < 0 \end{cases}$$

We first initialize  $\beta$  at  $\beta^{(0)} \in \mathbb{R}^p$

Iterate until convergence:  $\beta^{(k)} = S(A, \frac{1}{2c} \lambda)$

Stop if  $\|\beta^{(k)} - \beta^{(k-1)}\|_2 < \varepsilon$  for some pre-defined stepping threshold  $\varepsilon$ .



(f)  $R_n$  is misclassification rate of rule in (e).

Suppose we can show  $\hat{\beta} \xrightarrow{p} \Sigma^{-1}(\mu_0 - \mu_1)$  as  $n \rightarrow \infty$ . Show  $R_n \xrightarrow{p} R^*$

From (b), we know  $R^* = \Phi\left(-\frac{1}{2}(\delta^T \Sigma^{-1} \delta)^{1/2}\right) = \Phi\left(-\frac{1}{2}((\mu_0 - \mu_1)^T \Sigma^{-1} (\mu_0 - \mu_1))^{1/2}\right)$

The rule in (e) assigns  $X$  to class  $\theta=0$  iff  $\hat{\beta}^T(X - \hat{\mu}) \geq 0$

$$\begin{aligned} \text{So, } R_n &= 1 \cdot P(\hat{\beta}^T(X - \hat{\mu}) \geq 0 | \theta=1) P(\theta=1) + 1 \cdot P(\hat{\beta}^T(X - \hat{\mu}) < 0 | \theta=0) P(\theta=0) \\ &= \frac{1}{2} P(\hat{\beta}^T(X - \hat{\mu}) \geq 0 | \theta=1) + \frac{1}{2} P(\hat{\beta}^T(X - \hat{\mu}) < 0 | \theta=0) \end{aligned}$$

$$P(\hat{\beta}^T(X - \hat{\mu}) \geq 0 | \theta=1) = P(\hat{\beta}^T(X - \mu_1 - (\hat{\mu} - \mu_1)) \geq 0 | \theta=1)$$

$$= P(\hat{\beta}^T(X - \mu_1) \geq \hat{\beta}^T(\hat{\mu} - \mu_1) | \theta=1) ; \quad X | (\theta=1) \sim N_p(\mu_1, \Sigma)$$

$$\Rightarrow \hat{\beta}^T(X - \mu_1) | (\theta=1) \sim N_p(0, \hat{\beta}^T \Sigma \hat{\beta})$$

$$= P\left(\frac{\hat{\beta}^T(X - \mu_1)}{\sqrt{\hat{\beta}^T \Sigma \hat{\beta}}} \geq \frac{\hat{\beta}^T(\hat{\mu} - \mu_1)}{\sqrt{\hat{\beta}^T \Sigma \hat{\beta}}} \mid \theta=1\right) \leftarrow \Rightarrow \frac{\hat{\beta}^T(X - \mu_1)}{\sqrt{\hat{\beta}^T \Sigma \hat{\beta}}} \mid \theta=1 \sim N_p(0, 1)$$

$$= \Phi\left(\frac{-\hat{\beta}^T(\hat{\mu} - \mu_1)}{\sqrt{\hat{\beta}^T \Sigma \hat{\beta}}}\right)$$

Similarly,  $P(\hat{\beta}^T(X - \hat{\mu}) < 0 | \theta=0) = \Phi\left(\frac{\hat{\beta}^T(\hat{\mu} - \mu_0)}{\sqrt{\hat{\beta}^T \Sigma \hat{\beta}}}\right)$

$$\text{So, } R_n = \frac{1}{2} \Phi\left(\frac{-\hat{\beta}^T(\hat{\mu} - \mu_1)}{\sqrt{\hat{\beta}^T \Sigma \hat{\beta}}}\right) + \frac{1}{2} \Phi\left(\frac{\hat{\beta}^T(\hat{\mu} - \mu_0)}{\sqrt{\hat{\beta}^T \Sigma \hat{\beta}}}\right)$$

Since  $\hat{\beta} \xrightarrow{p} \Sigma^{-1}(\mu_0 - \mu_1)$ , then  $R_n \rightarrow \frac{1}{2} \Phi\left(\frac{(\mu_0 - \mu_1)^T \Sigma^{-1} (\hat{\mu} - \mu_1)}{((\mu_0 - \mu_1)^T \Sigma^{-1} (\mu_0 - \mu_1))^{1/2}}\right) + \frac{1}{2} \Phi\left(\frac{(\mu_0 - \mu_1)^T \Sigma^{-1} (\hat{\mu} - \mu_0)}{((\mu_0 - \mu_1)^T \Sigma^{-1} (\mu_0 - \mu_1))^{1/2}}\right)$

$$\hat{\mu} - \mu_1 = \frac{\mu_0 + \mu_1}{2} - \mu_1 = \frac{1}{2} \delta \quad \Rightarrow \quad = \frac{1}{2} \Phi\left(\frac{-\frac{1}{2}(\mu_0 - \mu_1)^T \Sigma^{-1} \hat{\delta}}{((\mu_0 - \mu_1)^T \Sigma^{-1} (\mu_0 - \mu_1))^{1/2}}\right) + \frac{1}{2} \Phi\left(\frac{-\frac{1}{2}(\mu_0 - \mu_1)^T \Sigma^{-1} \hat{\delta}}{((\mu_0 - \mu_1)^T \Sigma^{-1} (\mu_0 - \mu_1))^{1/2}}\right)$$

$$\hat{\mu} - \mu_0 = -\frac{1}{2} \delta \quad = \frac{1}{2} \Phi\left(\frac{-\frac{1}{2}(\mu_0 - \mu_1)^T \Sigma^{-1} (\mu_0 - \mu_1)}{((\mu_0 - \mu_1)^T \Sigma^{-1} (\mu_0 - \mu_1))^{1/2}}\right) + \frac{1}{2} \Phi\left(\frac{-\frac{1}{2}(\mu_0 - \mu_1)^T \Sigma^{-1} (\mu_0 - \mu_1)}{((\mu_0 - \mu_1)^T \Sigma^{-1} (\mu_0 - \mu_1))^{1/2}}\right)$$

$$\hat{\mu}_1 \rightarrow \mu_1, \hat{\mu}_0 \rightarrow \mu_0$$

$$\text{so } \hat{\delta} \rightarrow (\mu_0 - \mu_1)$$

$$= \Phi\left(-\frac{1}{2}((\mu_0 - \mu_1)^T \Sigma^{-1} (\mu_0 - \mu_1))^{1/2}\right) = R^* \quad \checkmark$$