

1. (25 points) Consider a binary classification problem that $\theta \in \{0, 1\}$ denotes the class label, $\mathbf{X} | (\theta = 0) \sim N_p(\boldsymbol{\mu}_0, \boldsymbol{\Sigma})$ and $\mathbf{X} | (\theta = 1) \sim N_p(\boldsymbol{\mu}_1, \boldsymbol{\Sigma})$, where N_p denotes the p -dimensional multivariate normal distribution. Suppose 0-1 loss is used, and the prior distribution of θ is $P(\theta = 0) = 1/2$ and $P(\theta = 1) = 1/2$.

- (i) (4 points) Derive the Bayes rule for a classifying a new observation $\mathbf{x} \in \mathcal{R}^p$.
- (ii) (4 points) Derive the misclassification rate R^* of the Bayes rule.
- (iii) (4 points) Let \mathbf{X}_{0i} ($i = 1, \dots, n_0$) be independent and identically distributed (i.i.d) samples from the class of $\theta = 0$ and \mathbf{X}_{1i} ($i = 1, \dots, n_1$) be i.i.d samples from the class of $\theta = 1$, and \mathbf{X}_{0i} is independent of \mathbf{X}_{1i} . Derive the maximum likelihood estimators $(\hat{\boldsymbol{\mu}}_0, \hat{\boldsymbol{\mu}}_1, \hat{\boldsymbol{\Sigma}})$ of $(\boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \boldsymbol{\Sigma})$.
- (iv) (4 points) If we replace $(\boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \boldsymbol{\Sigma})$ in the Bayes rule with $(\hat{\boldsymbol{\mu}}_0, \hat{\boldsymbol{\mu}}_1, \hat{\boldsymbol{\Sigma}})$, prove that the misclassification rate of the resulting rule, i.e., the probability of classifying \mathbf{x} to a wrong class given the training data $\{\mathbf{X}_{0i}\}_{i=1}^{n_0}$ and $\{\mathbf{X}_{1i}\}_{i=1}^{n_1}$, is given by

$$\frac{1}{2} \Phi \left(\frac{\hat{\boldsymbol{\delta}}^T \hat{\boldsymbol{\Sigma}}^{-1} (\boldsymbol{\mu}_1 - \hat{\boldsymbol{\mu}})}{\sqrt{\hat{\boldsymbol{\delta}}^T \hat{\boldsymbol{\Sigma}}^{-1} \boldsymbol{\Sigma} \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\delta}}}} \right) + \frac{1}{2} \Phi \left(-\frac{\hat{\boldsymbol{\delta}}^T \hat{\boldsymbol{\Sigma}}^{-1} (\boldsymbol{\mu}_0 - \hat{\boldsymbol{\mu}})}{\sqrt{\hat{\boldsymbol{\delta}}^T \hat{\boldsymbol{\Sigma}}^{-1} \boldsymbol{\Sigma} \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\delta}}}} \right),$$

where $\hat{\boldsymbol{\delta}} = \hat{\boldsymbol{\mu}}_0 - \hat{\boldsymbol{\mu}}_1$ and $\hat{\boldsymbol{\mu}} = (\hat{\boldsymbol{\mu}}_0 + \hat{\boldsymbol{\mu}}_1)/2$.

- (v) (5 points) We propose another classification rule that assigns \mathbf{x} to the class of $\theta = 0$ if and only if $\hat{\boldsymbol{\beta}}^T (\mathbf{x} - \hat{\boldsymbol{\mu}}) \geq 0$, where $\hat{\boldsymbol{\mu}} = (\hat{\boldsymbol{\mu}}_0 + \hat{\boldsymbol{\mu}}_1)/2$ and $\hat{\boldsymbol{\beta}}$ solves the following problem

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta} \in \mathcal{R}^p}{\operatorname{argmin}} \frac{1}{2} \boldsymbol{\beta}^T \hat{\boldsymbol{\Sigma}} \boldsymbol{\beta} - (\hat{\boldsymbol{\mu}}_0 - \hat{\boldsymbol{\mu}}_1)^T \boldsymbol{\beta} + \lambda \sum_{j=1}^p |\beta_j|.$$

Derive the Majorization-Minimization algorithm for solving $\hat{\boldsymbol{\beta}}$. Give an explicit choice of step size and closed-form expressions on how iterations need to be done.

- (vi) (4 points) Let R_n denote the misclassification rate of the rule described in (v). Suppose we can show that $\hat{\boldsymbol{\beta}} \xrightarrow{P} \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)$ as $n \rightarrow \infty$. Using this result to prove $R_n \xrightarrow{P} R^*$.

You may use the following facts.

- (a) The density of $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is $\{(2\pi)^p |\boldsymbol{\Sigma}| \}^{-1/2} \exp\{-(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})/2\}$;
- (b) For symmetric matrices \mathbf{A} and \mathbf{M} ,

$$\frac{\partial \operatorname{tr}(\mathbf{A}\mathbf{M})}{\partial \mathbf{M}} = \frac{\partial \operatorname{tr}(\mathbf{M}\mathbf{A})}{\partial \mathbf{M}} = \mathbf{A}.$$

$$\frac{\partial \log |\mathbf{M}|}{\partial \mathbf{M}} = \mathbf{M}^{-1}.$$

Solution:

- (i) Under 0-1 loss, the Bayes rule is the posterior mode. Therefore, the Bayes rule assigns \mathbf{x} to $\theta = 0$ iff $f(\theta = 0|\mathbf{x}) > f(\theta = 1|\mathbf{x})$. That is equivalent as $-(1/2)(\mathbf{x} - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_0) > -(1/2)(\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_1)$. After some algebra, we find that the Bayes rule assigns \mathbf{x} to $\theta = 0$ iff $\boldsymbol{\delta}^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) > 0$, where $\boldsymbol{\delta} = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1$ and $\boldsymbol{\mu} = (\boldsymbol{\mu}_0 + \boldsymbol{\mu}_1)/2$.
- (ii) The misclassification rate of Bayes rule is given by

$$\begin{aligned}
R^* &= \frac{1}{2}P(\boldsymbol{\delta}^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) > 0 | \theta = 1) + \frac{1}{2}P(\boldsymbol{\delta}^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \leq 0 | \theta = 0) \\
&= \frac{1}{2}P\left(\frac{\boldsymbol{\delta}^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) - \boldsymbol{\delta}^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu})}{\sqrt{\boldsymbol{\delta}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\delta}}} > -\frac{\boldsymbol{\delta}^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu})}{\sqrt{\boldsymbol{\delta}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\delta}}}\right) \\
&\quad + \frac{1}{2}P\left(\frac{\boldsymbol{\delta}^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) - \boldsymbol{\delta}^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_0 - \boldsymbol{\mu})}{\sqrt{\boldsymbol{\delta}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\delta}}} \leq -\frac{\boldsymbol{\delta}^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_0 - \boldsymbol{\mu})}{\sqrt{\boldsymbol{\delta}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\delta}}}\right) \\
&= \frac{1}{2}P\left(Z > -\frac{\boldsymbol{\delta}^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu})}{\sqrt{\boldsymbol{\delta}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\delta}}}\right) + \frac{1}{2}P\left(Z \leq -\frac{\boldsymbol{\delta}^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_0 - \boldsymbol{\mu})}{\sqrt{\boldsymbol{\delta}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\delta}}}\right) \\
&= \Phi\left(-\frac{1}{2}\sqrt{\boldsymbol{\delta}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\delta}}\right),
\end{aligned}$$

where $Z \sim N(0, 1)$ and $\Phi(\cdot)$ is the c.d.f of $N(0, 1)$.

- (iii) The log-likelihood

$$\begin{aligned}
\ell(\boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}) &= -\frac{1}{2} \sum_{i=1}^{n_0} (\mathbf{X}_{0i} - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1} (\mathbf{X}_{0i} - \boldsymbol{\mu}_0) - (n_0 p / 2) \log(2\pi) - \frac{n_0}{2} \log |\boldsymbol{\Sigma}| \\
&\quad - \frac{1}{2} \sum_{i=1}^{n_1} (\mathbf{X}_{1i} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\mathbf{X}_{1i} - \boldsymbol{\mu}_1) - (n_1 p / 2) \log(2\pi) - \frac{n_1}{2} \log |\boldsymbol{\Sigma}|
\end{aligned}$$

We set

$$\frac{\partial \ell(\boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \boldsymbol{\Sigma})}{\partial \boldsymbol{\mu}_0} = \sum_{i=1}^{n_0} \boldsymbol{\Sigma}^{-1} (\mathbf{X}_{0i} - \boldsymbol{\mu}_0) = \mathbf{0}.$$

Then, $\hat{\boldsymbol{\mu}}_0 = \bar{\mathbf{X}}_0 = (1/n_0) \sum_{i=1}^{n_0} \mathbf{X}_{0i}$. Similarly, $\hat{\boldsymbol{\mu}}_1 = \bar{\mathbf{X}}_1 = (1/n_1) \sum_{i=1}^{n_1} \mathbf{X}_{1i}$. Note that,

$$\begin{aligned}
\sum_{i=1}^{n_0} (\mathbf{X}_{0i} - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1} (\mathbf{X}_{0i} - \boldsymbol{\mu}_0) &= \sum_{i=1}^{n_0} (\mathbf{X}_{0i} - \bar{\mathbf{X}}_0)^T \boldsymbol{\Sigma}^{-1} (\mathbf{X}_{0i} - \bar{\mathbf{X}}_0) + n_0 (\bar{\mathbf{X}}_0 - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{X}}_0 - \boldsymbol{\mu}_0) \\
&= n_0 \text{tr}\{\boldsymbol{\Sigma}^{-1}(\mathbf{S}_0 + \mathbf{d}_0 \mathbf{d}_0^T)\},
\end{aligned}$$

where $\mathbf{S}_0 = (1/n_0) \sum_{i=1}^{n_0} (\mathbf{X}_{0i} - \bar{\mathbf{X}}_0)(\mathbf{X}_{0i} - \bar{\mathbf{X}}_0)^T$ and $\mathbf{d}_0 = \bar{\mathbf{X}}_0 - \boldsymbol{\mu}_0$. Let $\boldsymbol{\Omega} = \boldsymbol{\Sigma}^{-1}$, set

$$\frac{\partial \ell(\boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \boldsymbol{\Omega})}{\partial \boldsymbol{\Omega}} = -\frac{1}{2} n_0 (\mathbf{S}_0 + \mathbf{d}_0 \mathbf{d}_0^T) + \frac{1}{2} n_0 \boldsymbol{\Omega} - \frac{1}{2} n_1 (\mathbf{S}_1 + \mathbf{d}_1 \mathbf{d}_1^T) + \frac{1}{2} n_1 \boldsymbol{\Omega} = \mathbf{0}.$$

Insert $\hat{\boldsymbol{\mu}}_0$ and $\hat{\boldsymbol{\mu}}_1$ into the above equation, we have $\hat{\boldsymbol{\Sigma}} = (1/n)(n_0 \mathbf{S}_0 + n_1 \mathbf{S}_1)$, where $n = n_0 + n_1$, $\mathbf{S}_1 = (1/n_1) \sum_{i=1}^{n_1} (\mathbf{X}_{1i} - \bar{\mathbf{X}}_1)(\mathbf{X}_{1i} - \bar{\mathbf{X}}_1)^T$.

(iv) Use similar argument as in (ii), the misclassification rate is given by

$$\frac{1}{2}\Phi\left(\frac{\hat{\boldsymbol{\delta}}^T \hat{\boldsymbol{\Sigma}}^{-1}(\boldsymbol{\mu}_1 - \hat{\boldsymbol{\mu}})}{\sqrt{\hat{\boldsymbol{\delta}}^T \hat{\boldsymbol{\Sigma}}^{-1} \boldsymbol{\Sigma} \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\delta}}}}\right) + \frac{1}{2}\Phi\left(-\frac{\hat{\boldsymbol{\delta}}^T \hat{\boldsymbol{\Sigma}}^{-1}(\boldsymbol{\mu}_0 - \hat{\boldsymbol{\mu}})}{\sqrt{\hat{\boldsymbol{\delta}}^T \hat{\boldsymbol{\Sigma}}^{-1} \boldsymbol{\Sigma} \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\delta}}}}\right),$$

where $\hat{\boldsymbol{\delta}} = \hat{\boldsymbol{\mu}}_0 - \hat{\boldsymbol{\mu}}_1$ and $\hat{\boldsymbol{\mu}} = (\hat{\boldsymbol{\mu}}_0 + \hat{\boldsymbol{\mu}}_1)/2$.

(v) Suppose $\tilde{\boldsymbol{\beta}}$ is the solution of $\boldsymbol{\beta}$ at the current iteration. Let $L(\boldsymbol{\beta}) = (1/2)\boldsymbol{\beta}^T \hat{\boldsymbol{\Sigma}} \boldsymbol{\beta} - \hat{\boldsymbol{\delta}}^T \boldsymbol{\beta}$, we have

$$\begin{aligned} L(\boldsymbol{\beta}) &= L(\tilde{\boldsymbol{\beta}}) + \{\nabla L(\tilde{\boldsymbol{\beta}})\}^T (\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}) + \frac{1}{2}(\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}})^T \nabla^2 L(\tilde{\boldsymbol{\beta}}) (\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}) \\ &= L(\tilde{\boldsymbol{\beta}}) + (\hat{\boldsymbol{\Sigma}} \tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\delta}})^T (\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}) + \frac{1}{2}(\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}})^T \hat{\boldsymbol{\Sigma}} (\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}) \\ &\leq L(\tilde{\boldsymbol{\beta}}) + (\hat{\boldsymbol{\Sigma}} \tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\delta}})^T (\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}) + \frac{c}{2} \|\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}\|_2^2, \end{aligned}$$

where $c \geq \lambda_{\max}(\hat{\boldsymbol{\Sigma}})$ and $\tilde{\boldsymbol{\beta}}$ is a vector on the line segment connecting $\boldsymbol{\beta}$ and $\tilde{\boldsymbol{\beta}}$. Therefore, we solve

$$\tilde{\boldsymbol{\beta}}^{(\text{new})} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \frac{c}{2} \|\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}\|_2^2 + (\hat{\boldsymbol{\Sigma}} \tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\delta}})^T (\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}) + \lambda \sum_{j=1}^p |\beta_j|.$$

The KKT condition is given by

$$\begin{aligned} c(\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}})_j + (\hat{\boldsymbol{\Sigma}} \tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\delta}})_j + \operatorname{sign}(\beta_j) &= 0, \text{ for } \beta_j \neq 0; \\ |c(\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}})_j + (\hat{\boldsymbol{\Sigma}} \tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\delta}})_j| &< \lambda, \text{ for } \beta_j = 0. \end{aligned}$$

Hence, the solution is given by $\tilde{\beta}_j^{(\text{new})} = s(\tilde{\beta}_j - (1/c)(\hat{\boldsymbol{\Sigma}} \tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\delta}})_j, \lambda/c)$, where $s(x, \lambda) = \operatorname{sign}(x)(|x| - \lambda)_+$ is the soft-thresholding function. With stepsize being $1/c$ for any $c \geq \lambda_{\max}(\hat{\boldsymbol{\Sigma}})$, the algorithm is summarized as follows.

Step 1: Initialize $\boldsymbol{\beta}$ at $\boldsymbol{\beta}^{(0)}$.

Step 2: At the k th iteration, let

$$\boldsymbol{\beta}^{(k)} = s\left(\boldsymbol{\beta}^{(k-1)} - \frac{1}{c}(\hat{\boldsymbol{\Sigma}} \boldsymbol{\beta}^{(k-1)} - \hat{\boldsymbol{\delta}}), \frac{\lambda}{c}\right)$$

where s is the soft-thresholding function defined above.

Update the gradient vector $\hat{\boldsymbol{\Sigma}} \boldsymbol{\beta}^{(k-1)} - \hat{\boldsymbol{\delta}}$ with $\hat{\boldsymbol{\Sigma}} \boldsymbol{\beta}^{(k)} - \hat{\boldsymbol{\delta}}$.

Step 3: Iterate until convergence.

(vi) Similar as in (iv), the misclassification rate

$$R_n = \frac{1}{2}\Phi\left(\frac{\hat{\boldsymbol{\beta}}^T(\boldsymbol{\mu}_1 - \hat{\boldsymbol{\mu}})}{\sqrt{\hat{\boldsymbol{\beta}}^T \boldsymbol{\Sigma} \hat{\boldsymbol{\beta}}}}\right) + \frac{1}{2}\Phi\left(-\frac{\hat{\boldsymbol{\beta}}^T(\boldsymbol{\mu}_0 - \hat{\boldsymbol{\mu}})}{\sqrt{\hat{\boldsymbol{\beta}}^T \boldsymbol{\Sigma} \hat{\boldsymbol{\beta}}}}\right).$$

By Law of Large Numbers, $\hat{\boldsymbol{\mu}} \xrightarrow{P} \boldsymbol{\mu}$. This together with $\hat{\boldsymbol{\beta}} \xrightarrow{P} \boldsymbol{\beta}$ imply $R_n \xrightarrow{P} R^*$.