

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

- (a) (4 points) Derive the Bayes rule for classifying a new observation  $\mathbf{x} \in \mathcal{R}^p$ .
- (b) (4 points) Derive the misclassification rate  $R^*$  of the Bayes rule.
- (c) (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  $(\hat{\boldsymbol{\mu}}_0, \hat{\boldsymbol{\mu}}_1, \hat{\Sigma})$  of  $(\boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \Sigma)$ .
- (d) (4 points) If we replace  $(\boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \Sigma)$  in the Bayes rule with  $(\hat{\boldsymbol{\mu}}_0, \hat{\boldsymbol{\mu}}_1, \hat{\Sigma})$ , 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

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

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

- (e) (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

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

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.

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

You may use the following facts: (i) The density of  $N_p(\boldsymbol{\mu}, \Sigma)$  is  $\{(2\pi)^p |\Sigma|\}^{-1/2} \exp\{-(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})/2\}$ ; (ii) For symmetric matrices  $\mathbf{A}$  and  $\mathbf{M}$ ,

$$\frac{\partial \operatorname{tr}(\mathbf{AM})}{\partial \mathbf{M}} = \frac{\partial \operatorname{tr}(\mathbf{MA})}{\partial \mathbf{M}} = \mathbf{A}, \quad \frac{\partial \log |\mathbf{M}|}{\partial \mathbf{M}} = \mathbf{M}^{-1}.$$

2. Consider a binary classification problem

2018 Theory 1

$\theta \in \{0, 1\}$  denotes class label

$$X | \theta=0 \sim N_p(\mu_0, \Sigma) \quad X | \theta=1 \sim N_p(\mu_1, \Sigma) \quad \Sigma \text{ is positive definite}$$

Suppose 0-1 loss is used & the prior distribution of  $\theta$  is  $P(\theta=0) = P(\theta=1) = 1/2$

(a) Derive the Bayes rule for classifying a new observation  $x \in \mathbb{R}^p$

Under 0-1 loss, the Bayes rule is the posterior MODE.

$$\therefore \theta=0 \text{ if } f(\theta=0|x) \geq f(\theta=1|x)$$

$$\theta=1 \text{ if } f(\theta=1|x) > f(\theta=0|x)$$

$$P(\theta|x) \propto P(x|\theta) \pi(\theta)$$

$$f(\theta=1|x) = P(x|\theta=1) \pi(\theta=1) = \{(2\pi)^p |\Sigma|^{1/2} \exp\{-\frac{1}{2}(x-\mu_1)' \Sigma^{-1} (x-\mu_1)\}\} \cdot \frac{1}{2}$$

$$f(\theta=0|x) = P(x|\theta=0) \pi(\theta=0) = \{(2\pi)^p |\Sigma|^{1/2} \exp\{-\frac{1}{2}(x-\mu_0)' \Sigma^{-1} (x-\mu_0)\}\} \cdot \frac{1}{2}$$

$$f(\theta=0|x) \geq f(\theta=1|x) \Rightarrow \frac{f(\theta=0|x)}{f(\theta=1|x)} \geq 1 \Rightarrow \log\left(\frac{f(\theta=0|x)}{f(\theta=1|x)}\right) \geq 0$$

$$\frac{f(\theta=0|x)}{f(\theta=1|x)} = \frac{\exp\{-\frac{1}{2}(x-\mu_0)' \Sigma^{-1} (x-\mu_0)\}}{\exp\{-\frac{1}{2}(x-\mu_1)' \Sigma^{-1} (x-\mu_1)\}}$$

$$\log\left(\frac{f(\theta=0|x)}{f(\theta=1|x)}\right) = -\frac{1}{2}(x-\mu_0)' \Sigma^{-1} (x-\mu_0) + \frac{1}{2}(x-\mu_1)' \Sigma^{-1} (x-\mu_1) \geq 0$$

$$(x-\mu_1)' \Sigma^{-1} (x-\mu_1) - (x-\mu_0)' \Sigma^{-1} (x-\mu_0) \geq 0$$

$$x' \cancel{\Sigma^{-1} x} - x' \Sigma^{-1} \mu_1 - \mu_1' \Sigma^{-1} x + \mu_1' \Sigma^{-1} \mu_1 - x' \cancel{\Sigma^{-1} x} + x' \Sigma^{-1} \mu_0 + \mu_0' \Sigma^{-1} x - \mu_0' \Sigma^{-1} \mu_0 \geq 0$$

$$\mu_1' \Sigma^{-1} \mu_1 - \mu_0' \Sigma^{-1} \mu_0 + 2x' \Sigma^{-1} (\mu_1 - \mu_0) \geq 0$$

$$(\mu_1 + \mu_0)' \Sigma^{-1} (\mu_1 - \mu_0) + 2x' \Sigma^{-1} (\mu_1 - \mu_0) \geq 0$$

$$[(\mu_1 + \mu_0)' - 2x'] \Sigma^{-1} (\mu_1 - \mu_0) \geq 0$$

$$[(\frac{\mu_1 + \mu_0}{2})' - x'] \Sigma^{-1} (\mu_1 - \mu_0) \geq 0$$

$$[x' - (\frac{\mu_1 + \mu_0}{2})'] \Sigma^{-1} (\mu_1 - \mu_0) \geq 0$$

$$(\mu_1 - \mu_0)' \Sigma^{-1} (x - \frac{\mu_1 + \mu_0}{2}) \geq 0$$

$$\text{let } \delta = (\mu_1 - \mu_0)$$

$$\frac{\delta}{2} = \bar{\mu} \quad \delta' \Sigma^{-1} (x - \bar{\mu}) \geq 0$$

∴ the Bayes Rule is

$$\begin{cases} \theta=0 & \text{if } \delta' \Sigma^{-1} (x - \bar{\mu}) \geq 0 \\ \theta=1 & \text{if } \delta' \Sigma^{-1} (x - \bar{\mu}) < 0 \end{cases}$$

2.1(b) Derive the misclassification rate  $R^*$  of the Bayes Rule

2018 Theory 1

$$R^* = \left[ P(S' \Sigma^{-1} (x - \bar{\mu}) \geq 0 | \theta=1) + P(S' \Sigma^{-1} (x - \bar{\mu}) < 0 | \theta=0) \right] \frac{1}{2}$$

$$\begin{aligned} P(S' \Sigma^{-1} (x - \mu) \geq 0 | \theta=1) &= P(S' \Sigma^{-1} (x - \frac{\mu_0 + \mu_1}{2}) \geq 0 | \theta=1) && \text{when } \theta=1, x - \mu_i \sim N(0, \Sigma) \\ &= P(S' \Sigma^{-1} (x - \mu_1 - \frac{\mu_0 - \mu_1}{2}) \geq 0 | \theta=1) && S' \Sigma^{-1} (x - \mu_1) \sim N(0, S' \Sigma^{-1} S) \\ &= P(S' \Sigma^{-1} (x - \mu_1 - \frac{S}{2}) \geq 0 | \theta=1) && \frac{S' \Sigma^{-1} (x - \mu_1)}{\sqrt{S' \Sigma^{-1} S}} \sim N(0, 1) \\ &= P(S' \Sigma^{-1} (x - \mu_1) - \frac{1}{2} S' \Sigma^{-1} S > 0 | \theta=1) \\ &= P(S' \Sigma^{-1} (x - \mu_1) > \frac{1}{2} S' \Sigma^{-1} S | \theta=1) && \text{let } Z \equiv N(0, 1) \\ &= P\left(\frac{S' \Sigma^{-1} (x - \mu_1)}{\sqrt{S' \Sigma^{-1} S}} > \frac{S' \Sigma^{-1} S}{2\sqrt{S' \Sigma^{-1} S}} | \theta=1\right) \\ &= P\left(Z > \frac{1}{2}\sqrt{S' \Sigma^{-1} S}\right) = \Phi\left(-\frac{1}{2}\sqrt{S' \Sigma^{-1} S}\right) \end{aligned}$$

let  $\Phi(\cdot)$  be the CDF of a Standard Normal.

$$\begin{aligned} P(S' \Sigma^{-1} (x - \mu) < 0 | \theta=0) &= P(S' \Sigma^{-1} (x - \frac{\mu_0 + \mu_1}{2}) < 0 | \theta=0) && \text{similarly when } \theta=0 \\ &= P(S' \Sigma^{-1} (x - \mu_0 + \frac{S}{2}) < 0 | \theta=0) && x - \mu_0 \sim N(0, \Sigma) \\ &= P(S' \Sigma^{-1} (x - \mu_0) + \frac{S' \Sigma^{-1} S}{2} < 0 | \theta=0) && \frac{S' \Sigma^{-1} (x - \mu_0)}{\sqrt{S' \Sigma^{-1} S}} \sim N(0, 1) \\ &= P(S' \Sigma^{-1} (x - \mu_0) < -\frac{1}{2} S' \Sigma^{-1} S | \theta=0) \\ &= P(Z < -\frac{1}{2}\sqrt{S' \Sigma^{-1} S}) \\ &= \Phi\left(-\frac{1}{2}\sqrt{S' \Sigma^{-1} S}\right) \end{aligned}$$

$$\therefore R^* = \frac{1}{2} \Phi\left(-\frac{1}{2}\sqrt{S' \Sigma^{-1} S}\right) + \frac{1}{2} \Phi\left(-\frac{1}{2}\sqrt{S' \Sigma^{-1} S}\right) = \Phi\left(-\frac{1}{2}\sqrt{S' \Sigma^{-1} S}\right)$$

2. (c) Let  $X_{0i}$  ( $i=1, \dots, n_0$ ) be iid samples from  $\theta=0$

and  $X_{1i}$  ( $i=1, \dots, n_1$ ) be iid samples from  $\theta=1$   $X_{0i} \perp\!\!\!\perp X_{1i}$

Derive the MLES  $(\hat{\mu}_0, \hat{\mu}_1, \hat{\Sigma})$  of  $(\mu_0, \mu_1, \Sigma)$

Both involve  $\Sigma$ , so we need the joint likelihood of  $(\bar{X}_0, \bar{X}_1)$

$$L(\mu_0, \mu_1, \Sigma) = \prod_{i=1}^{n_0} \left\{ (2\pi)^{-p/2} |\Sigma|^{-1/2} \exp \left\{ -\frac{1}{2} (x_{0i} - \mu_0) \Sigma^{-1} (x_{0i} - \mu_0) \right\} \right\} \prod_{j=1}^{n_1} \left\{ (2\pi)^{-p/2} |\Sigma|^{-1/2} \exp \left\{ -\frac{1}{2} (x_{1j} - \mu_1) \Sigma^{-1} (x_{1j} - \mu_1) \right\} \right\}$$

$$\begin{aligned} l(\mu_0, \mu_1, \Sigma) &= -\frac{n_0}{2} [p \log(2\pi) + \log(|\Sigma|)] - \frac{1}{2} \sum_{i=1}^{n_0} (x_{0i} - \mu_0) \Sigma^{-1} (x_{0i} - \mu_0) \\ &\quad - \frac{n_1}{2} [p \log(2\pi) + \log(|\Sigma|)] - \frac{1}{2} \sum_{j=1}^{n_1} (x_{1j} - \mu_1) \Sigma^{-1} (x_{1j} - \mu_1) \end{aligned}$$

$$\frac{\partial l(\mu_0, \mu_1, \Sigma)}{\partial \mu_0} = \sum_{i=1}^{n_0} \Sigma^{-1} (x_{0i} - \mu_0) \stackrel{\text{set}}{=} 0 \quad \Rightarrow \quad \sum_{i=1}^{n_0} x_{0i} - n_0 \mu_0 = 0 \Rightarrow \hat{\mu}_0 = \frac{1}{n_0} \sum_{i=1}^{n_0} x_{0i} \stackrel{\text{let}}{=} \bar{x}_0$$

$$\sum_{i=1}^{n_0} (x_{0i} - \mu_0) = 0$$

$$\text{similarly, } \hat{\mu}_1 = \frac{1}{n_1} \sum_{j=1}^{n_1} x_{1j} \stackrel{\text{let}}{=} \bar{x}_1$$

$$\text{let } \Omega = \Sigma^{-1}$$

$$\begin{aligned} \frac{\partial l(\mu_0, \mu_1, \Sigma)}{\partial \Omega} &= \frac{n_0}{2} \Sigma + \frac{n_1}{2} \Sigma - \frac{1}{2} \sum_{i=1}^{n_0} (x_{0i} - \mu_0)(x_{0i} - \mu_0)' - \frac{1}{2} \sum_{j=1}^{n_1} (x_{1j} - \mu_1)(x_{1j} - \mu_1)' \stackrel{\text{set}}{=} 0 \\ &= \frac{n_0 + n_1}{2} \Sigma - \frac{1}{2} \left( \sum_{i=1}^{n_0} (x_{0i} - \hat{\mu}_0)(x_{0i} - \hat{\mu}_0)' + \sum_{j=1}^{n_1} (x_{1j} - \hat{\mu}_1)(x_{1j} - \hat{\mu}_1)' \right) \stackrel{\text{set}}{=} 0 \\ \Rightarrow \hat{\Sigma} &= \frac{1}{n_0 + n_1} \left( \sum_{i=1}^{n_0} (x_{0i} - \hat{\mu}_0)(x_{0i} - \hat{\mu}_0)' + \sum_{j=1}^{n_1} (x_{1j} - \hat{\mu}_1)(x_{1j} - \hat{\mu}_1)' \right) \end{aligned}$$

Derivation of MVN MLES  $\hat{\Sigma}$  using  $\text{tr}(\cdot)$ . redo this!

2(d) If we replace  $(\mu_0, \mu_1, \Sigma)$  in the Bayes Rule w/  $(\hat{\mu}_0, \hat{\mu}_1, \hat{\Sigma})$  prove the misclassification rate of the resulting rule is:

$$\frac{1}{2} \Phi \left( \frac{\hat{s}^T \hat{\Sigma}^{-1} (\mu_1 - \hat{\mu})}{\sqrt{\hat{s}^T \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} \hat{s}}} \right) + \frac{1}{2} \Phi \left( -\frac{\hat{s}^T \hat{\Sigma}^{-1} (\mu_0 - \hat{\mu})}{\sqrt{\hat{s}^T \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} \hat{s}}} \right) \quad \text{where } \hat{s} = \hat{\mu}_0 - \hat{\mu}_1 \text{ and } \hat{\mu} = \frac{\hat{\mu}_0 + \hat{\mu}_1}{2}$$

$$\text{let } \tilde{R} = \frac{1}{2} P(\hat{s}' \hat{\Sigma}^{-1} (x - \hat{\mu}) \geq 0 | \theta=1) + \frac{1}{2} P(\hat{s}' \hat{\Sigma}^{-1} (x - \hat{\mu}) < 0 | \theta=0)$$

$$P(\hat{s}' \hat{\Sigma}^{-1} (x - \hat{\mu}) \geq 0 | \theta=1) = P(\hat{s}' \hat{\Sigma}^{-1} (x - \mu_1 + (\mu_1 - \hat{\mu}_1)) \geq 0 | \theta=1) \quad \begin{matrix} \text{under } \theta=1 \\ x - \mu_1 \sim N(0, \Sigma) \end{matrix}$$

$$= P(\hat{s}' \hat{\Sigma}^{-1} (x - \mu_1) + \hat{s}' \hat{\Sigma}^{-1} (\mu_1 - \hat{\mu}_1) \geq 0 | \theta=1)$$

$$\hat{s}' \hat{\Sigma}^{-1} (x - \mu_1) \sim N(0, \hat{s}' \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} \hat{s})$$

$$= P(\hat{s}' \hat{\Sigma}^{-1} (x - \mu_1) \geq -\hat{s}' \hat{\Sigma}^{-1} (\mu_1 - \hat{\mu}_1) | \theta=1)$$

$$= P(z \geq -\frac{\hat{s}' \hat{\Sigma}^{-1} (\mu_1 - \hat{\mu}_1)}{\sqrt{\hat{s}' \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} \hat{s}}}) = \Phi \left( \frac{\hat{s}' \hat{\Sigma}^{-1} (\mu_1 - \hat{\mu}_1)}{\sqrt{\hat{s}' \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} \hat{s}}} \right)$$

$$P(\hat{s}' \hat{\Sigma}^{-1} (x - \hat{\mu}) < 0 | \theta=0) = P(\hat{s}' \hat{\Sigma}^{-1} (x - \mu_0 + \mu_0 - \hat{\mu}) < 0 | \theta=0)$$

$$\text{under } \theta=0$$

$$x - \mu_0 \sim N(0, \Sigma)$$

$$\hat{s}' \hat{\Sigma}^{-1} (x - \mu_0) \sim N(0, \hat{s}' \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} \hat{s})$$

$$= P(\hat{s}' \hat{\Sigma}^{-1} (x - \mu_0) + \hat{s}' \hat{\Sigma}^{-1} (\mu_0 - \hat{\mu}) < 0 | \theta=0)$$

$$= P(\hat{s}' \hat{\Sigma}^{-1} (x - \mu_0) < -\hat{s}' \hat{\Sigma}^{-1} (\mu_0 - \hat{\mu}) | \theta=0)$$

$$= P(z < -\frac{\hat{s}' \hat{\Sigma}^{-1} (\mu_0 - \hat{\mu})}{\sqrt{\hat{s}' \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} \hat{s}}}) = \Phi \left( -\frac{\hat{s}' \hat{\Sigma}^{-1} (\mu_0 - \hat{\mu})}{\sqrt{\hat{s}' \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} \hat{s}}} \right)$$

$$\therefore \tilde{R} = \frac{1}{2} \Phi \left( \frac{\hat{s}' \hat{\Sigma}^{-1} (\mu_1 - \hat{\mu})}{\sqrt{\hat{s}' \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} \hat{s}}} \right) + \frac{1}{2} \Phi \left( -\frac{\hat{s}' \hat{\Sigma}^{-1} (\mu_0 - \hat{\mu})}{\sqrt{\hat{s}' \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} \hat{s}}} \right)$$

2(e), We propose another classification rule that assigns  $x$  to  $\theta=0$  if & only if  $\hat{\beta}^T(x-\hat{\mu}) \geq 0$  where  $\hat{\mu} = (\frac{\hat{\mu}_0 + \hat{\mu}_1}{2})$  and  $\hat{\beta}$  solves

2018 Theory P

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^P} \frac{1}{2} \beta' \hat{\Sigma} \beta - (\hat{\mu}_0 - \hat{\mu}_1)' \beta + \lambda \sum_{j=1}^P |\beta_j|$$

Derive the majorization-minimization algorithm for solving  $\hat{\beta}$ .

as the penalty is  $\lambda \sum_{j=1}^P |\beta_j|$  we know we will have a soft-threshold in our soln.

$$\begin{aligned} l(\beta) &= l(\tilde{\beta}) + \nabla l(\tilde{\beta})'(\beta - \tilde{\beta}) + \frac{1}{2}(\beta - \tilde{\beta}) \nabla^2 l(\tilde{\beta})(\beta - \tilde{\beta}) \\ &\leq l(\tilde{\beta}) + \nabla l(\tilde{\beta})'(\beta - \tilde{\beta}) + \frac{c}{2}(\beta - \tilde{\beta})(\beta - \tilde{\beta}) \end{aligned}$$

$$\text{where } c \geq \sup_{\beta} \lambda \max(\nabla^2 l(\beta))$$

$$\text{let } l(\beta) = \frac{1}{2} \beta' \hat{\Sigma} \beta - (\hat{\mu}_0 - \hat{\mu}_1)' \beta$$

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

$$\therefore l(\beta) = l(\tilde{\beta}) + (\hat{\Sigma} \beta - (\hat{\mu}_0 - \hat{\mu}_1))'(\beta - \tilde{\beta}) + \frac{c}{2}(\beta - \tilde{\beta})(\beta - \tilde{\beta})$$

$$\therefore \tilde{\beta}^{(\text{new})} = \arg \min_{\beta \in \mathbb{R}^P} \frac{c}{2} \|\beta - \tilde{\beta}\|_2^2 + \hat{\Sigma} \beta - (\hat{\mu}_0 - \hat{\mu}_1) + \lambda \sum_{j=1}^P |\beta_j|$$

The KKT condition is given by:

$$\left. \begin{aligned} c(\beta - \tilde{\beta})_j + (\hat{\Sigma} \beta - (\hat{\mu}_0 - \hat{\mu}_1))_j + \text{sign}(\beta_j) &= 0 \quad \text{for } \beta_j \neq 0 \\ |c(\beta - \tilde{\beta})_j + (\hat{\Sigma} \beta - (\hat{\mu}_0 - \hat{\mu}_1))_j| &\leq \lambda \quad \text{for } \beta_j = 0 \end{aligned} \right\} \begin{matrix} \text{maybe not write out} \\ \text{on exam} \end{matrix}$$

$$\Rightarrow \tilde{\beta}_j^{(\text{new})} = s(\tilde{\beta}_j - \frac{1}{c}(\hat{\Sigma} \beta - (\hat{\mu}_0 - \hat{\mu}_1)), \lambda/c) \quad \text{with stepsize } 1/c. \quad c \geq \sup_{\beta} \lambda \max(\nabla^2 l(\beta))$$

The algorithm is as follows

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

2. At the  $k^{\text{th}}$  iteration, let  $\beta^{(k)} = s(\beta^{(k-1)} - \frac{1}{c}(\hat{\Sigma} \beta^{(k-1)} - \hat{\beta}), \lambda/c)$

3. Update  $\nabla l(\beta)$  with  $\nabla l(\beta^{(k)})$

Iterate steps 2 & 3 until convergence.

2(f). Let  $R_n$  denote the misclassification rate of (e)

2018 Theory I

Suppose we can show  $\hat{\beta} \xrightarrow{P} \Sigma^{-1}(\mu_0 - \mu_1)$  as  $n \rightarrow \infty$ . Use this result to prove  $R_n \xrightarrow{P} R^*$

$$R_n = \frac{1}{2} P(\hat{\beta}'(x - \hat{\mu}) \geq 0 | \theta=1) + \frac{1}{2} P(\hat{\beta}'(x - \hat{\mu}) < 0 | \theta=0)$$

$$P(\hat{\beta}'(x - \hat{\mu}) \geq 0 | \theta=1) = P(\hat{\beta}'(x - \mu_1 + \mu_1 - \hat{\mu}) \geq 0 | \theta=1)$$

$$= P(\hat{\beta}'(x - \mu_1) + \hat{\beta}'(\mu_1 - \hat{\mu}) \geq 0 | \theta=1)$$

$$= P(\hat{\beta}'(x - \mu_1) \geq -\hat{\beta}'(\mu_1 - \hat{\mu}) | \theta=1)$$

under  $\theta=0 \quad x - \mu_1 \sim N(0, \Sigma)$

$\hat{\beta}'(x - \mu_1) \sim N(0, \hat{\beta}' \Sigma \hat{\beta})$

$$= P\left(Z \geq -\frac{\hat{\beta}'(\mu_1 - \hat{\mu})}{\sqrt{\hat{\beta}' \Sigma \hat{\beta}}}\right) = \Phi\left(-\frac{\hat{\beta}'(\mu_1 - \hat{\mu})}{\sqrt{\hat{\beta}' \Sigma \hat{\beta}}}\right)$$

$\frac{\hat{\beta}'(x - \mu_1)}{\sqrt{\hat{\beta}' \Sigma \hat{\beta}}} \sim N(0, 1)$

$$\text{Similarly } P(\hat{\beta}'(x - \hat{\mu}) < 0 | \theta=0) = \Phi\left(-\frac{\hat{\beta}'(\mu_0 - \hat{\mu})}{\sqrt{\hat{\beta}' \Sigma \hat{\beta}}}\right)$$

$$R_n = \frac{1}{2} \Phi\left(\frac{\hat{\beta}'(\mu_1 - \hat{\mu})}{\sqrt{\hat{\beta}' \Sigma \hat{\beta}}}\right) + \frac{1}{2} \Phi\left(-\frac{\hat{\beta}'(\mu_0 - \hat{\mu})}{\sqrt{\hat{\beta}' \Sigma \hat{\beta}}}\right)$$

$$\text{recall } R^* = \Phi\left(-\frac{1}{2} \sqrt{S' \Sigma S}\right)$$

$$\equiv \Phi\left(-\frac{1}{2} \sqrt{(\mu_0 - \mu_1) \Sigma^{-1} (\mu_0 - \mu_1)}\right)$$

By LLN  $\hat{\mu} \xrightarrow{P} \mu \quad \hat{\beta}' \xrightarrow{P} (\mu_0 - \mu_1)' \Sigma^{-1}$

$$R_n \rightarrow \frac{1}{2} \Phi\left(\frac{(\mu_0 - \mu_1)' \Sigma^{-1} (\mu_1 - \mu)}{\sqrt{S' \Sigma S}}\right) + \frac{1}{2} \Phi\left(-\frac{(\mu_0 - \mu_1)' \Sigma^{-1} (\mu_0 - \mu)}{\sqrt{S' \Sigma S}}\right)$$

$$\frac{1}{2} \Phi\left(\frac{\frac{1}{2} S' \Sigma S}{\sqrt{S' \Sigma S}}\right) + \frac{1}{2} \Phi\left(-\frac{\frac{1}{2} S' \Sigma S}{\sqrt{S' \Sigma S}}\right) = \Phi\left(-\frac{1}{2} \sqrt{S' \Sigma S}\right)$$

$$R_n \xrightarrow{P} R^*$$

(a) Use lots of conditioning to get through it

Factorial moments  $\stackrel{(1)}{n} \stackrel{(2)}{n(n-1)} \stackrel{(3)}{n(n-1)(n-2)}$

possibly use the mgf?

$X_i^2 | N$  is  $\chi^2$

(b) This works for all modes of convergence

$X_n \xrightarrow{\text{a.s.}} X \quad Y_n \xrightarrow{\text{a.s.}} Y$   $g(x, y)$  is cont on a set  $A: P_{x,y}(A)=1$   
 $\therefore$  the function is continuous almost everywhere  
 WRT  $x$  and  $y$

$\therefore$  By CMT,  $g(X_n, Y_n) \xrightarrow{\text{a.s.}} g(X, Y)$

$\frac{X_n}{Y_n} \xrightarrow{\text{a.s.}} \frac{X}{Y}$  when  $P(Y=0)=0$

$X_n \xrightarrow{\text{d}} X \quad Y_n \xrightarrow{\text{d}} Y$  if  $X$  or  $Y$  are constant  
 then  $\begin{pmatrix} X_n \\ Y_n \end{pmatrix} \xrightarrow{d} \begin{pmatrix} X \\ Y \end{pmatrix}$

(c) make sure when using CLT state  
 the 2nd moment is bounded.

(e) to show  $\hat{p}^2$  is consistent, use CMT.

note, we want an asymptotic pivot "classic Slutsky's"

$$\frac{\sqrt{n}(\hat{W}_n - \sigma^2)}{\hat{P}^n} \xrightarrow{d} N(0, 1) \quad \text{as } \frac{\sqrt{n}(\hat{W}_n - \sigma^2)}{\sigma^2 \hat{P}}$$

$$\text{let } \hat{P}^2 = T_{2,22}^2$$

If  $X_n \rightarrow X$   $Y_n \rightarrow Y$  marginally,  
 then they converge jointly  
 automatically when  
 converging a.s. & in p.  
 But not necessarily in dist.  
 in distribution, marginal  
 conv.  $\not\Rightarrow$  joint convergence  
 unless independent or  $\oplus$

1. (25 points) Let  $N$  be Poisson distributed with parameter  $0 < \lambda < \infty$ , and let  $Z_1, Z_2, \dots$  be an i.i.d. sequence of  $N(0, \sigma^2)$  random variables, independent of  $N$ , with  $0 < \sigma^2 < \infty$ . Let

$$X = 1\{N > 0\} \sum_{j=1}^N Z_j,$$

where  $1\{A\}$  is the indicator of  $A$ . Let  $X_1, \dots, X_n$  be i.i.d. realizations of  $X$ , and let  $U_i = 1\{X_i = 0\}$ ,  $1 \leq i \leq n$ . Do the following:

- (a) (4 points) Show that  $EU_i = e^{-\lambda}$ ,  $EX_i = 0$ ,  $EX_i^2 = \lambda\sigma^2$ , and  $EX_i^4 = 3(\lambda + \lambda^2)\sigma^4$ .
- (b) (5 points) Show that  $\hat{T}_n = -\log(n^{-1} \sum_{i=1}^n U_i)$  is almost surely consistent for  $\lambda$ , and that  $\hat{W}_n = n^{-1} \sum_{i=1}^n X_i^2 / \hat{T}_n$  is almost surely consistent for  $\sigma^2$ , as  $n \rightarrow \infty$ .
- (c) (5 points) Show that

$$\sqrt{n} \begin{pmatrix} n^{-1} \sum_{i=1}^n U_i - e^{-\lambda} \\ n^{-1} \sum_{i=1}^n X_i^2 - \lambda\sigma^2 \end{pmatrix} \xrightarrow{d} N(0, \tau_1^2),$$

as  $n \rightarrow \infty$ , and give the form of  $\tau_1^2$ .

- (d) (6 points) Show that

$$\sqrt{n} \begin{pmatrix} \hat{T}_n - \lambda \\ \hat{W}_n - \sigma^2 \end{pmatrix} \xrightarrow{d} N(0, \tau_2^2),$$

as  $n \rightarrow \infty$ , where

$$\tau_2^2 = \begin{pmatrix} e^\lambda - 1 & -(e^\lambda - \lambda - 1)\sigma^2/\lambda \\ -(e^\lambda - \lambda - 1)\sigma^2/\lambda & \left(\frac{e^\lambda - 1}{\lambda^2} + 2 + \frac{1}{\lambda}\right)\sigma^4 \end{pmatrix}.$$

- (e) (5 points) Show that  $\hat{W}_n \pm z_{1-\alpha/2} \hat{\rho}_n / \sqrt{n}$ , where

$$\hat{\rho}^2 = \left( \frac{e^{\hat{T}_n} - 1}{\hat{T}_n^2} + 2 + \frac{1}{\hat{T}_n} \right) \hat{W}_n^2$$

and  $z_q$  is the  $q$ th-quantile of a standard normal, is an asymptotically valid  $1 - \alpha$  level confidence interval for  $\sigma^2$ .

Let  $N$  be Poisson( $\lambda$ ) and  $Z_1, Z_2, \dots \stackrel{iid}{\sim} N(0, \sigma^2)$   $0 < \lambda < \infty$   $0 < \sigma^2 < \infty$  [2018 Theory 1]

Let  $X = I\{N>0\} \sum_{j=1}^N Z_j$  and  $x_1, \dots, x_n$  be realizations of  $X$   
 let  $U_i = I\{x_i=0\}$   $1 \leq i \leq n$

$$(a) \text{ Show that } E[U_i] = e^{-\lambda}, E[X_i] = 0, E[X_i^2] = \lambda\sigma^2, \text{ and } E[X_i^4] = 3(\lambda + \lambda^2)\sigma^2$$

$$E[X_i] = E[I\{N>0\} \sum_{j=1}^N Z_j] = E[E[I\{N>0\} \sum_{j=1}^N Z_j | N]] = E[I\{N>0\} \sum_{j=1}^N E[Z_j | N]]$$

$$= E_N[I\{N>0\} \sum_{j=1}^N 0] = E_N[0] = 0$$

$$E[U_i] = E[I\{x_i=0\}] = P(x_i=0) = P(I\{N>0\} \sum_{j=1}^N Z_j = 0)$$

$$= P(I\{N>0\} = 0 \cup \sum_{j=1}^N Z_j = 0) = P(N=0) \cup P(\sum_{j=1}^N Z_j = 0)$$

$$\sum_{j=1}^N Z_j \sim N(0, N\sigma^2)$$

$$= f_N(0) + P(\sum_{j=1}^N Z_j = 0) - P(N=0 \text{ AND } \sum_{j=1}^N Z_j = 0)$$

$$= \frac{e^{-\lambda}}{0!} + 0 - 0 = e^{-\lambda}$$

$$\frac{1}{\sqrt{N\sigma^2}} \sum_{j=1}^N Z_j \sim N(0, 1)$$

$$X_i^2 = [I\{N>0\} \sum_{j=1}^N Z_j]^2 = I\{N>0\} \left(\sum_{j=1}^N Z_j\right)^2$$

$$E[X_i^2] = E[I\{N>0\} \left(\sum_{j=1}^N Z_j\right)^2] = E_N[E_{X|N}[I\{N>0\} \left(\sum_{j=1}^N Z_j\right)^2 | N]]$$

$$= E_N[I\{N>0\} E_{X|N}\left[\left(\sum_{j=1}^N Z_j\right)^2 | N\right]] = E_N[N\sigma^2] = \sigma^2 E[N] = \lambda\sigma^2$$

$$\frac{1}{N\sigma^2} (\sum_{j=1}^N Z_j)^2 \sim \chi^2_N$$

$$\text{var}(\frac{1}{N\sigma^2} \sum_{j=1}^N Z_j)^N = 2$$

$$\frac{1}{N^2\sigma^4} (\sum_{j=1}^N Z_j^2)^N = 2$$

$$\text{var}((\sum_{j=1}^N Z_j)^2) = 2N^2\sigma^4$$

$$X_i^4 = [I\{N>0\} \sum_{j=1}^N Z_j]^4 = I\{N>0\} \left(\sum_{j=1}^N Z_j\right)^4$$

$$E[X_i^4] = E[I\{N>0\} \left(\sum_{j=1}^N Z_j\right)^4] = E_N[E_{X|N}[I\{N>0\} \left(\sum_{j=1}^N Z_j\right)^4 | N]]$$

$$= E_N[I\{N>0\} E_{X|N}\left[\left(\sum_{j=1}^N Z_j\right)^4 | N\right]] = E_N[I\{N>0\} (\text{var}_{Z|N}(\sum_{j=1}^N Z_j^2) + E[(\sum_{j=1}^N Z_j^2)^2] | N)]$$

$$= E_N[2N^2(\sigma^2)^2 + (N\sigma^2)^2] = E_N[3N^2\sigma^4] = 3E_N[N^2]\sigma^4$$

$$= 3\sigma^4 [\text{var}(N) + E(N)^2] = 3\sigma^4 (\lambda + \lambda^2) \checkmark$$

1(b). Show that  $\hat{T}_n = -\log \left( \frac{1}{n} \sum_{i=1}^n U_i \right) \xrightarrow{\text{a.s.}} \lambda$

and  $\hat{W}_n = \frac{1}{n} \sum_{i=1}^n \frac{x_i^2}{\hat{T}_n} \xrightarrow{\text{a.s.}} \sigma^2$  as  $n \rightarrow \infty$

$U_i$  are iid w/ mean  $e^{-\lambda}$

$\therefore$  by SLLN,  $\frac{1}{n} \sum_{i=1}^n U_i \xrightarrow{\text{a.s.}} e^{-\lambda}$

$\therefore$  by CMT  $\log \left( \frac{1}{n} \sum_{i=1}^n U_i \right) \xrightarrow{\text{a.s.}} -\lambda$

$$-\log \left( \frac{1}{n} \sum_{i=1}^n U_i \right) \xrightarrow{\text{a.s.}} \lambda \Rightarrow \hat{T}_n \xrightarrow{\text{a.s.}} \lambda$$

$X_i^2$  are iid w/ mean  $\lambda \sigma^2$

$\therefore$  by SLLN  $\frac{1}{n} \sum_{i=1}^n X_i^2 \xrightarrow{\text{a.s.}} \lambda \sigma^2$

by CMT,  $\frac{1}{\hat{T}_n} \xrightarrow{\text{a.s.}} \frac{1}{\lambda}$

because marginal  $\Rightarrow$  joint convergence a.s.

$$\therefore \left( \frac{\frac{1}{n} \sum X_i^2}{\hat{T}_n} \right) \xrightarrow{\text{a.s.}} \left( \frac{\lambda \sigma^2}{\lambda} \right) \Rightarrow \text{by CMT}$$

~~$\therefore \hat{W}_n = \frac{\frac{1}{n} \sum X_i^2}{\hat{T}_n} \xrightarrow{\text{a.s.}} \frac{\lambda \sigma^2}{\lambda} = \sigma^2 \Rightarrow \hat{W}_n \xrightarrow{\text{a.s.}} \sigma^2$~~

when  $P(\lambda=0)=0$  which is true as  $0 < \lambda < \infty$

1(c). Show that  $\sqrt{n} \begin{pmatrix} \frac{1}{n} \sum u_i - e^{-\lambda} \\ \frac{1}{n} \sum x_i^2 - \lambda \sigma^2 \end{pmatrix} \xrightarrow{d} N(0, \Sigma_{T_1})$  as  $n \rightarrow \infty$ , give the form of  $\Sigma_{T_1}$

2018 Theory 1

by CLT we know  $\sqrt{n} \left( \frac{1}{n} \sum u_i - e^{-\lambda} \right) \xrightarrow{d} N(0, \text{var}(u_i))$

$\sqrt{n} \left( \frac{1}{n} \sum x_i^2 - \lambda \sigma^2 \right) \xrightarrow{d} N(0, \text{var}(x_i^2))$

$u_i \sim \text{Bernoulli}(e^{-\lambda})$

$$\text{var}(u_i) = E[u_i^2] - E[u_i]^2$$

$$= E[I(X_i=0)^2] - (e^{-\lambda})^2 = E[I(X_i=0)] - (e^{-\lambda})^2 = (e^{-\lambda}) - (e^{-\lambda})^2$$

$$= e^{-\lambda} - e^{-2\lambda} = e^{-\lambda}(1 - e^{-\lambda})$$

$$\text{var}(x_i^2) = E[x_i^4] - E[x_i^2]^2 = 3(\lambda + \lambda^2)\sigma^4 - (\lambda \sigma^2)^2$$

$$= 3\lambda\sigma^4 + 3\lambda^2\sigma^4 - \lambda^2\sigma^4 = 3\lambda\sigma^4 + 2\lambda^2\sigma^4 = \lambda\sigma^4(3 + 2\lambda)$$

$$\text{cov}(u_i, x_i^2) = E[u_i x_i^2] - E[u_i]E[x_i^2]$$

$E[u_i x_i^2] = E[I(X_i=0)x_i^2] = 0$  because the indicator will make any value of  $x_i$  to

$u_i x_i^2 = I(X_i=0)x_i^2 = \begin{cases} x_i^2 & \text{when } X_i=0 \\ 0 & \text{where } X_i \neq 0 \end{cases}$  in the product, so if values of  $X_i$ ,  $u_i x_i^2 = 0$

$$\therefore \text{cov}(u_i, x_i^2) = -E[u_i]E[x_i^2] = -e^{-\lambda}\lambda\sigma^2$$

$$\therefore \sqrt{n} \begin{pmatrix} \frac{1}{n} \sum u_i - e^{-\lambda} \\ \frac{1}{n} \sum x_i^2 - \lambda \sigma^2 \end{pmatrix} \xrightarrow{d} N(0, \Sigma_{T_1})$$

$$\text{where } \Sigma_{T_1} = \begin{pmatrix} \text{var}(u_i) & \text{cov}(u_i, x_i^2) \\ \text{cov}(u_i, x_i^2) & \text{var}(x_i^2) \end{pmatrix} = \begin{pmatrix} e^{-\lambda}(1 - e^{-\lambda}) & -e^{-\lambda}\lambda\sigma^2 \\ -e^{-\lambda}\lambda\sigma^2 & \lambda\sigma^4(3 + 2\lambda) \end{pmatrix}$$

I(d). Show that  $\sqrt{n} \begin{pmatrix} \hat{T}_n - \lambda \\ \hat{W}_n - \sigma^2 \end{pmatrix} \xrightarrow{d} N(0, \Sigma_2)$  as  $n \rightarrow \infty$

$$\text{let } \frac{1}{n} \sum U_i = \bar{U} \quad \hat{T}_n = -\log \bar{U}$$

$$\text{let } \frac{1}{n} \sum X_i^2 = \bar{X} \quad \hat{W}_n = \frac{\bar{X}}{-\log \bar{U}} \quad \text{let } g(a, b) = \begin{pmatrix} -\log(a) \\ -\frac{b}{\log(a)} \end{pmatrix}$$

$$\nabla g(a, b) = \begin{pmatrix} -\frac{1}{a} & 0 \\ \frac{b}{a \log(a)^2} & -\frac{1}{\log(a)} \end{pmatrix}$$

from (c) we know  $\sqrt{n} \begin{pmatrix} \bar{U} - e^{-\lambda} \\ \bar{X} - \lambda \sigma^2 \end{pmatrix} \xrightarrow{d} N(0, \Sigma_1)$

$\therefore$  by Delta method  $\sqrt{n} \begin{pmatrix} \hat{T}_n - \lambda \\ \hat{W}_n - \sigma^2 \end{pmatrix} \xrightarrow{d} \nabla g(e^{-\lambda}, \lambda \sigma^2) N(0, \Sigma_1)$   
 $\equiv N(0, \Sigma_2)$

$$\nabla g(e^{-\lambda}, \lambda \sigma^2) \begin{pmatrix} -e^{-\lambda} & 0 \\ \frac{\lambda \sigma^2}{e^{-\lambda} \lambda^2} & \frac{1}{\lambda} \end{pmatrix} = \begin{pmatrix} -e^{-\lambda} & 0 \\ \frac{e^{-\lambda} \sigma^2}{\lambda} & \frac{1}{\lambda} \end{pmatrix}$$

$$\Sigma_2 = \begin{pmatrix} -e^{-\lambda} & 0 \\ \frac{e^{-\lambda} \sigma^2}{\lambda} & \frac{1}{\lambda} \end{pmatrix} \begin{pmatrix} e^{-\lambda}(1-e^{-\lambda}) & -e^{-\lambda} \lambda \sigma^2 \\ -e^{-\lambda} \lambda \sigma^2 & \lambda \sigma^4(3+2\lambda) \end{pmatrix} \begin{pmatrix} -e^{-\lambda} & \frac{e^{-\lambda} \sigma^2}{\lambda} \\ 0 & \frac{1}{\lambda} \end{pmatrix}$$

$$= \begin{pmatrix} e^{-\lambda}-1 & \lambda \sigma^2 \\ \frac{1-e^{-\lambda} \sigma^2}{\lambda} - e^{-\lambda} \sigma^2 & \sigma^4 + \sigma^4(3+2\lambda) \end{pmatrix} \begin{pmatrix} -e^{-\lambda} & \frac{e^{-\lambda} \sigma^2}{\lambda} \\ 0 & \frac{1}{\lambda} \end{pmatrix}$$

$$= \begin{pmatrix} -(1-e^{-\lambda}) & (e^{-\lambda}-1) \frac{e^{-\lambda} \sigma^2}{\lambda} + \sigma^2 \\ \sigma^2 - e^{-\lambda} \left( \frac{1-e^{-\lambda} \sigma^2}{\lambda} \right) & \frac{e^{-\lambda} \sigma^2}{\lambda} \left( \frac{1-e^{-\lambda} \sigma^2}{\lambda} - e^{-\lambda} \sigma^2 \right) + \frac{1}{\lambda} (\sigma^4 + \sigma^4(3+2\lambda)) \end{pmatrix}$$

$$= \begin{pmatrix} e^{-\lambda}-1 & \frac{\sigma^2}{\lambda} - \frac{e^{-\lambda} \sigma^2}{\lambda} + \sigma^2 \\ \sigma^2 - \frac{e^{-\lambda} \sigma^2}{\lambda} + \frac{\sigma^2}{\lambda} & \frac{e^{-\lambda} \sigma^2 - \sigma^4}{\lambda^2} - \frac{\sigma^4}{\lambda} + \frac{\sigma^4}{\lambda} (1+3+2\lambda) \end{pmatrix} = \begin{pmatrix} e^{-\lambda}-1 & -\frac{\sigma^2}{\lambda} (e^{-\lambda}-1-\lambda) \\ -\frac{\sigma^2}{\lambda} (e^{-\lambda}-1-\lambda) & \end{pmatrix}$$

1.(e). Show that  $\hat{W}_n \pm Z_{1-\alpha/2} \hat{\rho}_n / \sqrt{n}$  is an asymptotically valid  $1-\alpha$  level CI for  $\hat{\sigma}^2$

Go back to Kushal's Solution