

**Q1**

(i). We have :  $X_1, \dots, X_n \sim \text{Unif}[0, \theta]$ .

$$L(X_1, \dots, X_n) = \prod_{j=1}^n \frac{1}{\theta} \mathbf{1}_{X_j \leq \theta} = \frac{1}{\theta^n} \mathbf{1}_{\max\{X_j\} \leq \theta}$$

$$\mathbf{E}[X] = \frac{\theta}{2}$$

So,  $\hat{\theta} = \max\{X_j\}$  and  $\tilde{\theta} = 2\bar{X}$ .

(ii). Clearly, for  $\mu_2 = \mathbf{E}[\tilde{\theta}]$  and  $\sigma_2^2 = \text{Var}[\tilde{\theta}]$

$$\begin{aligned} \hat{F}(x) &= \mathbb{P}\{\max\{X_j\} \leq x\} = \prod_{j=1}^n \mathbb{P}\{X_j \leq x\} \\ &= \left[\frac{x}{\theta}\right]^n \mathbf{1}_{\{0 < x < \theta\}} + \mathbf{1}_{\{x \geq \theta\}} \end{aligned}$$

And since,  $\mu_2 = \theta$  and  $\sigma_2^2 = \frac{1}{n^2} \left[n \frac{\theta^2}{12}\right] = \frac{\theta^2}{12n}$ , we can use the Central Limit Theorem to conclude that for large  $n$ ,

$$\begin{aligned} \tilde{F}(x) &= \mathbb{P}\{2\bar{X} \leq x\} = \mathbb{P}\left\{\frac{\bar{X} - \mu_2}{\sigma_2} \leq \sqrt{3} \frac{x - \theta}{\theta}\right\} \\ &\simeq \left[ \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\sqrt{3n}(\frac{x-\theta}{\theta})} e^{-x^2/2} dx \right] \mathbf{1}_{\{0 < x < \theta\}} + \mathbf{1}_{\{x \geq \theta\}} \end{aligned}$$

(iii). Clearly, for  $\mu_1 = \mathbf{E}[\hat{\theta}]$  and  $\sigma_1^2 = \text{Var}[\hat{\theta}]$

$$\begin{aligned} \mu_1 &= \int_0^M x \cdot \frac{n}{\theta} \left[\frac{x}{\theta}\right]^{n-1} dx = \frac{n}{n+1} \theta \\ \mu_1^2 + \sigma_1^2 &= \int_0^\theta x^2 \cdot \frac{n}{\theta} \left[\frac{x}{\theta}\right]^{n-1} dx = \frac{n}{n+2} \theta^2 \\ \sigma_1^2 &= \frac{n}{(n+1)^2(n+2)} \theta^2 \end{aligned}$$

So, the bias is  $\frac{\theta}{n+1}$  and the variance is as above.

(iv). Let  $\check{\theta} = (1 + \frac{1}{n}) \max\{X_j\}$ . Then  $\mathbf{E}[\check{\theta}] = \theta$  and  $\text{Var}[\check{\theta}] = \frac{1}{n(n+2)} \theta^2$ .

(v). We calculate the Fisher information as follows :

$$\mathcal{I}_n(\theta) = \mathbf{E} \left[ \left( \frac{\partial}{\partial \theta} \log \left[ \frac{nx^{n-1}}{\theta^n} \right] \right)^2 \middle| \theta \right] = \mathbf{E} \left[ \frac{n}{\theta^2} \middle| \theta \right] = \frac{n}{\theta^2}$$

So,

$$\sqrt{n}(\hat{\theta} - \theta) \rightarrow \mathcal{N}(0, \theta^2)$$

By delta method we get,

$$\sqrt{n}(\check{\theta} - \theta) \rightarrow \mathcal{N}(0, \theta^2)$$

**Q2**

(a)

$$\mathbb{E} \hat{\rho} = \frac{1}{n} \sum_{i=1}^n \mathbb{E} X_i Y_i = \mathbb{E} XY = \rho,$$

so we have unbiasedness. For consistency, define  $Z_i := X_i Y_i$ . Then  $\hat{\rho} = \frac{1}{n} \sum_{i=1}^n Z_i$ , thus by SLLN for the sample  $Z_1, \dots, Z_n$  we have  $\hat{\rho} \rightarrow \mathbb{E} Z_1 = \rho$ , where the convergence is almost sure, which gives us consistency.

(b) This depends on which formula for CLT students have in their cheat-sheet. I think it will be easy for them to read the formula wrong. For example,

$$\frac{\hat{\rho} - \rho}{\sqrt{\text{var}(\hat{\rho})}} \rightarrow \mathcal{N}(0, 1).$$

Using the hint,  $\text{var}(\hat{\rho}) = \frac{1}{n} \text{var}(X_1, Y_1) = \frac{1}{n} ((1 + 2\rho^2) - \rho^2) = \frac{1}{n} (1 + \rho^2)$ , leading to

$$\sqrt{n} \frac{\hat{\rho} - \rho}{\sqrt{1 + \rho^2}} \rightarrow \mathcal{N}(0, 1).$$

(c) Using the Slutsky's theorem, we can replace  $\rho$  in the denominator by  $\hat{\rho}$ , using the fact that  $\hat{\rho}$  is consistent (this is minimally sufficient justification). Then it is standard algebra to get the answer :

$$P(-q_{0.975} < \sqrt{n} \frac{\hat{\rho} - \rho}{\sqrt{1 + \hat{\rho}^2}} < -q_{0.975}) = 0.95$$

$$\dots CI_{0.95}(\rho) = (\hat{\rho} \pm q_{0.975} \frac{\sqrt{1 + \hat{\rho}^2}}{\sqrt{n}})$$

**Q3**

(a) Write  $F_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{[X_i \leq x]}$ . Then

$$\hat{f}(x) = \frac{1}{2nh} \sum_{i=1}^n \mathbf{1}_{[X_i \in (x-h, x+h)]} = \frac{1}{2nh} \sum_{i=1}^n \mathbf{1}_{[1 \leq \frac{x-X_i}{h} < h]} = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-X_i}{h}\right),$$

where  $K(y) = \frac{1}{2} \mathbf{1}_{[|y| \leq 1]}$ . So the kernel corresponds to  $U[-1, 1]$  distribution.

(b)

$$\mathbb{E} \hat{f}(x) = \frac{1}{2h} \left[ \mathbb{E} F_n(x+h) - \mathbb{E} F_n(x-h) \right] = \frac{1}{2h} \left[ F(x+h) - F(x-h) \right]$$

If  $h \rightarrow 0$ , then the previous expression is the definition for derivative of  $F$ , thus  $\mathbb{E} \hat{f}(x) \rightarrow f(x)$  for  $h \rightarrow 0$ . This leads to  $\text{bias}(\hat{f}(x)) \rightarrow 0$ .

Note that students will probably start by writing down the definition of bias. This leads to the same conclusion.

**Q4**

(a) As explained on slide 342 of the course, the uncentered coefficient of determination measures the proportion of the squared norm of  $\mathbf{Y}$  explained by the fitted values  $\hat{\mathbf{Y}}$ . As a ratio of two positive quantities,  $R^2$  is necessarily positive. Moreover, we have :

$$\mathbf{Y} = \mathbf{H}\mathbf{Y} + (\mathbf{I} - \mathbf{H})\mathbf{Y} = \hat{\mathbf{Y}} + (\mathbf{I} - \mathbf{H})\mathbf{Y}$$

where the hat matrix  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$  is an orthogonal projection operator. We can hence apply the Pythagorean theorem to obtain

$$\|\mathbf{Y}\|^2 = \|\hat{\mathbf{Y}}\|^2 + \|(\mathbf{I} - \mathbf{H})\mathbf{Y}\|^2,$$

and hence  $\|\mathbf{Y}\|^2 \geq \|\hat{\mathbf{Y}}\|^2 \Rightarrow R_0^2 \leq 1$ .

**Alternative solution :** Use spectral norm of  $\mathbf{H}$  :

$$\|\hat{\mathbf{Y}}\|^2 = \|\mathbf{H}\mathbf{Y}\|^2 \leq \|\mathbf{H}\|_2^2 \|\mathbf{Y}\|^2.$$

Since  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$  is an orthogonal projection it has maximal eigenvalue 1, yielding  $\|\mathbf{H}\|_2^2 = 1$  and finally  $R_0^2 \leq 1$ .

The inequality is saturated when  $p = n$  indeed, in this case  $\mathbf{X}$  is square and hence invertible from the full-rank assumption, yielding  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{X} \mathbf{X}^{-1} \mathbf{X}^{-T} \mathbf{X}^T = I$  and trivially  $\hat{\mathbf{Y}} = \mathbf{Y} \Rightarrow R_0^2 = 1$ .

(b) The new fitted values are given by

$$\begin{aligned} \tilde{\mathbf{Y}} &= \tilde{\mathbf{X}}(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{Y} \\ &= (\mathbf{X}, \mathbf{x}_{p+1}) \begin{pmatrix} (\mathbf{X}^T \mathbf{X})^{-1} & \mathbf{0} \\ 0 & 1/\|\mathbf{x}_{p+1}\|^2 \end{pmatrix} \begin{pmatrix} \mathbf{X}^T \\ \mathbf{x}_{p+1}^T \end{pmatrix} \mathbf{Y} \\ &= \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} + \frac{\mathbf{x}_{p+1}^T \mathbf{Y}}{\|\mathbf{x}_{p+1}\|^2} \mathbf{x}_{p+1} \\ &= \hat{\mathbf{Y}} + \frac{\mathbf{x}_{p+1}^T \mathbf{Y}}{\|\mathbf{x}_{p+1}\|^2} \mathbf{x}_{p+1}, \end{aligned}$$

since  $\mathbf{X}^T \mathbf{x}_{p+1} = 0$  by assumption. From Pythagore again we get

$$\|\tilde{\mathbf{Y}}\|^2 = \|\hat{\mathbf{Y}}\|^2 + |\mathbf{x}_{p+1}^T \mathbf{Y}|^2,$$

which yields :

$$\frac{\|\hat{\mathbf{Y}}\|^2}{\|\mathbf{Y}\|^2} \leq \frac{\|\tilde{\mathbf{Y}}\|^2}{\|\mathbf{Y}\|^2}.$$

A higher uncentered coefficient of determination  $R_0^2$  is not necessarily desirable. Indeed, when we add too many covariates to the model, we explain more and more of the variations in the data, including after a certain threshold the noise variations, which is of course undesirable. This phenomenon is called *overfitting*.

(c) We have  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$  and

$$\mathbf{Y} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n),$$

by assumption of the Gaussian linear model. From properties of Gaussian random vector, we know that  $\hat{\boldsymbol{\beta}}$  is again a Gaussian random vector (as linear transformation of a Gaussian rv), with exact distribution given by

$$\hat{\boldsymbol{\beta}} \sim \mathcal{N}(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}).$$

Similarly the distribution of  $\mathbf{c}^T \hat{\boldsymbol{\beta}}$  for some arbitrary vector  $\mathbf{c} \in \mathbb{R}^p$  is Gaussian and given by

$$\mathbf{c}^T \hat{\boldsymbol{\beta}} \sim \mathcal{N}(\mathbf{c}^T \boldsymbol{\beta}, \sigma^2 \mathbf{c}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{c}).$$

(d) Choosing  $\mathbf{c} = \mathbf{e}_i$  where  $\mathbf{e}_i$  is the  $i$ th canonical vector of  $\mathbb{R}^p$ , we get

$$\hat{\beta}_i = \mathbf{e}_i^T \hat{\boldsymbol{\beta}}.$$

From the previous question and the symmetry of the standard normal distribution we hence get that

$$\mathbb{P} \left\{ \frac{|\hat{\beta}_i - \beta_i|}{\sqrt{\sigma^2 (\mathbf{X}^T \mathbf{X})_{ii}^{-1}}} \leq \Phi_{\alpha/2} \right\} = 1 - \alpha,$$

where  $\Phi_{\alpha/2}$  denotes the quantile  $\alpha/2$  of the standard normal distribution. We get hence the following confidence intervals for each  $i = 1, \dots, p$ ,

$$IC_{1-\alpha}(\beta_i) = \left[ \hat{\beta}_i - \Phi_{\alpha/2} \sqrt{\sigma^2(\mathbf{X}^T \mathbf{X})_{ii}^{-1}}, \hat{\beta}_i + \Phi_{\alpha/2} \sqrt{\sigma^2(\mathbf{X}^T \mathbf{X})_{ii}^{-1}} \right].$$

(e) The width of the confidence interval is  $2\Phi_{\alpha/2} \sqrt{\sigma^2(\mathbf{X}^T \mathbf{X})_{ii}^{-1}}$ . Notice that it is proportional to  $\sqrt{(\mathbf{X}^T \mathbf{X})_{ii}^{-1}}$ . In the event of multicollinearity, the matrix  $(\mathbf{X}^T \mathbf{X})$  will be very ill-conditioned and hence the numerical inversion very unstable, yielding inaccurate and potentially very large coefficients  $\sqrt{(\mathbf{X}^T \mathbf{X})_{ii}^{-1}}$ , which may blow up the width of the confidence interval.

## Q5

(a) By plugging the given value for  $\mu_i$  into the Poisson density function one has :

$$L(y_i | \alpha, \beta) = c \cdot \exp(-e^{\alpha+\beta x_i}) \exp(\alpha y_i + \beta x_i y_i)$$

where  $c$  is a constant.

Therefore the loglikelihood is

$$l(\alpha, \beta) = - \sum_{i=1}^n e^{\alpha+\beta x_i} + \sum_{i=1}^n \alpha y_i + \beta x_i y_i + \text{constant}$$

(b) By factorization theorem  $t_1 = \sum y_i$  and  $t_2 = \sum x_i y_i$  are sufficient statistics for  $\alpha$  and  $\beta$ .

(c) Since

$$l(\alpha, \beta) = - \sum_{i=1}^n e^{\alpha+\beta x_i} + \alpha t_1 + \beta t_2 + \text{constant}$$

thus

$$\frac{\partial l}{\partial \alpha} = 0 \rightarrow t_1 = \sum_{i=1}^n e^{\alpha+\beta x_i}$$

and

$$\frac{\partial l}{\partial \beta} = 0 \rightarrow t_2 = \sum_{i=1}^n x_i e^{\alpha+\beta x_i}$$

(d) if  $\beta = 0$ ,  $l(\alpha, \beta) = - \sum_{i=1}^n e^\alpha + \alpha t_1$ . This is maximized with respect to  $\alpha$  with  $\alpha^* = \log(t_1/n)$

(e) One could refer to

$$2\{l(\hat{\alpha}, \hat{\beta}) - l(\alpha^*, 0)\}$$

having  $\chi_1^2$  distribution by Wilk's theorem.

**Bonus Q** Let  $Y_n \sim \text{Bernoulli}(1/n)$  and define  $X_n = nY_n$ . Let  $\varepsilon > 0$ . Then

$$\mathbb{P}[|X_n| > \varepsilon] = \mathbb{P}[|Y_n| > 0] = \frac{1}{n}.$$

for all  $n > \varepsilon$ . It follows that

$$X_n \xrightarrow{p} 0.$$

On the other hand, for all  $n$ ,

$$\mathbb{E}[X_n] = n \times \frac{1}{n} + 0 \times \left(1 - \frac{1}{n}\right) = 1.$$