

STAT210A - Homework 10

Hoang Duong

November 24, 2014

Problem 1. Read Chapter of Keener

Problem 2. Complete Sufficient Statistics

Proof. We have the density for Y is:

$$p_Y(y) = \frac{1}{(2\pi)^{-n/2} \sigma^n} \exp \left\{ -\frac{1}{2\sigma^2} (y - X\beta)^T (y - X\beta) \right\} \quad (1)$$

$$= \frac{1}{(2\pi)^{-n/2} \sigma^n} \exp \left\{ -\frac{1}{2\sigma^2} (y - X\hat{\beta} + X\hat{\beta} - X\beta)^T (y - X\hat{\beta} + X\hat{\beta} - X\beta) \right\} \quad (2)$$

$$= \frac{1}{(2\pi)^{-n/2} \sigma^n} \exp \left\{ -\frac{1}{2\sigma^2} \left[\|y - X\hat{\beta}\|_2^2 + \|X\hat{\beta} - X\beta\|_2^2 + 2(y - X\hat{\beta})^T (X\hat{\beta} - X\beta) \right] \right\} \quad (3)$$

$$= \frac{1}{(2\pi)^{-n/2} \sigma^n} \exp \left\{ -\frac{1}{2\sigma^2} \left[(n-p) S^2 + \|X\hat{\beta} - X\beta\|_2^2 + 2y^T (I_n - H)^T X (\hat{\beta} - \beta) \right] \right\} \quad (4)$$

For $H = X(X^T X)^{-1} X^T$ is the projection matrix. H is symmetric, thus $I_n - H$ is symmetric. So $(I_n - H)^T X = (I_n - H) X = X - HX = 0$. So the cross term in (4) is zero. Thus:

$$p_Y(y) = \frac{1}{(2\pi)^{-n/2} \sigma^n} \exp \left\{ -\frac{1}{2\sigma^2} \left[(n-p) S^2 + \|X\hat{\beta} - X\beta\|_2^2 \right] \right\}$$

By the property of exponential family, we have $(\hat{\beta}, S^2)$ is a complete sufficient statistics. \square

Problem 3. Hypothesis Testing

Proof. (a) Consider a test $\delta_b(X) = \mathbb{I}_{\{|\hat{\beta}| \geq b\}}$ of rejecting H_0 if $|\hat{\beta}| \geq b$, and fail to reject if not. In the setting of Keener Section 14.5, we have:

$$S^2 = \frac{1}{n-2} \sum e^2$$

$$\mathbb{P} \left[\hat{\beta}_2 - \frac{St_{\alpha/2, n-2}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} \leq \beta \leq \hat{\beta}_2 + \frac{St_{\alpha/2, n-2}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} \right] = 1 - \alpha$$

$$\mathbb{P} \left[|\hat{\beta} - \beta| \geq \frac{St_{\alpha/2, n-2}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} \right] = \alpha$$

So if we choose $b = \frac{St_{\alpha/2, n-2}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}$, we have $\delta_b(X)$ is a level- α test for $H_0 : \beta_2 = 0$ versus $H_1 : \beta_2 \neq 0$.
 (b) Idea taken from John P. Buonaccorsi's note at UMass, and Parker (2011).

Consider $\hat{\theta} = \bar{x} + (y_0 - \hat{\beta}_1) / \hat{\beta}_2$ as an estimate for θ . From Casella and Berger (2002), we note:

$$\begin{aligned}
\text{Var} \left[\frac{U}{V} \right] &\approx \frac{\text{Var} [U]}{\mathbb{E}^2 V} + \frac{\mathbb{E}^2 U}{\mathbb{E}^4 V} \text{Var} [V] - 2 \frac{\mathbb{E} U}{\mathbb{E}^3 V} \text{Cov} [U, V] \\
U &= y_0 - \hat{\beta}_1 \\
V &= \hat{\beta}_2 \\
\mathbb{E} U &= \mathbb{E} [y_0 - \hat{\beta}_1] = \mathbb{E} [\hat{\beta}_2 (\theta - \bar{x}) + e_0] \\
&= (\theta - \bar{x}) \mathbb{E} \hat{\beta}_2 = \beta_2 (\theta - \bar{x}) \\
\mathbb{E} V &= \beta_2 \\
\text{Var} U &= \text{Var} [y_0] + \text{Var} [\hat{\beta}_1] - 2 \text{Cov} [y_0, \hat{\beta}_1] \\
&= \sigma^2 + \frac{\sigma^2}{n} + \frac{\sigma^2 \bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \\
\text{Var} V &= \frac{\sigma^2}{S_{xx}} \\
\text{Cov} [U, V] &= \frac{\sigma^2 \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} \\
\Rightarrow \text{Var} \left[\frac{U}{V} \right] &\approx \frac{\sigma^2}{\beta_2^2} \left(1 + \frac{1}{n} + \frac{(\theta - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)
\end{aligned}$$

Based on the Delta Method, we have:

$$\mathbb{P} \left[\hat{\theta} - t_{1-\alpha/2, n-2} \frac{\hat{\sigma}^2}{\hat{\beta}_2^2} \sqrt{1 + \frac{1}{n} + \frac{(\theta - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} \leq \theta \leq \hat{\theta} + t_{1-\alpha/2, n-2} \frac{\hat{\sigma}^2}{\hat{\beta}_2^2} \sqrt{1 + \frac{1}{n} + \frac{(\theta - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} \right] = 1 - \alpha$$

Thus a level- α test of $H_0 : \theta = \theta_0$ versus $H_1 : \theta \neq \theta_0$ is to reject H_0 iff:

$$\begin{aligned}
\left| \hat{\theta} - \theta_0 \right| &\geq t_{1-\alpha/2, n-2} \frac{\hat{\sigma}^2}{\hat{\beta}_2^2} \sqrt{1 + \frac{1}{n} + \frac{(\theta_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} \\
\hat{\theta} &= \bar{x} + \frac{y_0 - \hat{\beta}_1}{\hat{\beta}_2} \\
\hat{\sigma}^2 &= \frac{1}{n-2} \sum_{i=1}^n \left(y_i - \hat{\beta}_1 - \hat{\beta}_2 (x_i - \bar{x}) \right)^2
\end{aligned}$$

(c) We have:

$$h(\theta) = \frac{y_0 - (\hat{\beta}_1 + \hat{\beta}_2 \theta)}{\left[\sigma^2 + \frac{\sigma^2}{n} + \theta^2 \frac{\sigma^2}{S_{xx}} + 2\theta \frac{\sigma^2 \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]^{1/2}}$$

is t-distributed with $n-2$ degrees of freedom. The set of θ where $h^2(\theta) \leq t_{1-\alpha/2, n-2}^2$ is a $1-\alpha$ confidence region for θ . \square

Problem 4. Confidence Interval

Proof. Rewrite the regression as a linear model:

$$\begin{bmatrix} Y_1 \\ \dots \\ Y_{n_1} \\ Y_{n_1+1} \\ \dots \\ Y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 & 0 & 0 \\ \dots & \dots & \dots & \dots \\ 1 & x_{n_1} & 0 & 0 \\ 0 & 0 & 1 & x_{n_1+1} \\ \dots & \dots & \dots & \dots \\ 0 & 0 & 1 & x_n \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \dots \\ \epsilon_{n_1} \\ \epsilon_{n_1+1} \\ \dots \\ \epsilon_n \end{bmatrix}$$

$$Y = X\beta + \epsilon$$

$$\epsilon \sim \mathcal{N}(0, \sigma^2 I)$$

The OLS estimator for this linear model is:

$$\begin{aligned} \hat{\beta} &= X(X^T X)^{-1} X^T y \\ \hat{y} &= X\hat{\beta} \\ \text{Var}\hat{\beta} &= \sigma^2 (X^T X)^{-1} \\ c &= \begin{bmatrix} 0 \\ -1 \\ 0 \\ 1 \end{bmatrix} \\ \mathbb{E}[c^T \hat{\beta}] &= c^T \beta = \beta_4 - \beta_2 \\ \Rightarrow \text{Var}[c^T \hat{\beta}] &= \sigma^2 c^T (X^T X)^{-1} c \\ \hat{\sigma}^2 &= \frac{1}{n-4} \|y - \hat{y}\|_2^2 \\ \Rightarrow \text{SE}(c^T \hat{\beta}) &= \hat{\sigma}^2 c^T (X^T X)^{-1} c \end{aligned}$$

We have $(c^T \hat{\beta} - c^T \beta) / \text{SE}(c^T \beta)$ follows a standard t distribution with $n - 4$ degree of freedom. Thus:

$$\begin{aligned} &\mathbb{P}\left[-t_{\alpha/2, n-4} \leq \frac{c^T \hat{\beta} - c^T \beta}{\text{SE}(c^T \hat{\beta})} \leq t_{\alpha/2, n-4}\right] = 1 - \alpha \\ \Rightarrow &\mathbb{P}\left[-c^T \hat{\beta} - t_{\alpha/2, n-4} \text{SE}(c^T \hat{\beta}) \leq -c^T \beta \leq -c^T \hat{\beta} + t_{\alpha/2, n-4} \text{SE}(c^T \hat{\beta})\right] = 1 - \alpha \\ \Rightarrow &\mathbb{P}\left[c^T \hat{\beta} - t_{\alpha/2, n-4} \text{SE}(c^T \hat{\beta}) \leq c^T \beta \leq c^T \hat{\beta} + t_{\alpha/2, n-4} \text{SE}(c^T \hat{\beta})\right] = 1 - \alpha \end{aligned}$$

So the $(1 - \alpha)$ confident interval for $\beta_4 - \beta_2$ is:

$$\begin{aligned} &\left[c^T \hat{\beta} - t_{\alpha/2, n-4} \text{SE}(c^T \hat{\beta}), c^T \hat{\beta} + t_{\alpha/2, n-4} \text{SE}(c^T \hat{\beta})\right] \\ c^T &= [0, -1, 0, 1] \\ \hat{\beta} &= (X^T X)^{-1} X^T y \\ \text{SE}(c^T \hat{\beta}) &= \hat{\sigma}^2 c^T (X^T X)^{-1} c \\ \hat{\sigma}^2 &= \frac{1}{n-4} \|y - \hat{y}\|_2^2 \\ \hat{y} &= X\hat{\beta} \end{aligned}$$

$t_{\alpha/2, n-4}$ is the value such that the CDF of standard t-distribution with $n - 4$ degree of freedom evaluate to $1 - \alpha/2$. \square

Problem 5. Log-normal Distribution

Proof. (a) Let $X = \log Y \sim \mathcal{N}(\mu, \sigma^2)$

$$\begin{aligned}
\mathbb{E}Y &= \mathbb{E} \exp X \\
&= \int \exp\{x\} \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\} dx \\
&= \int \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{x^2 - 2x\mu + \mu^2 - 2x\sigma^2}{2\sigma^2}\right\} dx \\
&= \int \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{\mu + \frac{\sigma^2}{2}\right\} \exp\left\{-\frac{x^2 - 2x(\mu + \sigma^2) + \mu^2 + 2\mu\sigma^2 + \sigma^4}{2\sigma^2}\right\} dx \\
&= \exp\left\{\mu + \frac{\sigma^2}{2}\right\} \int \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{x^2 - 2x(\mu + \sigma^2) + \mu^2 + 2\mu\sigma^2 + \sigma^4}{2\sigma^2}\right\} dx \\
&= \exp\left\{\mu + \frac{\sigma^2}{2}\right\} \\
\mathbb{E}Y^2 &= \mathbb{E} \exp 2X \\
&= \int \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{x^2 - 2x\mu + \mu^2 - 4x\sigma^2}{2\sigma^2}\right\} dx \\
&= \int \frac{1}{\sqrt{2\pi}\sigma} \exp\{2\mu + 2\sigma^2\} \exp\left\{-\frac{x^2 - 2x(\mu + 2\sigma^2) + \mu^2 + 4\mu\sigma^2 + 4\sigma^4}{2\sigma^2}\right\} dx \\
&= \exp\{2\mu + 2\sigma^2\} \\
\mathbb{E}Y^\alpha &= \exp\{\alpha\mu + \alpha^2\sigma^2/2\} \\
\Rightarrow \text{Var}Y &= \mathbb{E}Y^2 - \mathbb{E}^2Y = \exp\{2\mu + 2\sigma^2\} - \exp\{2\mu + \sigma^2\} \\
\mathbb{P}[Y \leq y] &= \mathbb{P}[\log Y \leq \log y] = \Phi(\log y) \\
\Rightarrow p_Y(y) &= \frac{\partial \Phi(\log y)}{\partial y} \\
&= \frac{\partial \Phi(\log y)}{\partial \log y} \frac{\partial \log y}{\partial y} \\
&= \phi(\log y) \frac{1}{y} \\
&= \frac{1}{\sqrt{2\pi}\sigma y} \exp\left\{-\frac{(\log y - \mu)^2}{2\sigma^2}\right\}
\end{aligned}$$

(b) Y_1, \dots, Y_n are i.i.d log-normal is equivalent to $\log Y_1, \log Y_2, \dots, \log Y_n$ are i.i.d normal μ, σ^2 . In the Gaussian setting, we know that sample mean and sample variance is the UMVU. Thus $(\sum_{i=1}^n \log Y_i)/n$ is the UMVU for μ .

(c) Consider $\hat{\eta}_1 = \exp\{(\sum_{i=1}^n \log Y_i)/n\}$ as an estimate for $\eta = \mathbb{E}Y_i$. We have:

$$\begin{aligned}
\mathbb{E}\hat{\eta}_1 &= \mathbb{E} \left[\prod_{i=1}^n Y_i^{1/n} \right] \\
&= \prod_{i=1}^n \mathbb{E} \left[Y_i^{1/n} \right] \quad Y_i \text{'s are independent} \\
&= \prod_{i=1}^n \exp \left\{ \frac{\mu}{n} + \frac{\sigma^2}{2n^2} \right\} \\
&= \exp \left\{ \mu + \frac{\sigma^2}{2n} \right\} \\
&= \eta \exp \left\{ \frac{\sigma^2}{2n} \right\} \\
&\Rightarrow \mathbb{E} \frac{\hat{\eta}_1}{\exp \{ \sigma^2 / (2n) \}} = \eta \\
&\Leftrightarrow \mathbb{E} \left[\exp \left\{ \frac{1}{n} \sum_{i=1}^n \log Y_i + \frac{1}{2n} \sigma^2 \right\} \right] = \eta \\
&\hat{\eta} = \exp \left\{ \frac{1}{n} \sum_{i=1}^n \log Y_i + \frac{1}{2n} \sigma^2 \right\}
\end{aligned}$$

So $\hat{\eta}$ is an unbiased estimator of η . Now consider the density of (Y_1, Y_2, \dots, Y_n) :

$$\begin{aligned}
p_{Y_1, \dots, Y_n}(y_1, \dots, y_n) &= \prod_{i=1}^n \left[\frac{1}{\sqrt{2\pi}\sigma y_i} \exp \left\{ -\frac{(\log y_i - \mu)^2}{2\sigma^2} \right\} \right] \\
&= \frac{1}{(2\pi)^{n/2} \sigma^n \prod_{i=1}^n y_i} \exp \left\{ -\frac{1}{2\sigma^2} \left[\sum_{i=1}^n \log^2 y_i - 2\mu \sum \log y_i + n\mu^2 \right] \right\}
\end{aligned}$$

With σ^2 a known constant, by the factor theorem for sufficient statistics, we have $\sum \log y_i$ is a sufficient statistic. $\hat{\eta}$ is a function of sufficient statistic, thus it is UMVU.

(d) I would log-transform Y_i into $\log Y_i$ then perform typical OLS on $\log Y_i \sim \beta_1 + \beta_2 x_i$. So the estimator should be $(X^T X)^{-1} X^T \log Y$. For $\log Y$ is the element-wise log of $Y = [Y_1, \dots, Y_n]$. This estimator is also the MLE estimator according to the multivariate log-normal density we have in (c).

References:

- Buonaccorsi, J.P., 2012. STAT505/ST697R: Regression Analysis, Fall 2012 Note. <http://people.math.umass.edu/~johnpb/>
Casella, G., Berger, R.L. (2002). Statistical Inference, 2nd ed., Duxbury, CA.
Parker, P.A., Vining, G.G., Wilson, S.R., Szarka III, J.L., Johnson, N.G., 2011. The Prediction Properties of Inverse and Reverse Regression for the Simple Linear Calibration Problem. <http://ntrs.nasa.gov/archive/nasa/casi.ntrs.nasa.gov/>

□