

Question 3

3(a)

We know that $\hat{\beta} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}$ is the OLS estimator for β in our model $y_i = \beta x_i + u_i$, $i = 1, \dots, n$ and $u_i \sim N(0, \sigma^2)$.

Firstly,

$$\begin{aligned}
 \hat{\beta} &= \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2} \\
 &= \frac{\sum_{i=1}^n x_i (\beta x_i + u_i)}{\sum_{i=1}^n x_i^2} && \text{substituting } y_i = \beta x_i + u_i \\
 &= \frac{\sum_{i=1}^n (\beta x_i^2 + x_i u_i)}{\sum_{i=1}^n x_i^2} \\
 &= \frac{\sum_{i=1}^n \beta x_i^2}{\sum_{i=1}^n x_i^2} + \frac{\sum_{i=1}^n x_i u_i}{\sum_{i=1}^n x_i^2} \\
 &= \beta \frac{\sum_{i=1}^n x_i^2}{\sum_{i=1}^n x_i^2} + \frac{\sum_{i=1}^n x_i u_i}{\sum_{i=1}^n x_i^2} \\
 &= \beta + \frac{\sum_{i=1}^n x_i u_i}{\sum_{i=1}^n x_i^2}
 \end{aligned}$$

Let $\lambda_i = \frac{x_i}{\sum_{i=1}^n x_i^2}$, so

$$\hat{\beta} = \beta + \sum_{i=1}^n \lambda_i u_i \tag{1}$$

Now, it follows that

$$\begin{aligned}
 E[\hat{\beta}] &= E[\beta + \sum_{i=1}^n \lambda_i u_i] \\
 &= E[\beta] + E[\sum_{i=1}^n \lambda_i u_i] && \text{by linearity of expectation} \\
 &= \beta + E[\sum_{i=1}^n \lambda_i u_i] \\
 &= \beta + E[E[\sum_{i=1}^n \lambda_i u_i \mid x_1, \dots, x_n]] && \text{by Law of Iterated Expectation} \\
 &= \beta + E[\sum_{i=1}^n \lambda_i E[u_i \mid x_1, \dots, x_n]] \\
 &= \beta + E[\sum_{i=1}^n \lambda_i \times 0] && \text{by Zero Conditional Mean assumption} \\
 &= \beta + E[0] \\
 &= \beta
 \end{aligned}$$

This proves that the OLS estimator β is unbiased.

3(b)

$$\begin{aligned}
Var[\hat{\beta} \mid x_1, \dots, x_n] &= Var[\beta + \sum_{i=1}^n \lambda_i u_i \mid x_1, \dots, x_n] && \text{from (1)} \\
&= Var[\beta \mid x_1, \dots, x_n] + Var[\sum_{i=1}^n \lambda_i u_i \mid x_1, \dots, x_n] \\
&= Var[\sum_{i=1}^n \lambda_i u_i \mid x_1, \dots, x_n] \\
&= Var[\frac{\sum_{i=1}^n x_i u_i}{\sum_{i=1}^n x_i^2} \mid x_1, \dots, x_n] && \text{substituting back in } \lambda_i \\
&= \frac{1}{(\sum_{i=1}^n x_i^2)^2} Var[\sum_{i=1}^n x_i u_i \mid x_1, \dots, x_n] \\
&= \frac{1}{(\sum_{i=1}^n x_i^2)^2} \sum_{i=1}^n Var[x_i u_i \mid x_1, \dots, x_n] \\
&= \frac{1}{(\sum_{i=1}^n x_i^2)^2} \sum_{i=1}^n x_i^2 Var[u_i \mid x_1, \dots, x_n] \\
&= \frac{1}{\sum_{i=1}^n x_i^2} Var[u_i \mid x_1, \dots, x_n] \\
&= \frac{1}{\sum_{i=1}^n x_i^2} \sigma^2 && \text{as } u_i \sim N(0, \sigma^2)
\end{aligned}$$

It follows immediately that

$$SE[\hat{\beta} \mid x_1, \dots, x_n] = \frac{\sigma}{\sqrt{\sum_{i=1}^n x_i^2}}$$

as $SE[\hat{\beta} \mid x_1, \dots, x_n] = \sqrt{Var[\hat{\beta} \mid x_1, \dots, x_n]}.$

3(c)

$$\begin{aligned}
Var[\bar{u}_j \mid \bar{x}_1, \dots, \bar{x}_J] &= Var[\frac{1}{n_j} \sum_{i=1}^{n_j} u_{ij} \mid \bar{x}_1, \dots, \bar{x}_J] \\
&= \frac{1}{n_j^2} Var[\sum_{i=1}^{n_j} u_{ij} \mid \bar{x}_1, \dots, \bar{x}_J] \\
&= \frac{1}{n_j^2} \sum_{i=1}^{n_j} Var[u_{ij} \mid \bar{x}_1, \dots, \bar{x}_J]
\end{aligned}$$

Now, we know that $u_i \sim N(0, \sigma^2)$. So, if individuals with such an error term distribution are divided into j groups, the consequent error term will likewise follow $u_{ij} \sim N(0, \sigma^2)$. Therefore,

$$\begin{aligned}
Var[\bar{u}_j \mid \bar{x}_1, \dots, \bar{x}_J] &= \frac{1}{n_j^2} \sum_{i=1}^{n_j} Var[u_{ij} \mid \bar{x}_1, \dots, \bar{x}_J] \\
&= \frac{1}{n_j^2} \sum_{i=1}^{n_j} \sigma^2 \\
&= \frac{1}{n_j^2} n_j \sigma^2 \\
&= \frac{\sigma^2}{n_j} \\
&\neq \sigma^2
\end{aligned}$$

Therefore, the error terms \bar{u}_j is heteroskedastic.

3(d)

We need the OLS estimator for $\tilde{\beta}$. This is derivable from finding the minimum of the residual sum of squares.

$$\begin{aligned}
SSR(\hat{\beta}) &= \sum_{j=1}^J (\bar{y}_j - \hat{y})^2 \\
&= \sum_{j=1}^J (\bar{y}_j - \hat{\beta}\bar{x}_j)^2 \\
\frac{\partial}{\partial \hat{\beta}} (SSR(\hat{\beta})) &= -2 \sum_{j=1}^J \bar{x}_j (\bar{y}_j - \hat{\beta}\bar{x}_j) = 0 \\
\implies \hat{\beta} &= \frac{\sum_{i=1}^n \bar{x}_i \bar{y}_i}{\sum_{i=1}^n \bar{x}_i^2} \quad \text{analogous to what is given in question 3(a)}
\end{aligned}$$

It follows from this that

$$\begin{aligned}
\hat{\beta} &= \frac{\sum_{i=1}^n \bar{x}_i \bar{y}_i}{\sum_{i=1}^n \bar{x}_i^2} \\
&= \frac{\sum_{j=1}^J \bar{x}_j (\tilde{\beta}\bar{x}_j + \bar{u}_j)}{\sum_{j=1}^J \bar{x}_j^2} \quad \text{substituting } \bar{y}_j = \tilde{\beta}\bar{x}_j + \bar{u}_j \\
&= \frac{\sum_{j=1}^J (\tilde{\beta}\bar{x}_j^2 + \bar{x}_j \bar{u}_j)}{\sum_{j=1}^J \bar{x}_j^2} \\
&= \frac{\sum_{j=1}^J \tilde{\beta}\bar{x}_j^2 + \sum_{j=1}^J \bar{x}_j \bar{u}_j}{\sum_{j=1}^J \bar{x}_j^2} \\
&= \tilde{\beta} \frac{\sum_{j=1}^J \bar{x}_j^2}{\sum_{j=1}^J \bar{x}_j^2} + \frac{\sum_{j=1}^J \bar{x}_j \bar{u}_j}{\sum_{j=1}^J \bar{x}_j^2} \\
&= \tilde{\beta} + \frac{\sum_{j=1}^J \bar{x}_j \bar{u}_j}{\sum_{j=1}^J \bar{x}_j^2}
\end{aligned}$$

which is analogous to question 3(a).

Let $\lambda_j = \frac{\bar{x}_j}{\sum_{j=1}^J \bar{x}_j^2}$, so

$$\hat{\beta} = \tilde{\beta} + \sum_{j=1}^J \lambda_j \bar{u}_j \quad (2)$$

Following, the expectation of $\hat{\beta}$ is

$$\begin{aligned} E[\hat{\beta}] &= E[\tilde{\beta} + \sum_{j=1}^J \lambda_j \bar{u}_j] \\ &= E[\tilde{\beta}] + E[\sum_{j=1}^J \lambda_j \bar{u}_j] && \text{by linearity of expectation} \\ &= \tilde{\beta} + E[\sum_{j=1}^J \lambda_j \bar{u}_j] \\ &= \tilde{\beta} + E[E[\sum_{j=1}^J \lambda_j \bar{u}_j \mid \bar{x}_1, \dots, \bar{x}_J]] && \text{by Law of Iterated Expectation} \\ &= \tilde{\beta} + E[\sum_{j=1}^J \lambda_j E[\bar{u}_j \mid \bar{x}_1, \dots, \bar{x}_J]] \\ &= \tilde{\beta} + E[\sum_{j=1}^J \lambda_j \times 0] && \text{by Zero Conditional Mean assumption} \\ &= \tilde{\beta} + E[0] \\ &= \tilde{\beta} \end{aligned}$$

Hence, the OLS estimator for $\tilde{\beta}$ is unbiased.

3(e)

$$\begin{aligned}
Var[\hat{\tilde{\beta}} \mid \bar{x}_1, \dots, \bar{x}_J] &= Var[\tilde{\beta} + \sum_{i=1}^n \lambda_j \bar{u}_i \mid \bar{x}_1, \dots, \bar{x}_J] && \text{from (2)} \\
&= Var[\tilde{\beta} \mid \bar{x}_1, \dots, \bar{x}_J] + Var[\sum_{i=1}^n \lambda_j \bar{u}_i \mid \bar{x}_1, \dots, \bar{x}_J] \\
&= Var[\sum_{i=1}^n \lambda_j \bar{u}_i \mid \bar{x}_1, \dots, \bar{x}_J] \\
&= Var[\frac{\sum_{i=1}^n \bar{x}_i \bar{u}_i}{\sum_{i=1}^n \bar{x}_i^2} \mid \bar{x}_1, \dots, \bar{x}_J] && \text{substituting back in } \lambda_j \\
&= \frac{1}{(\sum_{i=1}^n \bar{x}_i^2)^2} Var[\sum_{i=1}^n \bar{x}_i \bar{u}_i \mid \bar{x}_1, \dots, \bar{x}_J] \\
&= \frac{1}{(\sum_{i=1}^n \bar{x}_i^2)^2} \sum_{i=1}^n Var[\bar{x}_i \bar{u}_i \mid \bar{x}_1, \dots, \bar{x}_J] \\
&= \frac{1}{(\sum_{i=1}^n \bar{x}_i^2)^2} \sum_{i=1}^n \bar{x}_i^2 Var[\bar{u}_i \mid \bar{x}_1, \dots, \bar{x}_J] \\
&= \frac{1}{\sum_{i=1}^n \bar{x}_i^2} Var[\bar{u}_i \mid \bar{x}_1, \dots, \bar{x}_J] \\
&= \frac{\sigma^2}{n_j \sum_{i=1}^n \bar{x}_i^2} && \text{as } Var[\bar{u}_j \mid \bar{x}_1, \dots, \bar{x}_J] = \frac{\sigma^2}{n_j} \text{ from 3(c)} \\
&= \frac{\sigma^2}{n_j \sum_{i=1}^n (\sum_{i=1}^{n_j} \frac{x_{ij}}{n_j})^2} \\
&= \frac{\sigma^2}{\frac{n_j}{n_j^2} \sum_{i=1}^n (\sum_{i=1}^{n_j} x_{ij})^2} \\
&= \frac{n_j \sigma^2}{\sum_{i=1}^n (\sum_{i=1}^{n_j} x_{ij})^2}
\end{aligned}$$

It follows immediately that

$$SE[\hat{\tilde{\beta}} \mid \bar{x}_1, \dots, \bar{x}_J] = \frac{\sqrt{n_j} \sigma}{\sqrt{\sum_{i=1}^n (\sum_{i=1}^{n_j} x_{ij})^2}}$$

$$\text{as } SE[\hat{\tilde{\beta}} \mid \bar{x}_1, \dots, \bar{x}_J] = \sqrt{Var[\hat{\tilde{\beta}} \mid \bar{x}_1, \dots, \bar{x}_J]}.$$

3(f)

Given that the error term is heteroskedastic $Var[\bar{u}_j \mid \bar{x}_1, \dots, \bar{x}_J] = \frac{\sigma^2}{n_j}$, the OLS estimators for our regression model parameters are still unbiased. However, OLS estimates are no longer Best Linear Unbiased Estimator (BLUE). In other words, among all the unbiased estimators, OLS does not provide the estimate with the smallest variance. This is because the standard errors from our model will be biased when heteroskedasticity is present. This in turn leads to biased test statistics and confidence intervals. So, significance tests can be too high or too low.

For the t-test, this is because the test statistic, under the null hypothesis ($\tilde{\beta} = 0$), is calculated as $\frac{\hat{\tilde{\beta}}}{SE[\hat{\tilde{\beta}}]} = \frac{\hat{\tilde{\beta}}}{\sqrt{\frac{\sigma^2}{\sum_{i=1}^n (\sum_{i=1}^{n_j} x_{ij})^2}}}$ from question 3(e). This test statistic follows a t-distribution with $n - 2$ degrees of freedom.

But, if errors are homoskedastic then the test statistics would be $\frac{\hat{\beta}}{\sqrt{n_j \sum_{i=1}^n (\sum_{i=1}^{n_j} x_{ij})^2}}$. So, the test statistic is biased.

For the F-test, this is because the test statistic under the null hypothesis ($\hat{\beta} = 0$) is calculated as $\frac{SSR_R - SSR_{UR}/1}{SSR_R/(n-2)}$ which follows a F-distribution with (1, 2) degrees of freedom. The restricted model (indicated by the subscript R and unrestricted by UR) would come from null hypothesis restriction that $\hat{\beta} = 0$. Note that since $Var[\bar{u}_j] = \frac{SSR_{UR}}{n-2} \implies SSR_{UR} = (n-2)Var[\bar{u}_j]$, and $Var[\bar{u}_j] = \frac{\sigma^2}{n_j} \neq \sigma^2$ (heteroskedastic errors), our SSR_{UR} will be biased. Consequently, our test-statistic is biased.

3(g)

The functional form of the error variance for White's test is

$$Var[\bar{u}_j | \bar{x}_j] = \alpha_1 \bar{x}_j + \alpha_2 \bar{x}_j^2$$

The resulting auxiliary regression model is

$$\hat{u}_j = \alpha_1 \bar{x}_j + \alpha_2 \bar{x}_j^2 + v_i$$

where \hat{u}_j comes from the observed OLS residuals for the j th group.

We test the null and alternative hypothesis:

$$\begin{aligned} H_0 : Var[\bar{u}_j] &= E[\hat{u}_j^2 | \bar{x}_j] = \sigma^2 \implies \alpha_1 = \alpha_2 = 0 \\ H_1 : Var[\bar{u}_j | \bar{x}_j] &= \alpha_1 \bar{x}_j + \alpha_2 \bar{x}_j^2 \implies \alpha_1 \text{ and/or } \alpha_2 \neq 0 \end{aligned}$$

The White test statistic is

$$W = J \times R_{\hat{u}_j}^2$$

, where J is the total number of groups (i.e. observations) and $R_{\hat{u}_j}^2$ is the coefficient of determination from the estimated auxiliary regression model. W follows, asymptotically, a chi-squared distribution with 2 degrees of freedom:

$$W \stackrel{asy}{\sim} \chi^2(2)$$

Decision rule: We reject the null hypothesis H_0 , if the realised W from our data is greater than the $((1 - \alpha) \times 100)^{th}$ quantile of a $\chi^2(2)$ distribution. That is, if $W_{calc} > \chi^2(2)^{@1-\alpha}$.

3(h)

Given that our model is $\bar{y}_j = \tilde{\beta} \bar{x}_j + \bar{u}_j$, multiplying both sides by $\sqrt{n_j}$ yields:

$$\bar{y}_j \sqrt{n_j} = \tilde{\beta} \bar{x}_j \sqrt{n_j} + \bar{u}_j \sqrt{n_j}$$

Hence, the variance of our regression's error term is

$$\begin{aligned} Var[\bar{u}_j \sqrt{n_j} | \bar{x}_1, \dots, \bar{x}_J] &= n_j Var[\bar{u}_j | \bar{x}_1, \dots, \bar{x}_J] \\ &= n_j \frac{\sigma^2}{n_j} && \text{from 3(c)} \\ &= \sigma^2 \end{aligned}$$

Therefore, the error term is homoskedastic. That is, applying weight $\sqrt{n_j}$ to \bar{y}_j corrects for heteroskedasticity.

3(i)

Another than the issue of heteroskedasticity, we may prefer individual data over group data over the possible issue of autocorrelation existing. An assumption of OLS is that all observations are independent. With grouped data like this, you do not have J independent observations, you have J observations each taking n_j data points from the *same* pool of data. That is, we are effectively sampling with replacement. This means that there will be correlation between each observation. Consequently, we may overstate the true degree of freedom in our regression model and the reported standard errors may be artificially small (leading to biased results of significance). This is absent if we simply run a regression across all the individual data, as there is no (effective) sampling happening in this case.