# 資料分析方法

經研一 江彥亨 R11323040

1.a. $Cov(\hat{\beta_0}, \hat{\beta_1}) = Cov(\bar{y} - \hat{\beta_1}\bar{x}, \hat{\beta_1})$

$= \underbrace{Cov(\bar{y}, \hat{\beta_1})}_{=0} - \bar{x} Var(\hat{\beta_1})$

$= -\bar{x} Var\left(\frac{\sum_{i=1}^{n}(x_i - \bar{x})y_i}{\sum_{i=1}^{n}(x_i - \bar{x})^2}\right) = -\bar{x} Var\left(\frac{\sum_{i=1}^{n}(x_i - \bar{x})(\beta_0 + \beta_1 x_i + u_i)}{\sum_{i=1}^{n}(x_i - \bar{x})^2}\right)$

$= -\bar{x} Var\left(\beta_1 + \frac{\sum_{i=1}^{n}(x_i - \bar{x})u_i}{\sum_{i=1}^{n}(x_i - \bar{x})^2}\right)$

$= -\bar{x} \frac{\sum_{i=1}^{n}(x_i - \bar{x})^2}{(\sum_{i=1}^{n}(x_i - \bar{x})^2)^2} Var(u_i) = -\bar{x} \bar{\sigma}^2 \frac{n}{\sum_{i=1}^{n}(x_i - \bar{x})^2}$

$= -\bar{x}\bar{\sigma}^2 / Sxx \quad \#$

b. $Cov(\bar{y}, \hat{\beta_1}) = Cov(\hat{\beta_0} + \hat{\beta_1}\bar{x}, \hat{\beta_1})$

$= Cov(\hat{\beta_0}, \hat{\beta_1}) + \bar{x} Var(\hat{\beta_1})$

$= -\bar{x} Var(\hat{\beta_1}) + \bar{x} Var(\hat{\beta_1})$

$= 0 \quad \#$

2. $SSR = \sum_{i=1}^{n}(\hat{y_i} - \bar{y})^2 = \sum_{i=1}^{n} \hat{y_i}^2 - 2\hat{y_i}\bar{y} + \bar{y}^2$

$= \sum_{i=1}^{n} \hat{y_i}^2 - 2\bar{y}\sum_{i=1}^{n}\hat{y_i} + \sum_{i=1}^{n}\bar{y}^2$

$\sum_{i=1}^{n}\hat{y_i} = \sum_{i=1}^{n}\hat{\beta} x_i = \hat{\beta}\sum_{i=1}^{n}x_i = n \cdot \hat{\beta}\bar{x}$

Since least square regression curve must pass sample average $\bar{y}$ and $\bar{x}$ $\Rightarrow \bar{y} = \hat{\beta}\bar{x}$

$\Rightarrow \sum_{i=1}^{n}\hat{y_i} = n\hat{\beta}\bar{x} = n\bar{y}$

$\Rightarrow SSR = \sum_{i=1}^{n}\hat{y_i}^2 - 2n\bar{y} + n\bar{y} = \sum_{i=1}^{n}\hat{y_i}^2 - n\bar{y} \quad \#$

3.a $HH = X(X^TX)^{-1}X^T X(X^TX)^{-1}X^T = X(X^TX)^{-1}(X^TX)(X^TX)^{-1}X^T$

$= XI(X^TX)^{-1}X^T = X(X^TX)^{-1}X^T = H$

$(I-H)(I-H) = II - HI - IH + HH = I - 2H + H = I - H$

3.b $V(\hat{Y}) = V(X\hat{\beta}) = X \, Var(\hat{\beta}) X^T$

$Var(\hat{\beta}) = Var((X^TX)^{-1}X^Ty) = (X^TX)^{-1}X^T \, Var(Y) X(X^TX)^{-1}$

$= (X^TX)^{-1}X^T \, \sigma^2 I \, X (X^TX)^{-1}$

$\Rightarrow Var(\hat{Y}) = X(X^TX)^{-1}X^T \, \sigma^2 I \, X(X^TX)^{-1}X^T = H\sigma^2 I H$

$= \sigma^2 H H = \sigma^2 H$

4. By definition, $R^2 = \dfrac{\sum_{i=1}^{n}(\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^{n}(Y_i - \bar{Y})^2}$

and $\sum_{i=1}^{n}(Y_i - \bar{Y})^2$ can be divided into

$\sum_{i=1}^{n}(Y_i - \hat{Y}_i)^2 + \sum_{i=1}^{n}(\hat{Y}_i - \bar{Y})^2$

$\Rightarrow R^2 = \dfrac{\sum_{i=1}^{n}(\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^{n}(Y_i - \hat{Y}_i)^2 + \sum_{i=1}^{n}(\hat{Y}_i - \bar{Y})^2}$

Thus, $R^2$ must be smaller than one. In addition, $\sum_{i=1}^{n}(\hat{Y}_i - \bar{Y})^2$ must be larger or equal to zero, so $R^2$ must not be smaller than zero. Intuitively, $R^2$ means how much variance in the explained variable can be explained by the explanatory variables. Hence, $R^2$ must falls between zero and one.

5.

```python
y = data['mpg']
xvar = ['cylinders', 'displacement', 'horsepower', 'weight', 'acceleration', 'model year', 'origin']
X = data[xvar]
model = sm.OLS(y,X).fit()
VIFs = [variance_inflation_factor(X.values, i) for i in range(X.shape[1])]
for x,vif in zip(xvar,VIFs):
    print(f'the VIFs of {x} is {vif}')
✓ 0.0s
```

```
the VIFs of cylinders is 117.70854741316116
the VIFs of displacement is 96.90976248526793
the VIFs of horsepower is 67.07215430169894
the VIFs of weight is 139.45416214259953
the VIFs of acceleration is 69.69976893381958
the VIFs of model year is 115.7946144892812
the VIFs of origin is 8.469941669334569
```

由結果可知，除 origin 變數外，其餘變數之膨脹因子 VIFs 皆非常大，故這些

變數具有高度共線性問題，因去除其中幾項較不顯著的變數。