**Question 1**

(a) Prove that for any two random variables $X$ and $Y$,

$$\text{Cov}\left(X, Y\right) = \text{E}\left(XY\right) - \text{E}\left(X\right)\text{E}\left(Y\right)$$

(b) Prove that for random variables $X_1, X_2, \ldots, X_n$ and $Y_1, Y_2, \ldots, Y_m$ that

$$\text{Cov}\left(\sum_{i=1}^{n} a_i X_i, \sum_{j=1}^{m} b_j Y_j\right) = \sum_{i=1}^{n}\sum_{j=1}^{m} a_i b_j \text{Cov}\left(X_i, Y_j\right)$$

**Question 2**

Suppose the following 11 values are the transaction amounts (in £) of online purchases for a particular credit card customer in a given month.

$$45, \ 81, \ 52, \ 23, \ 147, \ 92, \ 76, \ 124, \ 287, \ 103, \ 65$$

Tukey's criterion states that, given the lower quartile $q_{0.25}$, the upper quartile $q_{0.75}$ and the interquartile range IQR, if a value $x$ is either $x < q_{0.25} - k\text{IQR}$ or $x > q_{0.75} + k\text{IQR}$, for $k = 1.5$, then $x$ is considered to be an outlier.

(a) Compute the lower and upper quartiles, and the interquartile range for this dataset.

(b) According to Tukey's criterion, are any of these transaction amounts outliers?

(c) If any of the transactions is an outlier, would you take any action? What could be the consequences of (i) inaction (doing nothing) or (ii) taking action (preventing the transaction from going through)?

(d) If you were designing your own fraud detector for this customer (not using Tukey's criterion) for the next month, how high would a value need to be for you to decide that a value is anomalous and potentially fraudulent? In other words, at what value would you set the threshold?

**Question 3**

Recall that given the random variables $Y_1, Y_2, \ldots, Y_n$ and the observations $x_1, x_2, \ldots, x_n$, the estimators $\widehat{\beta}_0$ and $\widehat{\beta}_1$ are defined by

$$\widehat{\beta}_0 = \sum_{i=1}^{n} \left( \frac{1}{n} - \frac{(x_i - \overline{x})\, \overline{x}}{S_{xx}} \right) Y_i$$

$$\widehat{\beta}_1 = \frac{1}{S_{xx}} \sum_{i=1}^{n} (x_i - \overline{x})\, Y_i$$

where $S_{xx} = \sum_{i=1}^{n} (x_i - \overline{x})^2$.

(a) Show that

$$\text{Var}(\widehat{\epsilon}_i) = \text{Var}(Y_i) + \text{Var}(\widehat{\beta}_0) + x_i^2 \text{Var}(\widehat{\beta}_1) + 2x_i \text{Cov}(\widehat{\beta}_0, \widehat{\beta}_1) - 2\text{Cov}(Y_i, \widehat{\beta}_0) - 2x_i \text{Cov}(Y_i, \widehat{\beta}_1).$$

(b) Show that

$$\text{Cov}(Y_i, \widehat{\beta}_0) = \left( \frac{1}{n} - \frac{(x_i - \overline{x})\, \overline{x}}{S_{xx}} \right) \sigma^2.$$

(c) Show that

$$\text{Cov}(Y_i, \widehat{\beta}_1) = \left( \frac{x_i - \overline{x}}{S_{xx}} \right) \sigma^2.$$

(d) Given that

$$\sum_{i=1}^{n} \text{Var}(\widehat{\epsilon}_i) = \sigma^2 \sum_{i=1}^{n} \left[ \frac{n-2}{n} + \frac{1}{S_{xx}} \left( \frac{1}{n} \sum_{j=1}^{n} x_j^2 + x_i^2 - 2x_i \overline{x} - 2(x_i - \overline{x})^2 \right) \right],$$

prove that

$$\frac{1}{n} \sum_{i=1}^{n} \text{Var}(\widehat{\epsilon}_i) = \left( \frac{n-2}{n} \right) \sigma^2.$$