

# Statistics Exercise 2

Madhukara S Holla

13th October 2023

## Question 1

Project members: Madhukara S Holla (Master of Science in Computer Science, 1st Year)

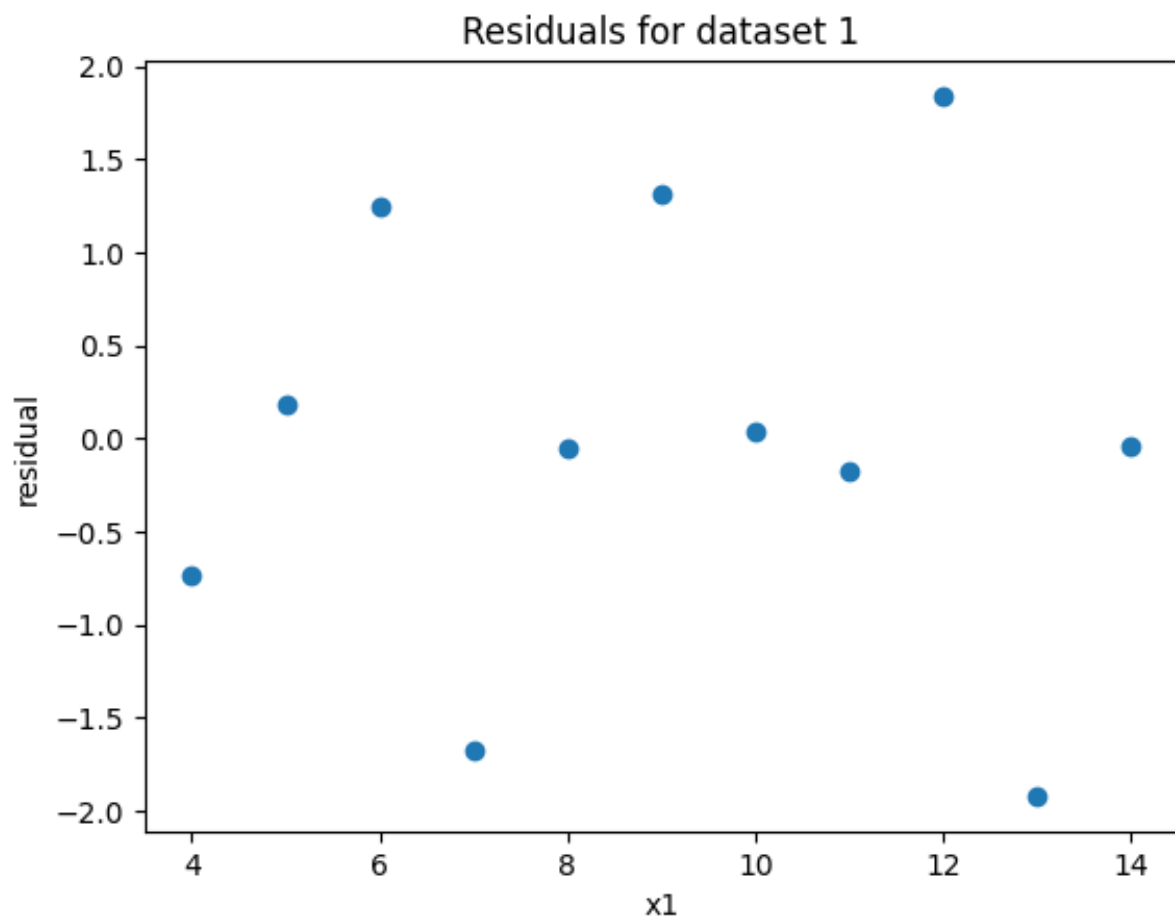
Project description: NA

## Question 2

### 2.a

#### Anscombe's Dataset 1

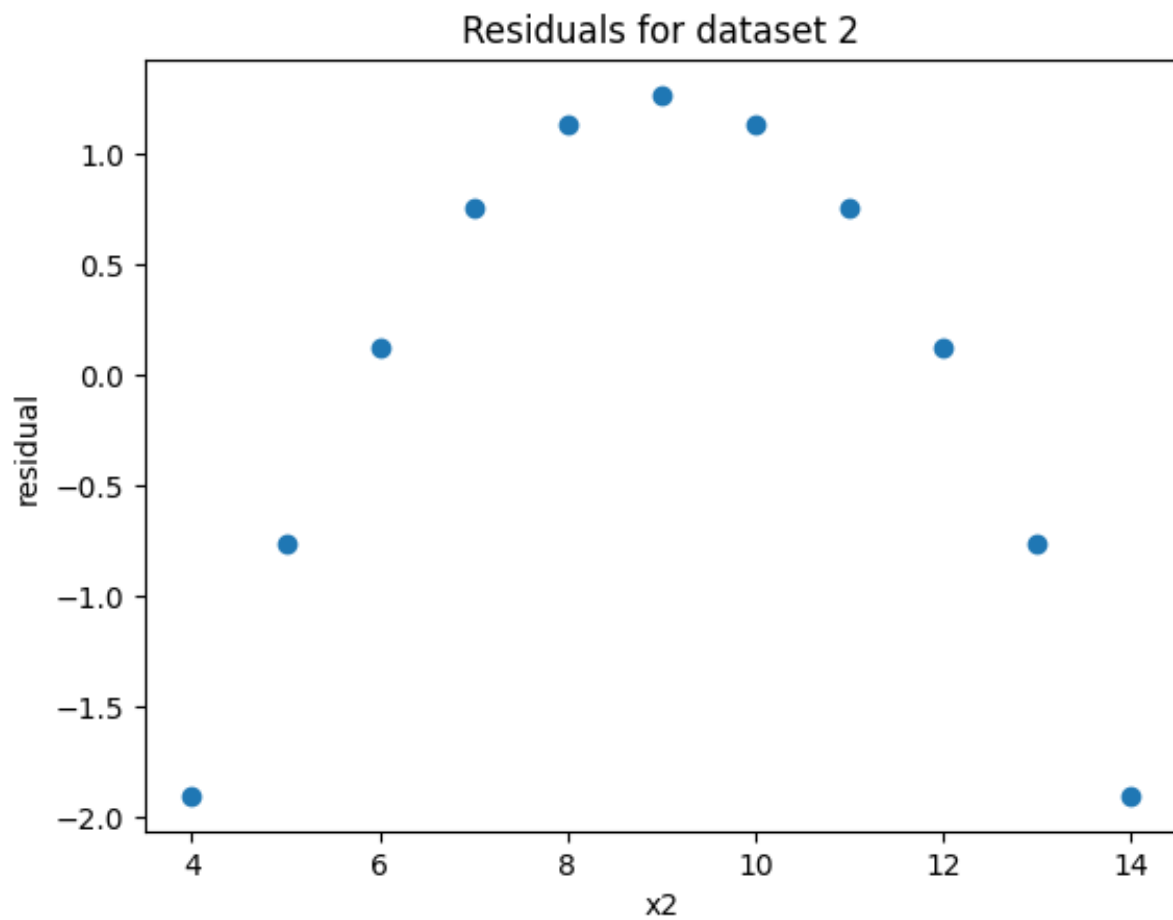
- Observed pearson co-efficient: 0.8164
- Observed co-efficient of multiple determination: 0.6665



## 2.a

### Anscombe's Dataset 2

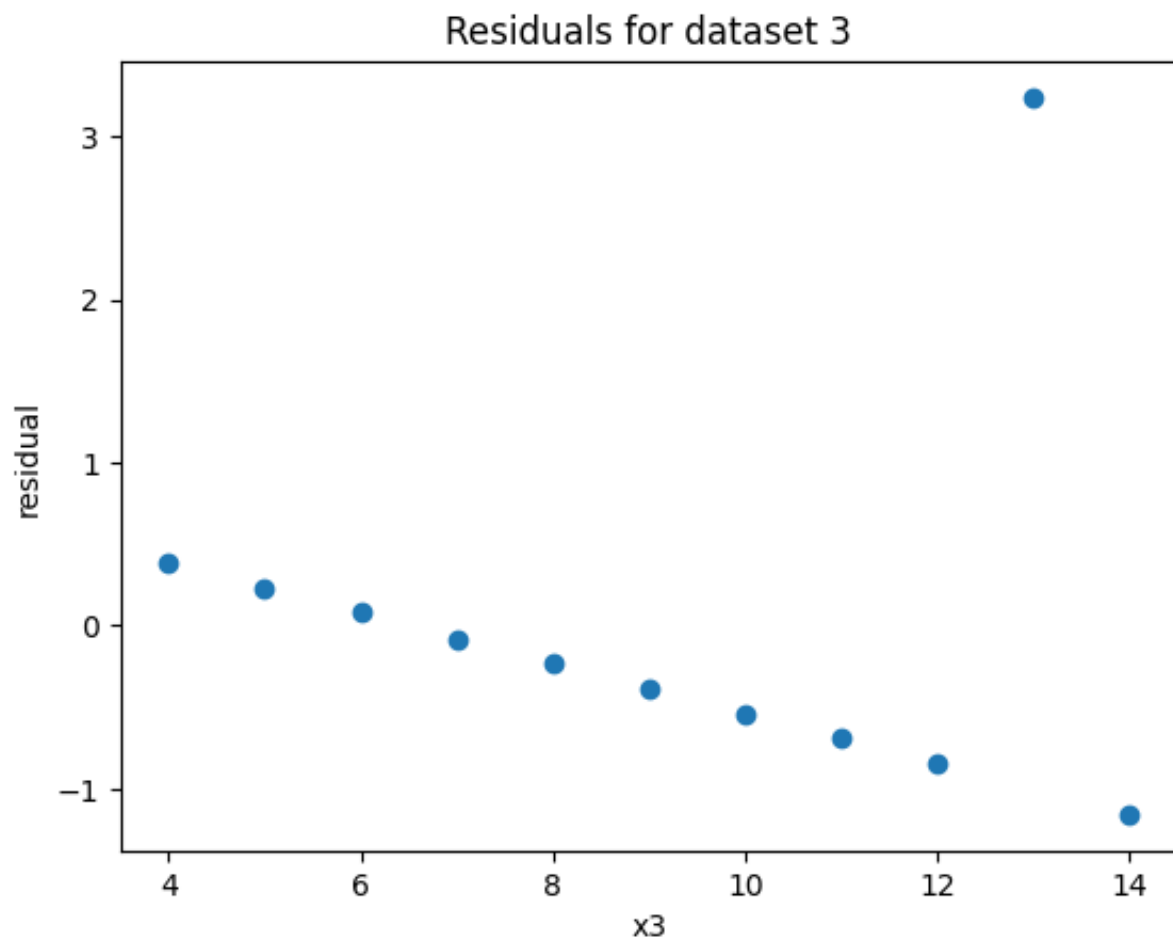
- Observed pearson co-efficient: 0.8162
- Observed co-efficient of multiple determination: 0.6662



## 2.a

### Anscombe's Dataset 3

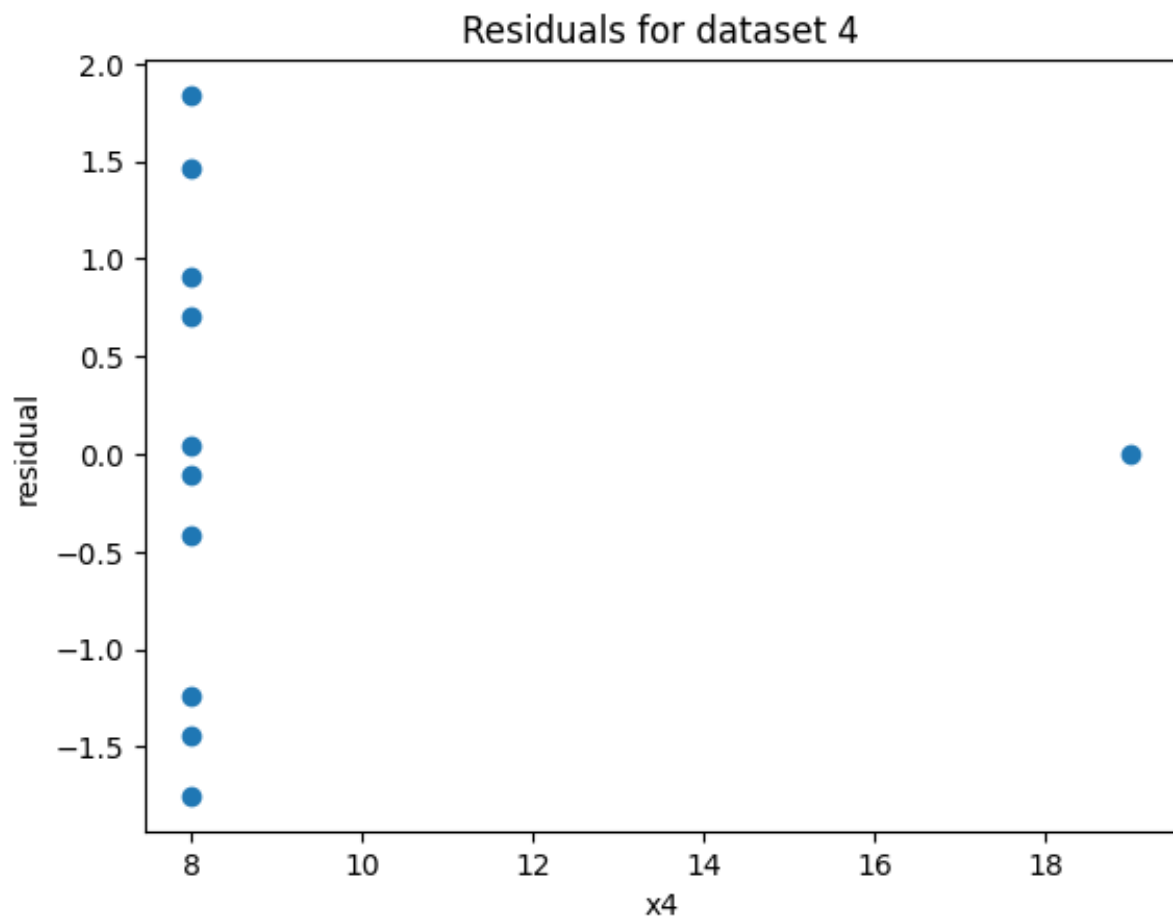
- Observed pearson co-efficient: 0.8163
- Observed co-efficient of multiple determination: 0.6663



## 2.a

### Anscombe's Dataset 4

- Observed pearson co-efficient: 0.8165
- Observed co-efficient of multiple determination: 0.6667



## 2.b

Lack of fit tests for Anscombe's datasets.

Datasets grouped by x values (2 values per group)

Null Hypothesis  $H_0$ : A simple linear model is adequate to explain the systematic variations in the data.

Alternate Hypothesis  $H_a$ : A linear model is not adequate and a nonlinear model is required to capture the systematic variations in the data.

### Anscombe's Dataset 1

P value: 0.86 - Fail to reject  $H_0$ .

No significant lack of fit. The dataset appears to be a simple linear relationship, and a linear regression model seems appropriate for this dataset.

### Anscombe's Dataset 2

P value: 0.03 - Reject  $H_0$  in favor of  $H_a$ .

Significant lack of fit. The data clearly follows a non-linear (quadratic) relationship, indicating that the linear model does not capture all systematic variations in the dataset.

### Anscombe's Dataset 3

P value: 0.83 - Fail to reject  $H_0$ .

The outlier is ignored when we group the data by x values in pairs of 2.

No significant lack of fit. Since the influence of the outlier is ignored while calculating lack of fit, the dataset appears to be a simple linear relationship.

### Anscombe's Dataset 4

P value: 0.04 - Reject  $H_0$  in favor of  $H_a$ .

The outlier is ignored when we group the data by x values in pairs of 2.

Significant lack of fit. Test is not appropriate due to the nature of the data. If forced, likely a significant lack of fit.

The lack of fit test assumes that there's some variation in the independent variable (x) that corresponds to variation in the dependent variable (y). In Dataset 4, for all but one observation, there's no variation in x. This goes against the fundamental premise of regression that we're trying to understand how y changes as x changes.

## 2.c

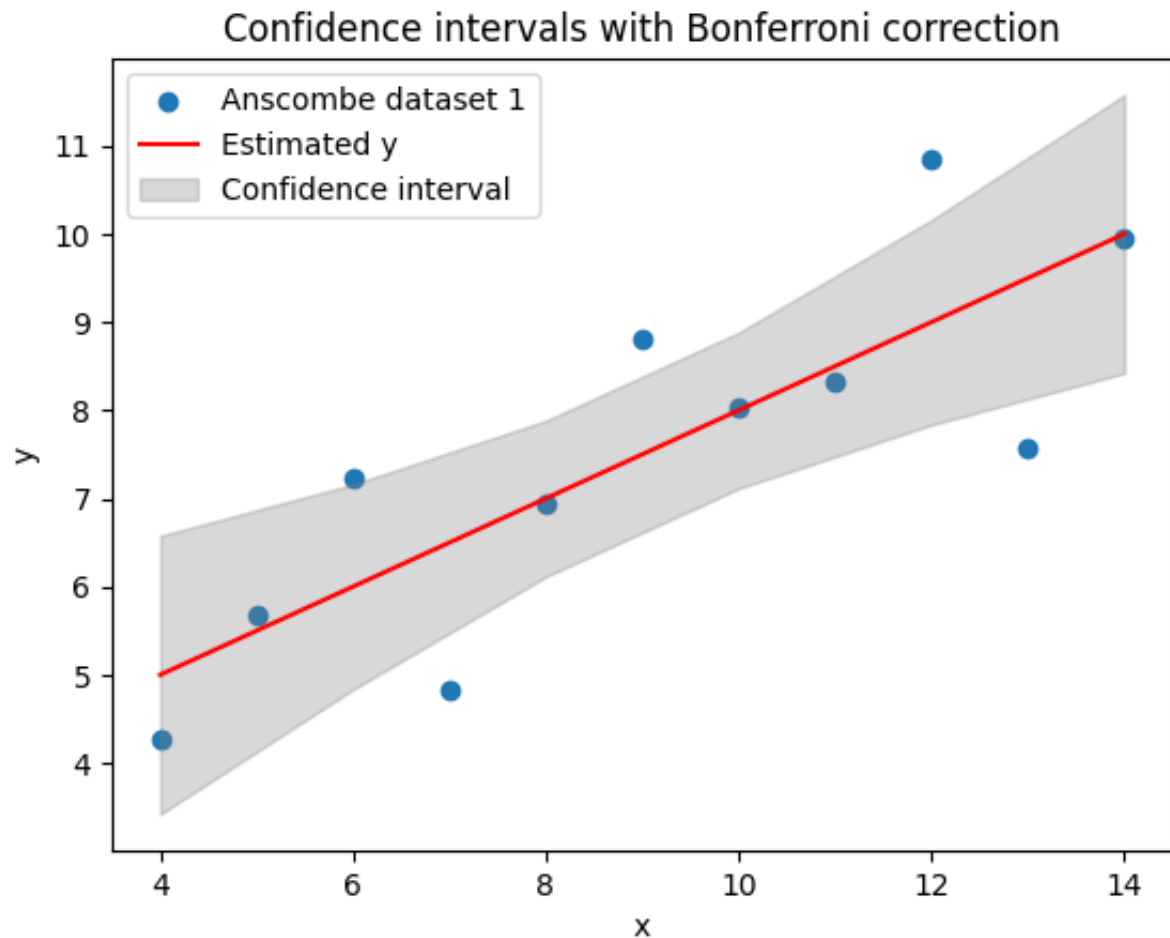
- Patterns in residual plots can indicate non-linearity and outliers, helping us to identify problems with the model.
- Lack of fit tests provide a formal statistical test to validate the model assumptions.
- But the lack of fit tests require replicate observations in the data which may not be available, or it may be inappropriate to run on datasets like Dataset 4.
- Both pearson coefficient and co-efficient of multiple determination do not clearly indicate the goodness of fit of the model. They just indicate the strength of the linear relationship and proportion of variance explained by the model respectively.

In conclusion, we need to use multiple methods such as visualization, lack of fit tests, and co-efficient determinations to validate the model assumptions.



## Question 3

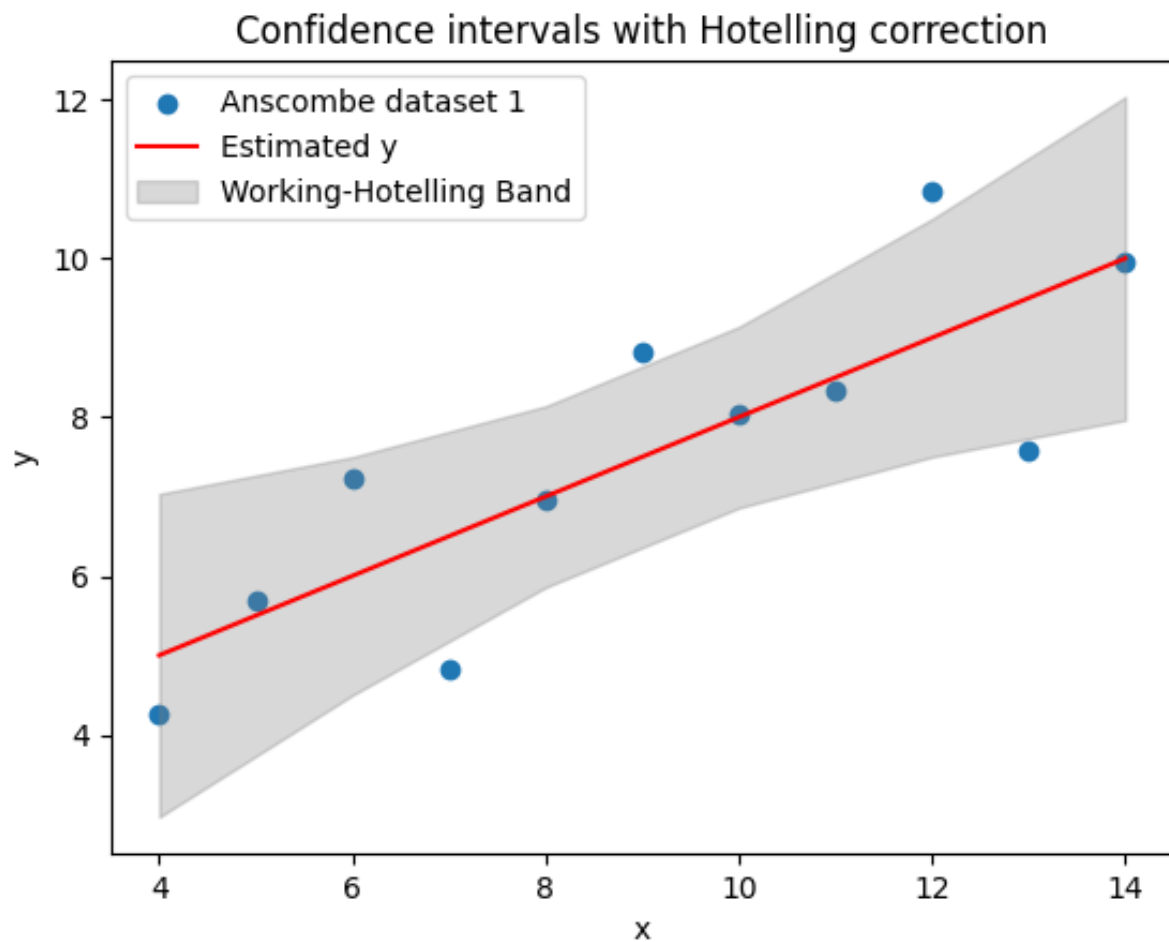
3.a



The Bonferroni adjustment controls the Familywise Error Rate (FWER). This is the probability of making at least one Type I error when performing multiple statistical tests.

By adjusting the significance level with the Bonferroni method, we ensure that our overall chance of incorrectly rejecting a true null hypothesis remains at the specified level (like 5%).

3.b



The Working-Hotelling confidence method is used to construct confidence intervals for the difference between the means of two multivariate data sets.

It controls for the precision of the estimate of the difference between the means. It provides a range of values within which we can be reasonably confident that the true difference between the means lies.

### 3.c

For a simple linear regression model using least squares estimation,

$$Y = \beta_0 + \beta_1.X + \epsilon$$

The estimates for  $\beta_0$  and  $\beta_1$  are given by:

$$b_1 = \frac{\Sigma(X_i - \bar{X})(Y_i - \bar{Y})}{\Sigma(X_i - \bar{X})^2}$$
$$b_0 = \bar{Y} - b_1.X$$

and  $\epsilon$  is the error term.

$$\text{cov}(b_0, b_1) = E[(b_0 - E[b_0])(b_1 - E[b_1])]$$
$$\text{cov}(b_0, b_1) = E[b_0 b_1] - E[b_0].E[b_1]$$

Substituting for  $b_0$ ,

$$E[b_0 b_1] = E[(\bar{Y} - b_1 \bar{X})(b_1)]$$
$$= E[\bar{Y} b_1 - b_1^2 \bar{X}]$$
$$= \bar{Y} E[b_1] - \bar{X} E[b_1^2]$$

We know  $E[b_1] = \beta_1$  and

$$E[b_0] = E[\bar{Y} - b_1 \bar{X}]$$
$$= \bar{Y} - \bar{X} E[b_1] = \bar{Y} - \bar{X}.\beta_1$$

Now we have,

$$E[b_0 b_1] = \beta_1 \bar{Y} - \bar{X} E[b_1^2]$$
$$E[b_0].E[b_1] = (\bar{Y} - \bar{X}.\beta_1)\beta_1$$

Substituting in  $\text{cov}(b_0, b_1)$  we have,

$$\text{cov}(b_0, b_1) = \beta_1 \bar{Y} - \bar{X} E[b_1^2] - (\bar{Y} - \bar{X}.\beta_1)\beta_1$$
$$= \beta_1 \bar{Y} - \bar{X} E[b_1^2] - \beta_1 \bar{Y} + \bar{X}.\beta_1^2$$
$$= \bar{X}.\beta_1^2 - \bar{X} E[b_1^2]$$

We need a formula for  $E[b_1^2]$

$$\text{Var}(b_1) = E[b_1^2] - (E[b_1])^2 = E[b_1^2] - \beta_1^2$$

From linear regression,

$$\text{Var}(b_1) = \frac{\sigma^2}{\Sigma(X_i - \bar{X})^2}$$
$$E[b_1^2] - \beta_1^2 = \frac{\sigma^2}{\Sigma(X_i - \bar{X})^2}$$
$$E[b_1^2] = \beta_1^2 + \frac{\sigma^2}{\Sigma(X_i - \bar{X})^2}$$

Substituting into  $cov(b_0, b_1)$  we have,

$$\begin{aligned} cov(b_0, b_1) &= \bar{X}.\beta_1^2 - \bar{X}E[b_1^2] \\ &= \bar{X}.\beta_1^2 - \bar{X}(\beta_1^2 + \frac{\sigma^2}{\Sigma(X_i - \bar{X})^2}) \\ &= -\frac{\sigma^2}{\Sigma(X_i - \bar{X})^2}\bar{X} \\ cov(b_0, b_1) &= -\bar{X}.Var(b_1) \end{aligned}$$