

# COMP9417: Homework Set #1

z5113817

University of New South Wales — June 20, 2021

## Question 1

a)

We know from the normal equations for a the gradient of a minimised linear regression is:

$$\hat{\beta}_1 = \frac{\bar{X}\bar{Y} - \bar{X}^2}{(\bar{X}^2) - (\bar{X})^2}$$

Which can expand into:

$$\hat{\beta}_1 = \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2}$$

So now substituting in the transformation, we get:

$$\hat{\beta}_1 = \frac{\frac{1}{n} \sum_{i=1}^n (\tilde{X}_i - \bar{X})(Y_i - \bar{Y})}{\frac{1}{n} \sum_{i=1}^n (\tilde{X}_i - \bar{X})^2} \quad (1)$$

$$\hat{\beta}_1 = \frac{\frac{1}{n} \sum_{i=1}^n (c(X_i + d) - \bar{X})(Y_i - \bar{Y})}{\frac{1}{n} \sum_{i=1}^n (c(X_i + d) - \bar{X})^2} \quad (2)$$

$$\dots \quad (3)$$

b)

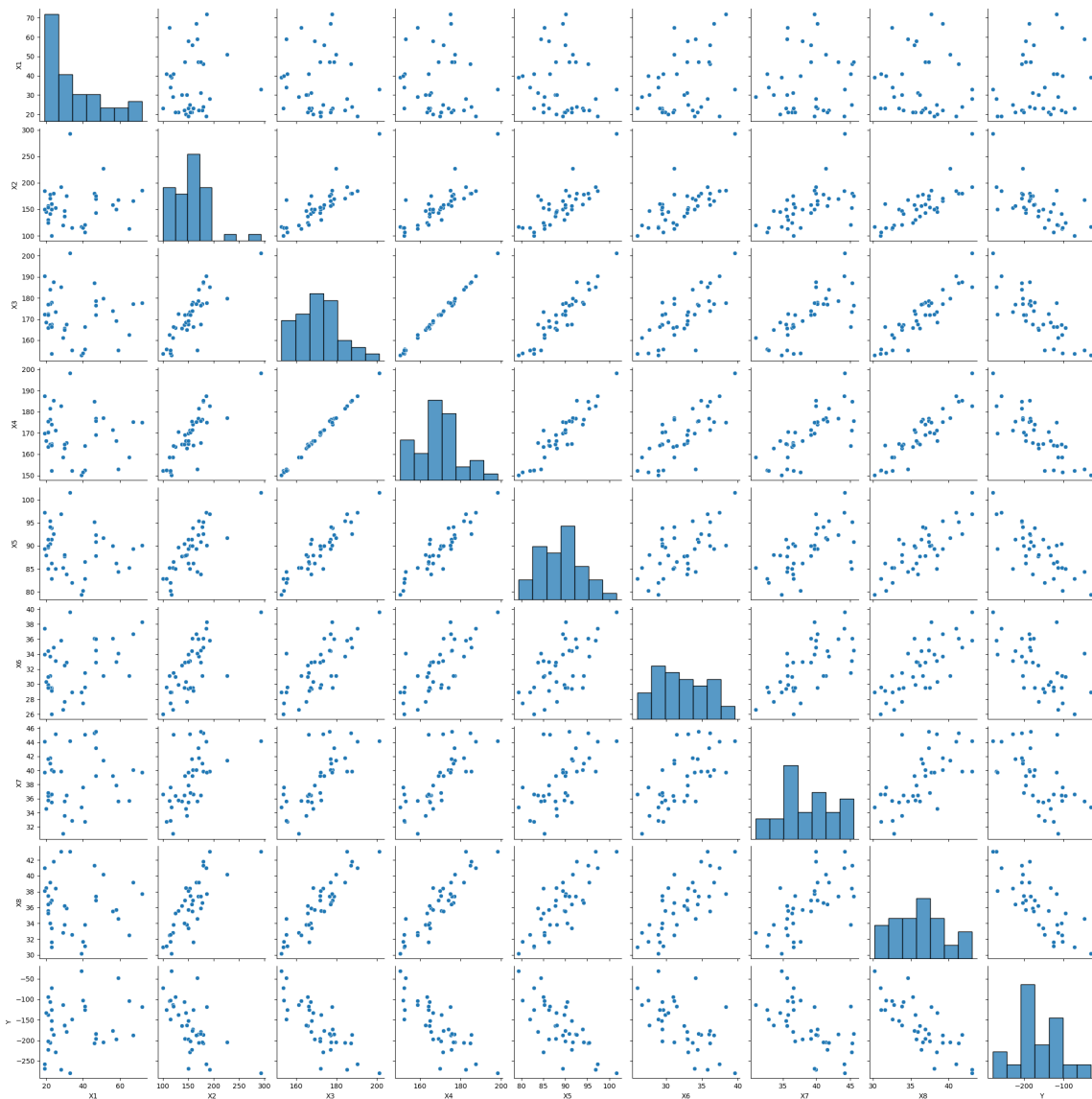
## Question 2

See Github repository [here](#) for all of the python code used in this question.

a)

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

df = pd.read_csv("../data/data.csv")
sns.pairplot(df)
plt.savefig("../outputs/q2a_pairplots.png")
```



The pairs plots shows a scatter plots between variables in the dataset, while the histogram on the diagonal shows the distribution of each variable. For example, the plot on the fourth row in the second column is the scatter plot of X3 on the y-axis and X2 on the x-axis. Such a matrix is useful for seeing correlations between variables. This is important in linear regression to prevent multicollinearity between the variables, which can skew the coefficients of the regression model and lead to unreliable statistical inferences.

**b)**

The code in **q2b.py** will produce:

Sum of squares for each transformed feature:

X1: 38.0000000000000014

X2: 38.0

X3: 38.0000000000000014

X4: 37.999999999999986

X5: 37.99999999999998

X6: 38.0

X7: 38.000000000000001

X8: 37.999999999999986

These should all be 38 exactly however they are not due to floating point errors.

**c)**