# 1 Data Preprocessing

1. There were no errors or missing values in the dataset.

2. Gradient descent algorithms tend to work better when the input features have similar scales. If the features have different scales, the optimization may take longer or get stuck in local minima.

   Rescaling the data can help prevent numerical issues. When features have different scales, the matrix of input data can become poorly fit, which can cause numerical instability during training.

   It also improves the interpretability of the model. If one feature is in meters and another feature is in kilometers, the coefficients of the model will be difficult to interpret.

3. Models that perform regularization according to the norm of the parameters are sensitive to the scale of the input features or variables. If the input features have very different scales, the regularization penalty may be dominated by the features with large scales and may not effectively control the model's complexity. Essentially, rescaling ensures that ll parameters have equal importance.

4. Rescaling ensures that each feature has equal importance. If the features have different scales, then some features may dominate the distance calculation, leading to biased predictions.

   Distance-based models are sensitive to the scale of the input features. Features with large scales may overshadow features with smaller scales.

   Rescaling the data can make the model interpretable. If one feature is measured in meters and another feature is measured in kilometers, the distance between two data points will be dominated by the kilometers feature.

5.
```
in sample error: 0.005300060606060606
out sample error: 0.029683672727272726
in sample error: 0.005304063299663301
out sample error: 0.029392521212121214
[0.02918048484848485, 0.0288448484848484848, 0.03139187878787879, 0.030262545454545455, 0.027228242424242426, 0.02850339393939394, 0.030491030303030305, 0.031527878787878785, 0.02883721212121212, 0.030569212121212122]
[0.028286424242424245, 0.028020000000000003, 0.029015030303030306, 0.028037575757575758, 0.02804169696969697, 0.027090181818181822, 0.03254569696969697, 0.03304278787878788, 0.02885260606060606, 0.030993212121212123]
in sample t-Test
TtestResult(statistic=-0.07971425097307931, pvalue=0.9382088806133239, df=9)
out sample t-Test
TtestResult(statistic=0.6087505128419045, pvalue=0.5577400530280223, df=9)
```

## 2 Feature Selection

1. The advantage of having many features is that there is more information at our disposal that can be used to learn a complex function. More features allow for more complex models that can capture more complex relationships between the input features and the target variable.
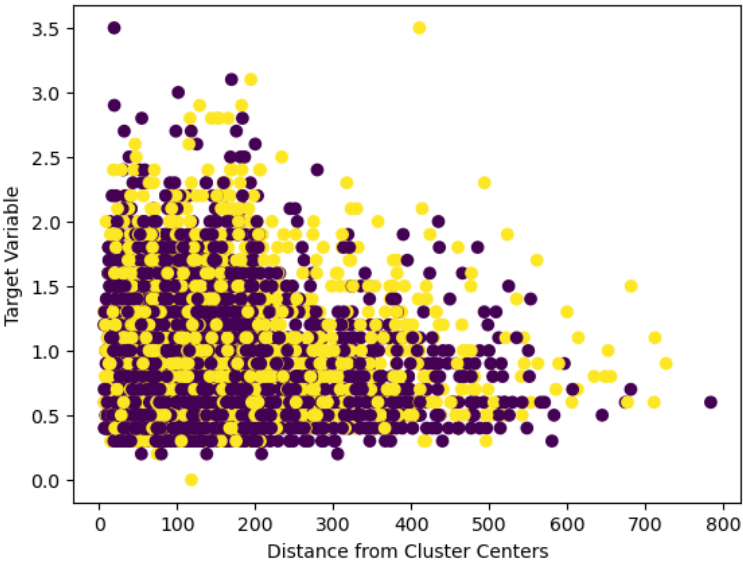
   The disadvantage of too many features can lead to a phenomenon known as the curse of dimensionality, where the complexity of the data and the number of possible feature combinations grow exponentially with the number of features. This can make it difficult to find the best model and lead to overfitting, where the model is too complex and fits the training data too well, but does not generalize well to new data.

2. KDE Their KDE plots showed a very low density and high overlap with other features, which further suggests that it may not be informative.

3. A feature that may be discarded is feature 33. It had a correlation score below a threshold of 0.01 with a value of 0.00287. A low correlation score indicates that the feature may not be strongly related to the target variable and not provide much information to the model.

   A feature that may be kept is feature 6. It had a correlation score of 0.705. A high correlation score between a feature and the target variable indicates that the feature is strongly related to the target variable. Changes in the feature value are highly correlated with changes in the target variable. This suggests that the feature may be informative and could potentially be useful for predicting the target variable.
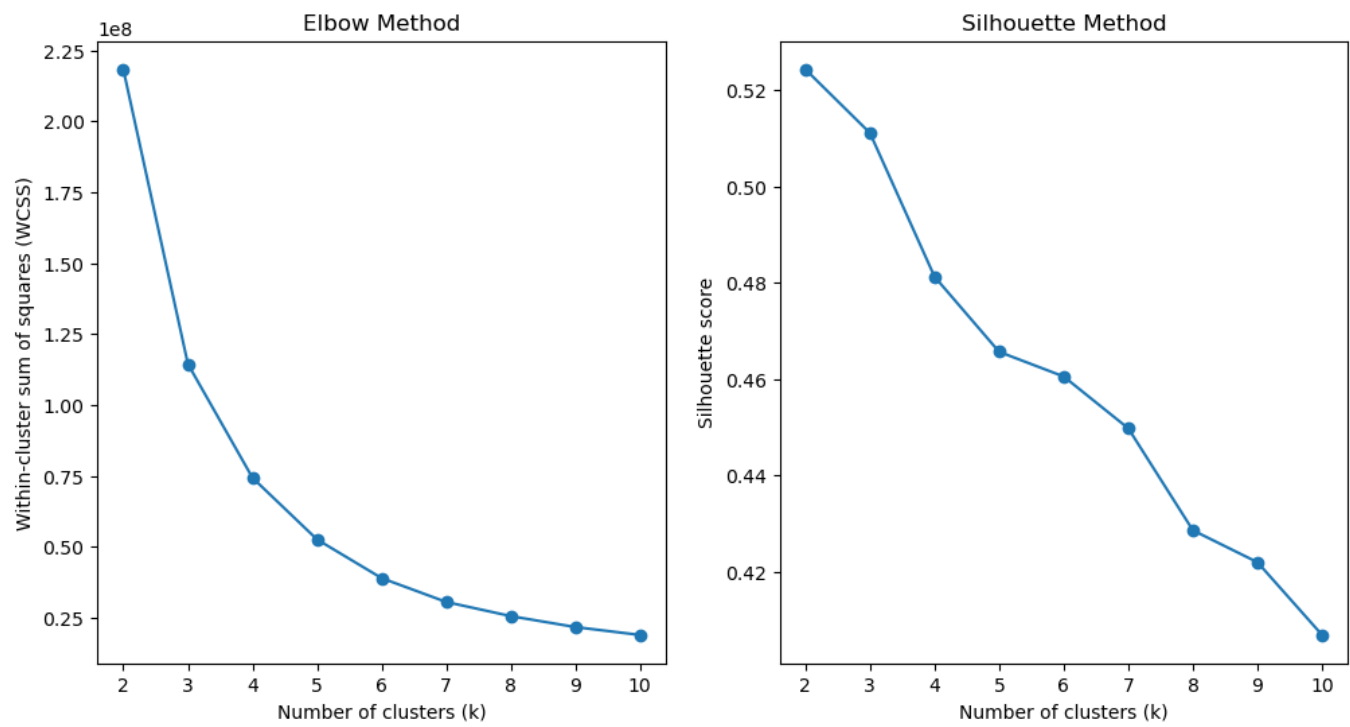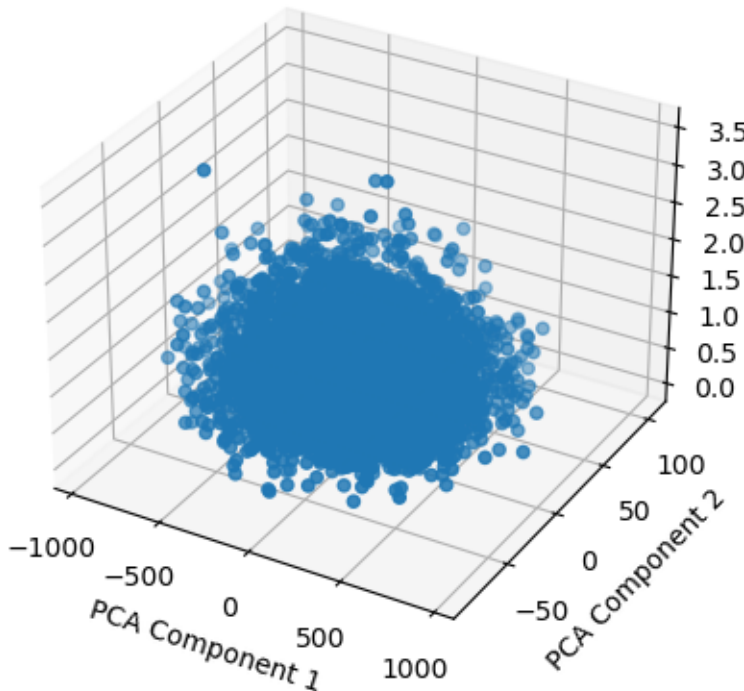
4.

# 3   Clustering



1.

2.

# 4   PCA



1.

2.



3. The negative $R^2$ scores suggest that the linear regression models are not fitting the data well, and in fact are performing worse than a horizontal line. This can be due to many factors, such as non-linear relationships between the features and the target variable, high variance in the data, or the need for more complex models.

   One possible explanation for the decreasing performance with increasing number of PCA components is that too many components are being used, which may result in overfitting to the training data and

poor generalization to new data. Alternatively, it could be due to a poor choice of hyperparameters in the PCA or linear regression models.

To further investigate the issue, you may want to try different combinations of hyperparameters and/or explore other models and performance metrics. Additionally, you could analyze the relationships between the features and the target variable more closely, for example by plotting the data and fitting non-linear models.

```
PCA with 10 components — Average R^2 score: −1.1513
PCA with 20 components — Average R^2 score: −18.7283
PCA with 30 components — Average R^2 score: −5694493951.6642
PCA with 40 components — Average R^2 score: −3552288062163501275152384.0000
Original standardized dataset — Average R^2 score: −11978076907850751832752128.0000
```