# Data Analysis (1)

*Presentation of Data Analysis Project*
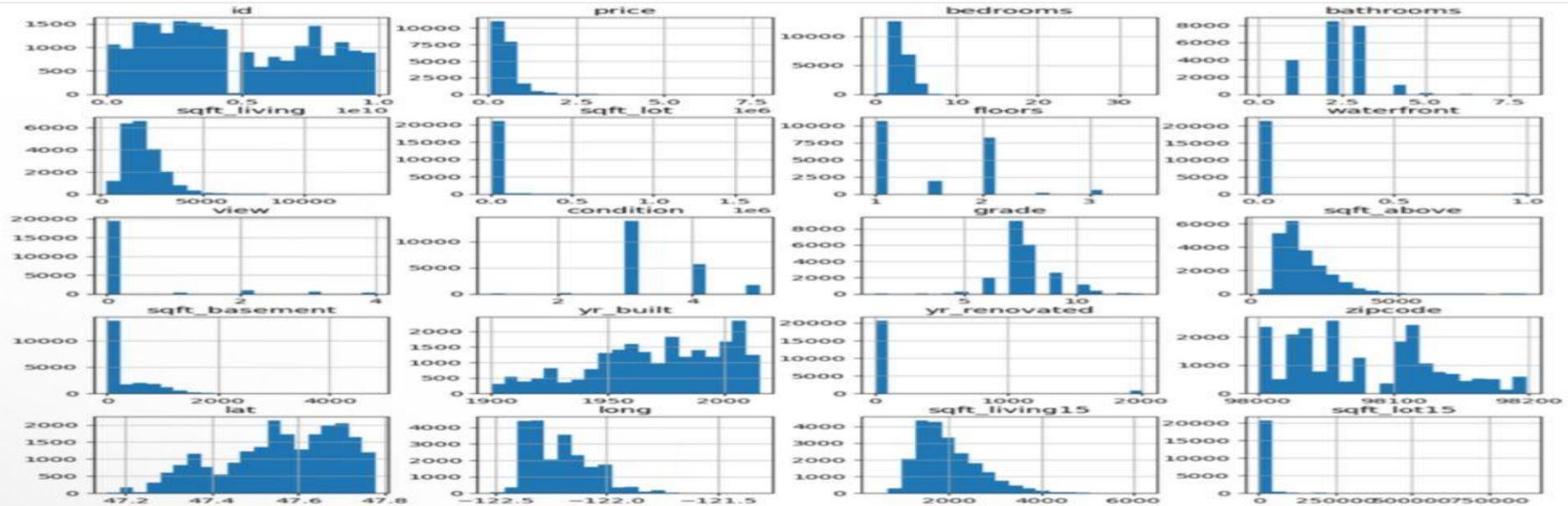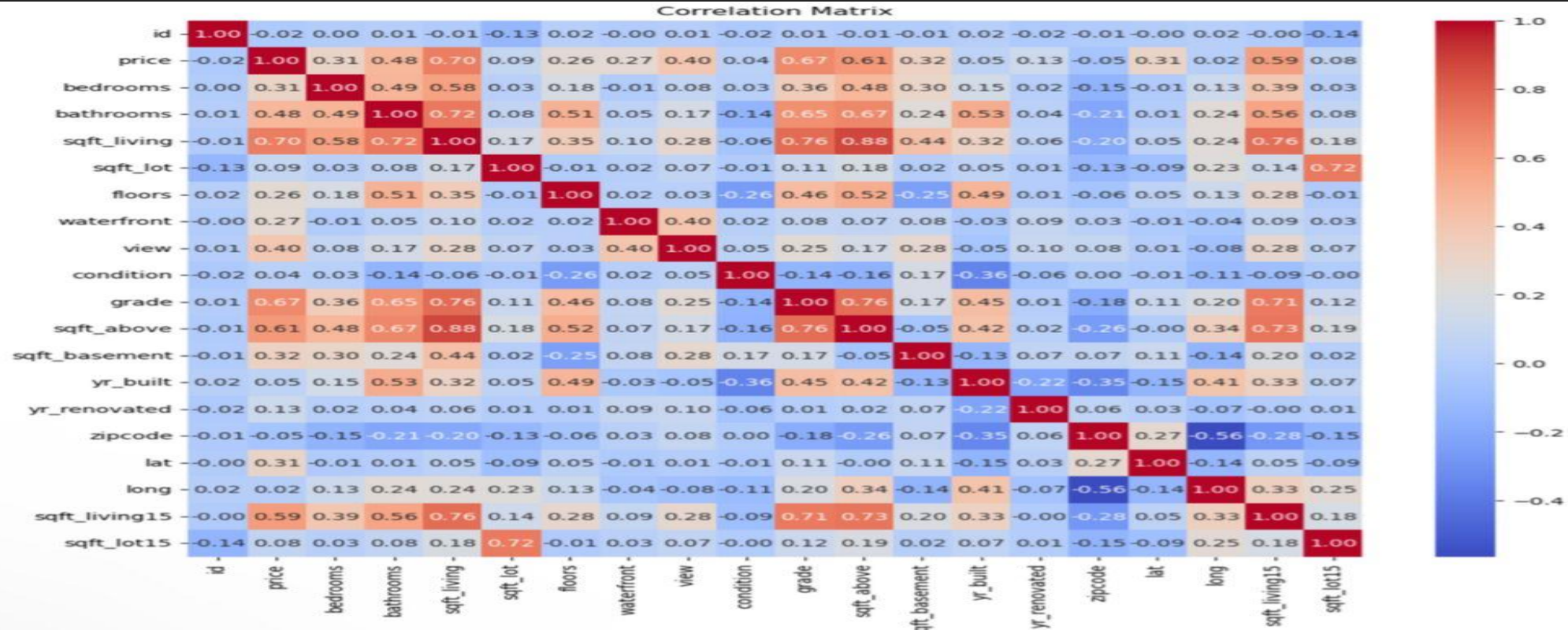
# Outline

# Introduction and Exploratory about Dataset

- First, the data set was chosen, and the choice fell on a data set related to home selling prices in King County. The data covers a period from May 2014 to May 2015. After the selection, the data was formatted, missing values were checked, and some methods were applied that make the data ready to build the model.

- there are some important features: Price, bedrooms, bathrooms, sqft_living, floors, etc…

# Exploratory about Dataset



*The histogram has been applied to all of our features in the dataset to make it easier to build the model so that we can better understand the distribution of the data and identify any potential issues or outliers that may need to be addressed before training the model.*
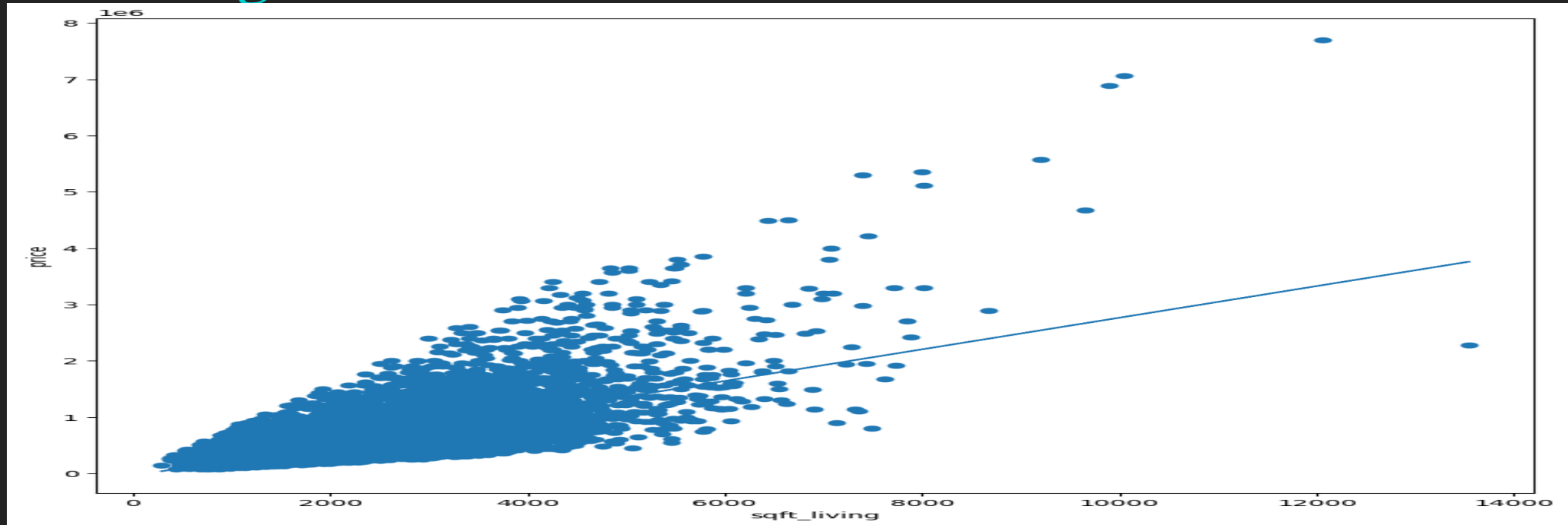
# Heatmap of our Dataset



Correlation Matrix

This correlation matrix expresses the strength of the relationship between each attribute. When the value approaches 1 and the color is red, the relationship between the attributes becomes stronger. This correlation matrix helps us in choosing the strongest relationships to use in the linear regression model.

# Outline

- ***Introduction and Exploratory about Dataset***
- ***Linear Regression model***
- ***Classification by (Decision tree)***
- ***Clustering***
- ***Anomaly Detection***
- ***PCA***

# Linear Regression model



Linear regression is a fundamental statistical method used in data science and machine learning.
Based on the linear regression analysis, there is a linear relationship between the living area (square feet) and the property price. As the living area per square foot increases or decreases, the price of the property tends to rise or fall accordingly.
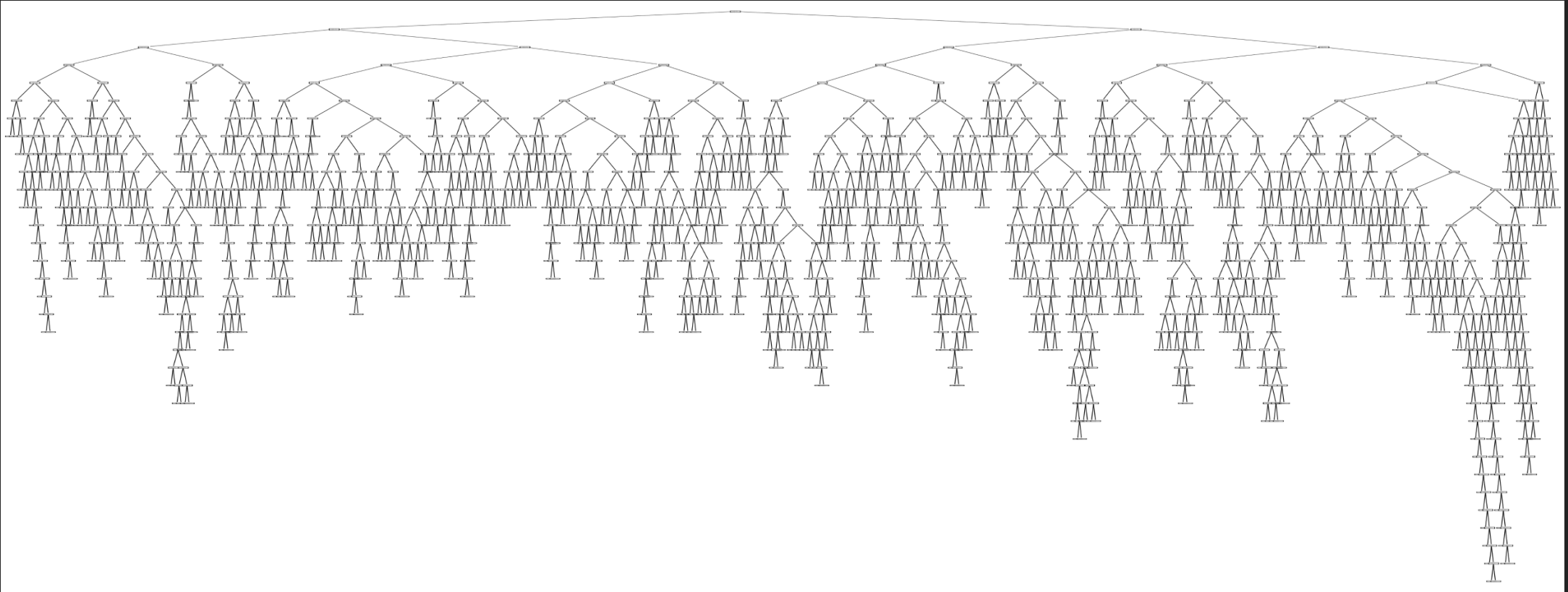The model was also tested and the results shown are:

RMSE=1.3933333680162217     and     R^2=06670504922774556

# Outline

- ***Introduction and Exploratory about Dataset***
- ***Linear Regression model***
- ***Classification by (Decision tree)***
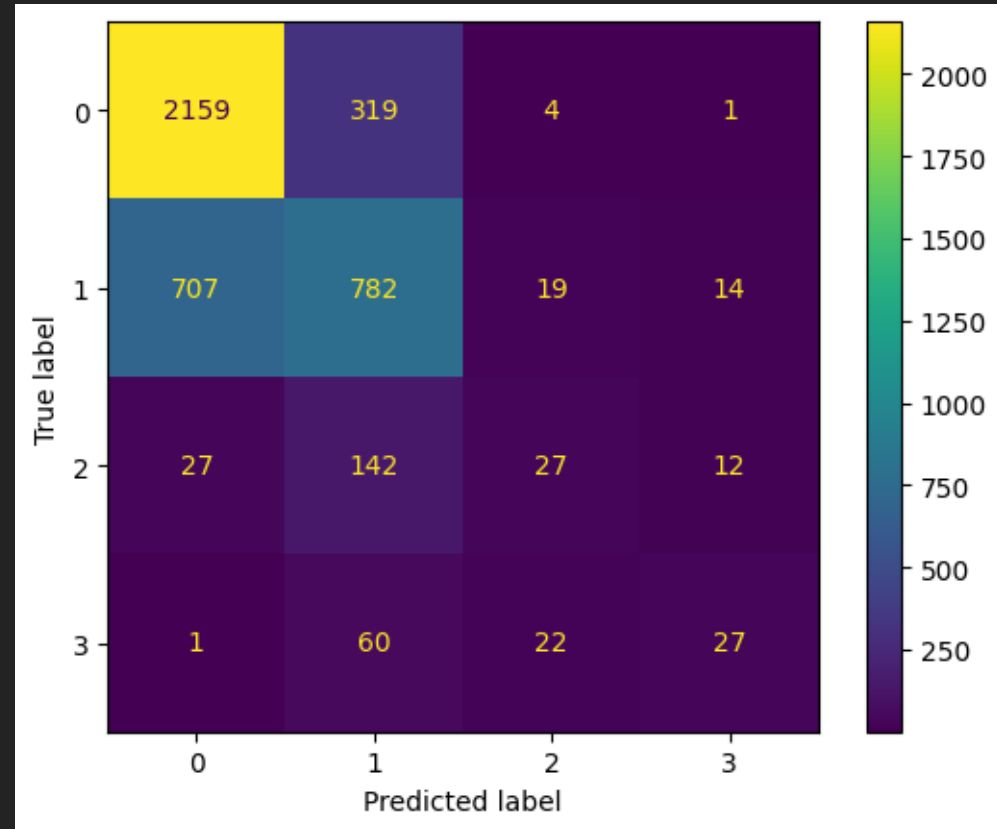- ***Clustering***
- ***Anomaly Detection***
- ***PCA***

# Classification by (Decision tree)



The Decision Tree model has been executed on the variable sqft_living. This allows us to understand how the number of sqft_living in a house impacts the predicted sale price.
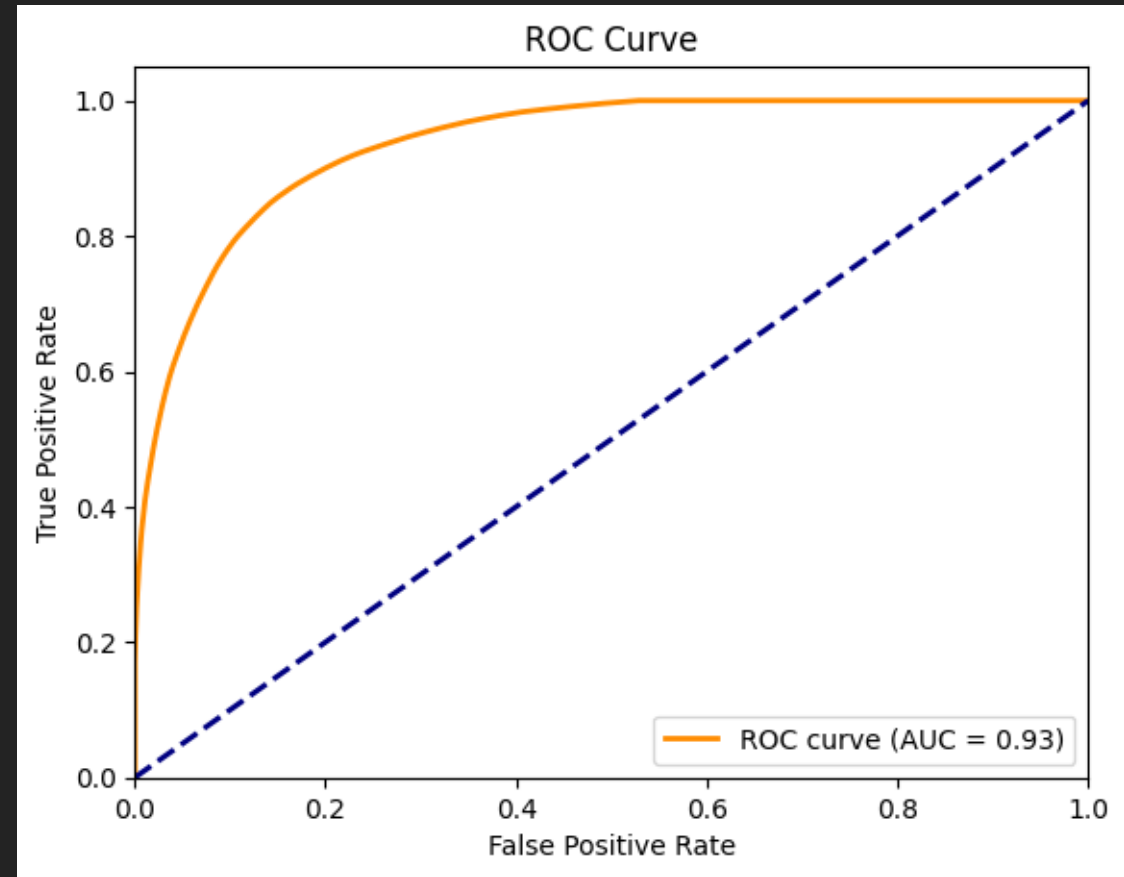
# Model Evaluation by Confusion matrix

*As we can see, most of the data is in the true positive position, which means that the model's predictive ability is high*

# Model Evaluation by ROC Method

*AUC (Area Under the Curve): 0.93 Indicates excellent performance.*

*The model performs significantly better than random guessing and is highly effective at classification.*
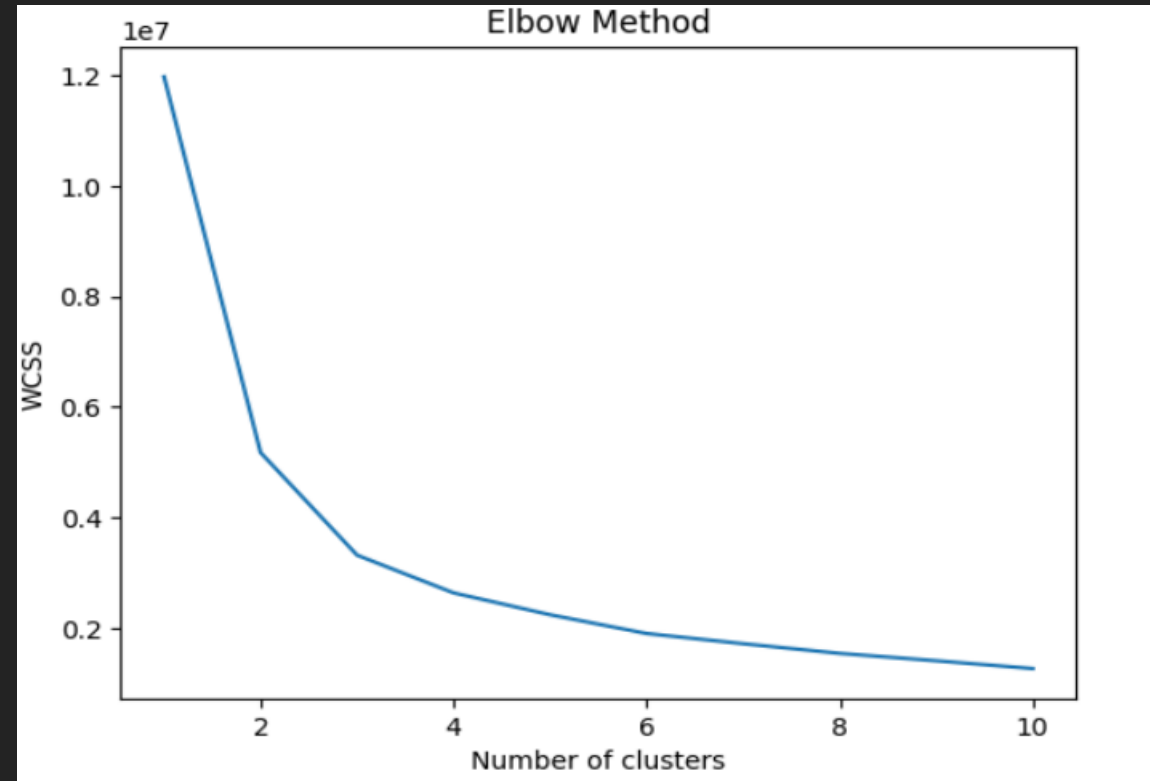
# Outline

- ***Introduction and Exploratory about Dataset***
- ***Linear Regression model***
- ***Classification by (Decision tree)***
- ***Clustering***
- ***Anomaly Detection***
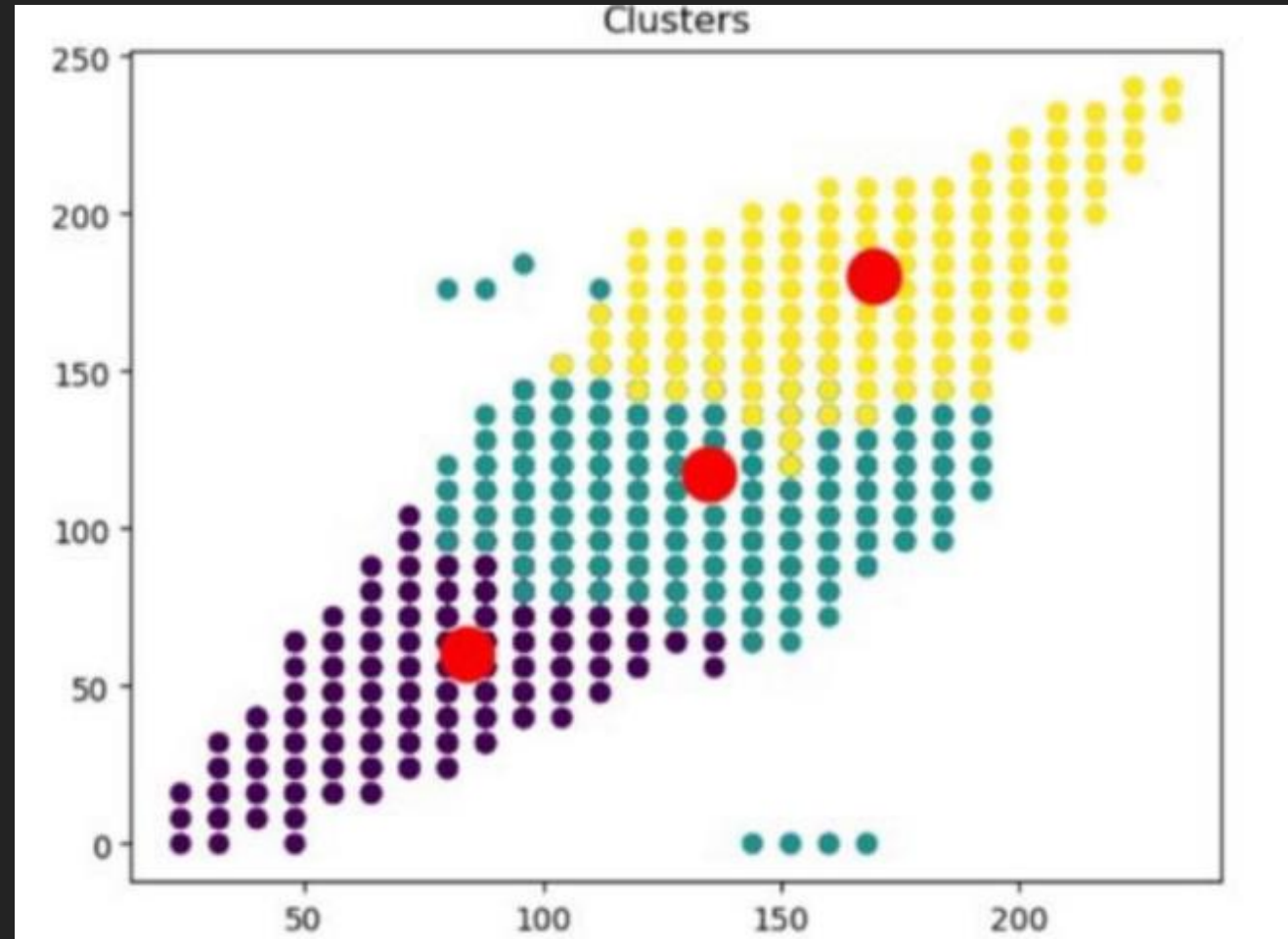- ***PCA***

# Clustering

*Elbow Method for Optimal Cluster Selection: Plot demonstrating the within-cluster sum of squares (WCSS) as a function of the number of clusters. The elbow point indicates an optimal number of clusters for K-means clustering.*

# Clustering(.cont)

Cluster Visualization: Scatter plot illustrating the clusters generated by K-means clustering algorithm. Each point represents a data instance colored by its assigned cluster, with centroids highlighted in red.
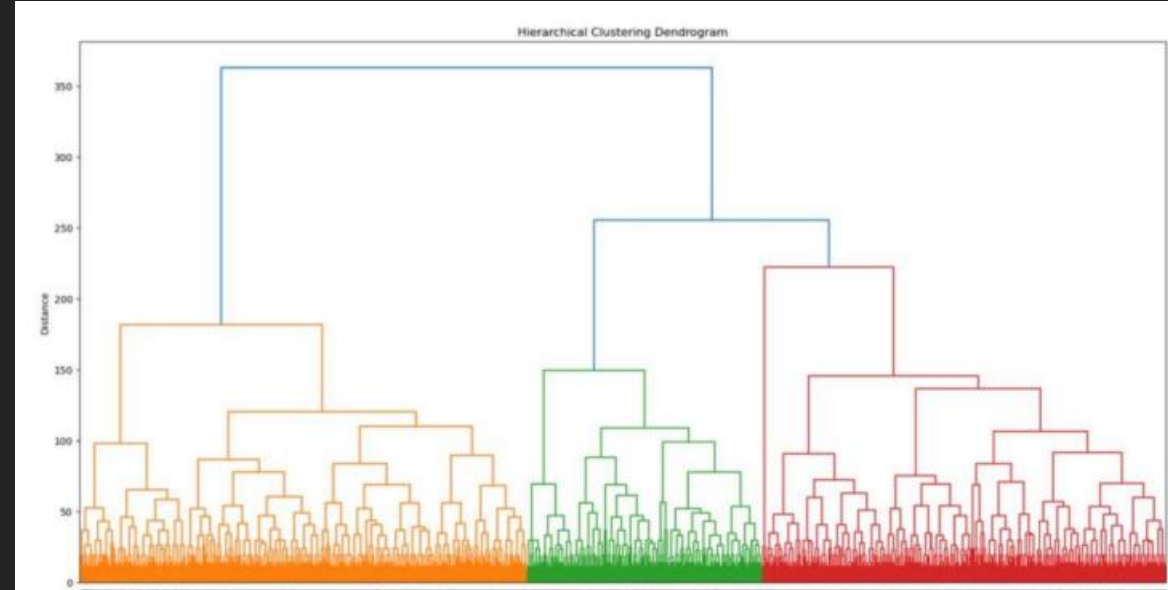Sum of Squared Error (SSE): 3318367.90223293

# Clustering(.cont)

Hierarchical clustering groups similar items together in a tree-like structure based on their similarity. It doesn't require predefining the number of clusters and is used to discover patterns in various fields. Dendrograms visualize the cluster hierarchy, Each tree represents a type of data in a specific order.
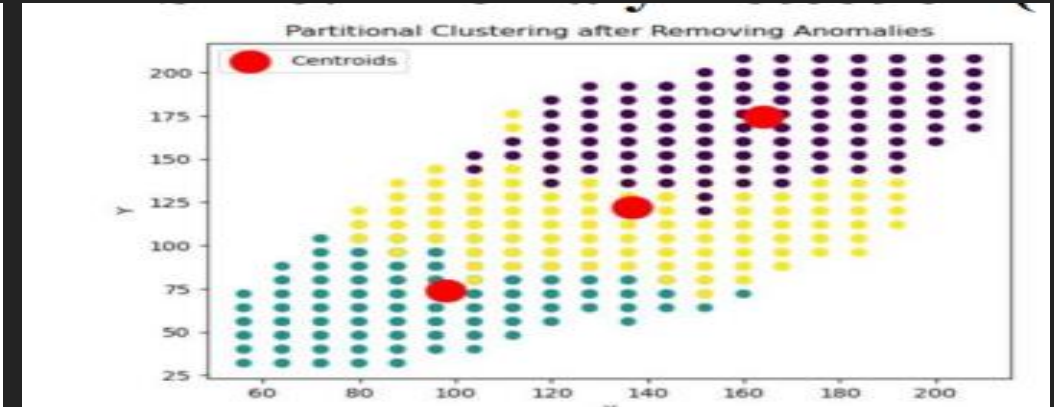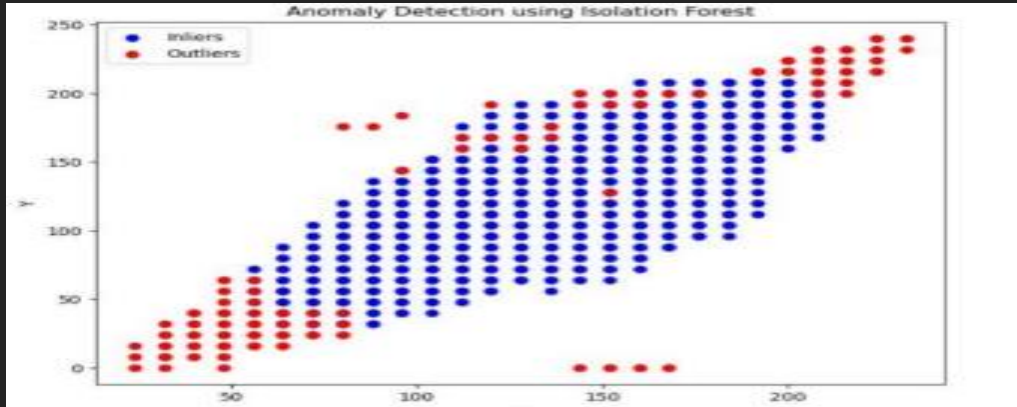
**Sum of Squared Error (SSE):**



```
SSE for test set with Euclidean distance: 24    242215.333776
0      216348.602204
40     205260.942476
dtype: float64
```

# Outline

- ***Introduction and Exploratory about Dataset***
- ***Linear Regression model***
- ***Classification by (Decision tree)***
- ***Clustering***
- ***Anomaly Detection***
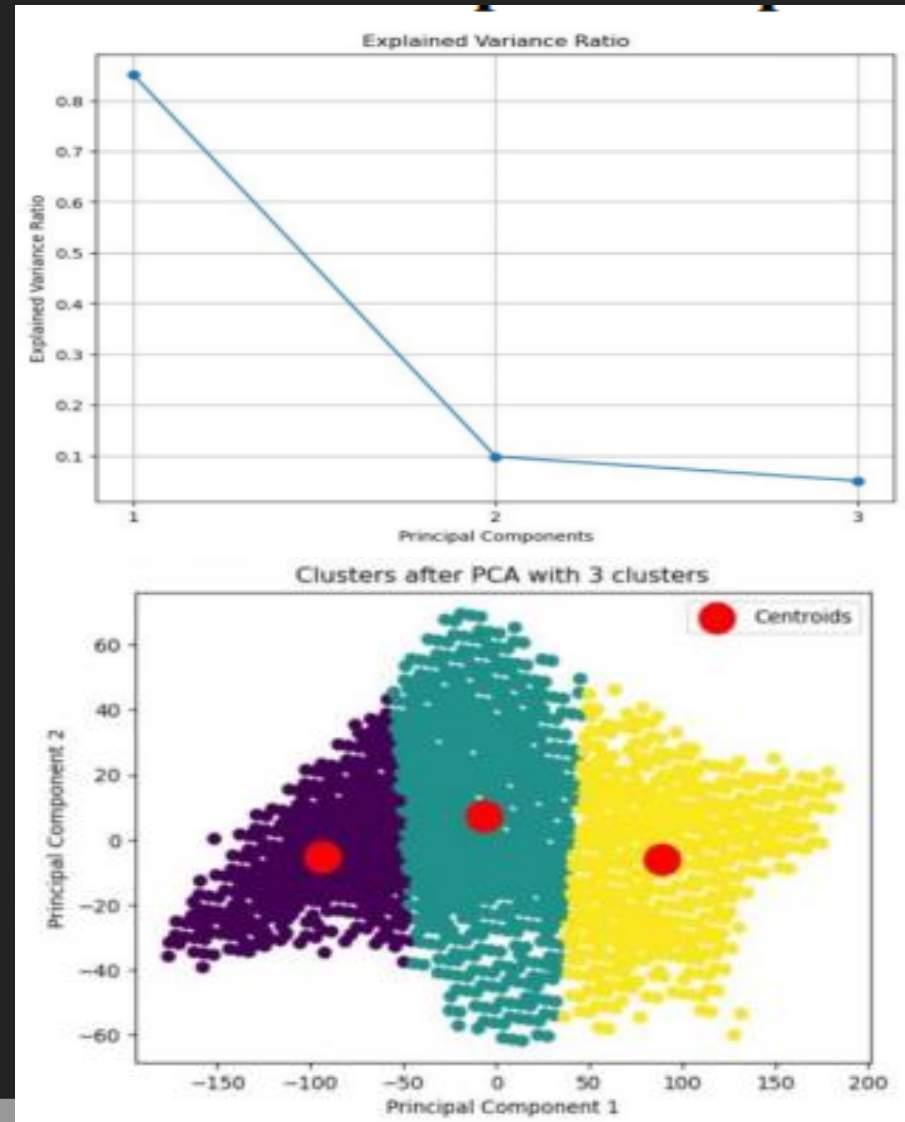- ***PCA***

# Anomaly Detection



Anomaly detection is finding unusual patterns in data. It helps to identify outliers or distortions that deviate from the norm. In this task, we used it well and were able to detect and recognize outliers, as shown in the graphics
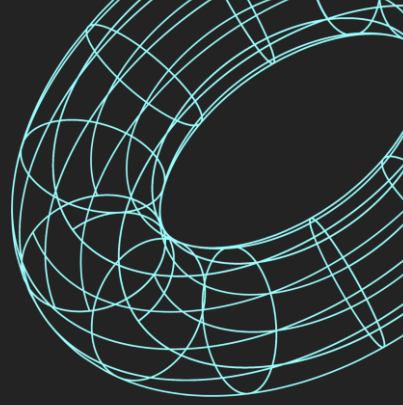
# Outline

- *__Introduction and Exploratory about Dataset__*
- *__Linear Regression model__*
- *__Classification by (Decision tree)__*
- *__Clustering__*
- *__Anomaly Detection__*
- *__PCA__*

# Principal Component Analysis (PCA)

PCA (Principal Component Analysis) is a technique used to reduce the dimensionality of data while preserving important information. It helps in finding the most important patterns in data by extracting key features. In this task, we used PCA to reduce the dimensions available in the data set, and we modified it and reduced the dimensions. This process facilitates the process of understanding the data set and making a specific decision
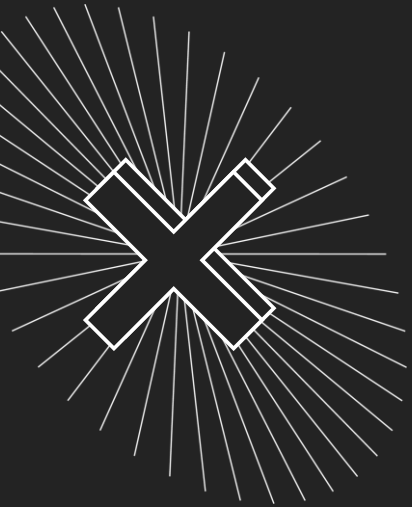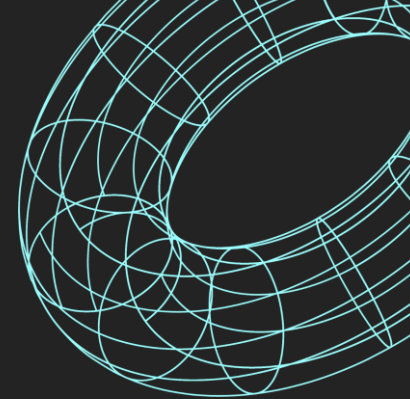
# SUBMITTED BY:

*Eyad Mohammed Alharthi:444000005*

*Mohammed Nedal Alshareef:444004824*

*Ziyad omar altalhi:444000039*

*Aseel Abdulfattah Almaghrabi:444004296*

*Raed Falah Alhelali:444006561*

# Thanks for listening!