# HOUSING PRICES PREDICTION

## USING DATA MINING

# OBJECTIVES

**1** **PROBLEM**

**2** **DATA**
data overview , and data graphs

**3** **DATA PREPROCESSING**
data cleaning , and data transformation

**4** **DATA MINING TASK**
data mining Techniques and findings

# PROBLEM:

Since housing is one of the most significant investments in a person's life, using data mining techniques can help predict housing prices and make informed decisions.

In our project, we analyzed a housing dataset, which provided valuable insights into the factors influencing housing prices. This data enabled us to identify key patterns and challenges, such as multicollinearity, and apply effective techniques to improve prediction accuracy.

# OUR DATA

- **We applied our data mining tasks on data set consisting of :**

- **Number of tupels : 545**
- **13 attibutes which are:**

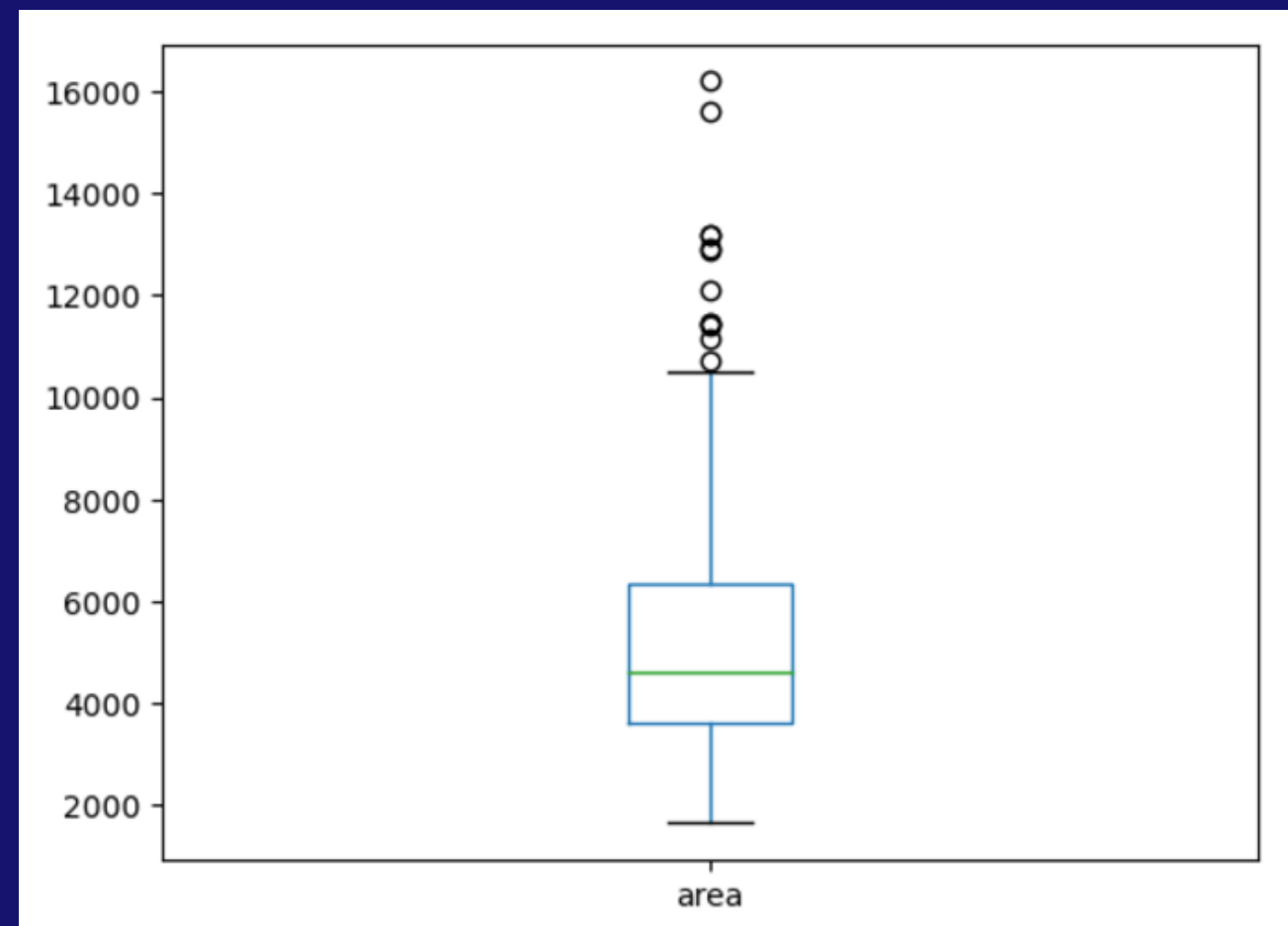| Price | Area | bedrooms | bathrooms | stories |
| mainroad | guestroom | basement | parking | prefarea |

| Hotwater heating | Air conditioning | Furnishing status |

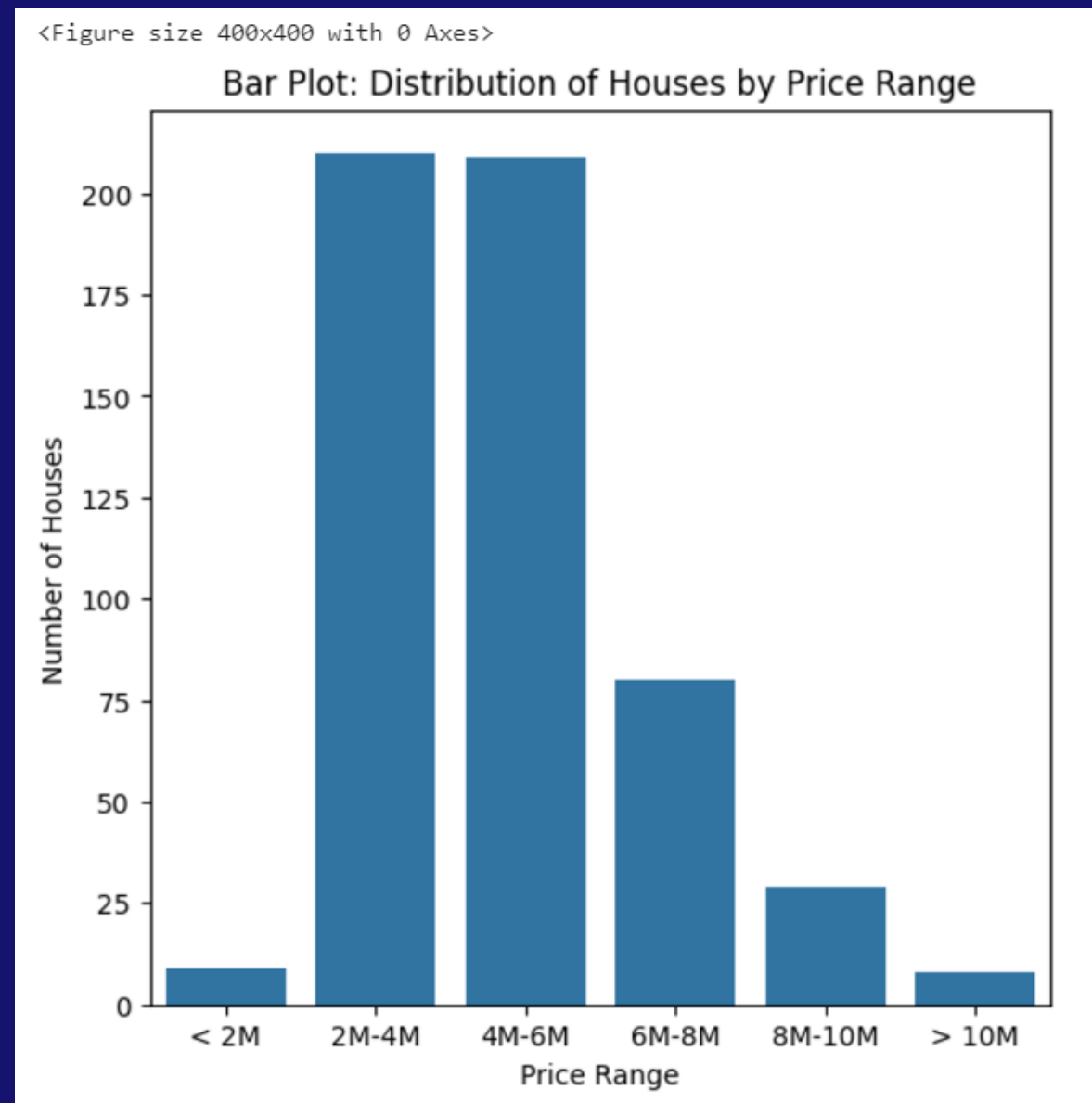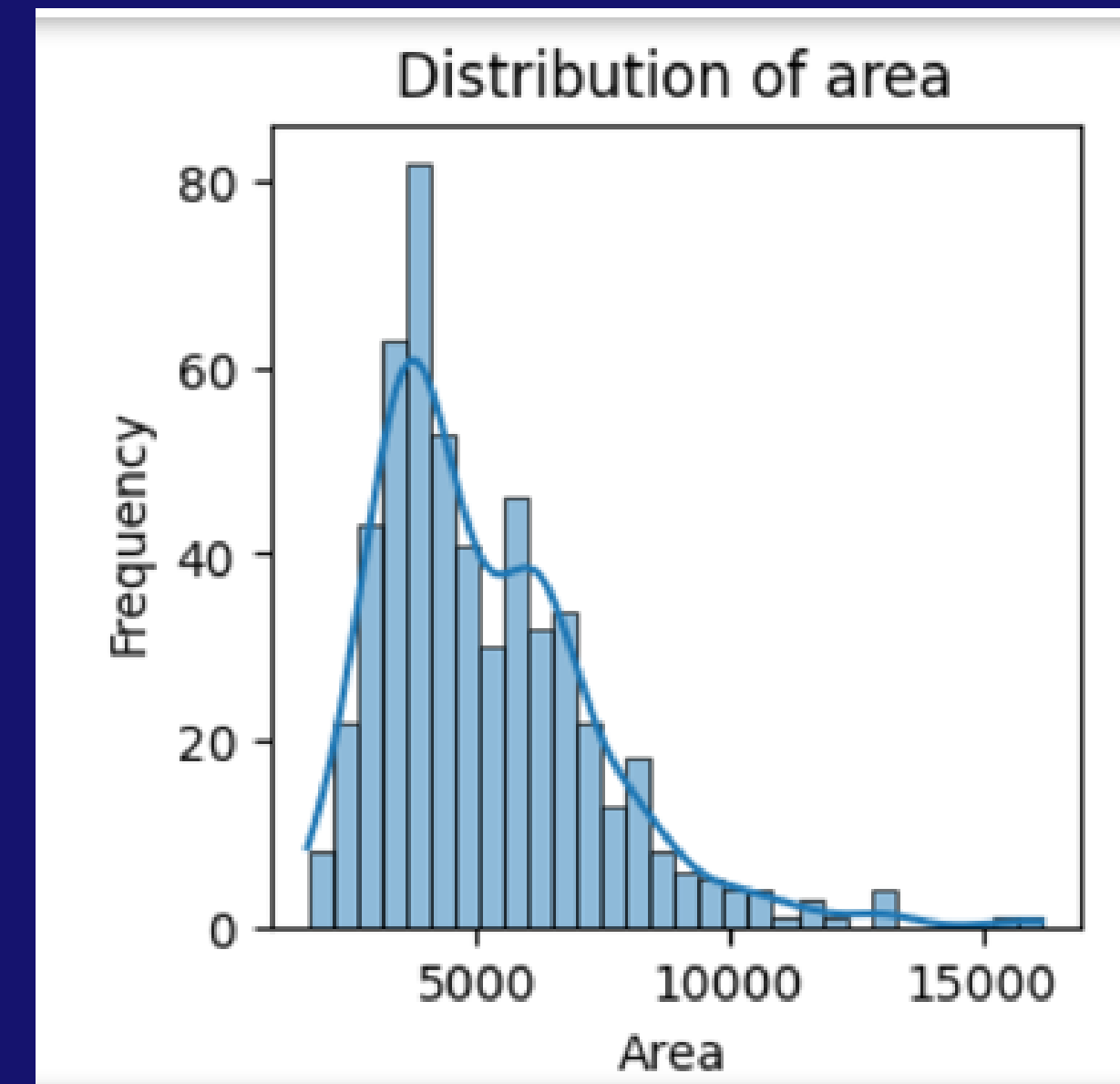# DATA GRAPHS:



**Box plot**



**Scatter plot**

# DATA GRAPHS:



**Bar plot**



**Histogram**

# DATA PREPROCESSING:

- **To achieve optimal accuracy, various preprocessing techniques were applied to enhance data efficiency, including:**

**1**

**DATA CLEANINIG**

**2**

**DATA TRANSFORMATION**

# Data cleaning:

- We checked our data for missing or null values and found none.

- We identified outliers in the numeric attributes and removed the rows containing these outliers.

# Min-Max Normalization:

- **The area attribute was scaled to a range between 0 and 1 to standardize its influence in modeling.**

- **This ensures all features have equal weight during analysis and prevents any feature from dominating due to its scale.**

# Min-Max Normalization:

| area | b |
|------|---|
| 7420 | |
| 8960 | |
| 9960 | |
| 7500 | |

→

| area |
|------|
| 0.491525 |
| 0.559322 |
| 0.333333 |
| 0.538983 |
| 0.301695 |

# Encoding:

- **In our dataset, we transformed categorical attributes into numerical representations to make them usable for machine learning models:**

- **Attributes like mainroad, guestroom were encoded as 0 for "No" and 1 for "Yes".these outliers.**

# Encoding:

| mainroad | guestroom | basement | hotwaterheating | airconditioning |
|---|---|---|---|---|
| yes | no | no | no | yes |
| yes | no | no | no | yes |
| yes | no | yes | no | no |
| yes | no | yes | no | yes |
| yes | yes | yes | no | yes |

| mainroad | guestroom | basement | hotwaterheating | airconditioning |
|---|---|---|---|---|
| 1 | 0 | 1 | 0 | 0 |
| 1 | 1 | 1 | 0 | 1 |
| 1 | 1 | 0 | 0 | 1 |
| 1 | 0 | 0 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 |

# DATA MINING TASK : CLASSIFICATION

Classification is a supervised learning technique, which means it requires a class label to classify objects.

We trained our model to be able to predict if the price was high or medium or low using (price) class labe

# DATA MINING TASK : CLASSIFICATION

To build our model We used a decision tree algorithm based on ( IG and gini index )which is a recursive algorithm produces a tree with a leaf nodes representing the final decisions.

The final decision for class price is either true or false. the prediction is made on the rest attributes.
('area', 'bedrooms', 'bathrooms', 'airconditioning', 'prefarea')

# DATA MINING TASK : CLASSIFICATION

- **This technique includes dividing the dataset into two sets:**

**1**

**TRAINING DATASET:**
USED FOR BUILDING
THE DECISION TREE.

**2**

**TESTING DATASET:**
USED TO EVALUATE THE
CONSTRUCTED MODEL.

- **We tried 3 different sizes of training and testing data to get the best result for construction and evaluation.**

**1**

90% TRAINING
10% TESTING

**2**

80% TRAINING
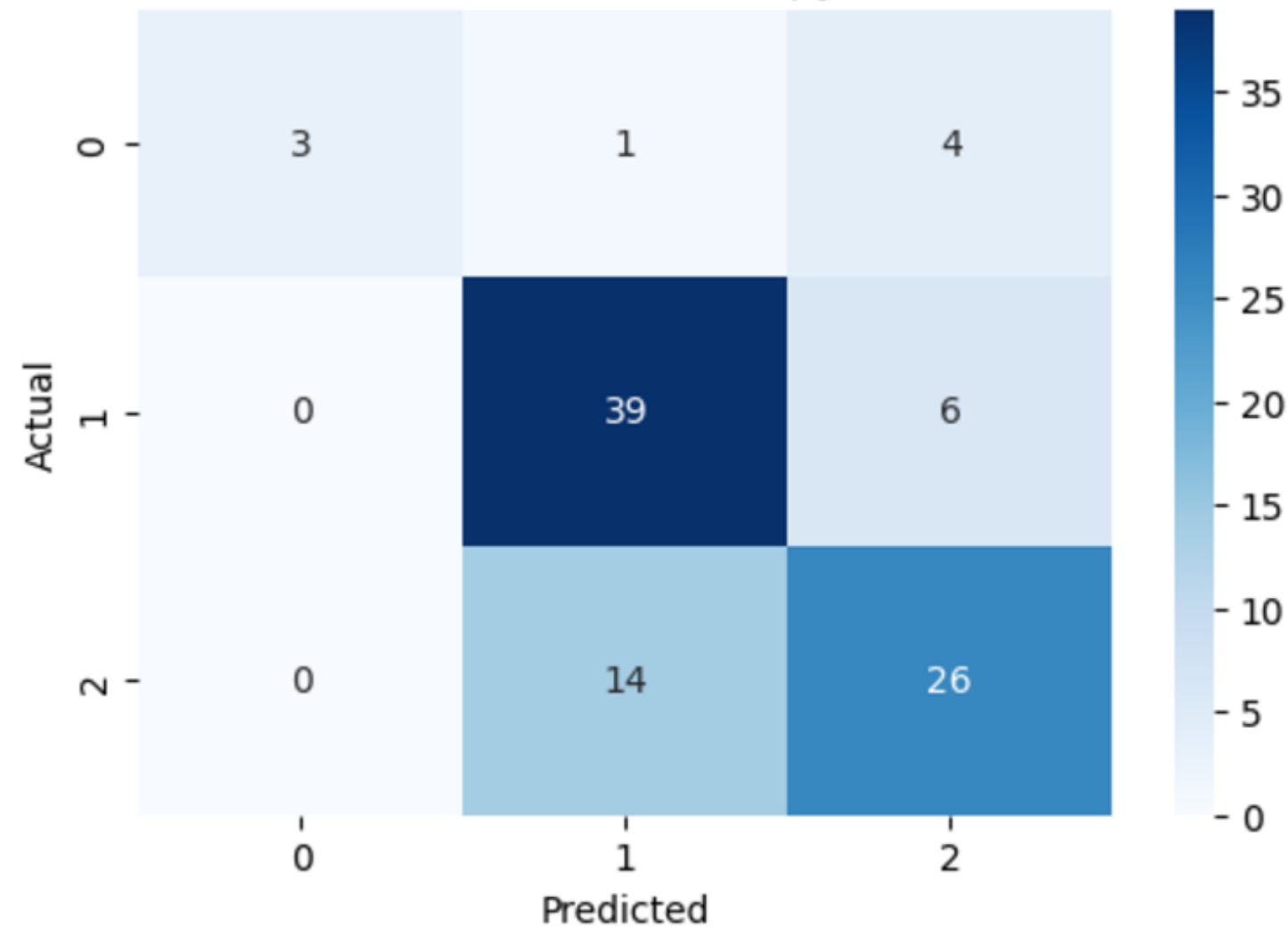20% TESTING
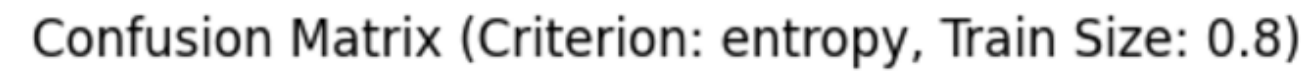
**3**

70% TRAINING
30% TESTING

# ACCURACY COMPARISON

| | 90-10 | | 80-20 | | 70-30 | |
|---|---|---|---|---|---|---|
| | IG | Gini | IG | Gini | IG | Gini |
| Accuracy | 0.72 | 0.68 | 0.73 | 0.68 | 0.58 | 0.66 |

| Train Size | Splitting Criterion | Class Label | Precision | Recall | Specificity |
|---|---|---|---|---|---|
| 0.7 | entropy | High | 30.8 | 28.6 | 92.8 |
| | | Low | 62.4 | 80.3 | 56.2 |
| | | Medium | 56.1 | 39 | 77.5 |
| | gini | High | 100 | 14.3 | 100 |
| | | Low | 67.1 | 83.3 | 63 |
| | | Medium | 63.6 | 59.3 | 75 |
| 0.8 | entropy | High | 100 | 37.5 | 100 |
| | | Low | 72.2 | 86.7 | 68.8 |
| | | Medium | 72.2 | 65 | 81.1 |
| | gini | High | 66.7 | 50 | 97.6 |
| | | Low | 66.7 | 84.4 | 60.4 |
| | | Medium | 70 | 52.5 | 83 |
| 0.9 | entropy | High | 100 | 33.3 | 100 |
| | | Low | 69.2 | 81.8 | 68 |
| | | Medium | 75 | 68.2 | 80 |
| | gini | High | 25 | 33.3 | 93.2 |
| | | Low | 66.7 | 90.9 | 60 |
| | | Medium | 84.6 | 50 | 92 |

# FINDINGS

We evaluate our model by measuring the accuracy, Precision, Recall and Specificity measures of the testing dataset

overall the best Best Train Size is 0.8 (80% training, 20% testing) and the Best Algorithm is Entropy Criterion (Information gain)

Confusion Matrix (Criterion: entropy, Train Size: 0.8)

# DATA MINING TASK : CLUESTRING

Clustering is an unsupervised machine learning technique used to group data points into clusters based on their similarity and dissimilarity.

Our model will create a set of clusters for houses with similar characteristics, These clusters reveal market trends and help predict characteristics or pricing for new properties.
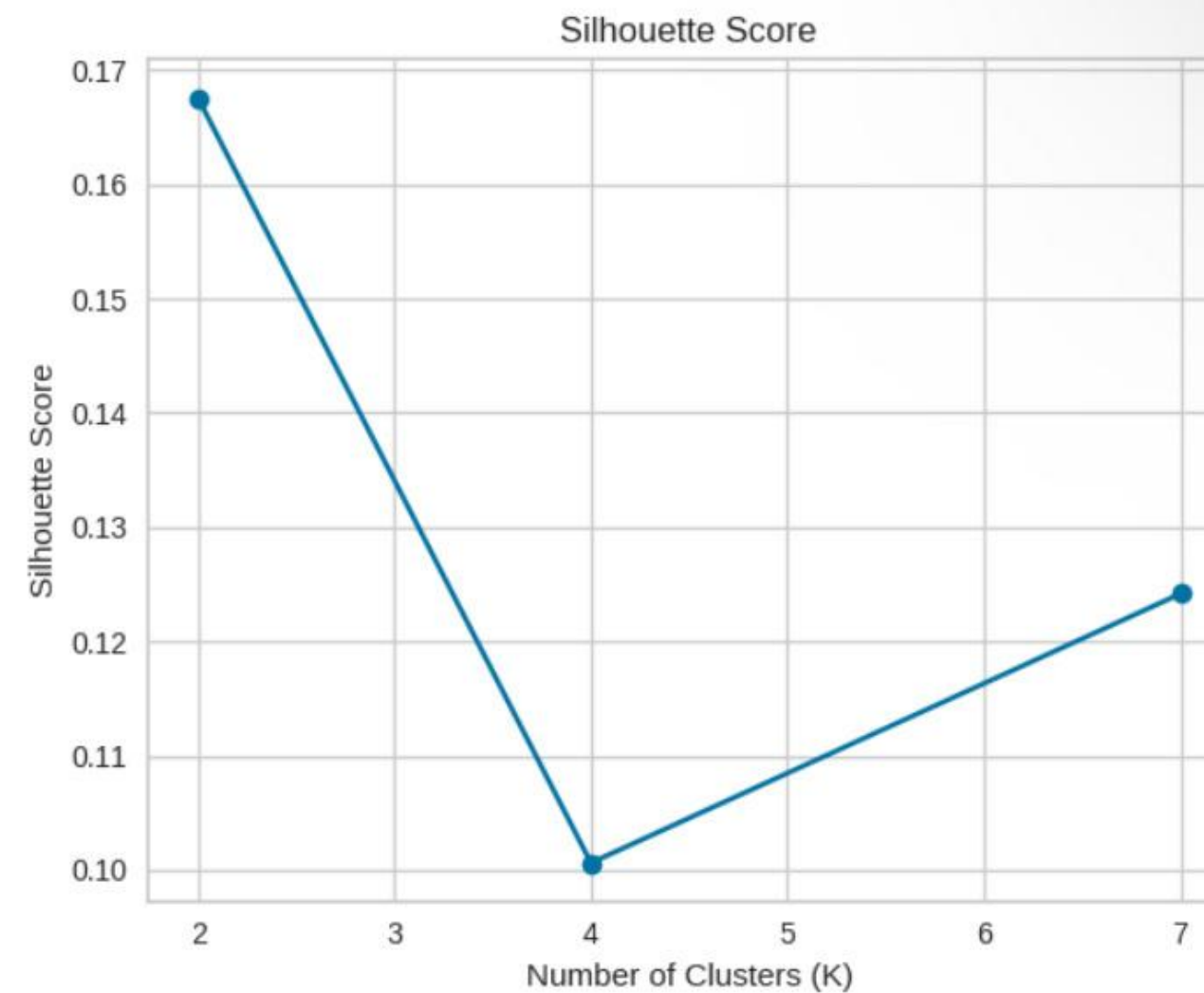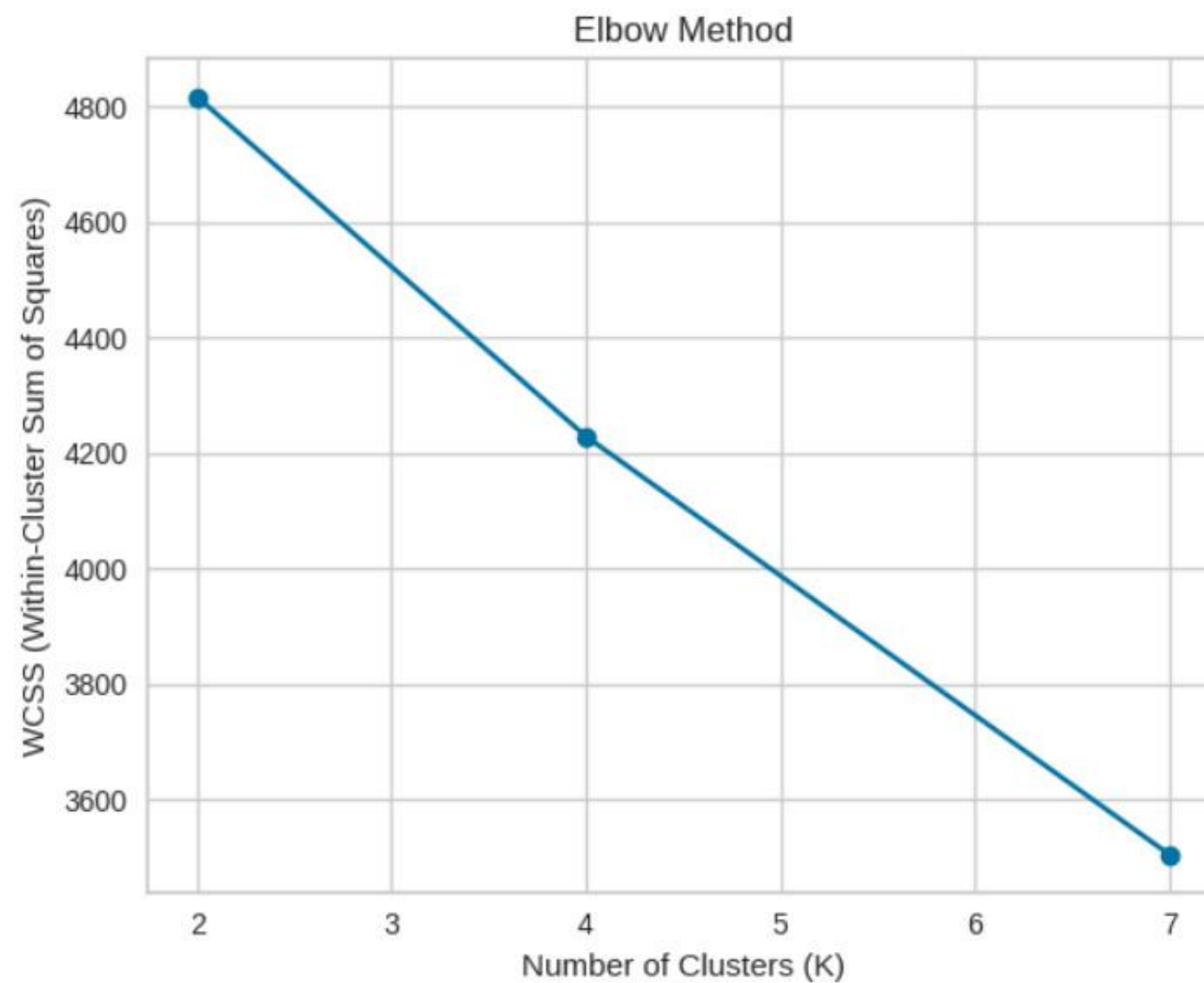
# DATA MINING TASK : CLUESTRING

We used the K-Means algorithm, which groups data points into K clusters, where each cluster is represented by a central point . The algorithm assigns each data point to the cluster with the nearest centroid.

then iteratively recalculates the center and reassigns data points until the centroids no longer change, indicating that the data points are correctly grouped.
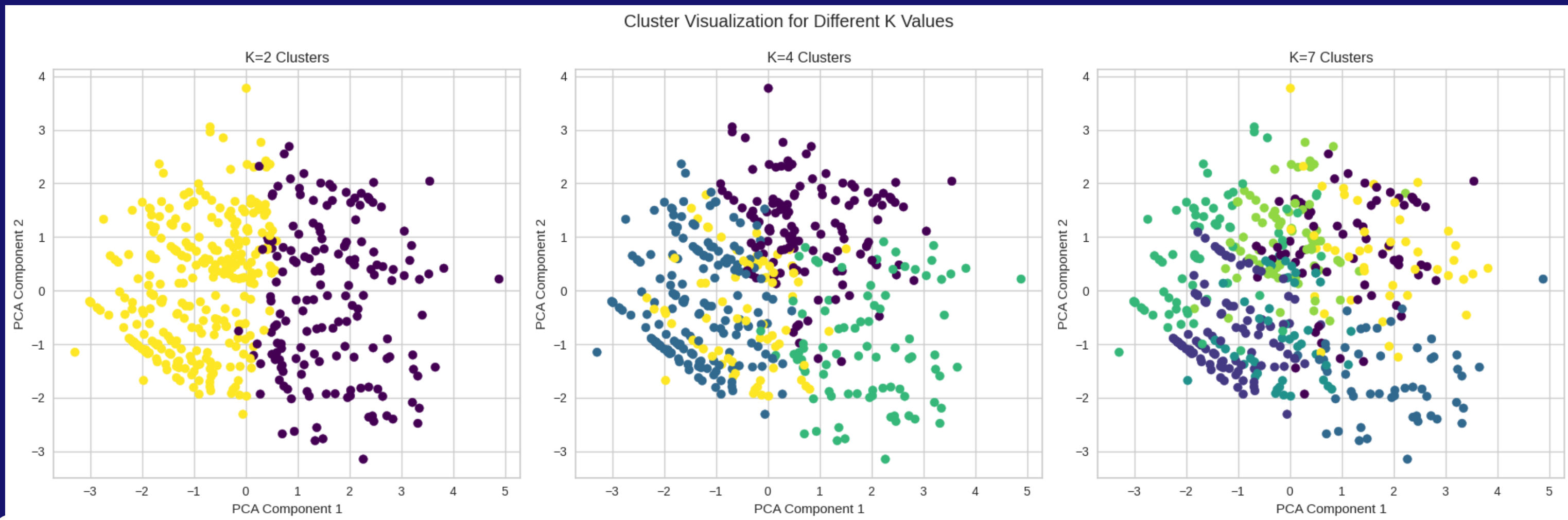
# DATA MINING TASK : CLUESTRIN G



K=2: WCSS=4816.887852648317, Silhouette Score=0.16758863832753448
K=4: WCSS=4229.139404394767, Silhouette Score=0.10061890960368892
K=7: WCSS=3505.3898179922026, Silhouette Score=0.1241750788552499

# FINDINGS

We tested three different values K, evaluating the models using the Silhouette Width and Within-Cluster Sum of Squares (WCSS) metrics.

The optimal model is K=4, providing a balance between compactness and meaningful segmentation of the data.



Cluster Visualization for Different K Values

# THANK YOU FOR LISTENING!

Team members:

Haifa Alsaif - Flwah Alrashed

Majd Alruways - Norah Alfaheed

Shouq Altamimi

Supervised by: Dr. Hessah Alsaaran