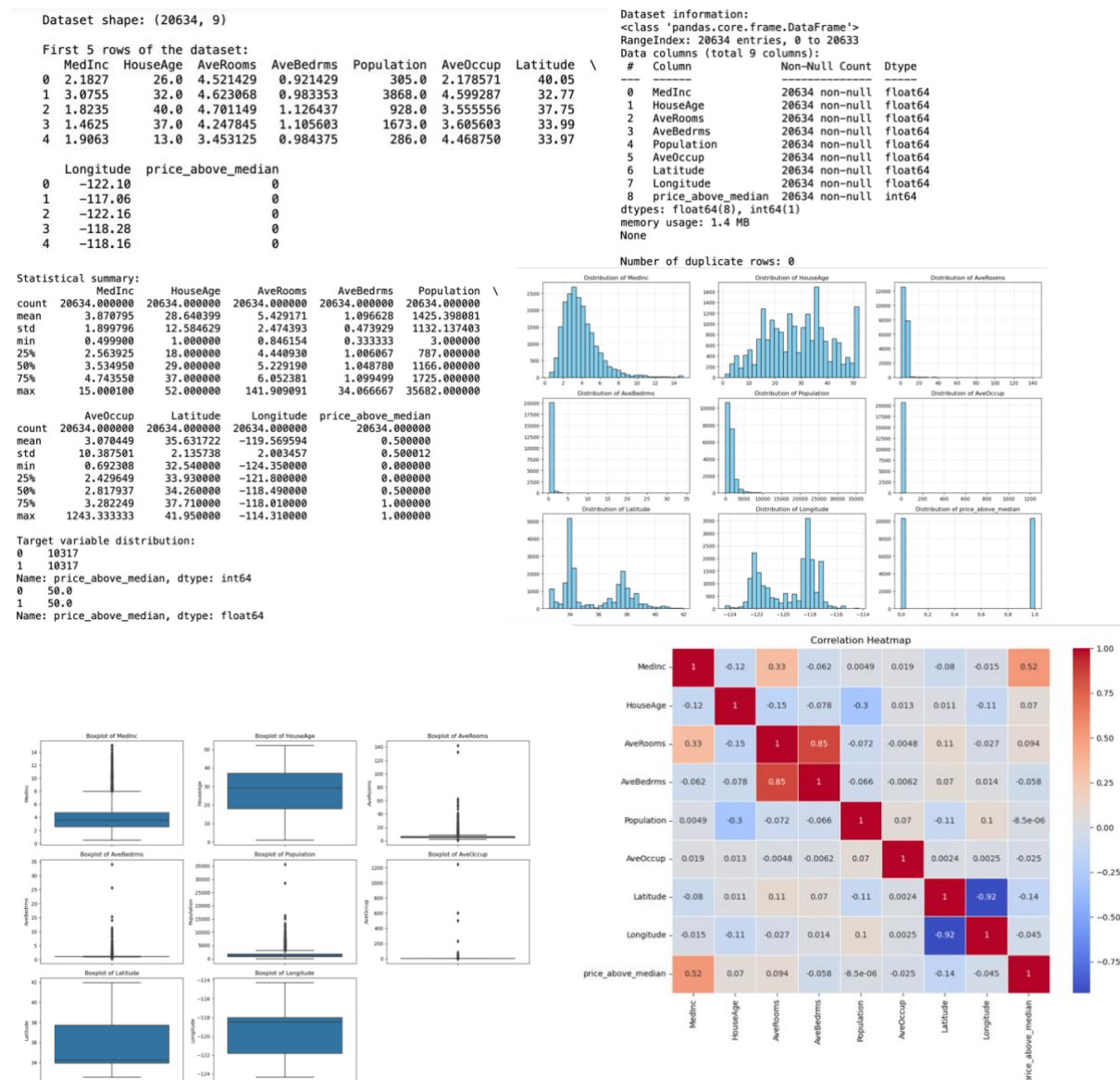


Part 1

The dataset contains 20,634 rows and 9 columns, including 8 features and 1 target variable. The features represent different characteristics of housing blocks in California, and the target variable indicates whether the house price is above the median.

All features in the dataset are of floating-point type, but the target variable is an integer type (0 or 1). The analysis shows no duplicate rows in the dataset, so no data cleaning was required.



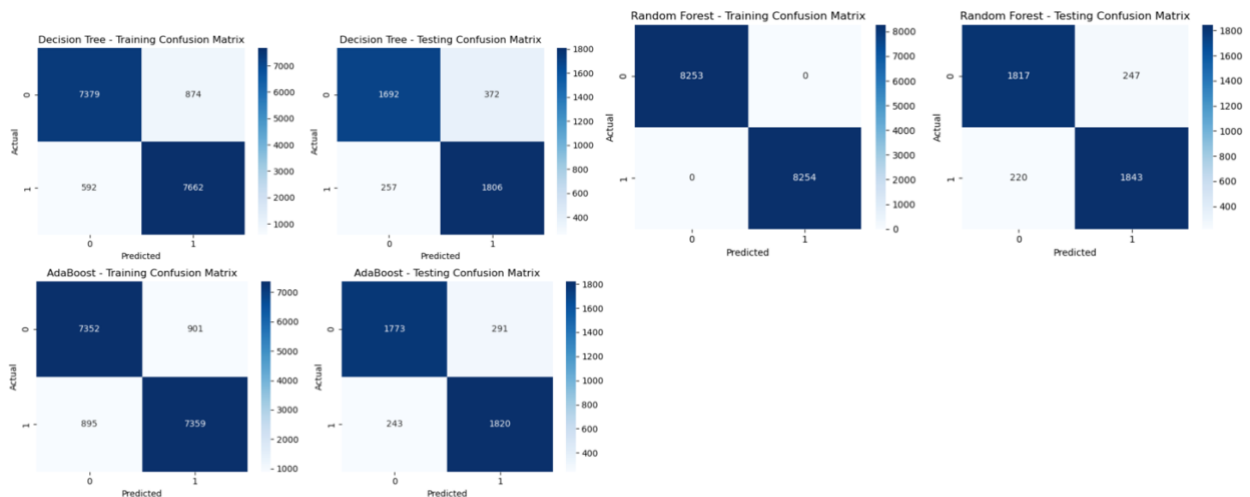
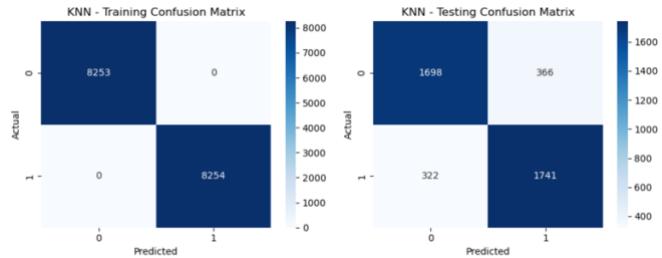
Part 2

The dataset was split into training and test sets. Training set: 16,506 instances (80%) and test set: 4,128 instances (20%). The class distribution was preserved in both sets, with exactly 50% of instances in each class. I used KNN, Decision Tree, Random Forest and AdaBoost classifier.

```

Training set shape: (16507, 8)
Testing set shape: (4127, 8)
Training set class distribution:
0    50.003029
1    49.996971
Name: price_above_median, dtype: float64
Testing set class distribution:
0    50.012115
1    49.987885
Name: price_above_median, dtype: float64

```



Model Performance Comparison:

Model Metrics Summary:

	Model	Train Accuracy	Test Accuracy	Precision	Recall	F1 Score
0	KNN	1.000000	0.833293	0.826293	0.843917	0.835012
1	Decision Tree	0.911189	0.847589	0.829201	0.875424	0.851686
2	Random Forest	1.000000	0.886843	0.881818	0.893359	0.887551
3	AdaBoost	0.891198	0.870608	0.862151	0.882210	0.872065

Best performing model based on test accuracy: Random Forest

Part 3

1. Which techniques did you use to train the models?

All models were trained using a stratified train-test split way to ensure balanced class representation in both sets.

KNN was used with distance-weighted voting and an optimal $k=11$ neighbors. Decision Tree classifier was used with parameters including $\text{max_depth}=10$ and $\text{min_samples_split}=10$ to control complexity. Random Forest combined 200 decision trees with a maximum depth of 20 to

improve predictive performance. AdaBoost was implemented with a learning rate of 1.0 and 200 estimators. And I asked chatgpt for helping me find the best ways to train the models.

2.Explain any techniques used to optimize model performance?

For hyperparameter tuning, I used Grid search CV with the help of chatgpt. And I also used Feature Standardization so that all features were standardized.

3.Compare the performance of all models to predict the dependent variable?

Random Forest has results with 88.68% test accuracy, AdaBoost has 87.06% , Decision Tree 84.76%, and KNN 83.33%. Random Forest also has 88.18% precision, 89.34% recall, and 88.76% F1 score.

Both KNN and Random Forest achieved 100% training accuracy, but Random Forest showed better generalization on test data. Decision Tree showed signs of overfitting with a notable gap between training and testing accuracies. AdaBoost has good generalization with minimal difference between training and testing performance. All models maintained balanced performance across both housing price classes.

4.Which model would you recommend to be used for this dataset?

Random Forest is recommended.

5.For this dataset, which metric is more important, why?

I think precision is more important because false classifying could lead to significant financial losses for buyers and investors.