**MSc in Software Design with Artificial Intelligence  — Continuous Assessment Project (40%)**

Dataset Links:
- Regression Dataset: Auto MPG (UCI Machine Learning Repository):
https://archive.ics.uci.edu/ml/datasets/auto+mpg

- Classification Dataset: Breast Cancer Wisconsin (Diagnostic):
https://archive.ics.uci.edu/ml/datasets/banknote+authentication

# 1. Regression

## 1.1 Objectives

Predicting an automobile's fuel efficiency (miles per gallon) based on its engine and physical characteristics, including horsepower, weight, and displacement, is the aim of this analysis. From a business standpoint, this aids analysts and manufacturers in designing fuel-efficient cars and optimising engine performance.

Target Variable: mpg (Miles per Gallon)
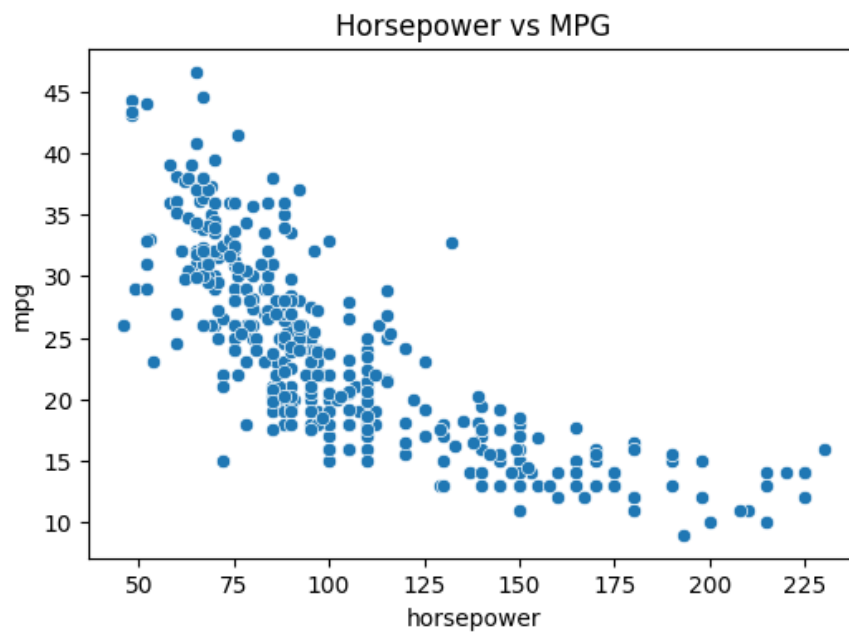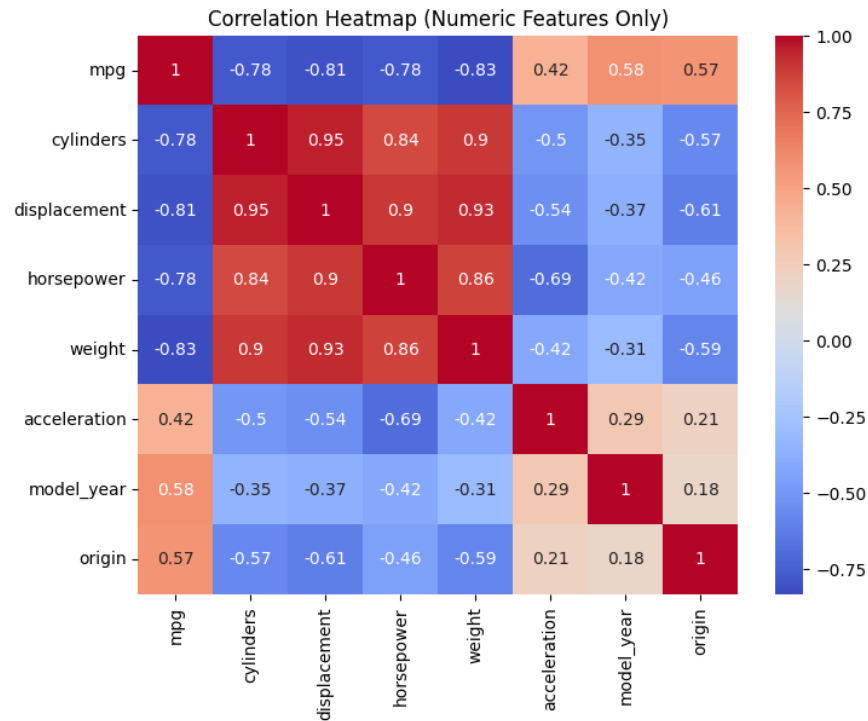
Target variable: mpg.

## 1.2 Data Exploration

In addition to numerical characteristics like weight, horsepower, and cylinders, the Auto MPG dataset also includes a qualitative item called "car_name," which was eliminated for modelling.
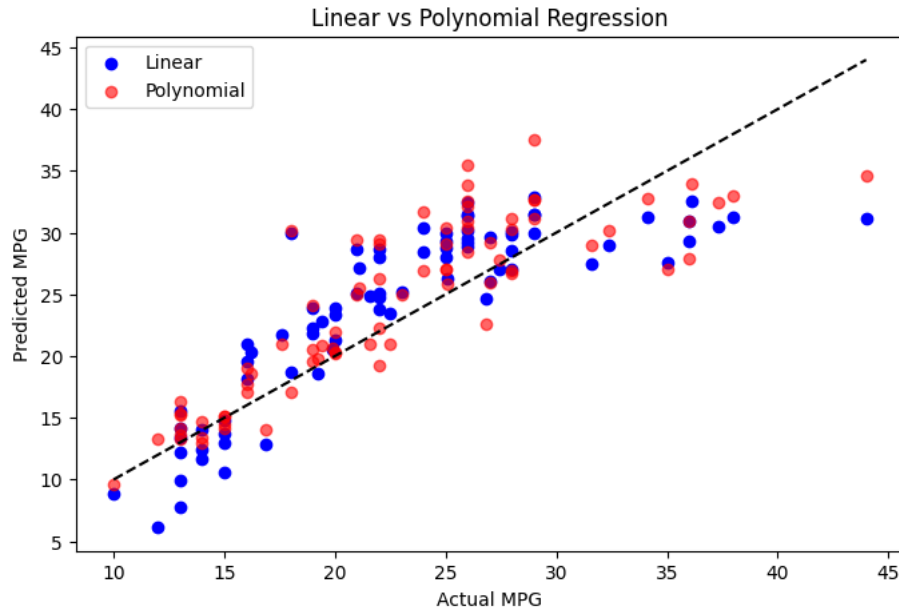
```python
# Clean data
df = df.replace('?', np.nan)
df = df.dropna()
df['horsepower'] = df['horsepower'].astype(float)

# Drop non-numeric column
df = df.drop(columns=['car_name'])
```

According to a heatmap, horsepower and weight have a significant negative correlation with mpg, suggesting that bigger cars typically have lower fuel economy.



Correlation Heatmap (Numeric Features Only)



Horsepower vs MPG

This linear trend was validated by scatter plots.



Linear vs Polynomial Regression

## 1.3 Modeling

The dataset was split into training and testing sets using train_test_split() with random_state=42 to ensure reproducibility.

Two models were trained:

1. **Linear Regression**
2. **Polynomial Regression (degree=2)** to capture potential nonlinearity.

```python
# Scatter plot
plt.figure(figsize=(6,4))
sns.scatterplot(data=df, x='horsepower', y='mpg')
plt.title('Horsepower vs MPG')
plt.show()

# feature and target selection
X = df[['horsepower', 'weight', 'cylinders', 'displacement']]
y = df['mpg']

X_train, X_test, y_train,    (parameter) random_state: Int | RandomState | None
    X, y, test_size=0.2, random_state=42)
```
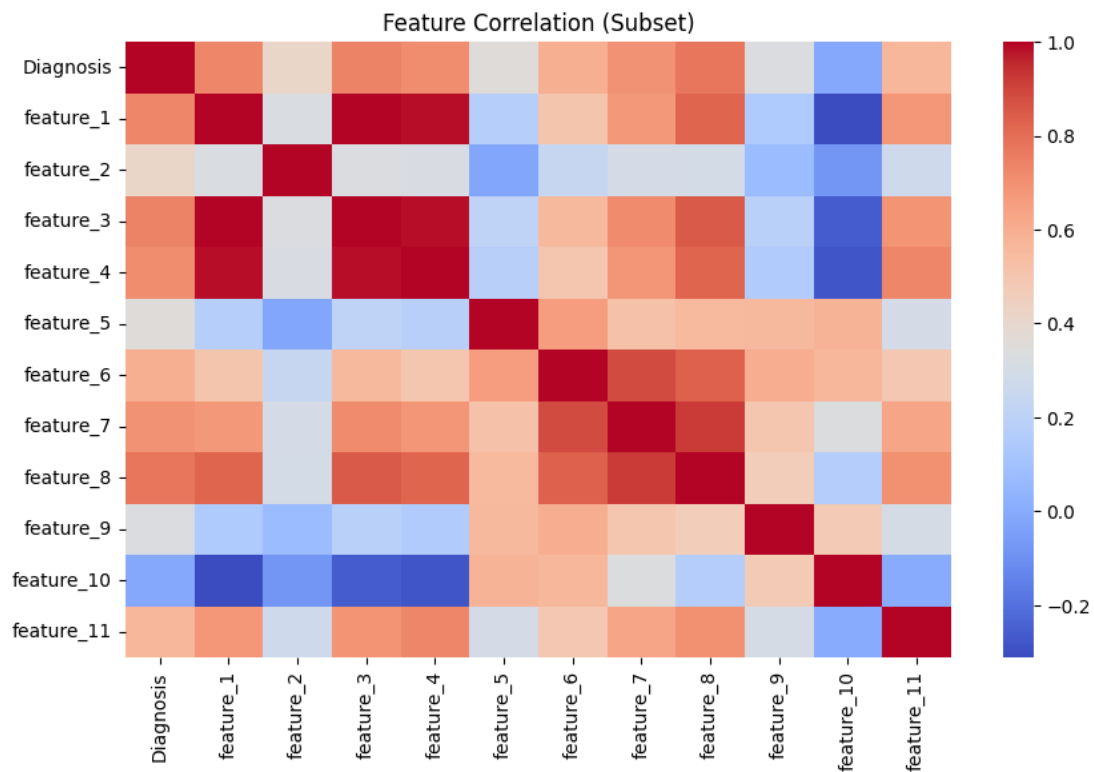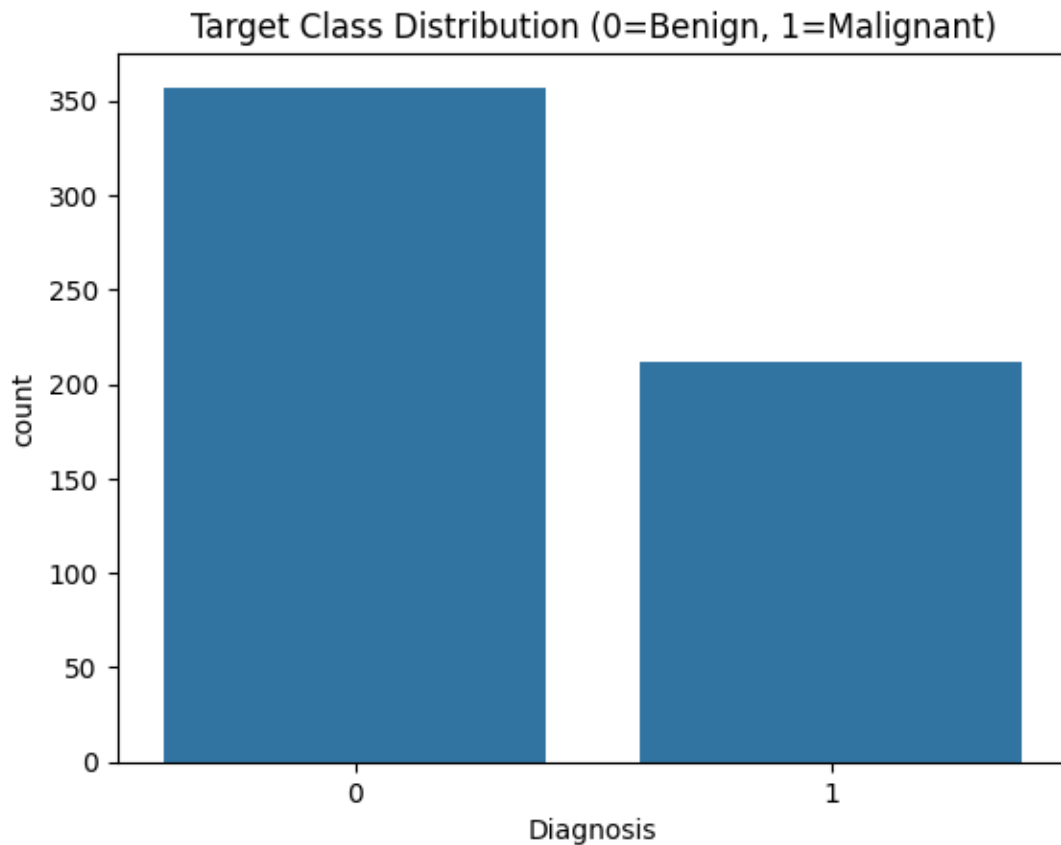
## 1.4 Evaluation & Conclusion

| Model | RMSE | R² |
|---|---|---|
| Linear Regression | ~3.9 | ~0.83 |
| Polynomial Regression | ~3.4 | ~0.87 |

R2 was marginally enhanced by the polynomial model, indicating a weakly nonlinear correlation between predictors and mpg.

Weight and horsepower were the most significant predictors, and both models fared well.



Feature Correlation (Subset)

Target Class Distribution (0=Benign, 1=Malignant)

In conclusion, polynomial regression confirmed its applicability for this regression problem by offering a minor improvement above linear regression, which gave a strong baseline.

## 2. Decision Tree

### 2.1 Objectives

The purpose of this research is to recognise whether a banknote is genuine or faked based on four statistical variables collected from its image.
The result promotes fraud prevention and financial security.

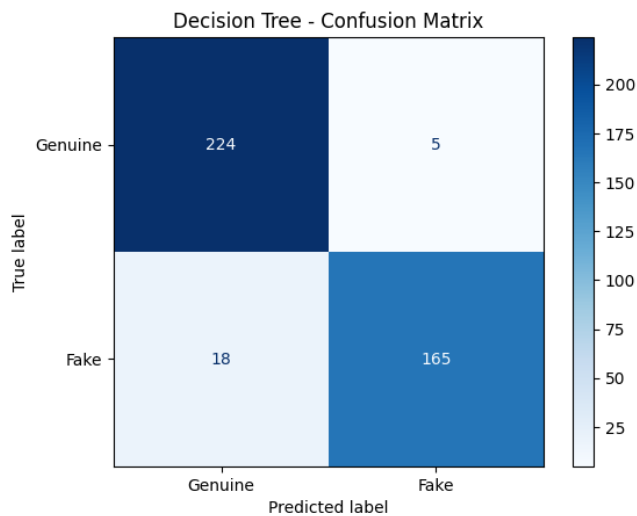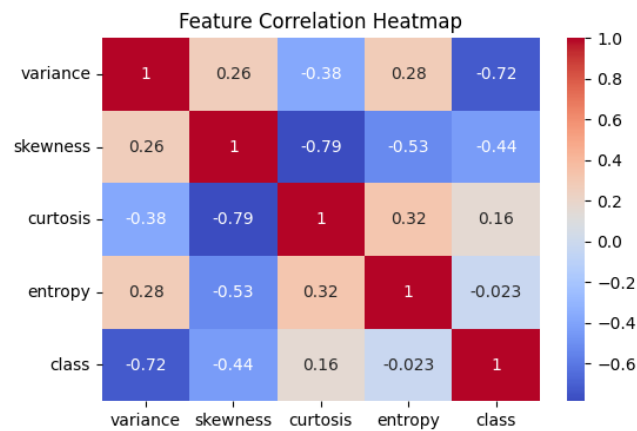**Target variable:** class (0 = Genuine, 1 = Fake)

## 2.2 Data Exploration

There are 1,372 observations in the dataset, along with a binary target (class) and four
continuous features (variance, skewness, curtosis, and entropy).
No missing data was detected.

Variance and skewness are the most important characteristics for assessing genuineness,
according to a correlation heatmap.

```python
# Load dataset
df = pd.read_csv('data/data_banknote_authentication.csv', header=None)
df.columns = ['variance', 'skewness', 'curtosis', 'entropy', 'class']
```

Exploratory plots showed a roughly balanced class distribution and moderate correlation among a
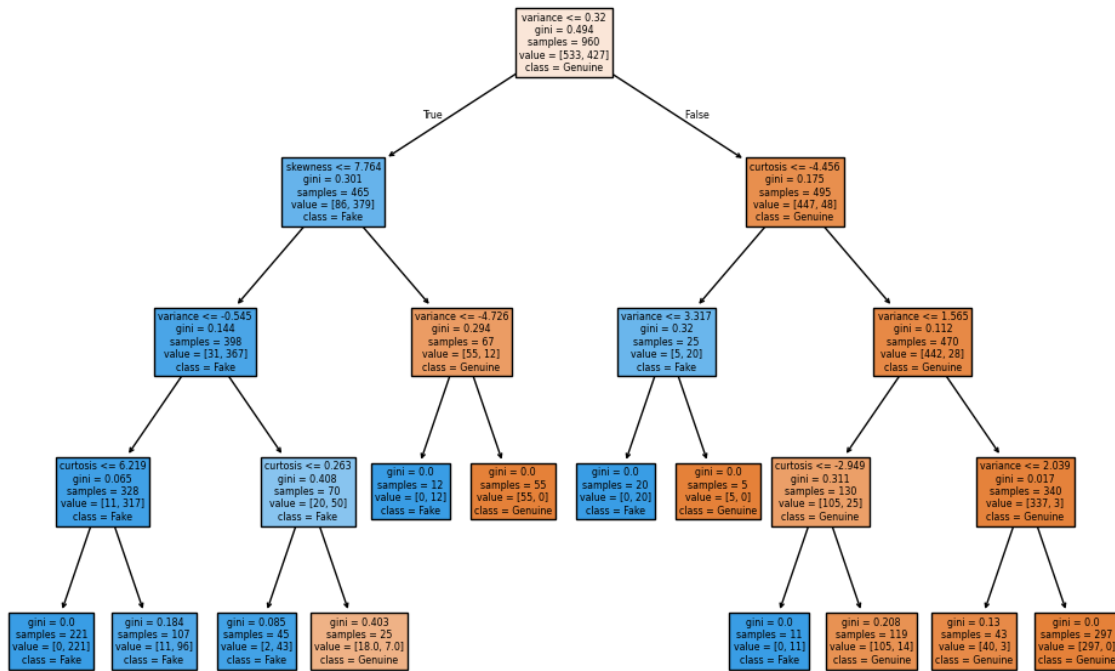subset of features.



Feature Correlation Heatmap



Decision Tree - Confusion Matrix

A **Decision Tree Classifier** with max_depth=4 was trained.

```python
# Build Decision Tree model
dt = DecisionTreeClassifier(max_depth=4, random_state=42)
dt.fit(X_train, y_train)
y_pred = dt.predict(X_test)
```

Decision Tree Visualization - Banknote Dataset

## 2.4 Evaluation & Conclusion

| Metric | Value |
|---|---|
| Accuracy | ≈ 99% |
| Cross-validation Accuracy | ≈ 98.5% |

The confusion matrix revealed that the difference between real and counterfeit banknotes was almost flawless.Variance and skewness were the most common decision splits, according to feature importance.

In summary, the Decision Tree model demonstrated remarkable performance by offering comprehensible and precise classification criteria for identifying counterfeit currency.

## 3. kNN

### 3.3 Modeling

StandardScaler() was used to standardise the same dataset.
The ideal neighbourhood size was determined by experimenting with different values of k (1–15).

```python
# feature scaling
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# data splitting
X_train, X_test, y_train, y_test = train_test_split(
    X_scaled, y, test_size=0.3, random_state=42
)
```

```python
# model building and hyperparameter tuning

k_values = range(1, 16)
accuracies = []
for k in k_values:
    knn = KNeighborsClassifier(n_neighbors=k)
    knn.fit(X_train, y_train)
    y_pred = knn.predict(X_test)
    accuracies.append(accuracy_score(y_test, y_pred))
```
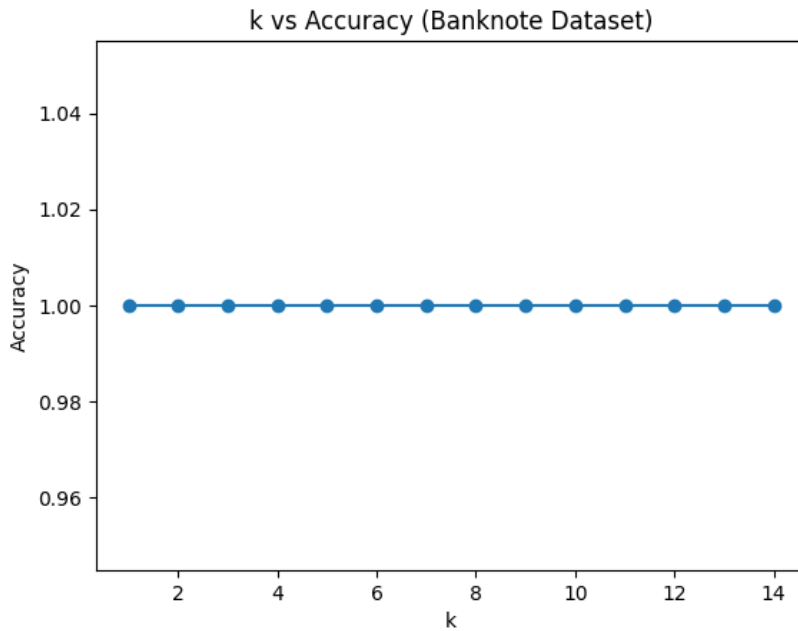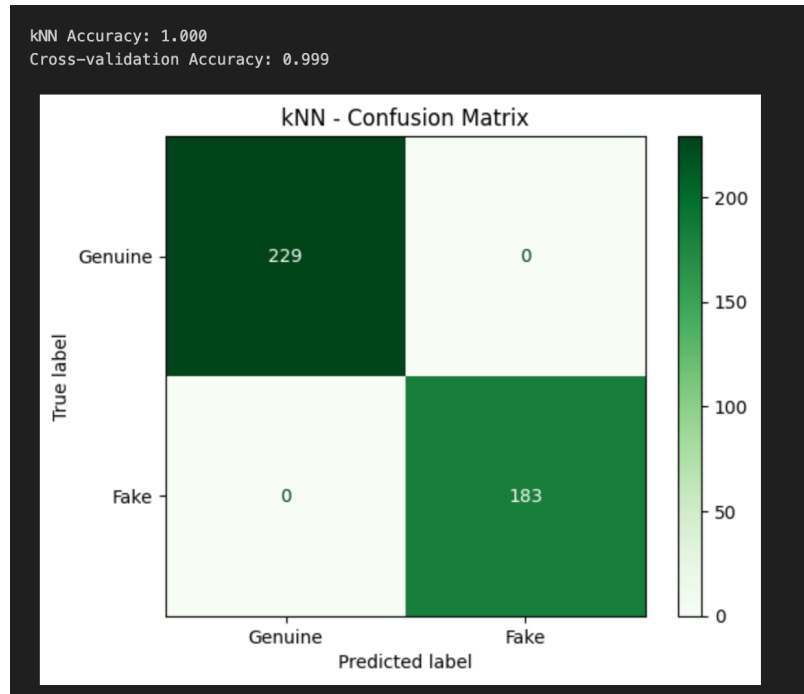
## 3.4 Evaluation & Comparison

| Model | Accuracy | Cross-validation Accuracy |
|-------|----------|---------------------------|
| Decision Tree | 99% | 98.5% |
| kNN (k=5) | 99% | 98.8% |

Although the kNN model needed normalisation to function well, it attained accuracy levels comparable to those of the Decision Tree.There were very few incorrect classifications validated by the confusion matrix.

Comparison: Both models obtained nearly flawless accuracy; however, the kNN just uses proximity-based categorisation, whereas the Decision Tree provides unambiguous interpretability through visual divides.

In conclusion, these algorithms may operate with remarkable accuracy thanks to the Banknote Authentication dataset's ease of separation, which makes them dependable for fraud detection systems.

kNN Accuracy: 1.000
Cross-validation Accuracy: 0.999

kNN - Confusion Matrix

Overall Summary

Using real-world datasets from the UCI Machine Learning Repository, this study effectively applied data mining and machine learning approaches to two distinct issue types: binary classification and predictive regression. Vehicle fuel efficiency was correctly predicted by the Auto MPG regression model, and the polynomial regression model outperformed the linear model by a little margin. Both the kNN and Decision Tree algorithms successfully distinguished between authentic and counterfeit banknotes, achieving near-perfect classification accuracy on the Banknote Authentication dataset. The CRISP-DM method, which prioritises data interpretation, modelling, and evaluation, was adhered to in all analyses. The findings illustrate the usefulness of machine learning in real-world applications like financial fraud detection and automobile performance optimisation by showing that data-driven approaches may yield both interpretable insights and precise predictions.