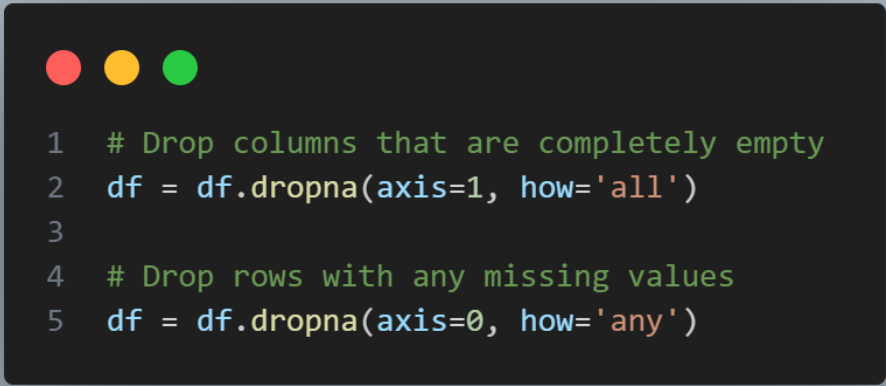


Data Science Lab1

Team nr: 16	Student 1: Antero Morgado	IST nr: 1119213
	Student 2: David Ferreira	IST nr: 1107077
	Student 3: José Fernandes	IST nr: 1103727
	Student 4: Olha Buts	IST nr: 1116276

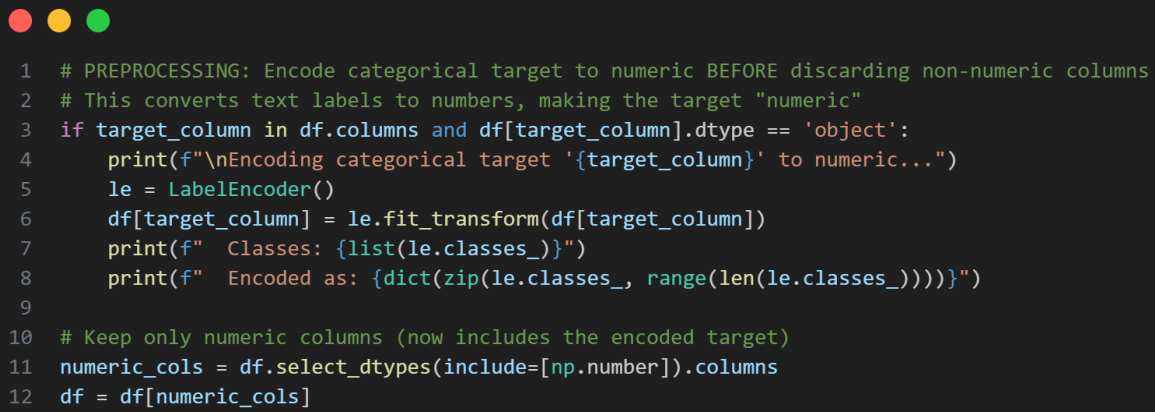
CLASSIFICATION

1 DATA CLEANING



```
1 # Drop columns that are completely empty
2 df = df.dropna(axis=1, how='all')
3
4 # Drop rows with any missing values
5 df = df.dropna(axis=0, how='any')
```

Figure 1: Dropping all empty variables and all records with missing values.



```

1  # PREPROCESSING: Encode categorical target to numeric BEFORE discarding non-numeric columns
2  # This converts text labels to numbers, making the target "numeric"
3  if target_column in df.columns and df[target_column].dtype == 'object':
4      print(f"\nEncoding categorical target '{target_column}' to numeric...")
5      le = LabelEncoder()
6      df[target_column] = le.fit_transform(df[target_column])
7      print(f"  Classes: {list(le.classes_)}")
8      print(f"  Encoded as: {dict(zip(le.classes_, range(len(le.classes_))))}")
9
10 # Keep only numeric columns (now includes the encoded target)
11 numeric_cols = df.select_dtypes(include=[np.number]).columns
12 df = df[numeric_cols]

```

Figure 2: Discard all non-numeric data excluding the target variable.

2 RESULTS

2.1 Model Performance Summary

Table 1: Performance Metrics - Traffic Accidents Dataset

Model	Accuracy	Precision	Recall	F1-Score
Naïve Bayes	0.823	0.863	0.823	0.813
Logistic Regression	0.821	0.864	0.821	0.811
KNN	0.826	0.822	0.811	0.806
Decision Tree	0.829	0.859	0.829	0.821
Multi-layer Perceptron	0.829	0.859	0.829	0.821

Table 2: Performance Metrics - Combined Flights 2022 Dataset

Model	Accuracy	Precision	Recall	F1-Score
Naïve Bayes	0.968	0.943	0.968	0.955
Logistic Regression	0.970	0.971	0.970	0.956
KNN	0.970	0.954	0.970	0.956
Decision Tree	0.971	0.962	0.971	0.960
Multi-layer Perceptron	0.970	0.970	0.970	0.956

2.2 Best Hyperparameters

Table 3: Optimal Hyperparameters - Traffic Accidents

Model	Best Hyperparameters
Logistic Regression	penalty: l1, solver: liblinear, max_iter: 500
KNN	n_neighbors: 10, metric: euclidean
Decision Tree	criterion: entropy, max_depth: 6
Multi-layer Perceptron	learning_rate_ini: 0.5, max_iter: 500, activation: logistic, solver: sgd, learning_rate: adaptive

Table 4: Optimal Hyperparameters - Combined Flights 2022

Model	Best Hyperparameters
Logistic Regression	penalty: l1, solver: liblinear, max_iter: 500
KNN	n_neighbors: 15, metric: manhattan
Decision Tree	criterion: gini, max_depth: 8
Multi-layer Perceptron	learning_rate_ini: 0.5, max_iter: 500, activation: logistic, solver: sgd, learning_rate: adaptive

2.3 Naïve Bayes

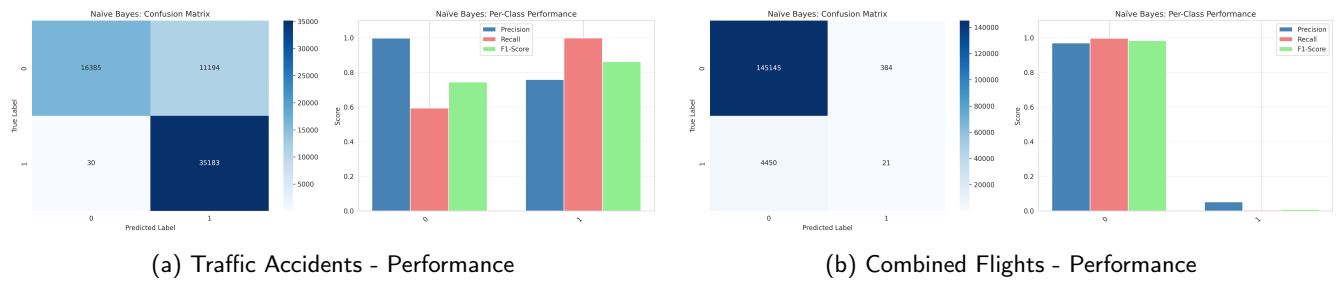


Figure 3: Naïve Bayes: Model Performance

2.4 Logistic Regression

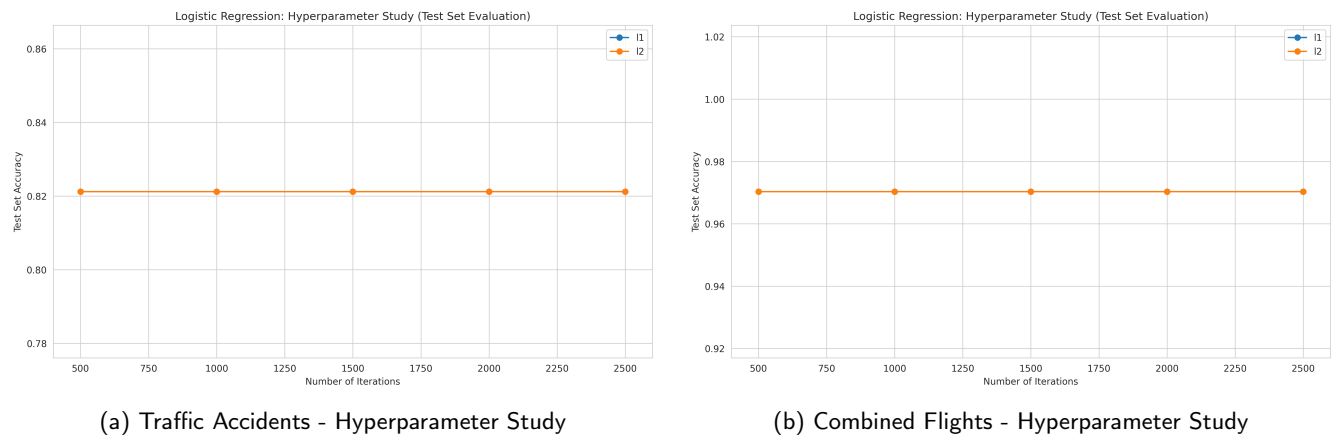


Figure 4: Logistic Regression: Hyperparameter Tuning Results

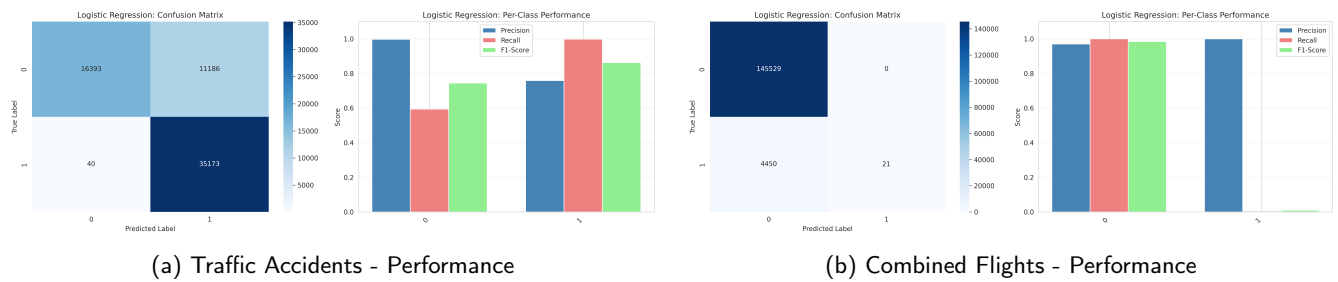
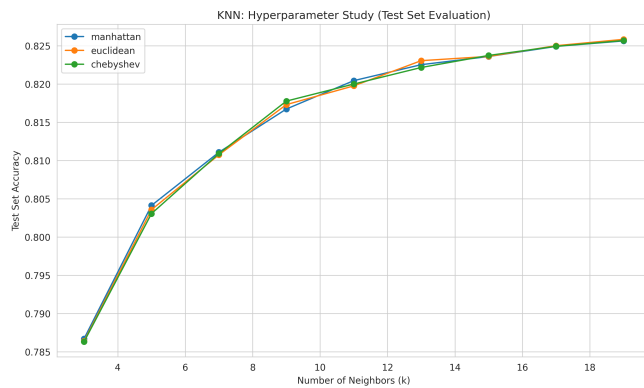
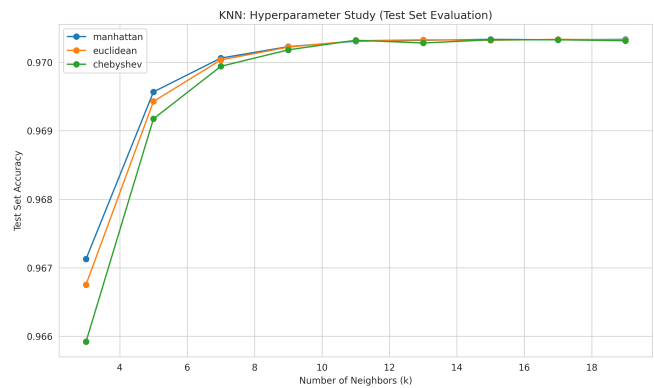


Figure 5: Logistic Regression: Model Performance

2.5 K-Nearest Neighbors (KNN)

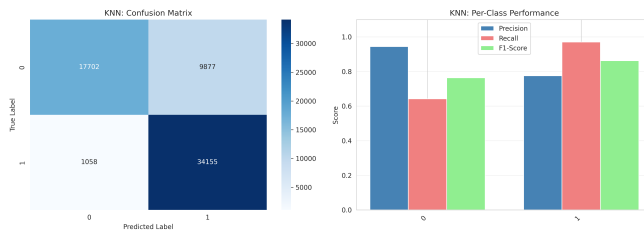


(a) Traffic Accidents - Hyperparameter Study

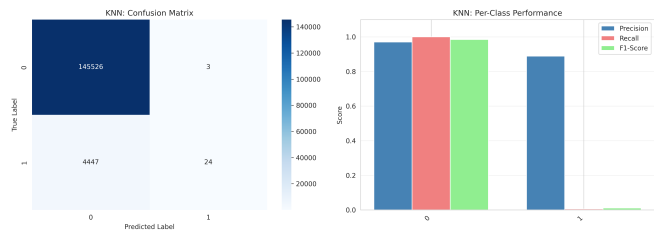


(b) Combined Flights - Hyperparameter Study

Figure 6: KNN: Hyperparameter Tuning Results



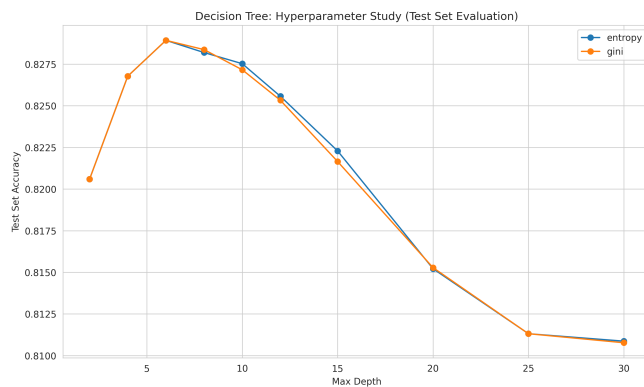
(a) Traffic Accidents - Performance



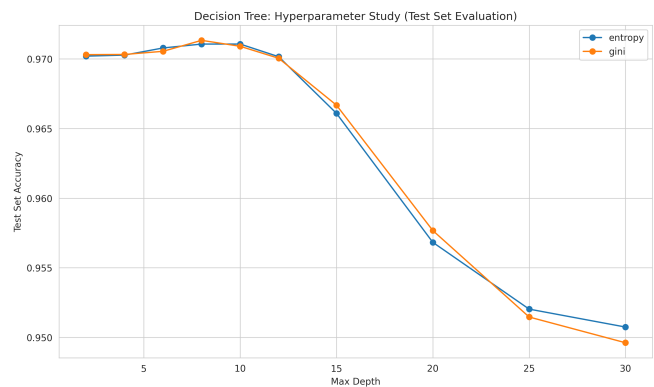
(b) Combined Flights - Performance

Figure 7: KNN: Model Performance

2.6 Decision Tree

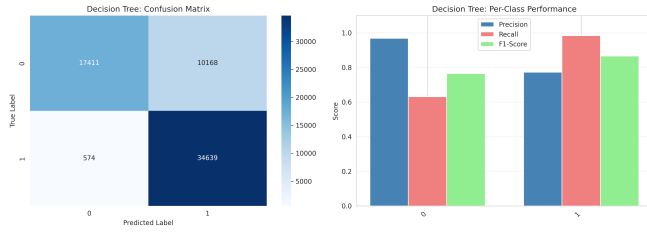


(a) Traffic Accidents - Hyperparameter Study

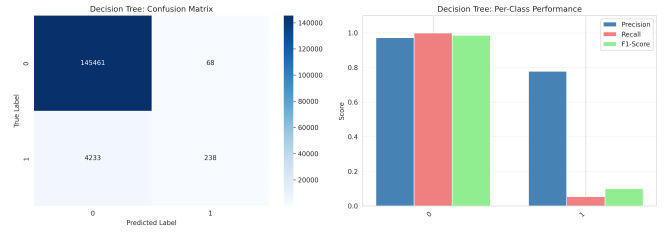


(b) Combined Flights - Hyperparameter Study

Figure 8: Decision Tree: Hyperparameter Tuning Results



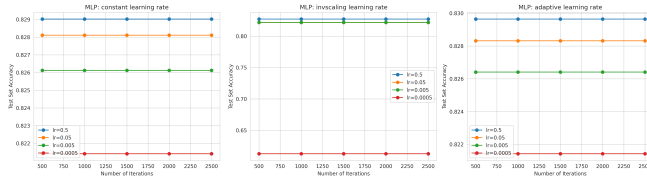
(a) Traffic Accidents - Performance



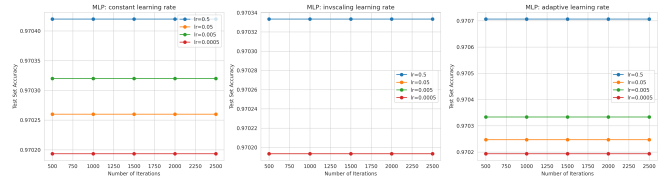
(b) Combined Flights - Performance

Figure 9: Decision Tree: Model Performance

2.7 Multi-layer Perceptron

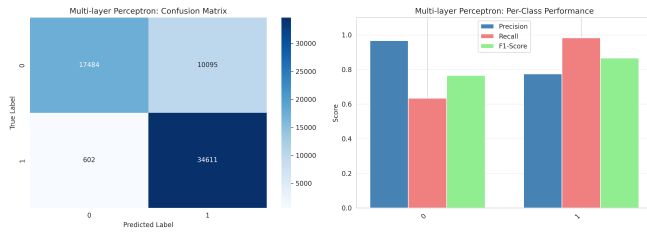


(a) Traffic Accidents - Hyperparameter Study

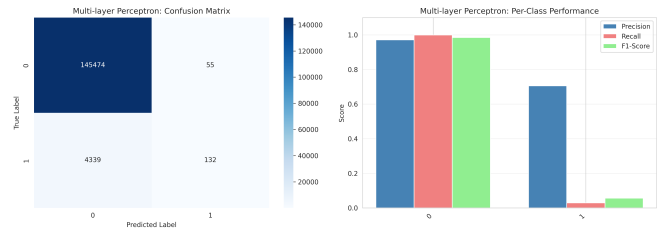


(b) Combined Flights - Hyperparameter Study

Figure 10: Multi-layer Perceptron: Hyperparameter Tuning Results



(a) Traffic Accidents - Performance



(b) Combined Flights - Performance

Figure 11: Multi-layer Perceptron: Model Performance