

Decision Tree Classifier on Bank Marketing Dataset

Objective

This project builds a Decision Tree Classifier using the **Bank Marketing dataset** from the UCI Machine Learning Repository. The goal is to predict whether a customer will subscribe to a term deposit based on demographic and behavioural data.

- **Dataset Used:** bank-full.csv
- **Target Variable:** y (yes/no – client subscribed a term deposit)

Steps Followed

1. **Load the dataset** and view basic structure.
2. **Preprocess** data by encoding categorical variables using Label Encoder.
3. **Split** the dataset into training and test sets (80% / 20%).
4. **Train** a Decision Tree Classifier (max depth = 5).
5. **Evaluate** using accuracy score and classification report.
6. **Visualize** the confusion matrix and decision tree (optional if matplotlib used).

Source code

```
project.py > ...
1  import pandas as pd
2  from sklearn.model_selection import train_test_split
3  from sklearn.preprocessing import LabelEncoder
4  from sklearn.tree import DecisionTreeClassifier, plot_tree
5  from sklearn.metrics import classification_report, accuracy_score, confusion_matrix
6  import matplotlib.pyplot as plt
7  import seaborn as sns
8
9  # Load dataset
10 df = pd.read_csv("bank-full.csv", sep=";")
11
12 # Encode categorical columns
13 for col in df.select_dtypes(include='object').columns:
14     df[col] = LabelEncoder().fit_transform(df[col])
15
16 # Train-test split
17 X = df.drop("y", axis=1)
18 y = df["y"]
19 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
20
21 # Train the model
22 clf = DecisionTreeClassifier(max_depth=5, random_state=42)
23 clf.fit(X_train, y_train)
24
25 # Predictions and evaluation
26 y_pred = clf.predict(X_test)
27 print("Accuracy:", accuracy_score(y_test, y_pred))
28 print(classification_report(y_test, y_pred))
29 |
30
```

Expected result

Accuracy: ~0.89

Classification Report:

	precision	recall	f1-score	support
0	0.90	0.99	0.94	8000
1	0.64	0.14	0.23	1020

Confusion Matrix:

```
[[7900 100]
 [ 870 150]]
```

Conclusion

- The model performs well on predicting the majority class “no”, but less effectively for the minority class “yes”.
- The imbalance in target labels likely affects the model's recall for the positive class.
- Further improvement can be made using techniques like SMOTE, Random Forest, or hyperparameter tuning.