

# Prodigy InfoTech Internship

## Task 3:

Build a decision tree classifier to predict whether a customer will purchase a product or service based on their demographic and behavioral data. Use a dataset such as the Bank Marketing dataset from the UCI Machine Learning Repository.

Sample Dataset: [Bank Marketing](#)

## Decision Tree Classifier of Bank Marketing Dataset

### Loading Libraries and Dataset:

```
[1] 1 import warnings
    2 warnings.filterwarnings('ignore')
    3 import numpy as np
    4 import pandas as pd
    5 import matplotlib.pyplot as plt
    6 import seaborn as sns
```

```
1 sns.set_theme(context='notebook', style='whitegrid', palette='muted')
```

```
[3] 1 data = pd.read_csv("/content/drive/MyDrive/Project_Datasets/Decision_Tree_Classifier/bank-full.csv", sep=';')
```

### Understanding the shape of the data:

```
1 data
```

|       | age | job          | marital  | education | default | balance | housing | loan | contact   | day | month | duration | campaign | pdays | previous | poutcome | y   |
|-------|-----|--------------|----------|-----------|---------|---------|---------|------|-----------|-----|-------|----------|----------|-------|----------|----------|-----|
| 0     | 58  | management   | married  | tertiary  | no      | 2143    | yes     | no   | unknown   | 5   | may   | 261      | 1        | -1    | 0        | unknown  | no  |
| 1     | 44  | technician   | single   | secondary | no      | 29      | yes     | no   | unknown   | 5   | may   | 151      | 1        | -1    | 0        | unknown  | no  |
| 2     | 33  | entrepreneur | married  | secondary | no      | 2       | yes     | yes  | unknown   | 5   | may   | 76       | 1        | -1    | 0        | unknown  | no  |
| 3     | 47  | blue-collar  | married  | unknown   | no      | 1506    | yes     | no   | unknown   | 5   | may   | 92       | 1        | -1    | 0        | unknown  | no  |
| 4     | 33  | unknown      | single   | unknown   | no      | 1       | no      | no   | unknown   | 5   | may   | 198      | 1        | -1    | 0        | unknown  | no  |
| ...   | ... | ...          | ...      | ...       | ...     | ...     | ...     | ...  | ...       | ... | ...   | ...      | ...      | ...   | ...      | ...      | ... |
| 45206 | 51  | technician   | married  | tertiary  | no      | 825     | no      | no   | cellular  | 17  | nov   | 977      | 3        | -1    | 0        | unknown  | yes |
| 45207 | 71  | retired      | divorced | primary   | no      | 1729    | no      | no   | cellular  | 17  | nov   | 456      | 2        | -1    | 0        | unknown  | yes |
| 45208 | 72  | retired      | married  | secondary | no      | 5715    | no      | no   | cellular  | 17  | nov   | 1127     | 5        | 184   | 3        | success  | yes |
| 45209 | 57  | blue-collar  | married  | secondary | no      | 668     | no      | no   | telephone | 17  | nov   | 508      | 4        | -1    | 0        | unknown  | no  |
| 45210 | 37  | entrepreneur | married  | secondary | no      | 2971    | no      | no   | cellular  | 17  | nov   | 361      | 2        | 188   | 11       | other    | no  |

45211 rows x 17 columns

By: Aman Kumar Jha  
[LinkedIn](#)

```
1 data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 45211 entries, 0 to 45210
Data columns (total 17 columns):
 #   Column      Non-Null Count  Dtype
---  -
 0   age         45211 non-null  int64
 1   job         45211 non-null  object
 2   marital     45211 non-null  object
 3   education   45211 non-null  object
 4   default     45211 non-null  object
 5   balance     45211 non-null  int64
 6   housing     45211 non-null  object
 7   loan        45211 non-null  object
 8   contact     45211 non-null  object
 9   day         45211 non-null  int64
10  month       45211 non-null  object
11  duration    45211 non-null  int64
12  campaign    45211 non-null  int64
13  pdays       45211 non-null  int64
14  previous    45211 non-null  int64
15  poutcome    45211 non-null  object
16  y           45211 non-null  object
dtypes: int64(7), object(10)
memory usage: 5.9+ MB
```

```
1 data.describe(include='object')
```

|        | job         | marital | education | default | housing | loan  | contact  | month | poutcome | y     |
|--------|-------------|---------|-----------|---------|---------|-------|----------|-------|----------|-------|
| count  | 45211       | 45211   | 45211     | 45211   | 45211   | 45211 | 45211    | 45211 | 45211    | 45211 |
| unique | 12          | 3       | 4         | 2       | 2       | 2     | 3        | 12    | 4        | 2     |
| top    | blue-collar | married | secondary | no      | yes     | no    | cellular | may   | unknown  | no    |
| freq   | 9732        | 27214   | 23202     | 44396   | 25130   | 37967 | 29285    | 13766 | 36959    | 39922 |

```
[7] 1 data.duplicated().sum()
```

```
0
```

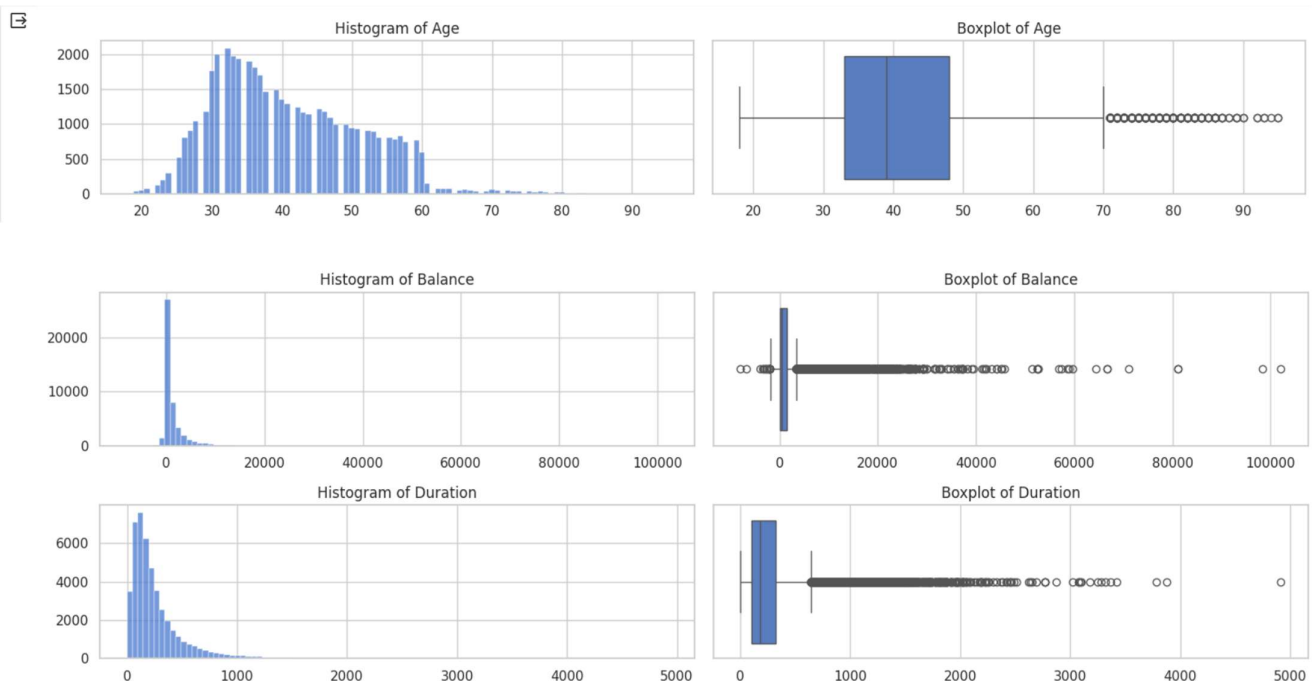
## Data Cleaning:

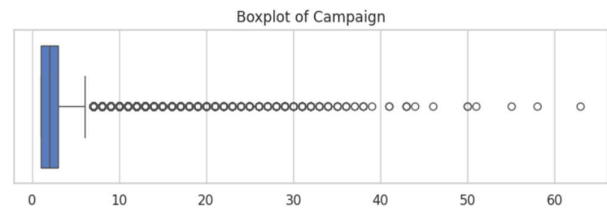
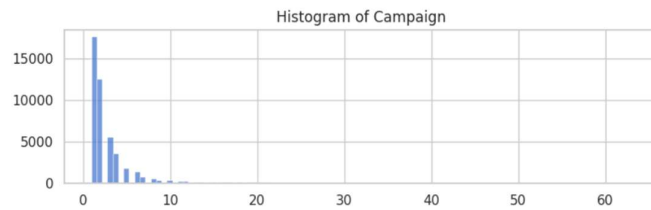
```
[8] 1 data = data.rename(columns={'y': 'subscribed'})
    2 data['subscribed'] = data['subscribed'].map({'yes': 'Subscribed', 'no': 'Not Subscribed'})
```

```
1 categorical_cols = ['job', 'marital', 'education', 'contact', 'month', 'outcome']
2 data[categorical_cols] = (data[categorical_cols].apply(lambda x: x.str.title()))
3 .astype('category')
4 binary_cols = ['default', 'housing', 'loan']
5 data[binary_cols] = data[binary_cols] == 'yes'
```

```
[10] 1 cols_with_outliers = ['age', 'balance', 'duration', 'campaign']
```

```
1 fig, axes = plt.subplots(4, 2, figsize=(15, 10))
2 for i, col in enumerate(cols_with_outliers):
3     hist_ax, box_ax = axes[i, :]
4     sns.histplot(data=data, x=col, bins=100, ax=hist_ax)
5     hist_ax.set_title(f'Histogram of {col.title()}')
6     hist_ax.set_xlabel('')
7     hist_ax.set_ylabel('')
8     sns.boxplot(data=data, x=col, ax=box_ax)
9     box_ax.set_title(f'Boxplot of {col.title()}')
10    box_ax.set_xlabel('')
11    box_ax.set_ylabel('')
12 plt.tight_layout()
13 plt.show();
```





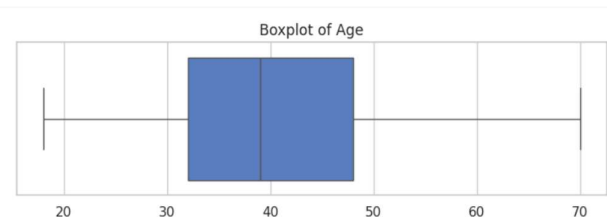
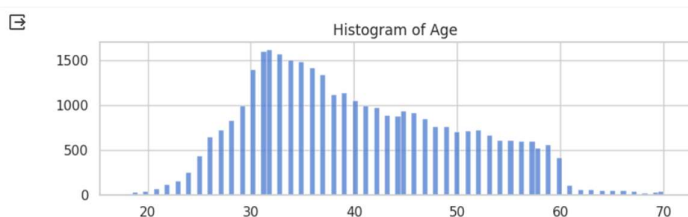
```
[12] 1 def remove_outliers(df, columns):
2     df_outliers_removed = data.copy()
3     for col in columns:
4         Q1 = df_outliers_removed[col].quantile(0.25)
5         Q3 = df_outliers_removed[col].quantile(0.75)
6         IQR = Q3 - Q1
7         lower_bound = Q1 - 1.5 * IQR
8         upper_bound = Q3 + 1.5 * IQR
9         df_outliers_removed = df_outliers_removed[
10             (df_outliers_removed[col] >= lower_bound) &
11             (df_outliers_removed[col] <= upper_bound)
12         ]
13     return df_outliers_removed
14 data = remove_outliers(data, cols_with_outliers)
```

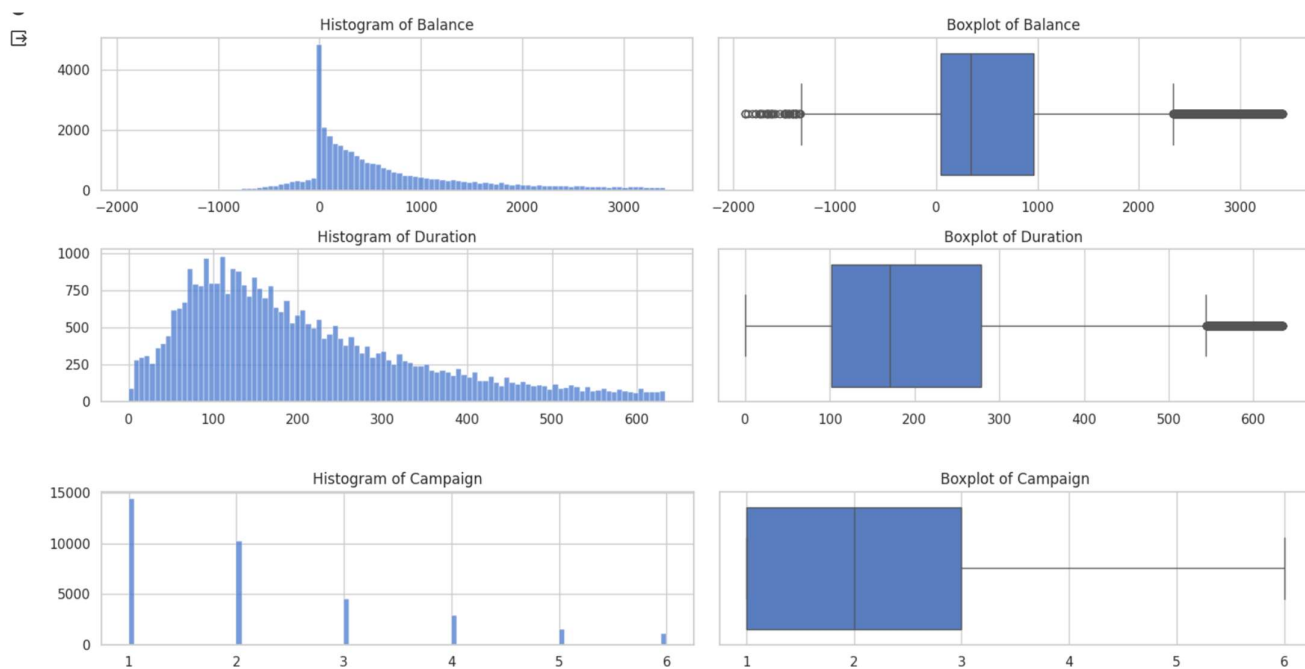
[13] 1 data

|       | age | job          | marital | education | default | balance | housing | loan  | contact   | day | month | duration | campaign | pdays | previous | poutcome | subscribed     |
|-------|-----|--------------|---------|-----------|---------|---------|---------|-------|-----------|-----|-------|----------|----------|-------|----------|----------|----------------|
| 0     | 58  | Management   | Married | Tertiary  | False   | 2143    | True    | False | Unknown   | 5   | May   | 261      | 1        | -1    | 0        | Unknown  | Not Subscribed |
| 1     | 44  | Technician   | Single  | Secondary | False   | 29      | True    | False | Unknown   | 5   | May   | 151      | 1        | -1    | 0        | Unknown  | Not Subscribed |
| 2     | 33  | Entrepreneur | Married | Secondary | False   | 2       | True    | True  | Unknown   | 5   | May   | 76       | 1        | -1    | 0        | Unknown  | Not Subscribed |
| 3     | 47  | Blue-Collar  | Married | Unknown   | False   | 1506    | True    | False | Unknown   | 5   | May   | 92       | 1        | -1    | 0        | Unknown  | Not Subscribed |
| 4     | 33  | Unknown      | Single  | Unknown   | False   | 1       | False   | False | Unknown   | 5   | May   | 198      | 1        | -1    | 0        | Unknown  | Not Subscribed |
| ...   | ... | ...          | ...     | ...       | ...     | ...     | ...     | ...   | ...       | ... | ...   | ...      | ...      | ...   | ...      | ...      | ...            |
| 45202 | 34  | Admin.       | Single  | Secondary | False   | 557     | False   | False | Cellular  | 17  | Nov   | 224      | 1        | -1    | 0        | Unknown  | Subscribed     |
| 45203 | 23  | Student      | Single  | Tertiary  | False   | 113     | False   | False | Cellular  | 17  | Nov   | 266      | 1        | -1    | 0        | Unknown  | Subscribed     |
| 45205 | 25  | Technician   | Single  | Secondary | False   | 505     | False   | True  | Cellular  | 17  | Nov   | 386      | 2        | -1    | 0        | Unknown  | Subscribed     |
| 45209 | 57  | Blue-Collar  | Married | Secondary | False   | 668     | False   | False | Telephone | 17  | Nov   | 508      | 4        | -1    | 0        | Unknown  | Not Subscribed |
| 45210 | 37  | Entrepreneur | Married | Secondary | False   | 2971    | False   | False | Cellular  | 17  | Nov   | 361      | 2        | 188   | 11       | Other    | Not Subscribed |

34563 rows × 17 columns

```
1 fig, axes = plt.subplots(4, 2, figsize=(15, 10))
2 for i, col in enumerate(cols_with_outliers):
3     hist_ax, box_ax = axes[i, :]
4     sns.histplot(data=data, x=col, bins=100, ax=hist_ax)
5     hist_ax.set_title(f'Histogram of {col.title()}')
6     hist_ax.set_xlabel('')
7     hist_ax.set_ylabel('')
8     sns.boxplot(data=data, x=col, ax=box_ax)
9     box_ax.set_title(f'Boxplot of {col.title()}')
10    box_ax.set_xlabel('')
11    box_ax.set_ylabel('')
12 plt.tight_layout()
13 plt.show();
```

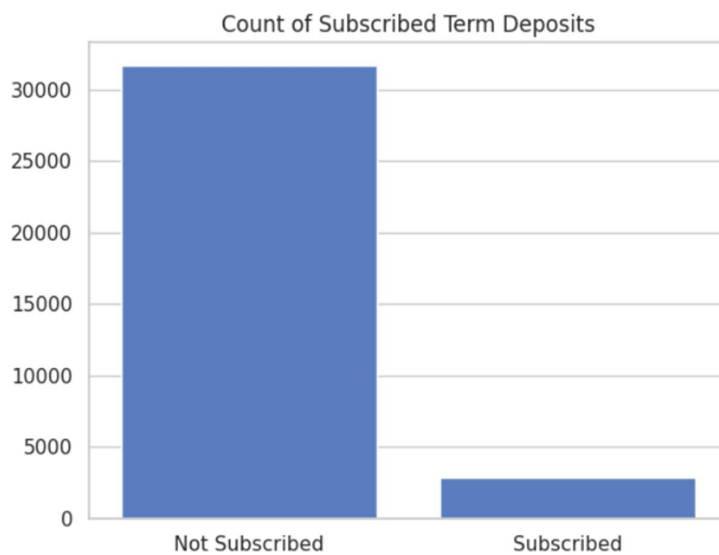




## Data Exploration:

```
[15] 1 num_cols = data.select_dtypes('number').columns.tolist()
      2 bool_cols = data.select_dtypes('bool').columns.tolist()
      3 cat_cols = data.select_dtypes('category').columns.tolist()
```

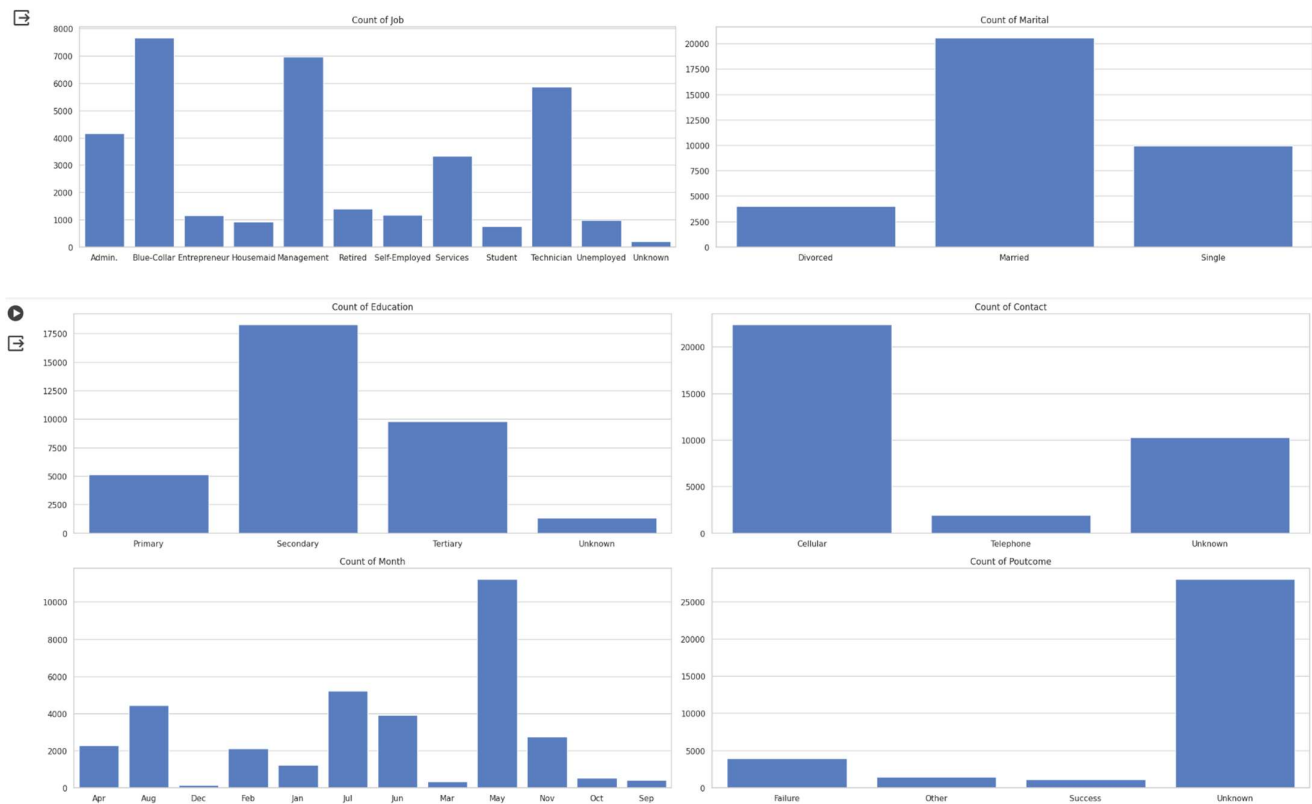
```
1 sns.countplot(data=data, x='subscribed');
2 plt.title('Count of Subscribed Term Deposits')
3 plt.xlabel('')
4 plt.ylabel('')
5 plt.show();
```



```

1 fig, axes = plt.subplots(3, 2, figsize=(25, 15))
2 for feature, ax in zip(cat_cols, axes.flatten()):
3     sns.countplot(data=data, x=feature, ax=ax)
4     ax.set_title(f'Count of {feature.title()}')
5     ax.set_xlabel('')
6     ax.set_ylabel('')
7 plt.tight_layout()
8 plt.show();

```



## Data Preprocessing for Model Training:

```

[18] 1 from sklearn.model_selection import train_test_split
2     from sklearn.preprocessing import OneHotEncoder, StandardScaler
3     from sklearn.compose import ColumnTransformer
4     from imblearn.over_sampling import RandomOverSampler

```

```

[19] 1 X = data.drop(columns='subscribed')
2     y = data['subscribed']
3     X_train, X_test, y_train, y_test = train_test_split(X,
4                                                         y,
5                                                         test_size=0.2,
6                                                         stratify=y,
7                                                         random_state=42)

```

```

[20] 1 num_vars = data.select_dtypes('number').columns.tolist()
2     cat_vars = data.select_dtypes('category').columns.tolist()

```

```

[21] 1 preprocessing_pipeline = ColumnTransformer([
2     ('numerical', StandardScaler(), num_vars),
3     ('categorical', OneHotEncoder(), cat_vars),
4 ])
5 X_train = preprocessing_pipeline.fit_transform(X_train)
6 X_test = preprocessing_pipeline.transform(X_test)

```



```
[22] 1 sampler = RandomOverSampler(random_state=42)
      2 X_train, y_train = sampler.fit_resample(X_train, y_train)
```

## Building Basic Model:

```
[23] 1 from sklearn.tree import DecisionTreeClassifier
      2 from sklearn.metrics import classification_report
```

```
[24] 1 %%time
      2 model = DecisionTreeClassifier(random_state=42)
      3 model.fit(X_train, y_train)
```

CPU times: user 2.24 s, sys: 0 ns, total: 2.24 s  
Wall time: 2.24 s

```
DecisionTreeClassifier
DecisionTreeClassifier(random_state=42)
```

```
[25] 1 y_pred = model.predict(X_test)
      2 accuracy = model.score(X_test, y_test)
      3 report = classification_report(y_test, y_pred)
      4 print(f'Accuracy: {accuracy:.2%}')
      5 print(f'Classification Report:\n{report}')
```

Accuracy: 90.16%

Classification Report:

|                | precision | recall | f1-score | support |
|----------------|-----------|--------|----------|---------|
| Not Subscribed | 0.95      | 0.95   | 0.95     | 6343    |
| Subscribed     | 0.41      | 0.41   | 0.41     | 570     |
| accuracy       |           |        | 0.90     | 6913    |
| macro avg      | 0.68      | 0.68   | 0.68     | 6913    |
| weighted avg   | 0.90      | 0.90   | 0.90     | 6913    |

## Fine Tuning the Model:

```
[26] 1 from sklearn.model_selection import GridSearchCV
      2 from sklearn.metrics import make_scorer, f1_score
```

```
[27] 1 param_grid = {
      2     'max_depth': [None, 10, 20],
      3     'min_samples_split': [2, 5, 10],
      4     'min_samples_leaf': [1, 2, 4],
      5 }
```

```
[28] 1 scorer = make_scorer(f1_score, pos_label='Subscribed')
```

```
[29] 1 base_model = DecisionTreeClassifier(random_state=42)
      2 grid_search = GridSearchCV(estimator=base_model,
      3                           param_grid=param_grid,
      4                           cv=5,
      5                           scoring=scorer,
      6                           verbose=1,
      7                           n_jobs=-1)
```

```
1 %%time
2 grid_search.fit(X_train, y_train)
```

Fitting 5 folds for each of 27 candidates, totalling 135 fits  
CPU times: user 4.3 s, sys: 230 ms, total: 4.53 s  
Wall time: 2min 18s

```
GridSearchCV
GridSearchCV(cv=5, estimator=DecisionTreeClassifier(random_state=42), n_jobs=-1,
             param_grid={'max_depth': [None, 10, 20],
                          'min_samples_leaf': [1, 2, 4],
                          'min_samples_split': [2, 5, 10]},
             scoring=make_scorer(f1_score, pos_label=Subscribed), verbose=1)
  estimator: DecisionTreeClassifier
    DecisionTreeClassifier
```

```
[31] 1 best_params = grid_search.best_params_
     2 best_model = grid_search.best_estimator_
     3 accuracy = best_model.score(X_test, y_test)
     4 print(f'Best Accuracy: {accuracy:.2%}')
     5 print(f'Best Parameters:\n{best_params}')
```

Best Accuracy: 90.16%  
Best Parameters:  
{'max\_depth': None, 'min\_samples\_leaf': 1, 'min\_samples\_split': 2}

```
[32] 1 y_pred = best_model.predict(X_test)
     2 report = classification_report(y_test, y_pred)
     3 print(f'Classification Report:\n{report}')
```

Classification Report:

|                | precision | recall | f1-score | support |
|----------------|-----------|--------|----------|---------|
| Not Subscribed | 0.95      | 0.95   | 0.95     | 6343    |
| Subscribed     | 0.41      | 0.41   | 0.41     | 570     |
| accuracy       |           |        | 0.90     | 6913    |
| macro avg      | 0.68      | 0.68   | 0.68     | 6913    |
| weighted avg   | 0.90      | 0.90   | 0.90     | 6913    |

## Testing the results:

```
[33] 1 from sklearn.metrics import confusion_matrix
```

```
[34] 1 conf_matrix = confusion_matrix(y_test, y_pred)
     2 labels = best_model.classes_
     3 sns.heatmap(conf_matrix,
     4             annot=True,
     5             fmt='d',
     6             cmap='Blues',
     7             cbar=False,
     8             xticklabels=labels,
     9             yticklabels=labels)
    10 plt.title('Confusion Matrix')
    11 plt.xlabel('Predicted')
    12 plt.ylabel('True')
    13 plt.show();
```





Confusion Matrix

| True           | Not Subscribed | Subscribed |
|----------------|----------------|------------|
|                | Predicted      |            |
| Not Subscribed | 5998           | 345        |
| Subscribed     | 335            | 235        |