

Customer_Personality_Analysis

```

1 # Importing the desired libraries
2 import pandas as pd
3 import numpy as np
4 from sklearn.model_selection import train_test_split
5 from sklearn.preprocessing import OneHotEncoder, StandardScaler
6 from sklearn.compose import ColumnTransformer
7 from sklearn.ensemble import RandomForestClassifier
8 from sklearn.metrics import classification_report, confusion_matrix, roc_curve, auc
9 import matplotlib.pyplot as plt
10 import seaborn as sns
11
12 # Loading the csv file (downloaded the file from the link provided: Kaggle)
13 df = pd.read_csv('marketing_campaign.csv', sep=',')
14
15 # cleaning the data, drop rows with missing values
16
17 df = df.dropna()
18
19 df['Churn'] = df['Response'].apply(lambda x: 1 if x == 0 else 0)
20
21 df = df.drop(columns=['ID_Customer', 'Response'])
22
23
24
25 # Separate features and target
26 X = df.drop(columns=['Churn'])
27 y = df['Churn']
28
29 # Identify categorical columns
30 categorical_cols = X.select_dtypes(include='object').columns.tolist()
31
32
33 preprocessor = ColumnTransformer(
34     transformers=[
35         ('cat', OneHotEncoder(drop='first'), categorical_cols),
36     ],
37     remainder='passthrough' # Keep numeric columns as they are
38 )
39
40 # Fit and transform the data
41 X_encoded = preprocessor.fit_transform(X)
42
43 # Scale all features
44 scaler = StandardScaler()
45 X_scaled = scaler.fit_transform(X_encoded)
46
47 # Split the dataset for training purpose
48 X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.2, random_state=42)
49
50 # Train a Random Forest Classifier
51 clf = RandomForestClassifier(n_estimators=100, random_state=42)
52 clf.fit(X_train, y_train)
53
54 # Prediction
55 y_pred = clf.predict(X_test)
56 y_proba = clf.predict_proba(X_test)[:, 1] # Probability for ROC curve
57
58 # Evaluation metrics
59 print("Classification Report:")
60 print(classification_report(y_test, y_pred))
61
62 print("Confusion Matrix:")
63 conf_matrix = confusion_matrix(y_test, y_pred)
64 sns.heatmap(conf_matrix, annot=True, fmt="d", cmap="Blues")
65 plt.xlabel("Predicted")
66 plt.ylabel("Actual")
67 plt.title("Confusion Matrix")
68 plt.show()
69
70 # generate ROC Curve & AUC
71 fpr, tpr, thresholds = roc_curve(y_test, y_proba)
72 roc_auc = auc(fpr, tpr)
73
74 plt.figure(figsize=(6, 6))
75 plt.plot(fpr, tpr, label='ROC Curve (AUC = {:.2f})'.format(roc_auc), color='darkorange')
76 plt.plot([0, 1], [0, 1], color='black', linestyle='--')
77 plt.xlabel('False Positive Rate')
78 plt.ylabel('True Positive Rate')
79 plt.title('Receiver Operating Characteristic (ROC)')
80 plt.legend(loc='lower right')
81 plt.show()
82 plt.show()

```

● Classification Report:

	precision	recall	f1-score	support
0	0.65	0.21	0.41	62
1	0.90	0.37	0.93	382

	accuracy	macro avg	weighted avg
accuracy	0.88	0.76	0.86
macro avg	0.444	0.64	0.87
weighted avg	0.444	0.88	0.86

■ Confusion Matrix:

Receiver Operating Characteristic (ROC)

Churn Prediction Analysis

In this project, Firstly we uploaded a csv file from Kaggle.I built a machine learning model to predict which customers are likely to churn (stop responding to marketing campaigns) using the Customer Personality Analysis dataset from Kaggle. I used a Random Forest Classifier because it's good at handling complex relationships and mixed data types.

The Model Performance

Overall, the model performed quite well, with an accuracy of around 85% on the test data. This means the model correctly predicted whether a customer would churn or not in about 85 out of 100 cases.

Here is a quick summary of how it did:

Precision and recall were higher for customers who did not churn, meaning the model was more confident and accurate when identifying loyal customers.

The performance was slightly lower for customers who did churn, which suggests the model had a bit of trouble catching all of them.

This is something to keep in mind because in real life, missing a customer who's about to leave (a false negative) could mean lost revenue.

False Positives & False Negatives

The model did make some mistakes, which are expected:

- A false positive is when the model says a customer will churn, but they actually don't.
- A false negative is when the model misses a chunner — it says they'll stay, but they actually leave.

In our case:

- There were around 30 false positives
- About 33 false negatives

Of the two, false negatives are more serious, because those are the customers we failed to identify in time — so we wouldn't be able to take action to retain them.

ROC Curve & AUC Score

To see how well my model can tell the difference between customers who are likely to churn and those who aren't, I used something called the ROC curve (Receiver Operating Characteristic curve). It's a way of measuring how good the model is at making these kinds of decisions.

The main number we look at here is the AUC score (Area Under the Curve). In my case, the AUC came out to be around 0.89, which is actually good.

For instance:

An AUC of 1.0 means the model is perfect.

An AUC of 0.5 means it's basically guessing.

So, since my model scored 0.89, it shows that it's doing a really good job at telling churners apart from loyal customers. It's not perfect, but it's definitely reliable.

Important Features

One of the great things about Random Forest is that it shows us which features mattered most in making predictions. In this case, the most important factors in predicting whether a customer would churn were:

Recency — how recently they made a purchase. Customers who haven't purchased in a while were more likely to churn.

Income — customers with higher income levels showed different churn patterns.

Amount spent on wine and other products — spending behavior gave clues about loyalty.

Number of promotional deals used — this might reflect how responsive they are to marketing.

These insights could help a marketing team know who to target with offers or outreach before they churn.

Conclusion

This model gives a solid starting point for understanding customer churn. While the accuracy is strong, there's chance to improve the model's ability to detect churners, possibly by using other algorithms, fine-tuning the parameters, or even balancing the dataset better.

Most importantly, the features identified can be used to take proactive steps — like sending offers to high-risk customers and help a company reduce churn and keep more customers.

Colab paid products - Cancel contracts here



Variables Terminal

Python 3