

```
In [37]: #Basic libraries  
import pandas as pd  
import numpy as np  
import matplotlib.pyplot as plt  
import seaborn as sns  
import math  
  
#Sampling methods  
from sklearn.model_selection import train_test_split  
from imblearn.under_sampling import RandomUnderSampler  
  
#feature engineering  
import datetime as dt  
import category_encoders as ce  
from sklearn.preprocessing import MinMaxScaler  
  
#Models  
from sklearn.ensemble import RandomForestClassifier  
from sklearn.neighbors import KNeighborsClassifier  
  
#Model evaluation  
from sklearn.metrics import f1_score  
from sklearn.metrics import roc_curve  
from sklearn.metrics import auc  
from sklearn.metrics import classification_report  
from sklearn import metrics  
import scipy.stats as stats  
from scipy.stats import skew
```

```
In [38]: train = pd.read_csv('fraudTrain.csv')
test = pd.read_csv('fraudTest.csv')
test.head()
```

Out[38]:

	Unnamed: 0	trans_date_trans_time	cc_num	merchant	category	amt	first	last	gender	street
0	0	2020-06-21 12:14:25	2291163933867244	fraud_Kirlin and Sons	personal_care	2.86	Jeff	Elliott	M	351 Darlene Green
1	1	2020-06-21 12:14:33	3573030041201292	fraud_Sporer-Keebler	personal_care	29.84	Joanne	Williams	F	3638 Marsh Union
2	2	2020-06-21 12:14:53	3598215285024754	fraud_Swaniawski, Nitzsche and Welch	health_fitness	41.28	Ashley	Lopez	F	9333 Valentine Point
3	3	2020-06-21 12:15:15	3591919803438423	fraud_Haley Group	misc_pos	60.05	Brian	Williams	M	32941 Krystal Mill Apt. 552
4	4	2020-06-21 12:15:17	3526826139003047	fraud_Johnston-Casper	travel	3.19	Nathan	Massey	M	5783 Evan Roads Apt. 465

```
In [39]: print(test.shape), print(train.shape)
```

```
(555719, 23)
(1296675, 23)
```

Out[39]: (None, None)

```
In [40]: train.isna().sum()
```

```
Out[40]: Unnamed: 0      0
trans_date_trans_time  0
cc_num                 0
merchant               0
category               0
amt                   0
first                  0
last                   0
gender                 0
street                 0
city                   0
state                  0
zip                    0
lat                    0
long                   0
city_pop               0
job                    0
dob                    0
trans_num              0
unix_time              0
merch_lat              0
merch_long             0
is_fraud               0
dtype: int64
```

```
In [41]: test.isnull().sum()
```

```
Out[41]: Unnamed: 0      0
trans_date_trans_time  0
cc_num                 0
merchant               0
category               0
amt                   0
first                  0
last                   0
gender                 0
street                 0
city                   0
state                  0
zip                    0
lat                    0
long                   0
city_pop               0
job                    0
dob                    0
trans_num              0
unix_time              0
merch_lat              0
merch_long             0
is_fraud               0
dtype: int64
```

In [42]: ▶ test.info(), train.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 555719 entries, 0 to 555718
Data columns (total 23 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Unnamed: 0                            555719 non-null  int64
1   trans_date_trans_time                 555719 non-null  object
2   cc_num                               555719 non-null  int64
3   merchant                             555719 non-null  object
4   category                             555719 non-null  object
5   amt                                   555719 non-null  float64
6   first                                555719 non-null  object
7   last                                 555719 non-null  object
8   gender                               555719 non-null  object
9   street                               555719 non-null  object
10  city                                 555719 non-null  object
11  state                                555719 non-null  object
12  zip                                  555719 non-null  int64
13  lat                                  555719 non-null  float64
14  long                                 555719 non-null  float64
15  city_pop                             555719 non-null  int64
16  job                                   555719 non-null  object
17  dob                                   555719 non-null  object
18  trans_num                             555719 non-null  object
19  unix_time                             555719 non-null  int64
20  merch_lat                             555719 non-null  float64
21  merch_long                             555719 non-null  float64
22  is_fraud                             555719 non-null  int64
dtypes: float64(5), int64(6), object(12)
```

memory usage: 97.5+ MB

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1296675 entries, 0 to 1296674
Data columns (total 23 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Unnamed: 0                            1296675 non-null  int64
1   trans_date_trans_time                 1296675 non-null  object
2   cc_num                               1296675 non-null  int64
3   merchant                             1296675 non-null  object
4   category                             1296675 non-null  object
5   amt                                   1296675 non-null  float64
6   first                                1296675 non-null  object
7   last                                 1296675 non-null  object
```

8	gender	1296675	non-null	object
9	street	1296675	non-null	object
10	city	1296675	non-null	object
11	state	1296675	non-null	object
12	zip	1296675	non-null	int64
13	lat	1296675	non-null	float64
14	long	1296675	non-null	float64
15	city_pop	1296675	non-null	int64
16	job	1296675	non-null	object
17	dob	1296675	non-null	object
18	trans_num	1296675	non-null	object
19	unix_time	1296675	non-null	int64
20	merch_lat	1296675	non-null	float64
21	merch_long	1296675	non-null	float64
22	is_fraud	1296675	non-null	int64

dtypes: float64(5), int64(6), object(12)

memory usage: 227.5+ MB

Out[42]: (None, None)

```
In [43]: train.drop("Unnamed: 0",axis=1,inplace=True)
test.drop("Unnamed: 0",axis=1,inplace=True)
train.head()
```

Out[43]:

	trans_date_trans_time	cc_num	merchant	category	amt	first	last	gender	street	city	st
0	2019-01-01 00:00:18	2703186189652095	fraud_Rippin, Kub and Mann	misc_net	4.97	Jennifer	Banks	F	561 Perry Cove	Moravian Falls	
1	2019-01-01 00:00:44	630423337322	fraud_Heller, Gutmann and Zieme	grocery_pos	107.23	Stephanie	Gill	F	43039 Riley Greens Suite 393	Orient	
2	2019-01-01 00:00:51	38859492057661	fraud_Lind- Buckridge	entertainment	220.11	Edward	Sanchez	M	594 White Dale Suite 530	Malad City	
3	2019-01-01 00:01:16	3534093764340240	fraud_Kutch, Hermiston and Farrell	gas_transport	45.00	Jeremy	White	M	9443 Cynthia Court Apt. 038	Boulder	
4	2019-01-01 00:03:06	375534208663984	fraud_Keeling- Crist	misc_pos	41.96	Tyler	Garcia	M	408 Bradley Rest	Doe Hill	

EDA

In [44]: `train.describe().T`

Out[44]:

	count	mean	std	min	25%	50%	75%	max
cc_num	1296675.0	4.171920e+17	1.308806e+18	6.041621e+10	1.800429e+14	3.521417e+15	4.642255e+15	4.992346e+18
amt	1296675.0	7.035104e+01	1.603160e+02	1.000000e+00	9.650000e+00	4.752000e+01	8.314000e+01	2.894890e+04
zip	1296675.0	4.880067e+04	2.689322e+04	1.257000e+03	2.623700e+04	4.817400e+04	7.204200e+04	9.978300e+04
lat	1296675.0	3.853762e+01	5.075808e+00	2.002710e+01	3.462050e+01	3.935430e+01	4.194040e+01	6.669330e+01
long	1296675.0	-9.022634e+01	1.375908e+01	-1.656723e+02	-9.679800e+01	-8.747690e+01	-8.015800e+01	-6.795030e+01
city_pop	1296675.0	8.882444e+04	3.019564e+05	2.300000e+01	7.430000e+02	2.456000e+03	2.032800e+04	2.906700e+06
unix_time	1296675.0	1.349244e+09	1.284128e+07	1.325376e+09	1.338751e+09	1.349250e+09	1.359385e+09	1.371817e+09
merch_lat	1296675.0	3.853734e+01	5.109788e+00	1.902779e+01	3.473357e+01	3.936568e+01	4.195716e+01	6.751027e+01
merch_long	1296675.0	-9.022646e+01	1.377109e+01	-1.666712e+02	-9.689728e+01	-8.743839e+01	-8.023680e+01	-6.695090e+01
is_fraud	1296675.0	5.788652e-03	7.586269e-02	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	1.000000e+00

In [45]: `train["amt"].describe()`

Out[45]:

```
count    1.296675e+06
mean      7.035104e+01
std       1.603160e+02
min       1.000000e+00
25%      9.650000e+00
50%      4.752000e+01
75%      8.314000e+01
max       2.894890e+04
Name: amt, dtype: float64
```

```
In [46]: ► donut = train["is_fraud"].value_counts().reset_index()

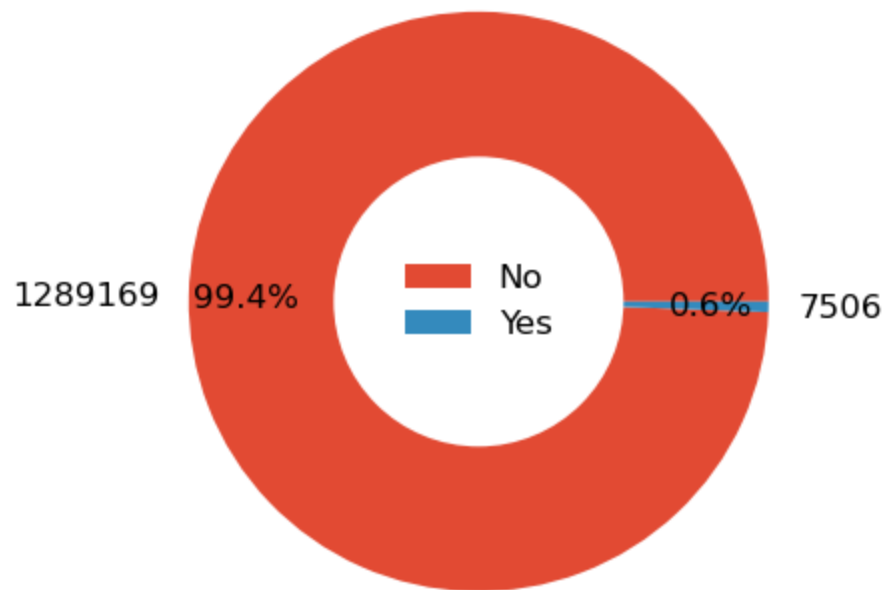
labels = ["No", "Yes"]
explode = (0, 0)

fig, ax = plt.subplots(dpi=120, figsize=(8, 4))
plt.pie(donut["is_fraud"],
        labels=donut["is_fraud"],
        autopct="%1.1f%%",
        pctdistance=0.8,
        explode=explode)

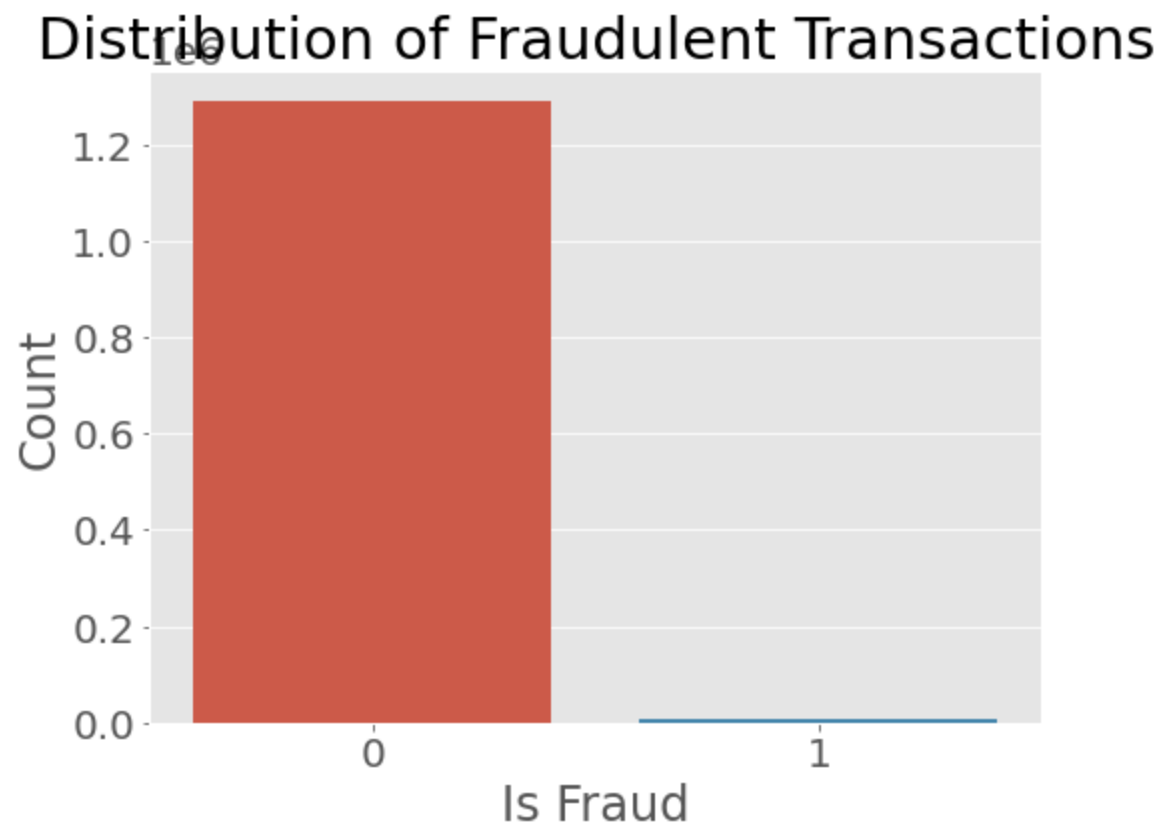
centre_circle = plt.Circle((0.0, 0.0), 0.5, fc='white')
fig = plt.gcf()
fig.gca().add_artist(centre_circle)

plt.title("Fraud proportion in Transactions")
plt.legend(labels, loc="center", frameon=False)
plt.show();
```

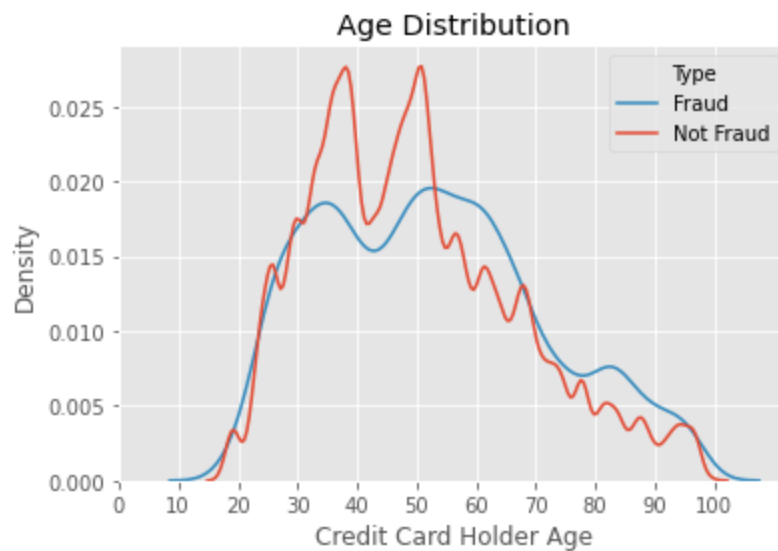
Fraud proportion in Transactions



```
In [57]: ▶ # Visualize the distribution of the target variable (fraudulent or not)
plt.figure(figsize=(8, 6))
sns.countplot(x='is_fraud', data=train)
plt.title('Distribution of Fraudulent Transactions')
plt.xlabel('Is Fraud')
plt.ylabel('Count')
plt.show()
```



```
In [47]: ▶ train['age'] = dt.date.today().year-pd.to_datetime(train['dob']).dt.year
ax = sns.kdeplot(x='age', data=train, hue='is_fraud', common_norm=False)
ax.set_xlabel('Credit Card Holder Age')
ax.set_ylabel('Density')
plt.xticks(np.arange(0, 110, 10))
plt.title('Age Distribution')
plt.legend(title='Type', labels=['Fraud', 'Not Fraud']);
```



```
In [50]: ► total = pd.concat([test,train])
total.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 1852394 entries, 0 to 1296674
Data columns (total 23 columns):
 #   Column                Dtype
---  -
 0   trans_date_trans_time  object
 1   cc_num                 int64
 2   merchant              object
 3   category              object
 4   amt                   float64
 5   first                 object
 6   last                  object
 7   gender                object
 8   street                object
 9   city                  object
10   state                 object
11   zip                   int64
12   lat                   float64
13   long                  float64
14   city_pop              int64
15   job                   object
16   dob                   object
17   trans_num             object
18   unix_time             int64
19   merch_lat             float64
20   merch_long            float64
21   is_fraud              int64
22   age                   float64
dtypes: float64(6), int64(5), object(12)
memory usage: 339.2+ MB
```

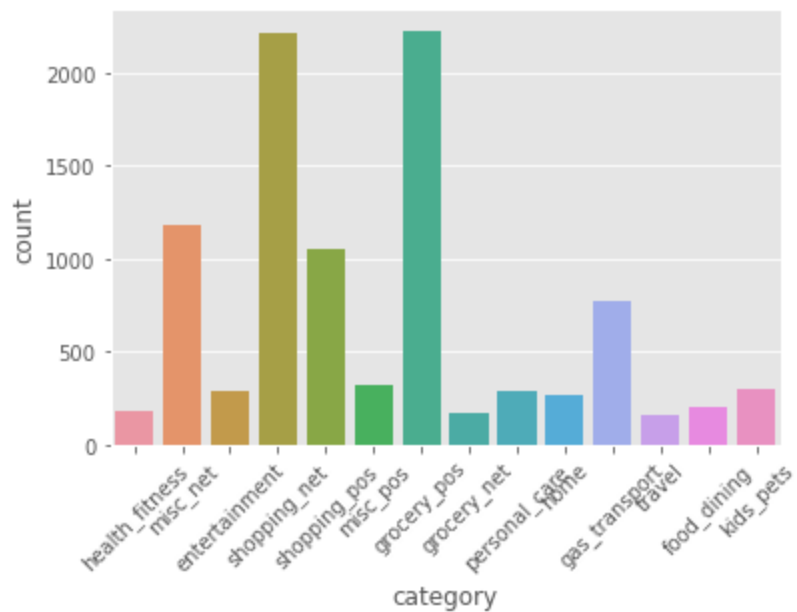
```
In [51]: ► total["is_fraud_cat"]=total.is_fraud.apply(lambda x: "T" if x==1 else "F")
total["is_fraud_cat"].astype("object")
```

```
Out[51]: 0      F
1      F
2      F
3      F
4      F
..
1296670  F
1296671  F
1296672  F
1296673  F
1296674  F
Name: is_fraud_cat, Length: 1852394, dtype: object
```

```
In [52]: ▶ sns.countplot(total[total['is_fraud_cat']=="T"].category)
plt.xticks(rotation=45)
plt.show()
```

C:\Users\dkond\anaconda3\lib\site-packages\seaborn_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

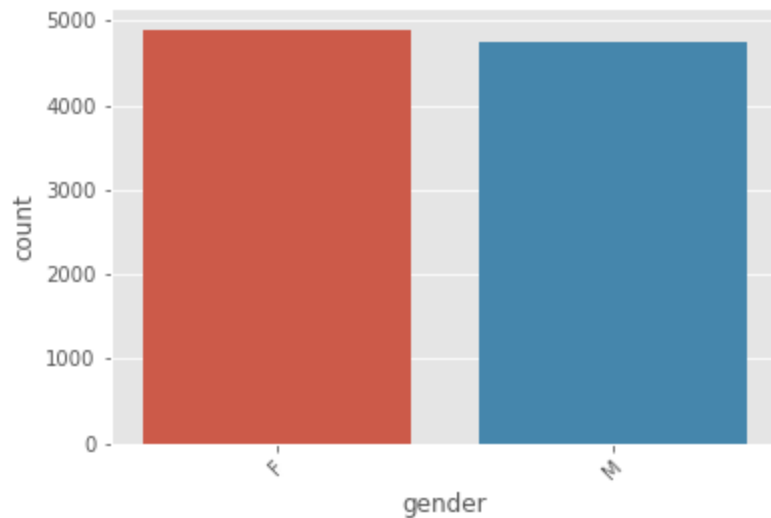
```
warnings.warn(
```




```
In [53]: ▶ sns.countplot(total[total['is_fraud_cat']=='T'].gender)
plt.xticks(rotation=45)
plt.show()
```

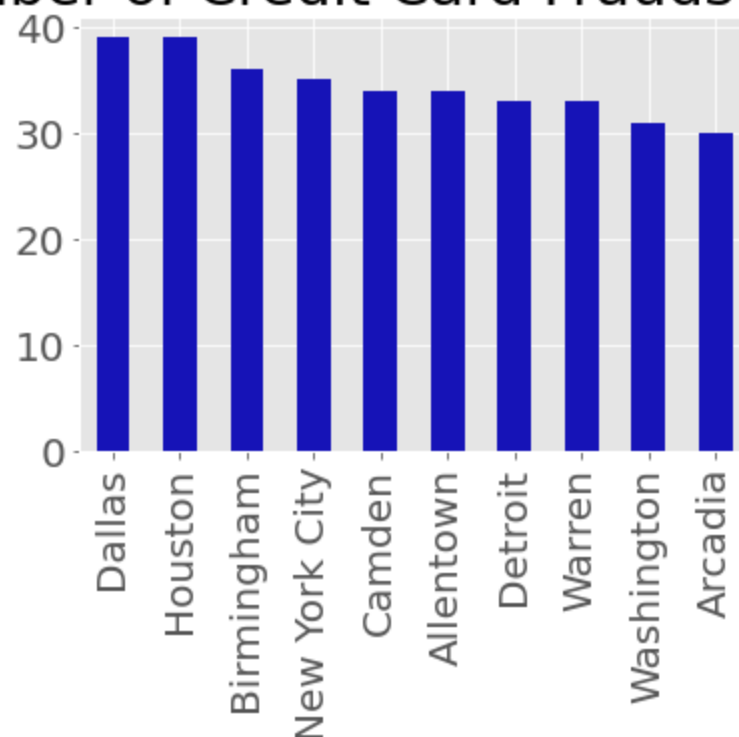
C:\Users\dkond\anaconda3\lib\site-packages\seaborn_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

```
warnings.warn(
```



```
In [54]: ▶ import random
def randomcolor():
    r = random.random()
    b = random.random()
    g = random.random()
    rgb = [r,g,b]
    return rgb
plt.rcParams.update({'font.size': 20})
total[total['is_fraud_cat']=="T"]['city'].value_counts(sort=True,ascending=False).head(10).plot(kind="bar")
plt.title("Number of Credit Card Frauds by City")
plt.show()
```

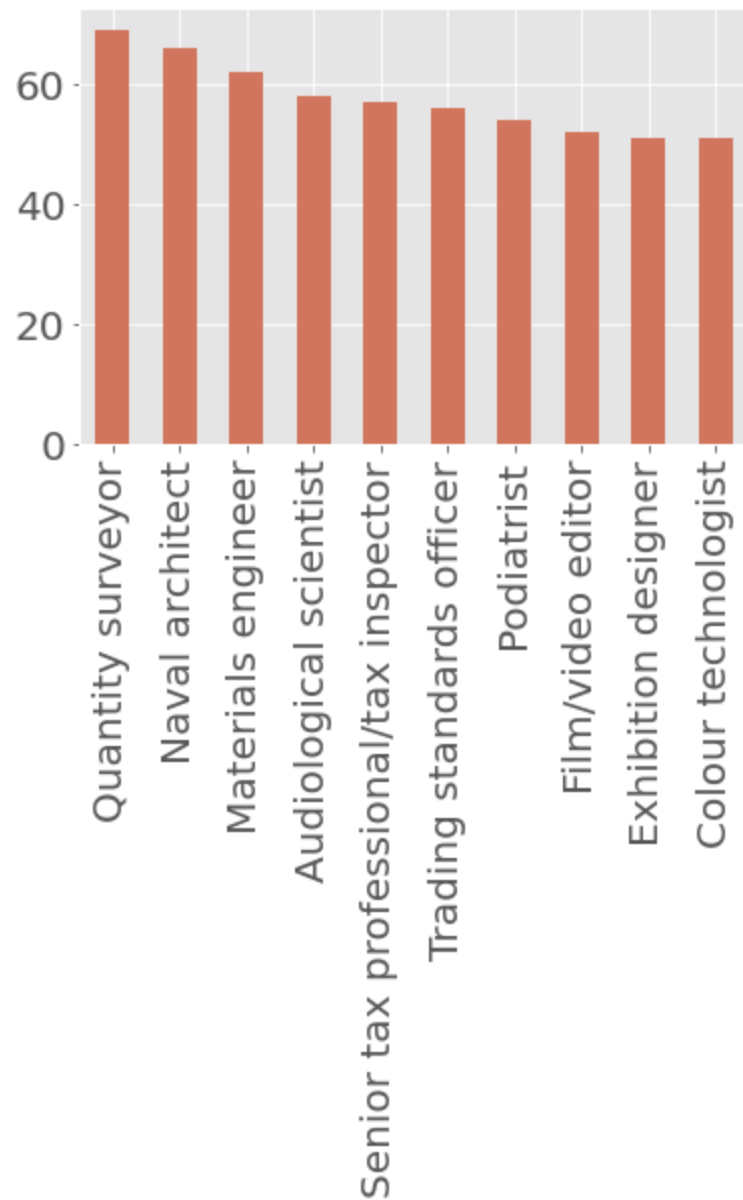
Number of Credit Card Frauds by City



```
In [55]: ▶ total[total['is_fraud_cat']=="T"]['job'].value_counts(sort=True,ascending=False).head(10).plot(kind="bar")  
plt.title("Number of Credit Card Frauds by Job")  
plt.show()
```



Number of Credit Card Frauds by Job



Imbalanced Data

Undersampling

```
In [12]: X = train.drop(columns=["is_fraud"])  
         y = train["is_fraud"]
```

```
In [13]: rus = RandomUnderSampler(sampling_strategy=0.1, random_state=23)
```

```
In [14]: X_undersampled, y_undersampled = rus.fit_resample(X, y)
```

```
In [15]: ▶ donut = y_undersampled.value_counts().reset_index()

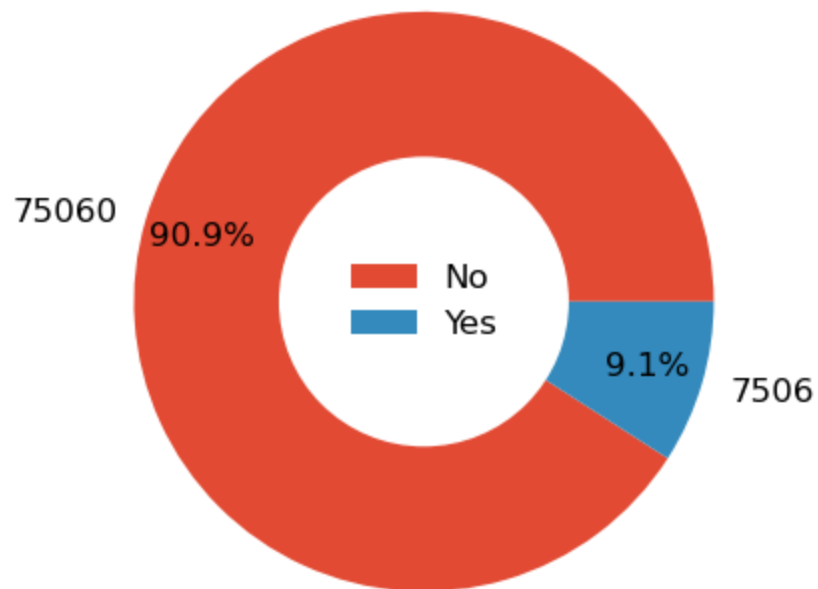
labels = ["No", "Yes"]
explode = (0, 0)

fig, ax = plt.subplots(dpi=120, figsize=(8, 4))
plt.pie(donut["is_fraud"],
        labels=donut["is_fraud"],
        autopct="%1.1f%%",
        pctdistance=0.8,
        explode=explode)

centre_circle = plt.Circle((0.0, 0.0), 0.5, fc='white')
fig = plt.gcf()
fig.gca().add_artist(centre_circle)

plt.title("Fraud Proportion with Undersampling")
plt.legend(labels, loc="center", frameon=False)
plt.show();
```

Fraud Proportion with Undersampling



Training the model

```
In [16]: ▶ from sklearn.preprocessing import StandardScaler, OneHotEncoder
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_curve, auc, conf
from imblearn.over_sampling import SMOTE
from sklearn.model_selection import train_test_split
from sklearn.utils import shuffle
```

```
In [17]: ▶ # Encode categorical variables
encoder = OneHotEncoder(drop='first')
categorical_cols = ['gender', 'category', 'state']
encoded_train_features = encoder.fit_transform(train[categorical_cols]).toarray()
encoded_test_features = encoder.transform(test[categorical_cols]).toarray()

# Feature scaling
scaler = StandardScaler()
numerical_cols = ['amt', 'lat', 'long', 'city_pop', 'unix_time', 'merch_lat', 'merch_long']
scaled_train_features = scaler.fit_transform(train[numerical_cols])
scaled_test_features = scaler.transform(test[numerical_cols])

# Concatenate encoded and scaled features for both train and test data
final_train_features = pd.concat([pd.DataFrame(encoded_train_features), pd.DataFrame(scaled_train_features)])
final_test_features = pd.concat([pd.DataFrame(encoded_test_features), pd.DataFrame(scaled_test_features)])

# Define target variables
train_target = train['is_fraud']
test_target = test['is_fraud']
```

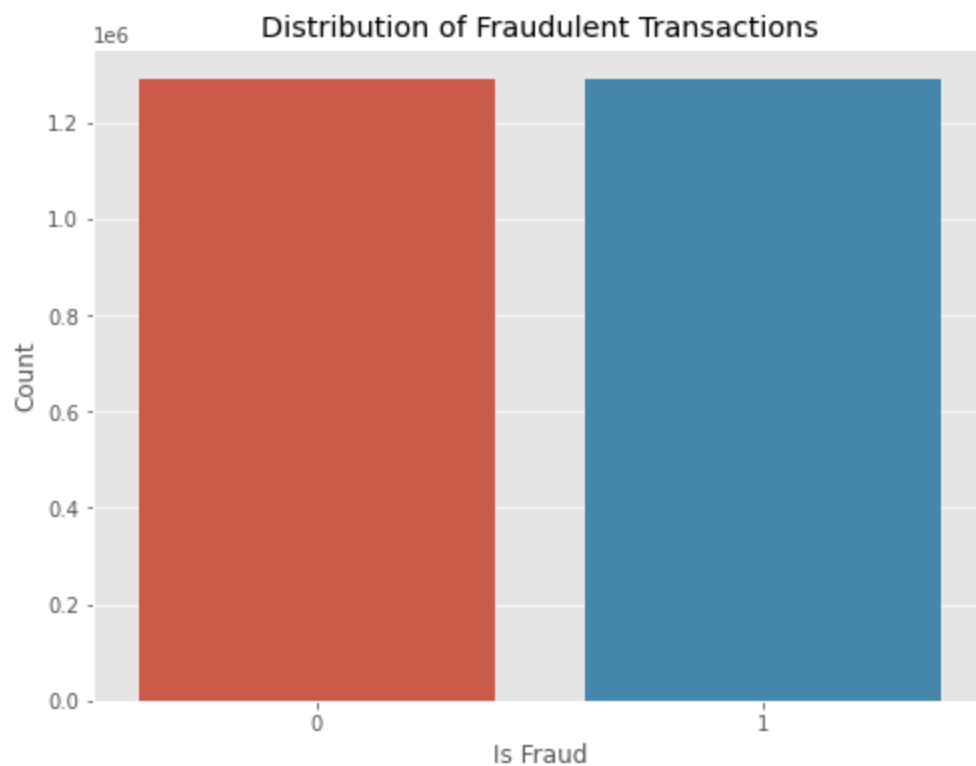
```
In [18]: ▶ # Generating synthetic data to balance the imbalanced dataset
smote = SMOTE(random_state=36)

x_train_resample, y_train_resample = smote.fit_resample(final_train_features, train_target)
```

```
In [19]: ▶ # checking newly created data
print('Current length of the training set: ', len(y_train_resample))
```

Current length of the training set: 2578338


```
In [20]: ▶ plt.figure(figsize=(8, 6))
sns.countplot(x=y_train_resample)
plt.title('Distribution of Fraudulent Transactions')
plt.xlabel('Is Fraud')
plt.ylabel('Count')
plt.show()
```



```
In [21]: ▶ X_shuffled, y_shuffled = shuffle(x_train_resample, y_train_resample, random_state=42)
```

```
In [22]: ▶ x_train, x_validation, y_train, y_validation = train_test_split(X_shuffled, y_shuffled, test_size=0.5)
```

```
In [23]: ▶ # for the initial selection process we will use a tiny
# portion of the actual training dataset
x_train_copy = x_train
y_train_copy = y_train

x_train = x_train[:10000]
y_train = y_train[:10000]
```

```
In [24]: ▶ # Train Logistic Regression model
lg_model = LogisticRegression()
lg_model.fit(x_train, y_train)

# Make predictions on test data
lg_predictions = lg_model.predict(x_validation)

# Calculate evaluation metrics on test data
lg_accuracy = accuracy_score(y_validation, lg_predictions)

# Print evaluation metrics with 3 decimal places, multiplied by 100
print("Logistic Regression Accuracy: {:.3f}%".format(lg_accuracy * 100))
```

C:\Users\dkond\anaconda3\lib\site-packages\sklearn\linear_model_logistic.py:814: ConvergenceWarning: lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html> (<https://scikit-learn.org/stable/modules/preprocessing.html>)

Please also refer to the documentation for alternative solver options:

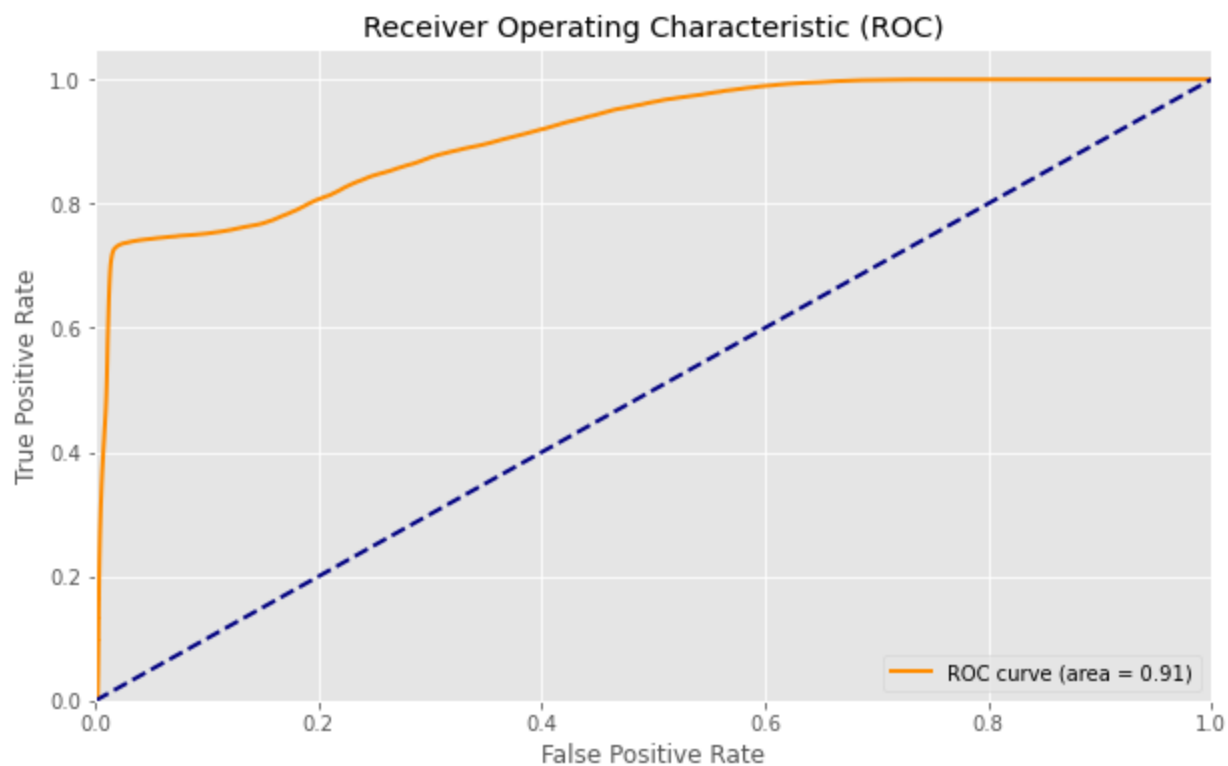
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression (https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)

```
n_iter_i = _check_optimize_result(
```

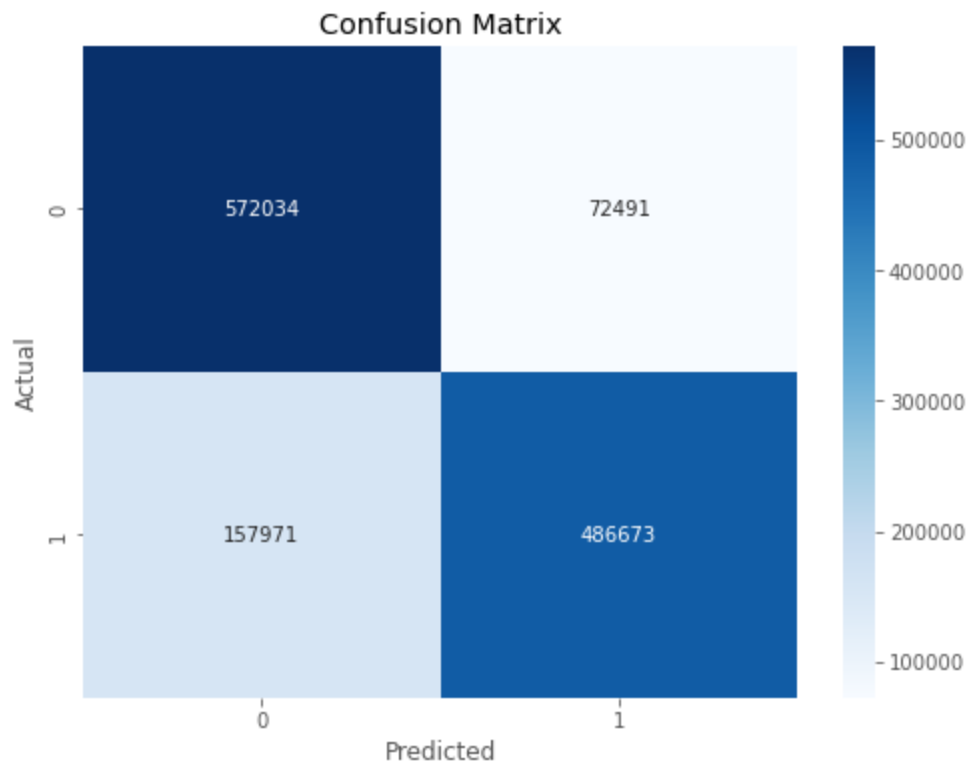
Logistic Regression Accuracy: 82.123%

```
In [25]: ▶ # Calculate ROC curve and AUC
probs = lg_model.predict_proba(x_validation)[: , 1]
fpr, tpr, thresholds = roc_curve(y_validation, probs)
roc_auc = auc(fpr, tpr)

# Plot ROC curve
plt.figure(figsize=(10, 6))
plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (area = %0.2f)' % roc_auc)
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic (ROC)')
plt.legend(loc="lower right")
plt.show()
```



```
In [26]: ▶ # Calculate and plot confusion matrix
conf_matrix = confusion_matrix(y_validation, lg_predictions)
plt.figure(figsize=(8, 6))
sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues')
plt.title('Confusion Matrix')
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.show()
```



```
In [27]: ► # Train SVM model
from sklearn.svm import SVC
svm_model = SVC(kernel='poly')
svm_model.fit(x_train, y_train)

# Make predictions on test data
svm_predictions = svm_model.predict(x_validation)

# Calculate evaluation metrics on test data
svm_accuracy = accuracy_score(y_validation, svm_predictions)

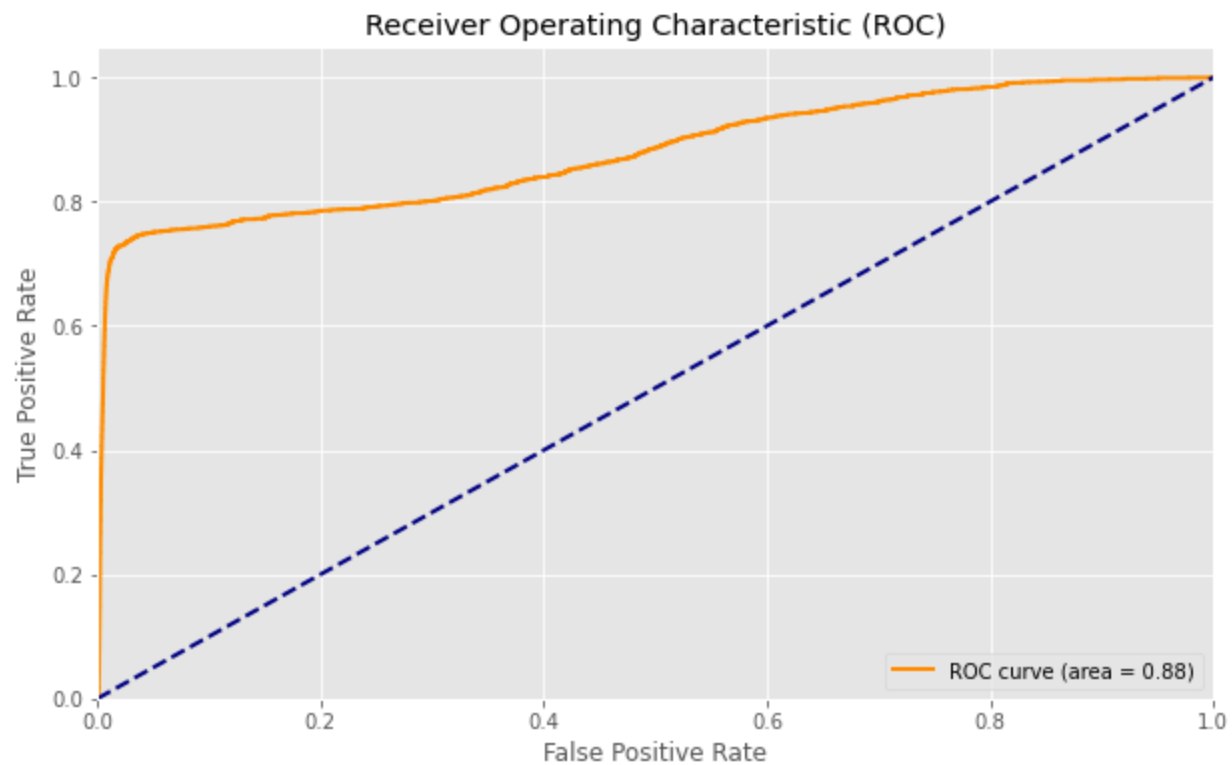
# Print evaluation metrics with 3 decimal places, multiplied by 100
print("SVM Accuracy: {:.3f}%".format(svm_accuracy * 100))
```

SVM Accuracy: 86.069%

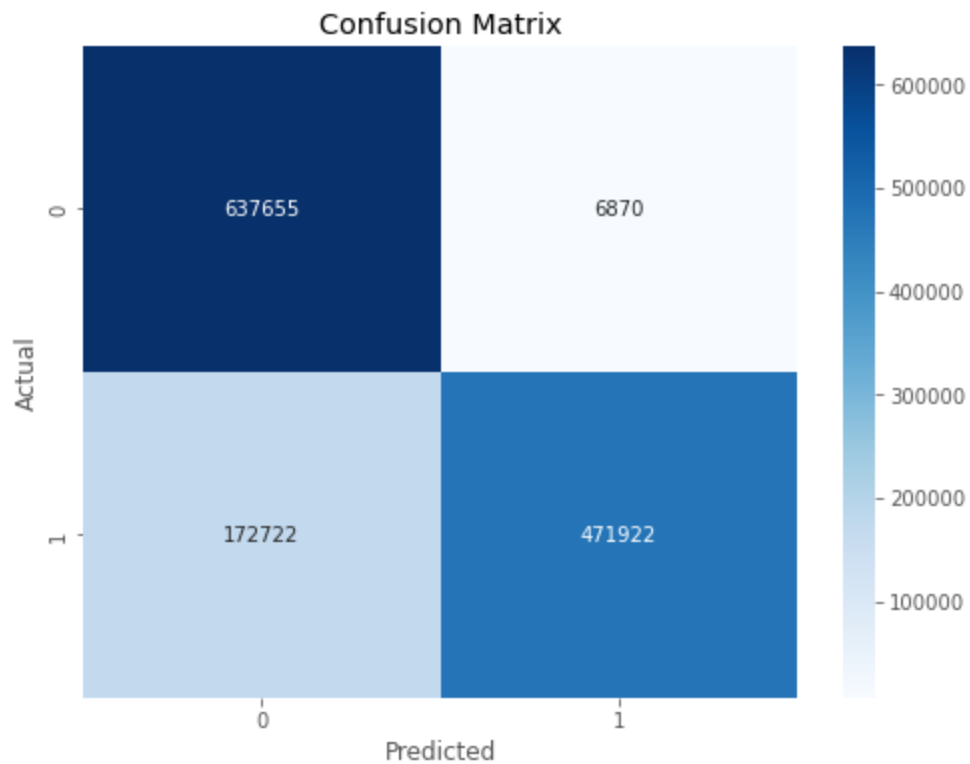
```
In [28]: ► # Calculate decision scores for the positive class
decision_scores = svm_model.decision_function(final_test_features)

# Calculate ROC curve and AUC
fpr, tpr, thresholds = roc_curve(test_target, decision_scores)
roc_auc = auc(fpr, tpr)

# Plot ROC curve
plt.figure(figsize=(10, 6))
plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (area = %0.2f)' % roc_auc)
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic (ROC)')
plt.legend(loc="lower right")
plt.show()
```



```
In [29]: ▶ # Calculate and plot confusion matrix
conf_matrix = confusion_matrix(y_validation, svm_predictions)
plt.figure(figsize=(8, 6))
sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues')
plt.title('Confusion Matrix')
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.show()
```




```
In [30]: ► # Train KNN model
from sklearn.neighbors import KNeighborsClassifier

knn_model = KNeighborsClassifier()
knn_model.fit(x_train, y_train)
# Make predictions on test data
knn_predictions = knn_model.predict(x_validation)

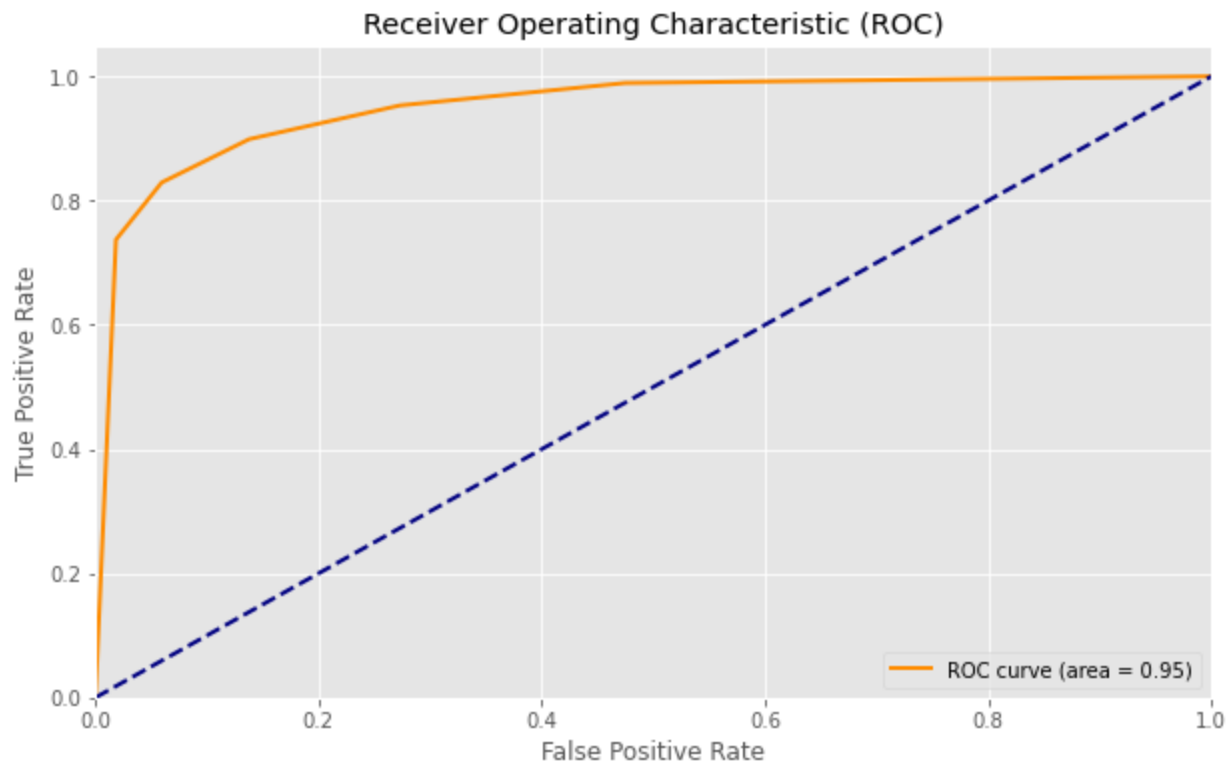
# Calculate evaluation metrics on test data
knn_accuracy = accuracy_score(y_validation, knn_predictions)

# Print evaluation metrics with 3 decimal places, multiplied by 100
print("KNN Accuracy: {:.3f}%".format(knn_accuracy * 100))
```

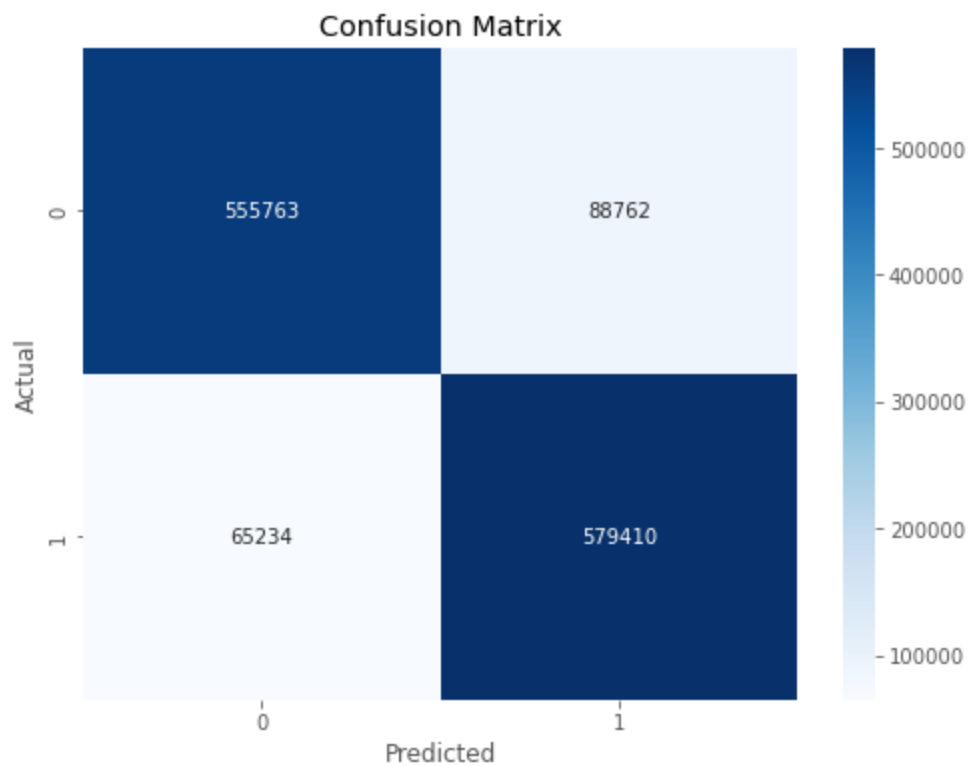
KNN Accuracy: 88.055%

```
In [31]: ▶ # Calculate ROC curve and AUC
probs = knn_model.predict_proba(x_validation)[: , 1]
fpr, tpr, thresholds = roc_curve(y_validation, probs)
roc_auc = auc(fpr, tpr)

# Plot ROC curve
plt.figure(figsize=(10, 6))
plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (area = %0.2f)' % roc_auc)
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic (ROC)')
plt.legend(loc="lower right")
plt.show()
```



```
In [32]: ▶ # Calculate and plot confusion matrix
conf_matrix = confusion_matrix(y_validation, knn_predictions)
plt.figure(figsize=(8, 6))
sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues')
plt.title('Confusion Matrix')
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.show()
```



```
In [33]: ► # Train Random Forest model
from sklearn.ensemble import RandomForestClassifier

rf_model = RandomForestClassifier()
rf_model.fit(x_train, y_train)
# Make predictions on test data
rf_predictions = rf_model.predict(x_validation)

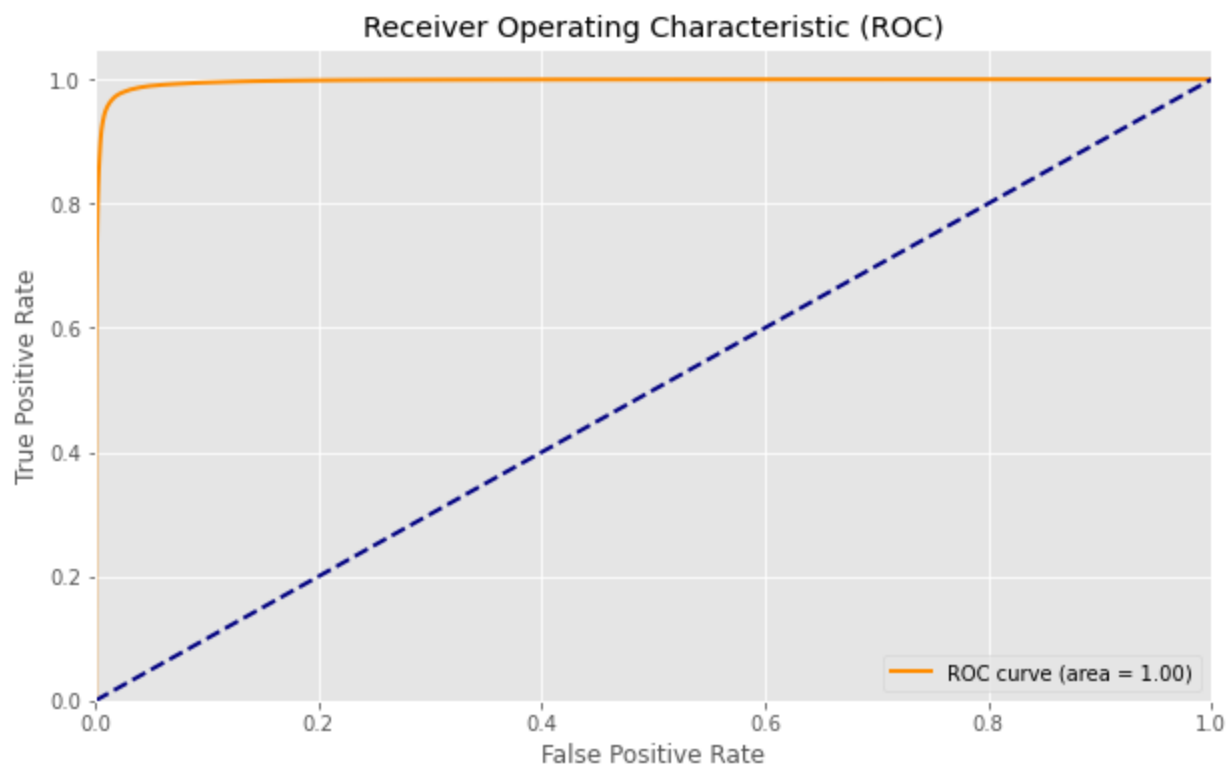
# Calculate evaluation metrics on test data
rf_accuracy = accuracy_score(y_validation, rf_predictions)

# Print evaluation metrics with 3 decimal places, multiplied by 100
print("Random Forest Accuracy: {:.3f}%".format(rf_accuracy * 100))
```

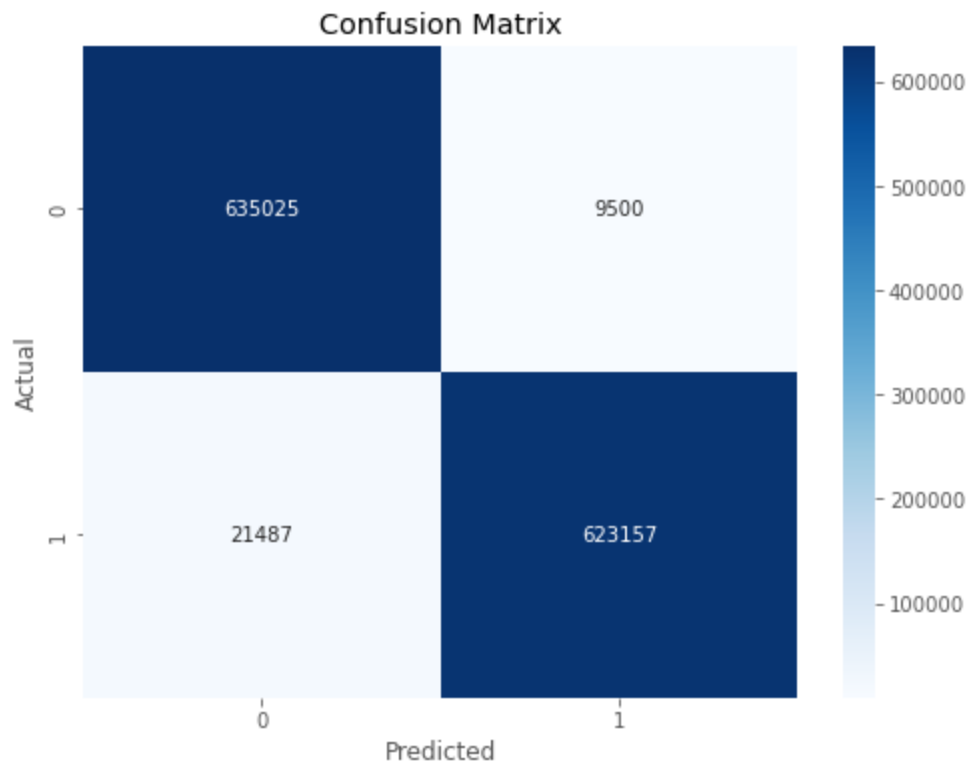
Random Forest Accuracy: 97.596%

```
In [34]: ▶ # Calculate ROC curve and AUC
probs = rf_model.predict_proba(x_validation)[: , 1]
fpr, tpr, thresholds = roc_curve(y_validation, probs)
roc_auc = auc(fpr, tpr)

# Plot ROC curve
plt.figure(figsize=(10, 6))
plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (area = %0.2f)' % roc_auc)
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic (ROC)')
plt.legend(loc="lower right")
plt.show()
```



```
In [35]: ▶ # Calculate and plot confusion matrix
conf_matrix = confusion_matrix(y_validation, rf_predictions)
plt.figure(figsize=(8, 6))
sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues')
plt.title('Confusion Matrix')
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.show()
```



```
In [36]: ▶ import pandas as pd
from sklearn.metrics import roc_auc_score, f1_score, precision_score, recall_score

# Define model names and instances
model_names = ['Logistic Regression', 'SVM', 'KNN', 'Random Forest']
model_instances = [lg_model, svm_model, knn_model, rf_model]

# Initialize lists to store accuracy and ROC scores
accuracy_scores = []
roc_scores = []
f1_scores = []
precision_scores = []
recall_scores = []

# Calculate accuracy and ROC scores for each model
for model in model_instances:
    predictions = model.predict(final_test_features)
    accuracy = accuracy_score(test_target, predictions)
    roc_score = roc_auc_score(test_target, predictions)
    accuracy_scores.append(accuracy)
    roc_scores.append(roc_score)
    f1_scores.append(f1_score(test_target, predictions))
    precision_scores.append(precision_score(test_target, predictions))
    recall_scores.append(recall_score(test_target, predictions))

# Create a DataFrame to compare results
results_df = pd.DataFrame({
    'Model': model_names,
    'Accuracy': accuracy_scores,
    'ROC Score': roc_scores,
    'F1 Score': f1_scores,
    'Precision Score': precision_scores,
    'Recall Score': recall_scores,
})

# Print the comparison table
print(results_df)
```

	Model	Accuracy	ROC Score	F1 Score	Precision Score	\
0	Logistic Regression	0.918964	0.823722	0.064832	0.033927	
1	SVM	0.971955	0.854268	0.168401	0.095083	
2	KNN	0.840139	0.806912	0.036004	0.018431	
3	Random Forest	0.985331	0.917637	0.308918	0.188789	

	Recall Score
0	0.727739
1	0.735664
2	0.773427
3	0.849417