



```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
```

```
train = pd.read_csv("fraudTrain.csv")
test = pd.read_csv("fraudTest.csv")
```


```
data = pd.concat([train, test])
data.describe()
```



	Unnamed: 0	cc_num	amt	zip	lat	l
count	1.852394e+06	1.852394e+06	1.852394e+06	1.852394e+06	1.852394e+06	1.852394e
mean	5.371934e+05	4.173860e+17	7.006357e+01	4.881326e+04	3.853931e+01	-9.022783e
std	3.669110e+05	1.309115e+18	1.592540e+02	2.688185e+04	5.071470e+00	1.374789e
min	0.000000e+00	6.041621e+10	1.000000e+00	1.257000e+03	2.002710e+01	-1.656723e
25%	2.315490e+05	1.800429e+14	9.640000e+00	2.623700e+04	3.466890e+01	-9.679800e
50%	4.630980e+05	3.521417e+15	4.745000e+01	4.817400e+04	3.935430e+01	-8.747690e
75%	8.335758e+05	4.642255e+15	8.310000e+01	7.204200e+04	4.194040e+01	-8.015800e
max	1.296674e+06	4.992346e+18	2.894890e+04	9.992100e+04	6.669330e+01	-6.795030e



```
print(train.shape)
print(test.shape)
```



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

```
display(data.head())
print(data.describe())
print(data.isnull().sum())
```



	Unnamed: 0	trans_date_trans_time	cc_num	merchant	category	amt	first	l
0	0	2019-01-01 00:00:18	2703186189652095	fraud_Rippin, Kub and Mann	misc_net	4.97	Jennifer	Ba
1	1	2019-01-01 00:00:44	630423337322	fraud_Heller, Gutmann and Zieme	grocery_pos	107.23	Stephanie	
2	2	2019-01-01 00:00:51	38859492057661	fraud_Lind-Buckridge	entertainment	220.11	Edward	Sanc
3	3	2019-01-01 00:01:16	3534093764340240	fraud_Kutch, Hermiston and	gas_transport	45.00	Jeremy	Wi

4

4

2019-01-01 00:03:06

375534208663984

fraud_Keeling-Crist

misc_pos

41.96

Tyler

Ga

5 rows × 23 columns

	Unnamed: 0	cc_num	amt	zip	lat \
count	1.852394e+06	1.852394e+06	1.852394e+06	1.852394e+06	1.852394e+06
mean	5.371934e+05	4.173860e+17	7.006357e+01	4.881326e+04	3.853931e+01
std	3.669110e+05	1.309115e+18	1.592540e+02	2.688185e+04	5.071470e+00
min	0.000000e+00	6.041621e+10	1.000000e+00	1.257000e+03	2.002710e+01
25%	2.315490e+05	1.800429e+14	9.640000e+00	2.623700e+04	3.466890e+01
50%	4.630980e+05	3.521417e+15	4.745000e+01	4.817400e+04	3.935430e+01
75%	8.335758e+05	4.642255e+15	8.310000e+01	7.204200e+04	4.194040e+01
max	1.296674e+06	4.992346e+18	2.894890e+04	9.992100e+04	6.669330e+01


	long	city_pop	unix_time	merch_lat	merch_long \
count	1.852394e+06	1.852394e+06	1.852394e+06	1.852394e+06	1.852394e+06
mean	-9.022783e+01	8.864367e+04	1.358674e+09	3.853898e+01	-9.022794e+01
std	1.374789e+01	3.014876e+05	1.819508e+07	5.105604e+00	1.375969e+01
min	-1.656723e+02	2.300000e+01	1.325376e+09	1.902742e+01	-1.666716e+02
25%	-9.679800e+01	7.410000e+02	1.343017e+09	3.474012e+01	-9.689944e+01
50%	-8.747690e+01	2.443000e+03	1.357089e+09	3.936890e+01	-8.744069e+01
75%	-8.015800e+01	2.032800e+04	1.374581e+09	4.195626e+01	-8.024511e+01
max	-6.795030e+01	2.906700e+06	1.388534e+09	6.751027e+01	-6.695090e+01

	is_fraud
count	1.852394e+06
mean	5.210015e-03
std	7.199217e-02
min	0.000000e+00
25%	0.000000e+00
50%	0.000000e+00
75%	0.000000e+00
max	1.000000e+00

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

test.info

 Show hidden output

train.info

 Show hidden output

```

from sklearn.preprocessing import LabelEncoder
label_encoders = {}

label_encode_cols = ['merchant', 'category', 'gender', 'state', 'job']
for col in label_encode_cols:
    le = LabelEncoder()
    data[col] = le.fit_transform(data[col])
    label_encoders[col] = le

    train[col] = le.fit_transform(train[col])
    label_encoders[col] = le

    test[col] = le.fit_transform(test[col])
    label_encoders[col] = le

data['trans_date_trans_time'] = pd.to_datetime(data['trans_date_trans_time'])
data['dob'] = pd.to_datetime(data['dob'])

data['transaction_year'] = data['trans_date_trans_time'].dt.year
data['transaction_month'] = data['trans_date_trans_time'].dt.month
data['transaction_day'] = data['trans_date_trans_time'].dt.day
data['transaction_hour'] = data['trans_date_trans_time'].dt.hour

data['birth_year'] = data['dob'].dt.year
data['birth_month'] = data['dob'].dt.month
data['birth_day'] = data['dob'].dt.day
data.drop(['trans_date_trans_time', 'dob'], axis=1, inplace=True)
train['trans_date_trans_time'] = pd.to_datetime(train['trans_date_trans_time'])
train['dob'] = pd.to_datetime(train['dob'])

train['transaction_year'] = train['trans_date_trans_time'].dt.year
train['transaction_month'] = train['trans_date_trans_time'].dt.month
train['transaction_day'] = train['trans_date_trans_time'].dt.day
train['transaction_hour'] = train['trans_date_trans_time'].dt.hour

train['birth_year'] = train['dob'].dt.year
train['birth_month'] = train['dob'].dt.month
train['birth_day'] = train['dob'].dt.day

train.drop(['trans_date_trans_time', 'dob'], axis=1, inplace=True)

test['trans_date_trans_time'] = pd.to_datetime(test['trans_date_trans_time'])
test['dob'] = pd.to_datetime(test['dob'])

test['transaction_year'] = test['trans_date_trans_time'].dt.year
test['transaction_month'] = test['trans_date_trans_time'].dt.month
test['transaction_day'] = test['trans_date_trans_time'].dt.day
test['transaction_hour'] = test['trans_date_trans_time'].dt.hour
test['birth_year'] = test['dob'].dt.year
test['birth_month'] = test['dob'].dt.month
test['birth_day'] = test['dob'].dt.day

test.drop(['trans_date_trans_time', 'dob'], axis=1, inplace=True)

data.drop(['first', 'last', 'street', 'city', 'trans_num'], axis=1, inplace=True)
train.drop(['first', 'last', 'street', 'city', 'trans_num'], axis=1, inplace=True)
test.drop(['first', 'last', 'street', 'city', 'trans_num'], axis=1, inplace=True)

```

```
print(train.shape)
print(test.shape)
print(data.shape)
```

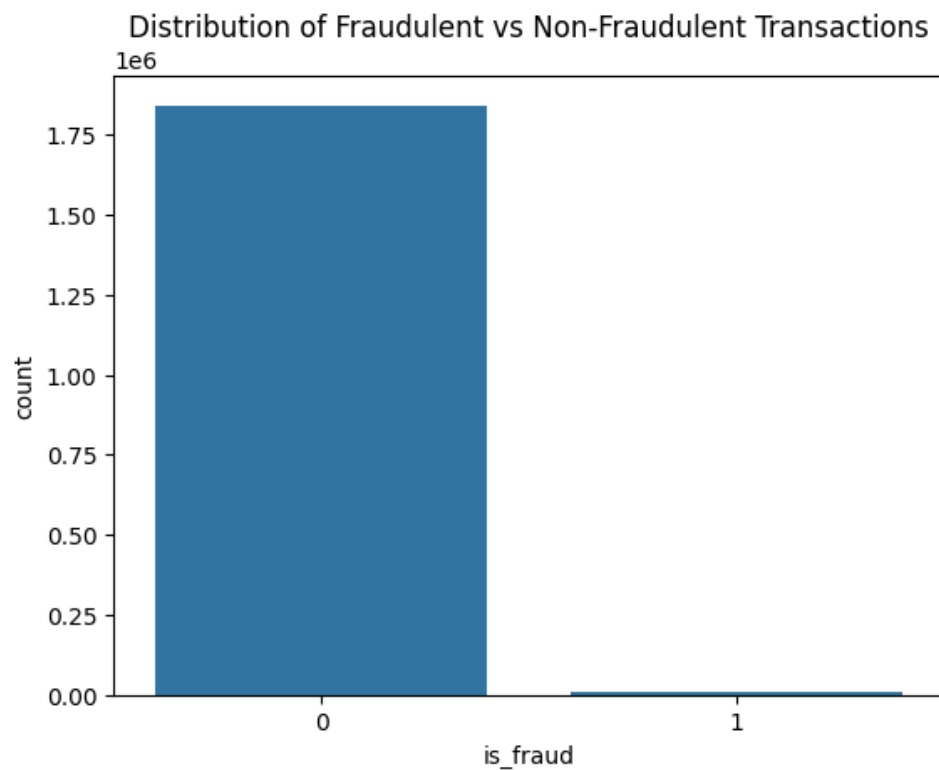
```
↗ (1296675, 23)
   (555719, 23)
   (1852394, 23)
```

```
print(data.head(0))
print(data.head())
print(data.describe())
print(data.isnull().sum())
```

```
↗ max    1.000000e+00  5.000000e+01  9.992100e+04  6.669330e+01 -6.795030e+01
    ...
count    ...  1.852394e+06  1.852394e+06  1.852394e+06    1.852394e+06
mean     ...  3.853898e+01 -9.022794e+01  5.210015e-03    2.019501e+03
```

```
import matplotlib.pyplot as plt
import seaborn as sns

sns.countplot(data=data, x='is_fraud')
plt.title('Distribution of Fraudulent vs Non-Fraudulent Transactions')
plt.show()
```

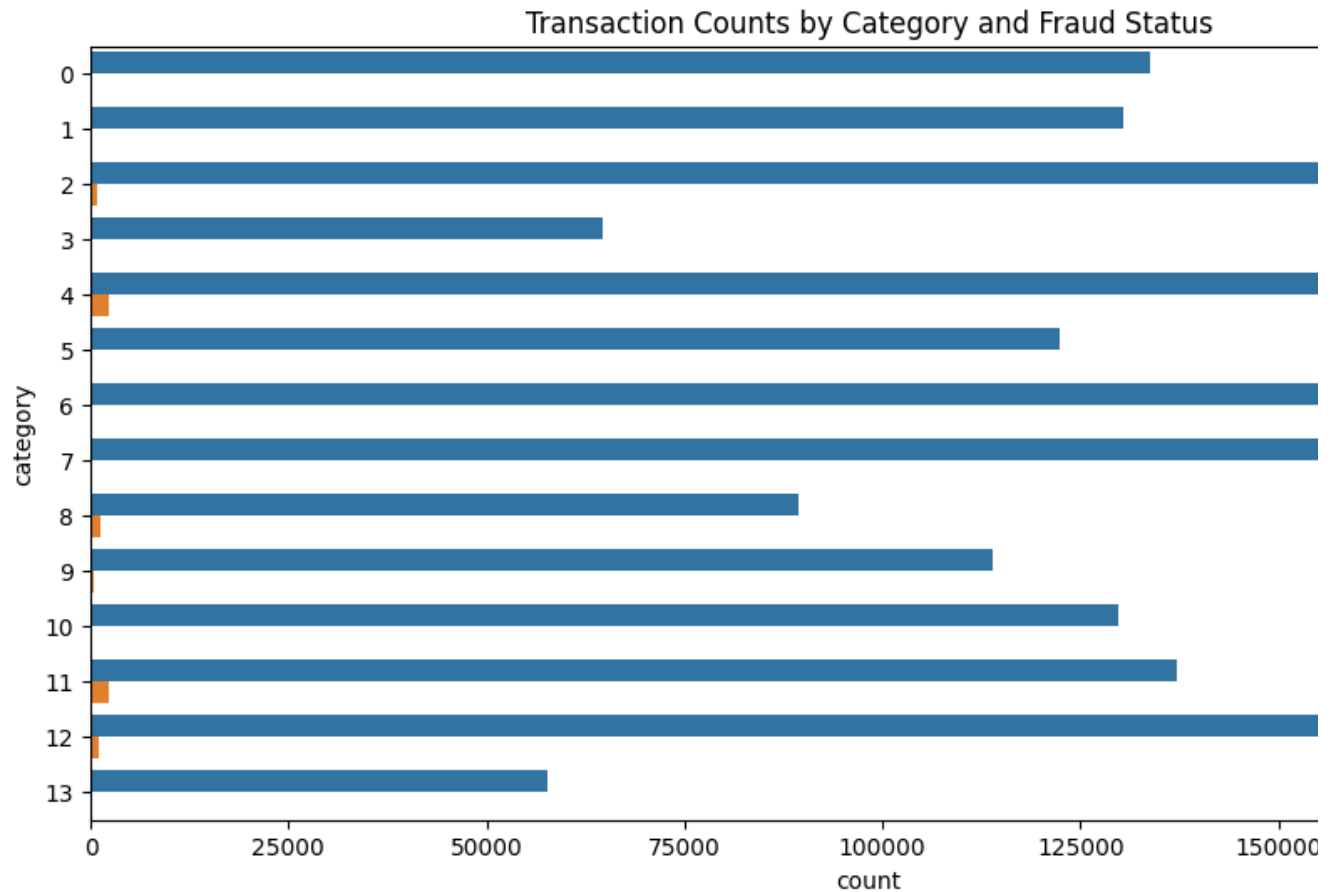


```
print(data.index.duplicated().sum())
data = data.reset_index(drop=True)
print(data.index.duplicated().sum())
```

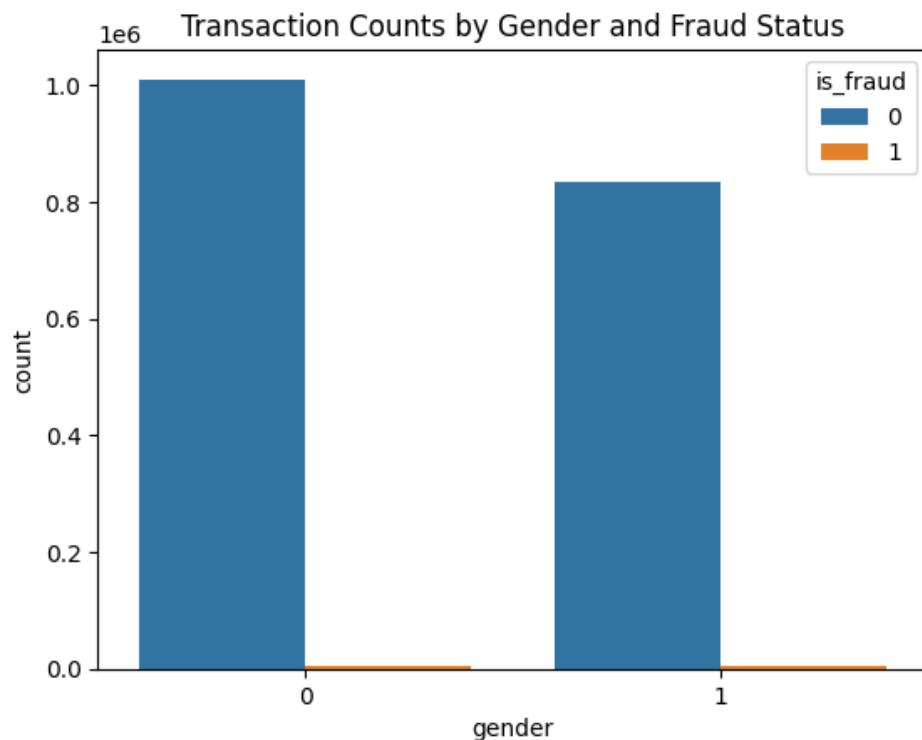


```
555719
0
```

```
plt.figure(figsize=(12, 6))
sns.countplot(data=data, y='category', hue='is_fraud')
plt.title('Transaction Counts by Category and Fraud Status')
plt.xticks(rotation=0)
plt.show()
```

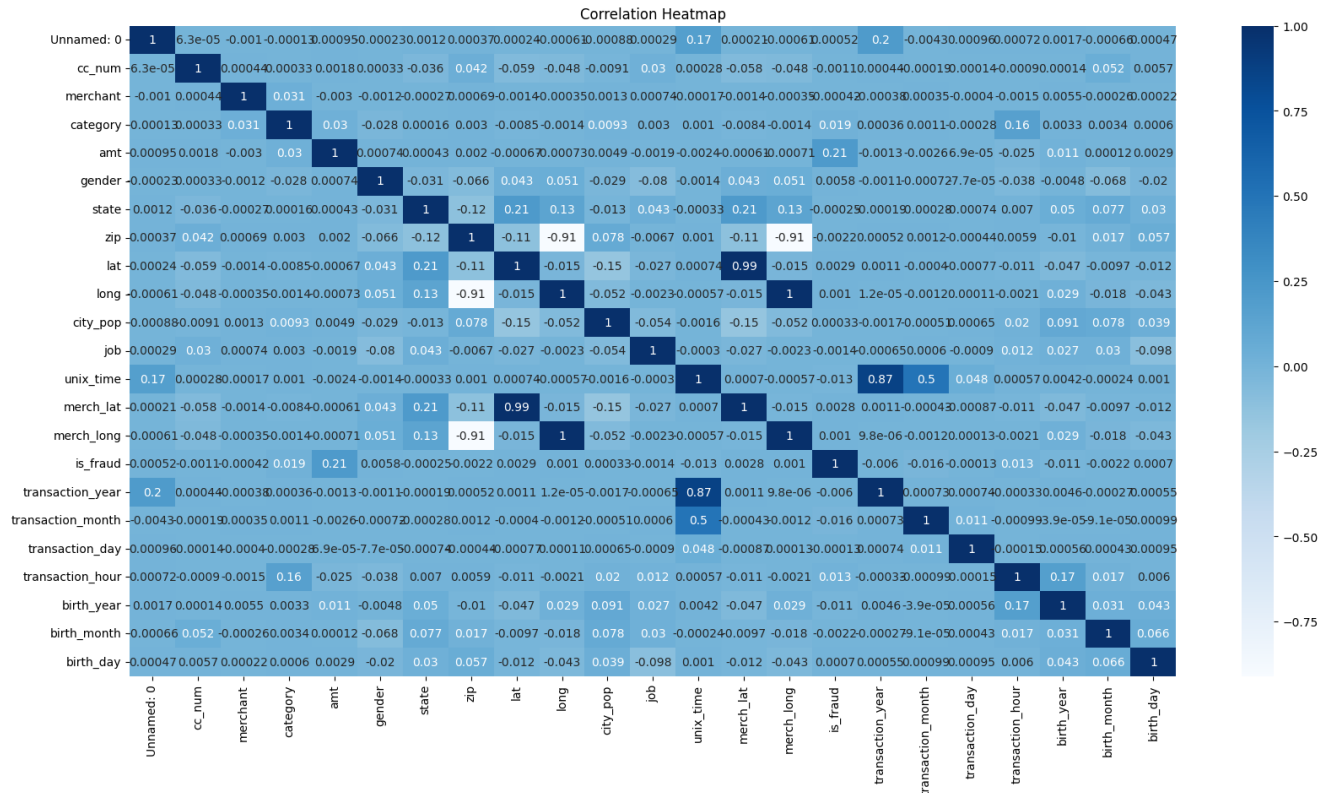


```
## The 0 represent male and 1 represent female
sns.countplot(data=data, x='gender', hue='is_fraud')
plt.title('Transaction Counts by Gender and Fraud Status')
plt.show()
```



```
plt.figure(figsize=(20,10))
sns.heatmap(data.corr(), annot=True, cmap='Blues')
plt.title('Correlation Heatmap')
```

plt.show()



```
X = data.drop('is_fraud', axis=1)
y = data['is_fraud']
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

```
# Training the model for logisticRegression
log_model = LogisticRegression(max_iter=1000)
log_model.fit(X_train, y_train)
y_pred = log_model.predict(X_test)
print(classification_report(y_test, y_pred))
print(confusion_matrix(y_test, y_pred))
print("Accuracy:", accuracy_score(y_test, y_pred))
```



```
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:460: ConvergenceWarning: lbfgs
ABNORMAL_TERMINATION_IN_LNSRCH.
```

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

<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

```
n_iter_i = _check_optimize_result(
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1471: UndefinedMetricWarning
_warn_prf(average, modifier, msg_start, len(result))
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1471: UndefinedMetricWarning
_warn_prf(average, modifier, msg_start, len(result))
```

	precision	recall	f1-score	support
0	0.99	1.00	1.00	368526
1	0.00	0.00	0.00	1953
accuracy			0.99	370479
macro avg	0.50	0.50	0.50	370479
weighted avg	0.99	0.99	0.99	370479

```
[[368526    0]
 [ 1953    0]]
```

Accuracy: 0.9947284461467452

```
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1471: UndefinedMetricWarning
  _warn_prf(average, modifier, msg_start, len(result))
```

```
# Training the model with DecisionTreeClassifier
tree_model = DecisionTreeClassifier()
tree_model.fit(X_train, y_train)
```

```
y_pred = tree_model.predict(X_test)
print(classification_report(y_test, y_pred))
print(confusion_matrix(y_test, y_pred))
print("Accuracy:", accuracy_score(y_test, y_pred))
```

```
↔
```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	368526
1	0.80	0.84	0.82	1953
accuracy			1.00	370479
macro avg	0.90	0.92	0.91	370479
weighted avg	1.00	1.00	1.00	370479

```
[[368112    414]
 [   306   1647]]
```

Accuracy: 0.9980565700080166

```
# Training the model with logistic Regression
lr_model = LogisticRegression(max_iter=1000)
lr_model.fit(X_train, y_train)
```

```
y_pred = lr_model.predict(X_test)
print(classification_report(y_test, y_pred))
print(confusion_matrix(y_test, y_pred))
print("Accuracy:", accuracy_score(y_test, y_pred))
```

```
↔ /usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:460: ConvergenceWarning: lbfgs
ABNORMAL_TERMINATION_IN_LNSRCH.
```

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

<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

```
n_iter_i = _check_optimize_result(
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1471: UndefinedMetricWarning
  _warn_prf(average, modifier, msg_start, len(result))
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1471: UndefinedMetricWarning
  _warn_prf(average, modifier, msg_start, len(result))
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1471: UndefinedMetricWarning
  _warn_prf(average, modifier, msg_start, len(result))
```

	precision	recall	f1-score	support
0	0.99	1.00	1.00	368526
1	0.00	0.00	0.00	1953
accuracy			0.99	370479
macro avg	0.50	0.50	0.50	370479
weighted avg	0.99	0.99	0.99	370479

```
[[368526    0]
 [ 1953    0]]
```

Accuracy: 0.9947284461467452