```python
import numpy as np
import pandas as pd
import seaborn as sns
import plotly.express as px
import matplotlib.pyplot as plt
from sklearn.model_selection import cross_val_score

from sklearn import metrics
from collections import Counter


train_df = pd.read_csv('/content/fraudTrain.csv')
test_df = pd.read_csv('/content/fraudTest.csv')


train_df.head()
```
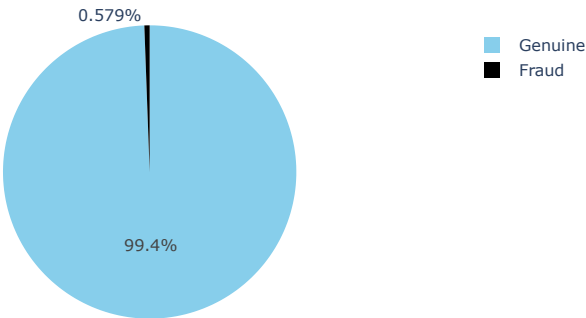
| | Unnamed: 0 | trans_date_trans_time | cc_num | merchant | category | amt |
|---|---|---|---|---|---|---|
| 0 | 0 | 2019-01-01 00:00:18 | 2703186189652095 | fraud_Rippin, Kub and Mann | misc_net | 4.97 |
| 1 | 1 | 2019-01-01 00:00:44 | 630423337322 | fraud_Heller, Gutmann and Zieme | grocery_pos | 107.23 |
| 2 | 2 | 2019-01-01 00:00:51 | 38859492057661 | fraud_Lind-Buckridge | entertainment | 220.11 |
| 3 | 3 | 2019-01-01 00:01:16 | 3534093764340240 | fraud_Kutch, Hermiston and Farrell | gas_transport | 45.00 |
| 4 | 4 | 2019-01-01 00:03:06 | 375534208663984 | fraud_Keeling-Crist | misc_pos | 41.96 |

5 rows × 23 columns

```python
fig = px.pie(values=train_df['is_fraud'].value_counts(), names=["Genuine","Fraud"] , width=700, height=400, color_discrete_sequence=["skyblue"
         ,title="Fraud vs Genuine transactions")
fig.show()
```
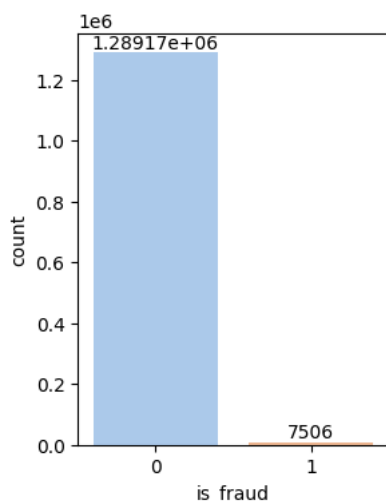
⤷

### Fraud vs Genuine transactions



```python
plt.figure(figsize=(3,4))
ax = sns.countplot(x='is_fraud',data=train_df,palette="pastel")
for i in ax.containers:
```

```
ax.bar_label(i,)
```



```
print('Genuine:', round(train_df['is_fraud'].value_counts()[0]/len(train_df) * 100,2), '% of the dataset')
print('Frauds:', round(train_df['is_fraud'].value_counts()[1]/len(train_df) * 100,2), '% of the dataset')
```

```
Genuine: 99.42 % of the dataset
Frauds: 0.58 % of the dataset
```

```
train_df.info(),test_df.info()
```

```
 #   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
<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
```

```
        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
       (None, None)
```

```
train_df.isnull().sum(),test_df.isnull().sum()
```

```
       (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,
        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)
```

```
drop_columns = ['Unnamed: 0','cc_num','merchant','trans_num','unix_time','first','last','street','zip']
train_df.drop(columns=drop_columns,inplace=True)
test_df.drop(columns=drop_columns,inplace=True)
```

```
print(train_df.shape)
print(test_df.shape)
```

```
       (1296675, 14)
       (555719, 14)
```

```
train_df['trans_date_trans_time']=pd.to_datetime(train_df['trans_date_trans_time'])
train_df['trans_date']=train_df['trans_date_trans_time'].dt.strftime('%Y-%m-%d')
train_df['trans_date']=pd.to_datetime(train_df['trans_date'])
train_df['dob']=pd.to_datetime(train_df['dob'])

test_df['trans_date_trans_time']=pd.to_datetime(test_df['trans_date_trans_time'])
test_df['trans_date']=test_df['trans_date_trans_time'].dt.strftime('%Y-%m-%d')
test_df['trans_date']=pd.to_datetime(test_df['trans_date'])
test_df['dob']=pd.to_datetime(test_df['dob'])
```

```python
train_df["age"] = train_df["trans_date"]-train_df["dob"]
train_df["age"]=train_df["age"].astype('timedelta64[Y]')

test_df["age"] = test_df["trans_date"]-test_df["dob"]
test_df["age"]=test_df["age"].astype('timedelta64[Y]')


train_df['Trans_month'] = pd.DatetimeIndex(train_df['trans_date']).month
train_df['Trans_year'] = pd.DatetimeIndex(train_df['trans_date']).year


train_df['Latitudinal_Distance'] = abs(round(train_df['merch_lat']-train_df['lat'],3))
train_df['Longitudinal_Distance'] = abs(round(train_df['merch_long']-train_df['long'],3))

test_df['Latitudinal_Distance'] = abs(round(test_df['merch_lat']-test_df['lat'],3))
test_df['Longitudinal_Distance'] = abs(round(test_df['merch_long']-test_df['long'],3))


drop_columns = ['trans_date_trans_time','city','lat','long','job','dob','merch_lat','merch_long','trans_date','state']
train_df.drop(columns=drop_columns,inplace=True)
test_df.drop(columns=drop_columns,inplace=True)


train_df.gender=train_df.gender.apply(lambda x: 1 if x=="M" else 0)
test_df.gender=test_df.gender.apply(lambda x: 1 if x=="M" else 0)


train_df = pd.get_dummies(train_df, columns=['category'], prefix='category')
test_df = pd.get_dummies(test_df, columns=['category'], prefix='category')

test_df = test_df.reindex(columns=train_df.columns, fill_value=0)


train_df.head()
```

|   | amt | gender | city_pop | is_fraud | age | Trans_month | Trans_year | Latitudinal_Distanc |
|---|-----|--------|----------|----------|-----|-------------|------------|---------------------|
| 0 | 4.97 | 0 | 3495 | 0 | 30.0 | 1 | 2019 | 0.06 |
| 1 | 107.23 | 0 | 149 | 0 | 40.0 | 1 | 2019 | 0.27 |
| 2 | 220.11 | 1 | 4154 | 0 | 56.0 | 1 | 2019 | 0.97 |
| 3 | 45.00 | 1 | 1939 | 0 | 51.0 | 1 | 2019 | 0.80 |
| 4 | 41.96 | 1 | 99 | 0 | 32.0 | 1 | 2019 | 0.25 |

5 rows × 23 columns

```python
test_df.head()
```

|   | amt | gender | city_pop | is_fraud | age | Trans_month | Trans_year | Latitudinal_Distance |
|---|-----|--------|----------|----------|-----|-------------|------------|----------------------|
| 0 | 2.86 | 1 | 333497 | 0 | 52.0 | 0 | 0 | 0.020 |
| 1 | 29.84 | 0 | 302 | 0 | 30.0 | 0 | 0 | 0.870 |
| 2 | 41.28 | 0 | 34496 | 0 | 49.0 | 0 | 0 | 0.177 |
| 3 | 60.05 | 1 | 54767 | 0 | 32.0 | 0 | 0 | 0.243 |
| 4 | 3.19 | 1 | 1126 | 0 | 64.0 | 0 | 0 | 0.706 |

5 rows × 23 columns

```python
X_train = train_df.drop('is_fraud', axis=1)
y_train = train_df['is_fraud']
X_test = test_df.drop('is_fraud', axis=1)
y_test = test_df['is_fraud']


from imblearn.over_sampling import SMOTE

smote = SMOTE(random_state=42)

X_train, y_train = smote.fit_resample(X_train, y_train)
```

```python
from sklearn.preprocessing import StandardScaler


scaler = StandardScaler()
scaler.fit(X_train)

X_train = scaler.transform(X_train)
X_test = scaler.transform(X_test)


from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import classification_report

clf = DecisionTreeClassifier(random_state=42)
clf.fit(X_train, y_train)
y_pred = clf.predict(X_test)
report = classification_report(y_test, y_pred)
print(report)
```

```
              precision    recall  f1-score   support

           0       1.00      0.99      1.00    553574
           1       0.33      0.72      0.45      2145

    accuracy                           0.99    555719
   macro avg       0.67      0.86      0.73    555719
weighted avg       1.00      0.99      0.99    555719
```