

```
In [1]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn import metrics

In [2]: train_df = pd.read_csv('Credit card/fraudTrain.csv.zip')
test_df = pd.read_csv('Credit card/fraudTest.csv.zip')
```

```
-----
FileNotFoundError                                Traceback (most recent call last)
Cell In[2], line 1
----> 1 train_df = pd.read_csv('fraudTrain.csv')
      2 test_df = pd.read_csv('fraudTest.csv')

File ~\AppData\Local\Programs\Python\Python312\Lib\site-packages\pandas\io\parsers\readers.py:1026, in read_csv(filepath_or_buffer, sep, delimiter, header, names, index_col, usecols, dtype, engine, converters, true_values, false_values, skipinitialspace, skiprows, skipfooter, nrows, na_values, keep_default_na, na_filter, verbose, skip_blank_lines, parse_dates, infer_datetime_for, rmat, keep_date_col, date_parser, date_format, dayfirst, cache_dates, iterator, chunksize, compression, thousands, decimal, lineterminator, quotechar, quoting, doublequote, escapechar, comment, encoding, encoding_errors, dialect, on_bad_lines, delim_whitespace, low_memory, memory_map, float_precision, storage_options, dtype_backend)
    1013 kwds_defaults = _refine_defaults_read(
    1014     dialect,
    1015     delimiter,
    1016     (...)
    1022     dtype_backend=dtype_backend,
    1023 )
    1024 kwds.update(kwds_defaults)
-> 1026 return _read(filepath_or_buffer, kwds)

File ~\AppData\Local\Programs\Python\Python312\Lib\site-packages\pandas\io\parsers\readers.py:620, in _read(filepath_or_buffer, kwds)
    617 _validate_names(kwds.get("names", None))
    619 # Create the parser.
-> 620 parser = TextFileReader(filepath_or_buffer, **kwds)
    622 if chunksize or iterator:
    623     return parser

File ~\AppData\Local\Programs\Python\Python312\Lib\site-packages\pandas\io\parsers\readers.py:1620, in TextFileReader.__init__(self, f, engine, **kwds)
    1617 self.options["has_index_names"] = kwds["has_index_names"]
    1619 self.handles: IOHandles | None = None
-> 1620 self._engine = self._make_engine(f, self.engine)

File ~\AppData\Local\Programs\Python\Python312\Lib\site-packages\pandas\io\parsers\readers.py:1880, in TextFileReader._make_engine(self, f, engine)
    1878     if "b" not in mode:
    1879         mode += "b"
-> 1880 self.handles = get_handle(
    1881     f,
    1882     mode,
    1883     encoding=self.options.get("encoding", None),
    1884     compression=self.options.get("compression", None),
    1885     memory_map=self.options.get("memory_map", False),
    1886     is_text=is_text,
    1887     errors=self.options.get("encoding_errors", "strict"),
    1888     storage_options=self.options.get("storage_options", None),
    1889 )
    1890 assert self.handles is not None
    1891 f = self.handles.handle

File ~\AppData\Local\Programs\Python\Python312\Lib\site-packages\pandas\io\common.py:873, in get_handle(path_or_buf, mode, encoding, compression, memory_map, is_text, errors, storage_options)
    868 elif isinstance(handle, str):
    869     # Check whether the filename is to be opened in binary mode.
    870     # Binary mode does not support 'encoding' and 'newline'.
    871     if ioargs.encoding and "b" not in ioargs.mode:
    872         # encoding
-> 873     handle = open(
    874         handle,
    875         ioargs.mode,
    876         encoding=ioargs.encoding,
    877         errors=errors,
    878         newline="",
    879     )
    880 else:
    881     # Binary mode
    882     handle = open(handle, ioargs.mode)

FileNotFoundError: [Errno 2] No such file or directory: 'fraudTrain.csv'
```

```
In [3]: train_df = pd.read_csv('Credit card/fraudTrain.csv.zip')
test_df = pd.read_csv('Credit card/fraudTest.csv.zip')
```

```
In [4]: train_df.head()
```

Out[4]:

| Unnamed: 0 | trans_date_trans_time | cc_num | merchant | category | amt | first | last | gender | street | ... | lat | long | city_pop | job | dob | trans_num |
|------------|-----------------------|---------------------|------------------|------------------------------------|---------------|--------|-----------|---------|--------|------------------------------|---------|-----------|----------|-----------------------------------|------------|----------------------------------|
| 0 | 0 | 2019-01-01 00:00:18 | 2703186189652095 | fraud_Rippin, Kub and Mann | misc_net | 4.97 | Jennifer | Banks | F | 561 Perry Cove | 36.0788 | -81.1781 | 3495 | Psychologist, counselling | 1988-03-09 | 0b242abb623afc578575680df30655b9 |
| 1 | 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 | 48.8878 | -118.2105 | 149 | Special educational needs teacher | 1978-06-21 | 1f76529f8574734946361c461b024d99 |
| 2 | 2 | 2019-01-01 00:00:51 | 38859492057661 | fraud_Lind-Buckridge | entertainment | 220.11 | Edward | Sanchez | M | 594 White Dale Suite 530 | 42.1808 | -112.2620 | 4154 | Nature conservation officer | 1962-01-19 | a1a22d70485983eac12b5b88dad1cf95 |
| 3 | 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 | 46.2306 | -112.1138 | 1939 | Patent attorney | 1967-01-12 | 6b849c168bdad6f867558c3793159a81 |
| 4 | 4 | 2019-01-01 00:03:06 | 375534208663984 | fraud_Keeling-Crist | misc_pos | 41.96 | Tyler | Garcia | M | 408 Bradley Rest | 38.4207 | -79.4629 | 99 | Dance movement psychotherapist | 1986-03-28 | a41d7549acf90789359a9a5346dc46 |

5 rows × 23 columns

```
In [5]: train_df.info()

<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

In [6]: test_df.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
```

```
In [7]: train_df.isnull().sum()
```

| Out[7]: | Unnamed: 0 | trans_date_trans_time | cc_num | merchant | category | amt | first | last | gender | street | city | state | zip | lat | long | city_pop | job | dob | trans_num | is_fraud | dtype |
|---------|------------|-----------------------|--------|----------|----------|-----|-------|------|--------|--------|------|-------|-----|-----|------|----------|-----|-----|-----------|----------|-------|
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | int64 |

```
In [8]: test_df.isnull().sum()
```

| Out[8]: | Unnamed: 0 | trans_date_trans_time | cc_num | merchant | category | amt | first | last | gender | street | city | state | zip | lat | long | city_pop | job | dob | trans_num | is_fraud | dtype |
|---------|------------|-----------------------|--------|----------|----------|-----|-------|------|--------|--------|------|-------|-----|-----|------|----------|-----|-----|-----------|----------|-------|
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | int64 |

```
In [9]: 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)
```

```
In [10]: print(train_df.shape)
print(test_df.shape)
```

(1296675, 14)
(555719, 14)

```
In [11]: 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))
```

```
In [12]: drop_columns = ['trans_date_trans_time', 'city', 'lat', 'long', 'job', 'dob', 'merch_lat', 'merch_long', 'state']
train_df.drop(columns=drop_columns, inplace=True)
test_df.drop(columns=drop_columns, inplace=True)
```

```
In [13]: # Convert categorical column gender into numerical
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)
```

```
In [14]: #One Hot Encoding of Category column
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)
```

```
In [15]: train_df.head()
```

| Out[15]: | amt | gender | city_pop | is_fraud | latitudinal_distance | longitudinal_distance | category_entertainment | category_food_dining | category_gas_transport | category_grocery_net | category_grocery_pos | category_health_fitness | category_home |
|----------|--------|--------|----------|----------|----------------------|-----------------------|------------------------|----------------------|------------------------|----------------------|----------------------|-------------------------|---------------|
| 0 | 4.97 | 0 | 3495 | 0 | 0.068 | 0.870 | False | False | False | False | False | False | False |
| 1 | 107.23 | 0 | 149 | 0 | 0.271 | 0.024 | False | False | False | False | True | False | False |
| 2 | 220.11 | 1 | 4154 | 0 | 0.970 | 0.108 | True | False | False | False | False | False | False |
| 3 | 45.00 | 1 | 1939 | 0 | 0.804 | 0.447 | False | False | True | False | False | False | False |
| 4 | 41.96 | 1 | 99 | 0 | 0.254 | 0.830 | False | False | False | False | False | False | False |

```
In [16]: test_df.head()
```

| Out[16]: | amt | gender | city_pop | is_fraud | latitudinal_distance | longitudinal_distance | category_entertainment | category_food_dining | category_gas_transport | category_grocery_net | category_grocery_pos | category_health_fitness | category_home |
|----------|-------|--------|----------|----------|----------------------|-----------------------|------------------------|----------------------|------------------------|----------------------|----------------------|-------------------------|---------------|
| 0 | 2.86 | 1 | 333497 | 0 | 0.020 | 0.265 | False | False | False | False | False | False | False |
| 1 | 29.84 | 0 | 302 | 0 | 0.870 | 0.476 | False | False | False | False | False | True | False |
| 2 | 61.28 | 0 | 34496 | 0 | 0.177 | 0.660 | False | False | False | False | False | False | False |
| 3 | 60.05 | 1 | 54767 | 0 | 0.243 | 0.064 | False | False | False | False | False | False | False |
| 4 | 3.19 | 1 | 1126 | 0 | 0.706 | 0.868 | False | False | False | False | False | False | False |

```
In [17]: 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']
```

```
In [18]: from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
```

```
In [19]: from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import confusion_matrix, accuracy_score
classifier= DecisionTreeClassifier(random_state=42)
classifier.fit(X_train, y_train)
y_pred1= classifier.predict(X_test)
cm = confusion_matrix(y_test, y_pred1)
print(cm)
print("accuracy:"+str(accuracy_score(y_test,y_pred1)))
```

[[552399 1265]
 [844 1301]]
accuracy:0.9962649165135617

```
In [20]: from sklearn.metrics import confusion_matrix, accuracy_score
from sklearn.linear_model import LogisticRegression
classifier = LogisticRegression(random_state = 0)
classifier.fit(X_train, y_train)
y_pred2=classifier.predict(X_test)
cm = confusion_matrix(y_test, y_pred2)
print(cm)
print("accuracy:"+str(accuracy_score(y_test,y_pred2)))
```

[[553251 323]
 [2145 8]]
accuracy:0.9955589656968872

```
In [21]: from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import confusion_matrix, accuracy_score
rf=RandomForestClassifier(n_estimators=10,criterion='entropy',random_state=0)
rf.fit(X_train,y_train)
y_pred3 = rf.predict(X_test)
cm = confusion_matrix(y_test,y_pred3)
print(cm)
print("accuracy:"+str(accuracy_score(y_test,y_pred3)))
```

[[553127 447]
 [904 1241]]
accuracy:0.9975689152251408

