

Credit Card Fraud Detection

Date

No.

→ ML with Python

→ Credit card data → Data Pre Processing (Imbalance) → Data Analysis

Evaluation ← Logistic Regression Model ← Train Test Splits

① Dependencies → Libraries

Swift + Enter = Run

Import the dependencies

1. import numpy as np → Making arrays
2. import pandas as pd → Pointers
3. from sklearn.model_selection import train-test-split → data frames
4. from sklearn.linear_model import LogisticRegression → Split data into training test data
5. from sklearn.metrics import accuracy_score

② Load dataset to Pandas dataframe

1. credit-card-data = pd.read_csv('...')
- ⇒ copy path of credit-data.csv → paste

③ First 5 rows of dataset
credit-card-data.head()

Label 0

Label 1

↑
Fraudulent

→ `credit-card-data.tail()`

After 2 days of transaction | Last 5 rows

(4) Dataset information
`credit-card-data.info()`

(5) Checking no. of missing values in each column
`credit-card-data.isnull().sum()`

(6) Distribution of legit transaction & fraudulent transaction

`credit-card-data['Class'].value_counts()`

↓
This dataset is highly unbalance.
0 - Normal transaction
1 - Fraudulent transaction

(7) Separating the data for analysis

```

ccd
legit = credit-card-data[Class == 0]
fraud = credit-card-data[Class == 1]
print(legit.shape)
print(fraud.shape)

```

(8) Statistical measures of the data

(1) `legit.Amount.describe()`

Count → no. of points

(2) `fraud.Amount.describe()`

Mean → More than legit.

- (9) Compare values for both transaction
 Credit-card-data.groupby('class').mean()
 0 - Normal
 1 - Fraudulent] - show
 wide diff. of both of mean.

(10) Undersampling

- ⇒ Build a sample dataset containing similar distribution of normal transaction and fraudulent transaction
 ⇒ No. of fraudulent transaction = 492

Ex: Legit (284315, 31) } we will take randomly
 Fraud (492, 31) } 492 data from legit
 and join to fraud.

Now 492 + 492
 (Fraud) (Legit)

Uniform dataset

legit-sample = legit.sample(n=492)

Concatenate 2 datasets/frames.

new-dataset = pd.concat([legit-sample,
 fraud], axis=0)

(Added 1 by 1)

axis=1 → Add column
 wise.

Add rows

MATRIKAS wise.

→ ① `new_dataset.head()`

first 5 rows.

② `new_dataset.tail()`

last 5 rows.

③ `new_dataset['class'].value_counts()`

1 492
0 492] show

④ `new_dataset.groupby('class').mean()`

↓

(nature of data not changed).

↓

Split the data into features & target

store
feature

`X = new_dataset.drop(columns='class', axis=1)`
`Y = new_dataset['class']`

↓

Don't have class column

`print(Y)`

||

split data into Training & Testing data

`X_train, X_test, Y_train, Y_test`
`= train_test_split`

`(X, Y, test_size=0.2,`
`stratify=Y, random_state=2)`

How to split 20% of your data to test data

Date

No.

```
print (X.shape , X_train.shape , X_test.  
shape)
```

(-, -) (-, -) (-, -)

Reference & Intelligence

Date

No.

Machine Learning \rightarrow Based on Predictions

Ex: for instance 100 Records

Types \rightarrow (1) Supervised (Predicted Models) \downarrow

(2) Unsupervised ID \rightarrow doesn't pay loan.

Why? \rightarrow Reasons \rightarrow (Bunch of reasons)

(1) Take Dataset

Customer C1 = 100K 20K (given) 80K (left)

find \rightarrow Avg. age, financial situations ... etc

then decide whether he/she is fraud or not.

(2) Make a characteristic to determine fraud

Ex: [Age: 30, salary = 50K, G = Male]

Bank issue credit card or not.

Not Paying
amount

Invalid
transacti-
ons.

Dataset

C1 C2 C3 C4 C5 ... CN

customer Amt. owed Age salary
ID to Bank

variables \rightarrow (dependent variable)
that you will predict

whether paying or
not on time
(✓)

(CN-1) \rightarrow independent variable

Dependent Variable \rightarrow (1) Categorical Problem

have fixed no. of values \rightarrow fraud, ontime, profit/loss \rightarrow ML Models

(1) Logistic, Decision tree, Neure.

Classification Models.

Random forest, KNN, SVM, LDA

MATRIKAS

② Numerical \rightarrow Sales, Profit, Salary? \rightarrow ML Models
Regression Models

\rightarrow Linear Regression, Decision Tree, RF \rightarrow Random Forest, KNN, SVM

Random Forest

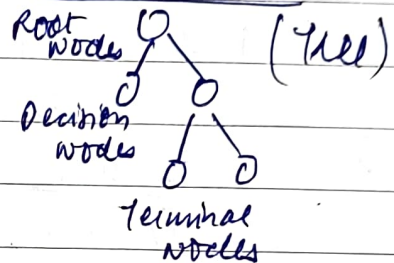
Tree-Based Model
(Flow Chart)

Ways to optimize tree

Ex: 5 trees - test all trees

and then want the prediction from each tree.

Decision Tree



Initially 1 tree prediction
then Forest of trees

GOOD MODEL CHARACTERISTICS

- ① Accuracy
- ② Stability

Dataset

Time	V_1	V_2	Amount	Class
	0 \rightarrow Legit
	1 \rightarrow Fraud