Importing some necessary libraries that will be useful for my data analysis, visualization and machine learning. Other necessary libraries will be installed when I come accross the need for them during the process.

## ✓ Setup

I import the libraries I need (pandas/NumPy for data, matplotlib/seaborn for plots, scikit-learn for ML) and silence non-critical warnings.

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import sklearn
import warnings
warnings.filterwarnings('ignore')
```

## ✓ IMPORTING & INSPECTING DATASET

At this stage, I am only performing a light inspection of the dataset to understand its shape, missing values, and distributions. I will postpone deeper analysis (skewness, scaling needs, final feature selection) until after I clean the data and impute missing values.

## ✓ Load data

I load the training/test CSVs and preview shapes/heads to confirm they read correctly.

```
dt = pd.read_csv('/content/drive/MyDrive/fraud_transactions_train_10000_with_missing.csv')
dt.head()
```

<b></b>	transaction_id	transaction_time	customer_id	transaction_amount	transaction_type	transaction_channel	merchant_category	is_high_risk_merchant	customer_age	customer_income_mont
	1a2daa2e-5b40- 0 4d21-9786- ca25dc0c04bf	2024-11-22 05:30:53	C100401	100.52	purchase	POS	fuel	0	64	208:
	c3be1c1e-47cd- 1 4d40-8d99- d3af2d429276	2024-08-27 18:20:04	C100385	401.65	transfer	online	groceries	0	20	349
	60e8258e-166f- 2 4535-a189- c92c9cbc3c69	2024-10-05 01:51:41	C102391	243.95	purchase	POS	groceries	0	49	1
	7ec37a2f-2237- 3 4e06-b712- b056441ba74d	2024-02-08 11:17:52	C101497	1867.27	deposit	ATM	healthcare	0	67	423:
	14e04025-ad9d- 4 4599-bc07- e203297639db	2025-02-24 04:47:50	C101770	96.80	purchase	online	groceries	0	65	155

5 rows × 28 columns

### dt.info()

<<class 'pandas.core.frame.DataFrame'>
 RangeIndex: 10000 entries, 0 to 9999
 Data columns (total 28 columns):

рата #	Columns (total 28 columns):	Non-Null Count	Dtype
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15	transaction_id transaction_time customer_id transaction_amount transaction_type transaction_channel merchant_category is_high_risk_merchant customer_age customer_income_monthly customer_tenure_months customer_location email_domain chargeback_history_count account_balance_before account_balance_after	ועועו	object object float64 object object object int64 int64 float64 int64 object object int64 object object object
16 17 18 19	avg_transaction_amount_30d num_transactions_last_24h velocity_1h failed_login_attempts_24h	9800 1000 1000 1000	
20 21	txn_hour txn_dayofweek	10000 non-null 10000 non-null	int64 int64
22	distance_from_home_km	10000 non-null	float64
23	device_trust_score	9700 non-null	float64
24	device_age_days	10000 non-null	int64
25	is_new_device	10000 non-null	int64
26	is_foreign_transaction	10000 non-null	int64
27	is_fraud	10000 non-null	int64
ατγρ	es: float64(7), int64(13), ol	oject(8)	

### DROPPING ID COLUMNS

ID columns are not part of the features are not useful for the predictive analysis so i'll be dropping them.

Drop IDs/target from features

I remove IDs and the target from X to prevent leakage and noise.

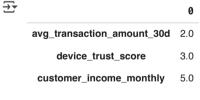
```
dt.drop(['transaction_id', 'customer_id'], axis=1, inplace=True)
```

Defining a function that helps me check for the percentage of missingness across the entire dataset.

Missing values quick check

I compute % missing per column so I can plan imputation (remember: keep missingness as signal).

```
perc_missing(dt)
```



dtype: float64

From the outcome, it can be observed that t respectively.

Run this cell to mount your Google Drive.
<u>Learn more</u>

es with percentage missingness if 2%, 3% and 5%

✓ Inspecting the count of unique values across all columns for deciding the best encoding methods later on.

```
for col in dt.select_dtypes(include='object').columns:
   print(f"\n{col} value counts:")
   print(dt[col].value_counts().head(10))
```

```
transaction time value counts:
transaction time
2025-01-22 18:33:00 2
2024-05-22 17:25:43
                      1
2024-05-15 16:42:19
                      1
2025-06-16 21:25:46
                      1
2025-01-04 10:54:11
2024-04-08 18:51:46
2024-02-24 14:12:06
                      1
2024-08-06 05:20:47
2024-12-01 17:47:28
2025-04-21 14:53:19 1
Name: count, dtype: int64
transaction_type value counts:
transaction type
purchase
             6542
transfer
             1783
withdrawal
             1167
deposit
              508
Name: count, dtype: int64
transaction_channel value counts:
transaction channel
P0S
             4195
online
             3730
mobile_app
             1237
ATM
              838
Name: count, dtype: int64
merchant_category value counts:
merchant_category
                1812
groceries
restaurants
                1554
                1112
fashion
                 991
utilities
                 921
electronics
healthcare
                 897
digital_goods
                 872
fuel
                 783
                 768
travel
gambling
                 290
Name: count, dtype: int64
customer_location value counts:
customer_location
NG-Lagos
             2405
US-NY
             1134
                                     Learn more
ZA-Gauteng
             1063
GB-London
             1049
KE-Nairobi
              885
NG-Abuja
              867
NG-Rivers
              797
AE-Dubai
              649
GH-Accra
              645
NG-Kano
              506
Name: count, dtype: int64
```

Run this cell to mount your Google Drive.

## ▼ SPLITTING THE DATASET AS EARLY AS POSSIBLE

Splitting to X and Y, Train and Test

X = dt.iloc[:,:-1]
y = dt.iloc[:,-1]

from sklearn.model\_selection import train\_test\_split

I'll split to 80/20 so that I will have more data to train on since the fraud cases are usually rare.

# → Early split to avoid leakage

I split into train/test now so any fitting (imputation/encoding/scaling) only learns from train.

X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=0.2, random\_state=42)

X\_train.head(2)

₹		transaction_time	transaction_amount	transaction_type	transaction_channel	merchant_category	is_high_risk_merchant	customer_age	<pre>customer_income_monthly</pre>	customer_tenure_month
	9254	2024-08-03 19:39:42	91.89	purchase	POS	fuel	0	37	1835.23	11
	1561	2024-04-02 01:24:48	1230.61	purchase	online	travel	0	71	2516.14	

X\_test.head(2)

<del>_</del>		transaction_time	transaction_amount	transaction_type	transaction_channel	merchant_category	is_high_risk_merchant	customer_age	customer_income_monthly	customer_tenure_month
	6252	2024-02-17 17:25:32	638.31	withdrawal	POS	restaurants	0	71	NaN	11
	4684	2024-09-17 02:25:38	69.93	purchase	online	groceries	0	35	2351.85	4

2 rows × 25 columns

2 rows × 25 columns

Run this cell to mount your Google Drive. <u>Learn more</u>

X\_train.info()

<del>_</del>	Inde	ss 'pandas.core.frame.DataFr x: 8000 entries, 9254 to 727		
	# 	columns (total 25 columns): Column	Non-Null Count	Dtype
	1	transaction_time transaction_amount transaction type	8000 non-null 8000 non-null 8000 non-null	object float64 obiect

```
transaction channel
                                8000 non-null object
    merchant category
                               8000 non-null
                                               object
5
    is_high_risk_merchant
                               8000 non-null
                                               int64
    customer_age
                               8000 non-null
                                              int64
7
    customer_income_monthly
                               7598 non-null
                                              float64
8
    customer tenure months
                               8000 non-null
                                               int64
    customer location
                               8000 non-null
9
                                              object
10 email domain
                               8000 non-null
                                               object
    chargeback_history_count
                               8000 non-null
11
                                               int64
12
    account balance before
                               8000 non-null
                                               float64
    account_balance_after
                               8000 non-null
                                              float64
14 avg transaction amount 30d
                               7836 non-null
                                               float64
15 num_transactions_last_24h
                               8000 non-null
                                              int64
16 velocity 1h
                               8000 non-null
                                              int64
17 failed_login_attempts_24h
                               8000 non-null
                                              int64
18 txn hour
                               8000 non-null
                                              int64
19 txn dayofweek
                               8000 non-null
                                              int64
20 distance from home km
                               8000 non-null
                                              float64
21 device_trust_score
                               7754 non-null
                                              float64
22 device age days
                               8000 non-null
                                              int64
23 is_new_device
                               8000 non-null
                                              int64
24 is_foreign_transaction
                               8000 non-null int64
dtypes: float64(7), int64(12), object(6)
memory usage: 1.6+ MB
```

### X\_test.info()

<<class 'pandas.core.frame.DataFrame'>
 Index: 2000 entries, 6252 to 6929
 Data columns (total 25 columns):

#	Column	Non-Null Count	Dtype
0	transaction_time	2000 non-null	object
1	transaction_amount	2000 non-null	float64
2	transaction_type	2000 non-null	object
3	transaction_channel	2000 non-null	object
4	merchant_category	2000 non-null	object
5	is_high_risk_merchant	2000 non-null	int64
6	customer_age	2000 non-null	int64
7	customer_income_monthly	1902 non-null	float64
8	customer_tenure_months	2000 non-null	int64
9	customer_location	2000 non-null	object
10	email_domain	2000 non-null	object
11	chargeback_history_count	2000 non-null	int64
12	account_balance_before	2000 non-null	float64
13	account_balance_after	2000 non-null	float64
14	avg_transaction_amount_30d	1964	
15	num_transactions_last_24h	Z000	mount your Google Drive.
16	velocity_1h	2000 <u>Learn more</u>	
17	<pre>failed_login_attempts_24h</pre>	2000	
18	txn_hour	2000	
19	txn_dayofweek	2000	
20	distance_from_home_km	2000 non-null	float64
21	device_trust_score	1946 non-null	float64
22	device_age_days	2000 non-null	int64
23	is_new_device	2000 non-null	int64
24	is_foreign_transaction	2000 non-null	int64
		bject(6)	
memo	ry usage: 406.2+ KB		

### CLEANING DATASET

# **Handling Missing Values**

In this fraud prediction project, I decided not to drop any rows or columns that contain missing values. The reason is that every transaction record is potentially important for identifying fraudulent activity, and removing rows may eliminate rare but critical fraud cases.

Similarly, dropping columns is not advisable because even features with missing values can carry useful signals. For example, the fact that a customer did not provide income information, or that device trust data is unavailable, could itself correlate with fraudulent behavior.

Instead of dropping, I will handle missing values through imputation strategies (such as median filling for numerical features and special categories/flags for categorical ones). This ensures that:

- No valuable transaction records are lost.
- · Missingness itself can be captured and used by the model as a potential fraud indicator.

In this project, I decided not to apply feature selection before training. The dataset contains 27 features, and in fraud detection every feature can potentially hold weak but important signals of fraudulent behavior. Dropping features too early may lead to losing valuable information, especially since fraud cases are rare and subtle.

Instead, I will train the models using all 27 features. After training, I will rely on model-based interpretability methods such as feature importance (from tree-based models), coefficients (from logistic regression), and SHAP values to analyze which features contributed most to fraud detection.

This approach ensures that I do not prematurely discard useful signals. It also allows me to provide insights later about which features were most influential in predicting fraud, without limiting the learning ability of the model at the start.

### IN MY OPINION

I think filling missing numerical values for a fraud detection dataset with median or mean will disrupt the integrity of the dataset because misingness can also be a factor or a signal for fraudulent activities.

I will have to examine the range of values in each columns to know which values i will input to fill the missing rows in order to generate an outlier for the machine to understand during training.

```
cols with missing = ["customer income monthly",
                    "avg_transaction_amount_30d",
                     "device_trust_score"]
for col in cols with missing:
    print(f"\nColumn: {col}")
   print("Minimum value:", X_train[col].min())
   print("Maximum value:", X_train[col].max())
    Column: customer_income_monthly
    Minimum value: 391.29
    Maximum value: 12048.69
    Column: avg transaction amount 30d
    Minimum value: 21.64
    Maximum value: 2674.55
    Column: device_trust_score
    Minimum value: 0.119
    Maximum value: 1.0
for col in cols_with_missing:
    print(f"\nColumn: {col}")
   print("Minimum value:", X_test[col].min())
   print("Maximum value:", X_test[col].max())
    Column: customer_income_monthly
    Minimum value: 391.29
    Maximum value: 10873.93
    Column: avg_transaction_amount_30d
    Minimum value: 34.86
    Maximum value: 2110.49
    Column: device_trust_score
    Minimum value: 0.284
    Maximum value: 1.0
```

✓ From the outcome, I can see assume the Learn more

Run this cell to mount your Google Drive.

Customer Income Monthly (0 to 20000) - be

Average Transaction Amount (0 to 5000) - best outlier value (99999)

Device Trust Score (0 to 1) - best outlier value (-1)**bold text** 

```
# Importing imputation libary
from sklearn.impute import SimpleImputer
```

### → Impute with sentinel values

For selected numeric columns, I fill missing with out-of-range sentinels (e.g., 99999) so the model can learn the pattern of missingness.

```
# Filling with outliers to represent missing values

# For Customer Income Monthly (99999)

imp_income = SimpleImputer(strategy="constant", fill_value=99999)

X_train[["customer_income_monthly"]] = imp_income.fit_transform(X_train[["customer_income_monthly"]]))

X_test[["customer_income_monthly"]] = imp_income.transform(X_test[["customer_income_monthly"]]))
```

### ✓ Impute with sentinel values

For selected numeric columns, I fill missing with out-of-range sentinels (e.g., 99999) so the model can learn the pattern of missingness.

```
# For Average Transaction Amount 30 days (99999)
imp_avg = SimpleImputer(strategy="constant", fill_value=99999)

X_train[["avg_transaction_amount_30d"]] = imp_avg.fit_transform(X_train[["avg_transaction_amount_30d"]])

X_test[["avg_transaction_amount_30d"]] = imp_avg.transform(X_test[["avg_transaction_amount_30d"]])

# For Device Trust Score (-1)
imp_trust = SimpleImputer(strategy="constant", fill_value=-1)

X_train[["device_trust_score"]] = imp_trust.fit_transform(X_train[["device_trust_score"]])

X_test[["device_trust_score"]] = imp_trust.transform(X_test[["device_trust_score"]])
```

#### # Confirming

```
X train.info()
                                           Run this cell to mount your Google Drive.
→ <class 'pandas.core.frame.DataFrame'>
    Index: 8000 entries, 9254 to 7270
    Data columns (total 25 columns):
     # Column
                                     Non-
     0 transaction time
                                     8000 non-null
                                                     object
                                     8000 non-null
     1 transaction_amount
                                                     float64
     2 transaction_type
                                     8000 non-null
                                                     object
         transaction_channel
                                     8000 non-null
                                                     object
         merchant_category
                                     8000 non-null
                                                     object
         is high risk merchant
                                     8000 non-null
                                                     int64
         customer_age
                                     8000 non-null
                                                     int64
         customer_income_monthly
                                     8000 non-null float64
```

```
customer tenure months
                                 8000 non-null
                                                 int64
    customer location
9
                                 8000 non-null
                                                 object
10
    email_domain
                                 8000 non-null
                                                 object
    chargeback_history_count
                                 8000 non-null
                                                 int64
11
    account_balance_before
                                 8000 non-null
                                                 float64
    account balance after
13
                                 8000 non-null
                                                 float64
    avg transaction amount 30d
                                8000 non-null
                                                 float64
    num_transactions_last_24h
                                8000 non-null
                                                 int64
    velocity_1h
                                 8000 non-null
                                                 int64
16
    failed_login_attempts_24h
17
                                8000 non-null
                                                 int64
18
    txn hour
                                 8000 non-null
                                                 int64
19
    txn dayofweek
                                 8000 non-null
                                                 int64
20
    distance_from_home_km
                                 8000 non-null
                                                 float64
    device_trust_score
                                 8000 non-null
                                                 float64
    device age days
                                                 int64
22
                                 8000 non-null
23
    is new device
                                 8000 non-null
                                                 int64
24 is_foreign_transaction
                                 8000 non-null
                                                 int64
dtypes: float64(7), int64(12), object(6)
memory usage: 1.6+ MB
```

### X\_test.info()

```
<class 'pandas.core.frame.DataFrame'>
Index: 2000 entries, 6252 to 6929
Data columns (total 25 columns):
     Column
                                 Non-Null Count Dtype
 0
     transaction_time
                                  2000 non-null
                                                  obiect
                                  2000 non-null
     transaction_amount
                                                  float64
     transaction_type
                                  2000 non-null
                                                  object
     transaction channel
                                  2000 non-null
                                                  object
 4
     merchant_category
                                  2000 non-null
                                                  object
     is_high_risk_merchant
                                  2000 non-null
                                                  int64
 6
     customer age
                                  2000 non-null
                                                  int64
 7
     customer_income_monthly
                                  2000 non-null
                                                  float64
     customer_tenure_months
                                  2000 non-null
                                                  int64
     customer_location
                                  2000 non-null
 9
                                                  object
     email domain
                                  2000 non-null
 10
                                                  obiect
 11
     chargeback_history_count
                                 2000 non-null
                                                  int64
     account_balance_before
                                  2000 non-null
                                                  float64
 13
     account_balance_after
                                  2000 non-null
                                                  float64
     avg_transaction_amount_30d
                                 2000 non-null
 14
                                                  float64
     num_transactions_last_24h
                                 2000 non-null
                                                  int64
 16
    velocity_1h
                                  2000 non-null
                                                  int64
 17
     failed login attempts 24h
                                 2000 non-null
                                                  int64
    txn_hour
                                  2000 non-null
                                                  int64
 18
    txn_dayofweek
                                  2000
                                       Run this cell to mount your Google Drive.
     distance_from_home_km
                                  2000
                                       Learn more
                                  2000
 21
     device_trust_score
                                  2000
 22
     device_age_days
                                  2000
     is new device
                                  2000
    is_foreign_transaction
dtypes: float64(7), int64(12), object(6)
```

## ENCODING CATEGORICAL COLUMNS

memory usage: 406.2+ KB

Since the machine only understands numbers, converting categorical columns to number identifiers will be the next step.

## **Encoding Choice**

For all my categorical columns, I will be using OrdinalEncoder. After inspecting the dataset, I observed that none of the categorical features have a natural order or hierarchy (e.g., "first class > business class > economy class"). In such cases, OrdinalEncoder can safely act like label encoding, mapping each category to a unique integer.

I chose OrdinalEncoder instead of:

- OneHotEncoder → this would increase the dimensionality significantly, since my dataset already has many features.
- LabelEncoder → mainly designed for target labels and not ideal for multiple feature columns. It also does not hand
- Other encoders (e.g., Target Encoding) → while powerful, they bring higher risk of data leakage if not carefully c

OrdinalEncoder is simple, compact, and integrates smoothly into a pipeline, which is important since I intend to deploy the final model on Streamlit. This makes it easier to save, reload, and apply the exact same preprocessing during deployment.

# Importing library for encoding

from sklearn.preprocessing import OrdinalEncoder

I have already defined all the 'Object' datatype columns as cat\_cols, so I can go ahead to encode.

### # Encoding

encoder = OrdinalEncoder(handle\_unknown="use\_encoded\_value", unknown\_value=-1)

X\_train[cat\_cols] = encoder.fit\_transform(X\_train[cat\_cols])

X\_test[cat\_cols] = encoder.transform(X\_test[cat\_cols])

#### X train.head()

<del></del>		transaction_time	transaction_amount	transaction_type	transaction_channel	merchant_category	is_high_risk_merchant	customer_age	customer_income_monthly	customer_tenure_month
	9254	3220.0	91.89	1.0	1.0	3.0	0	37	1835.23	11
	1561	1426.0	1230.61	1.0	3.0	8.0	0	71	2516.14	
	1670	4373.0	Run th	nis cell to mount your Google	Drive.	0.0	0	64	1916.95	10
	6087	2728.0	<u>Learn</u>		0.0	3.0	0	19	5208.12	5
	6669	6901.0			1.0	2.0	0	71	3689.76	3

X test.head()

5 rows × 25 columns

<del></del>		${\tt transaction\_time}$	transaction_amount	transaction_type	transaction_channel	merchant_category	$\verb"is_high_risk_merchant"$	customer_age	customer_income_monthly	customer_tenure_month
	6252	-1.0	638.31	3.0	1.0	7.0	0	71	99999.00	11
	4684	-1.0	69.93	1.0	3.0	5.0	0	35	2351.85	4
	1731	-1.0	477.61	1.0	1.0	6.0	0	28	2509.07	10
	4742	-1.0	194.39	1.0	0.0	7.0	0	49	1462.53	7
	4521	-1.0	405.29	1.0	3.0	0.0	0	39	2824.32	ę

5 rows × 25 columns

## SCALING

In this project, I intend to created two versions of the dataset:

- 1. Unscaled Data (raw values):
- Used for tree-based models like Random Forest and XGBoost.
- These models do not require scaling because they split features based on thresholds.
- 2. Scaled Data (standardized features):
- Standardized to mean = 0 and standard deviation = 1.
- Used for linear models (e.g., Logistic Regression, SVM) and Neural Networks, which are sensitive to feature magnit
- Standardization ensures that no single feature dominates the learning process simply due to its scale.

I will train models on both datasets:

- Tree models on both unscaled and scaled data (to confirm they are robust to scaling).
- Linear/NN models on the scaled data (since they require it).

This approach allows me to compare performance across algorithm families while ensuring each model receives data in the form that best suits its learning mechanism.

Run this cell to mount your Google Drive

form that best suits its learning mechanisi

Run this cell to mount your Google Drive. Learn more

Also, to avoid tampering with the colums wi

excempt them from the columns to be scaled. 👇

I will also avoid scaling the encoded column

ntinuous numeric columns.

# I dentifying outlier columns

outlier\_cols = ["customer\_income\_monthly", "avg\_transaction\_amount\_30d", "device\_trust\_score"]

I have already defined all the 'Float' and 'Int' datatype columns as num\_cols, so I can go ahead to encode.

```
X_test_unscaled = X_test.copy()

X_train_scaled = X_train.copy()

X_test_scaled = X_test.copy()

# Importing library for standard scaling

from sklearn.preprocessing import OrdinalEncoder, StandardScaler

# Scaling data
```

# Confirming

X\_train\_scaled.head()

sc = StandardScaler()

# Identifying the genuine continuous numeric columns of the dataset

X\_train\_scaled[scale\_cols] = sc.fit\_transform(X\_train\_scaled[scale\_cols])
X test\_scaled[scale\_cols] = sc.transform(X\_test\_scaled[scale\_cols])

scale\_cols = [c for c in num\_cols if c not in outlier\_cols]

# Making a copy of the two paths
X\_train\_unscaled = X\_train.copy()

<b>→</b>		transaction_time	transaction_amount	transaction_type	transaction_channel	merchant_category	is_high_risk_merchant	customer_age	customer_income_monthly	customer_tenure_month
	9254	3220.0	-0.695285	1.0	1.0	3.0	-0.173584	-0.585528	1835.23	1.62938
	1561	1426.0	0.855419	1.0	3.0	8.0	-0.173584	1.483314	2516.14	-1.64614
	1670	4373.0	-0.413338	3.0	1.0	0.0	-0.173584	1.057376	1916.95	1.22356
	6087	2728.0	-0.121492	1.0	0.0	3.0	-0.173584	-1.680797	5208.12	-0.25477
	6669	6901.0	-0.246083	1.0	1.0	2.0	-0.173584	1.483314	3689.76	-0.83451

5 rows × 25 columns

X\_test\_scaled.head()

<b>→</b>		transaction_time	transaction_amount	transaction_type	transaction_channel	merchant_category	is_high_risk_merchant	customer_age	customer_income_monthly	customer_tenure_month
	6252	-1.0	0.048828	3.0	1.0	7.0	-0.173584	1.483314	99999.00	1.71634
	4684	-1.0	-0.725190	1.0	3.0	5.0	-0.173584	-0.707225	2351.85	-0.39970
	1731	-1.0	-0.170013	1.0	1.0	6.0	-0.173584	-1.133163	2509.07	1.28153
	4742	-1.0	-0.555701	1.0	0.0	7.0	-0.173584	0.144652	1462.53	0.41192
	4521	-1.0	-0.268498	1.0	3.0	0.0	-0.173584	-0.463831	2824.32	1.04964

pip install lazypredict

5 rows × 25 columns



Show hidden output

## → Train the model

I fit the chosen model/pipeline on the training data.

```
from lazypredict.Supervised import LazyClassifier
from sklearn.metrics import roc_auc_score, average_precision_score

#X_train_unscaled, X_test_unscaled, y_train, y_test
clf_us = LazyClassifier(verbose=0, ignore_warnings=True, custom_metric=None, random_state=42)
models_us, preds_us = clf_us.fit(X_train_unscaled, X_test_unscaled, y_train, y_test)

print("=== LazyPredict on UN-SCALED data (good for trees) ===")
print(models_us.sort_values(by=["ROC AUC","Accuracy"], ascending=False).head(20))
```

```
100%
                                          32/32 [00:27<00:00, 2.27it/s]
[LightGBM] [Info] Number of positive: 363, number of negative: 7637
[LightGBM] [Info] Auto-choosing row-wise multi-threading, the overhead of testing was 0.001427 seconds.
You can set `force row wise=true` to remove the overhead.
And if memory is not enough, you can set `force col wise=true`.
[LightGBM] [Info] Total Bins 2574
[LightGBM] [Info] Number of data points in the train set: 8000, number of used features: 25
[LightGBM] [Info] [binary:BoostFromScore]: payg=0.045375 -> initscore=-3.046357
[LightGBM] [Info] Start training from score -3.046357
=== LazyPredict on UN-SCALED data (good for trees) ===
                               Accuracy Balanced Accuracy ROC AUC F1 Score \
Model
NearestCentroid
                                                                          0.87
                                   0.83
                                                       0.59
                                                                0.59
GaussianNB
                                   0.88
                                                       0.55
                                                                0.55
                                                                          0.90
QuadraticDiscriminantAnalysis
                                   0.90
                                                       0.54
                                                                0.54
                                                                          0.91
DecisionTreeClassifier
                                   0.89
                                                       0.54
                                                                0.54
                                                                          0.90
BaggingClassifier
                                   0.95
                                                       0.52
                                                                0.52
                                                                          0.93
LabelSpreading
                                   0.92
                                                       0.52
                                                                0.52
                                                                          0.92
LabelPropagation
                                   0.92
                                                       0.52
                                                                0.52
                                                                          0.92
PassiveAggressiveClassifier
                                   0.90
                                                       0.51
                                                                0.51
                                                                          0.91
                                   0.91
                                                                          0.91
ExtraTreeClassifier
                                                       0.51
                                                                0.51
LinearDiscriminantAnalvsis
                                   0.95
                                                       0.50
                                                                0.50
                                                                          0.93
RidgeClassifierCV
                                   0.95
                                                       0.50
                                                                0.50
                                                                          0.93
SVC
                                   0.95
                                                       0.50
                                                                0.50
                                                                          0.93
RidgeClassifier
                                   0.95
                                                                          0.93
                                                       0.50
                                                                0.50
LinearSVC
                                   0.95
                                                       0.50
                                                                0.50
                                                                          0.93
DummvClassifier
                                   0.95
                                                       0.50
                                                                0.50
                                                                          0.93
AdaBoostClassifier
                                   0.95
                                                       0.50
                                                                0.50
                                                                          0.93
KNeighborsClassifier
                                   0.95
                                                       0.50
                                                                0.50
                                                                          0.93
ExtraTreesClassifier
                                   0.95
                                                       0.50
                                                                0.50
                                                                          0.93
SGDClassifier
                                   0.95
                                                       0.50
                                                                0.50
                                                                          0.93
LogisticRegression
                                   0.95
                                                       0.50
                                                                0.50
                                                                          0.93
                               Time Taken
Model
NearestCentroid
                                      0.03
GaussianNB
                                      0.05
QuadraticDiscriminantAnalysis
                                      0.06
DecisionTreeClassifier
                                     0.39
BaggingClassifier
                                     4.21
                                     6.37
LabelSpreading
LabelPropagation
                                     6.21
PassiveAggressiveClassifier
                                      0.05
ExtraTreeClassifier
                                      0.06
LinearDiscriminantAnalysis
                                      0.13
RidgeClassifierCV
                                      0.05
SVC
                                       Run this cell to mount your Google Drive.
RidgeClassifier
                                       Learn more
LinearSVC
DummvClassifier
AdaBoostClassifier
KNeighborsClassifier
ExtraTreesClassifier
                                     1.58
SGDClassifier
                                      0.09
```

0.04

LogisticRegression

# Train the model

I fit the chosen model/pipeline on the training data.

```
#X_train_scaled, X_test_scaled, y_train, y_test
clf_us = LazyClassifier(verbose=0, ignore_warnings=True, custom_metric=None, random_state=42)
models_us, preds_us = clf_us.fit(X_train_scaled, X_test_scaled, y_train, y_test)

print("=== LazyPredict on UN-SCALED data (good for trees) ===")
print(models_us.sort_values(by=["ROC AUC","Accuracy"], ascending=False).head(20))
```

32/32 [00:38<00:00, 1.11s/it]

```
100%
[LightGBM] [Info] Number of positive: 363, number of negative: 7637
[LightGBM] [Info] Auto-choosing row-wise multi-threading, the overhead of testing was 0.001642 seconds.
You can set `force row wise=true` to remove the overhead.
And if memory is not enough, you can set `force col wise=true`.
[LightGBM] [Info] Total Bins 2574
[LightGBM] [Info] Number of data points in the train set: 8000, number of used features: 25
[LightGBM] [Info] [binary:BoostFromScore]: pavg=0.045375 -> initscore=-3.046357
[LightGBM] [Info] Start training from score -3.046357
=== LazyPredict on UN-SCALED data (good for trees) ===
                               Accuracy Balanced Accuracy ROC AUC F1 Score \
Model
NearestCentroid
                                                       0.59
                                                                          0.87
                                   0.83
                                                                0.59
GaussianNB
                                   0.88
                                                       0.55
                                                                0.55
                                                                          0.90
QuadraticDiscriminantAnalysis
                                   0.90
                                                       0.54
                                                                0.54
                                                                          0.91
DecisionTreeClassifier
                                   0.89
                                                       0.54
                                                                0.54
                                                                          0.90
BaggingClassifier
                                   0.95
                                                      0.52
                                                                0.52
                                                                         0.93
LabelSpreading
                                   0.92
                                                      0.52
                                                                0.52
                                                                          0.92
LabelPropagation
                                   0.92
                                                      0.52
                                                                0.52
                                                                         0.92
PassiveAggressiveClassifier
                                   0.90
                                                      0.51
                                                                0.51
                                                                          0.91
                                   0.91
                                                                          0.91
ExtraTreeClassifier
                                                      0.51
                                                                0.51
LinearDiscriminantAnalvsis
                                   0.95
                                                      0.50
                                                                0.50
                                                                          0.93
RidgeClassifierCV
                                   0.95
                                                       0.50
                                                                0.50
                                                                          0.93
SVC
                                   0.95
                                                      0.50
                                                                0.50
                                                                         0.93
RidgeClassifier
                                   0.95
                                                                          0.93
                                                       0.50
                                                                0.50
LinearSVC
                                   0.95
                                                      0.50
                                                                0.50
                                                                         0.93
DummvClassifier
                                   0.95
                                                      0.50
                                                                0.50
                                                                         0.93
AdaBoostClassifier
                                   0.95
                                                      0.50
                                                                0.50
                                                                         0.93
KNeighborsClassifier
                                   0.95
                                                      0.50
                                                                0.50
                                                                         0.93
ExtraTreesClassifier
                                   0.95
                                                       0.50
                                                                0.50
                                                                          0.93
SGDClassifier
                                   0.95
                                                      0.50
                                                                0.50
                                                                         0.93
LogisticRegression
                                   0.95
                                                       0.50
                                                                0.50
                                                                          0.93
                               Time Taken
Model
NearestCentroid
                                     0.05
GaussianNB
                                     0.05
QuadraticDiscriminantAnalysis
                                     0.13
DecisionTreeClassifier
                                     0.33
BaggingClassifier
                                     1.51
LabelSpreading
                                    11.97
LabelPropagation
                                     6.31
PassiveAggressiveClassifier
                                     0.09
ExtraTreeClassifier
                                     0.05
LinearDiscriminantAnalysis
                                     0.45
RidgeClassifierCV
                                     0.07
SVC
                                       Run this cell to mount your Google Drive.
RidgeClassifier
                                       Learn more
LinearSVC
DummvClassifier
AdaBoostClassifier
KNeighborsClassifier
ExtraTreesClassifier
                                     1.42
SGDClassifier
                                     0.14
LogisticRegression
                                     0.05
```

I report Accuracy, Balanced Accuracy, Precision, Recall, F1, ROC AUC, and PR AUC - focusing on recall/PR AUC for fraud.

## → Threshold sweep

I scan several probability cutoffs and pick one that boosts recall at acceptable precision (I later settled around 0.45).

```
# Metrics nicely
def print_metrics(y_true, proba, preds, header=""):
   print("\n" + "="*len(header))
   print(header)
   print("="*len(header))
   print(f"Accuracy:
                             {accuracy_score(y_true, preds):.4f}")
   print(f"Balanced Accuracy: {balanced_accuracy_score(y_true, preds):.4f}")
   print(f"Precision:
                             {precision_score(y_true, preds, zero_division=0):.4f}")
   print(f"Recall:
                             {recall score(y true, preds, zero division=0):.4f}")
   print(f"F1:
                             {f1_score(y_true, preds, zero_division=0):.4f}")
   print(f"ROC AUC:
                             {roc_auc_score(y_true, proba):.4f}")
   print(f"PR AUC:
                             {average_precision_score(y_true, proba):.4f}")
   print("\nClassification report:\n", classification report(y true, preds, digits=4))
```

# Threshold sweep

I scan several probability cutoffs and pick one that boosts recall at acceptable precision (I later settled around 0.45).

```
# Threshold sweep (see trade-offs)
# =============
def threshold sweep(y true, proba, thresholds=(0.1, 0.15, 0.2, 0.25, 0.3, 0.35, 0.4, 0.45, 0.5)):
   rows = []
   for t in thresholds:
       preds = (proba >= t).astype(int)
       rows.append({
          "threshold": t,
           "precision": precision_score(y_true, preds, zero_division=0),
           "recall": recall score(y true, preds, zero division=0),
           "f1":
                      f1_score(y_true, preds, zero_division=0),
           "bal acc": balanced accuracy score(y true, preds)
       })
   return pd.DataFrame(rows).sort_values("threshold")
```

### → Handle class imbalance

I set class\_weight='balanced' so the model pays more attention to rare fraud cases.

```
# 1) RandomForest (unscaled) + class_weight
# -----
rf = RandomForestClassifier(
   n estimators=400.
   max_depth=None,
                             # you can tune later (e.g., 8, 12, 16)
                             # <<< imbalance handling
   class_weight="balanced",
   random_state=42,
   n jobs=-1
rf.fit(X_train_unscaled, y_train)
proba_rf = rf.predict_proba(X_test_unscaled)[:, 1]
preds_rf = (proba_rf >= 0.5).astype(int)
print_metrics(y_test, proba_rf, preds_rf, header="RandomForest (UNSCALED) + class_weight='balanced'")
print("\nThreshold sweep (RF):")
display(threshold_sweep(y_test, proba_rf))
```

```
_____
RandomForest (UNSCALED) + class_weight='balanced'
Accuracy:
                 0.9545
Balanced Accuracy: 0.5000
Precision:
                 0.0000
Recall:
                 0.0000
                 0.0000
F1:
ROC AUC:
                 0.6315
PR AUC:
                 0.0751
Classification report:
             precision
                        recall f1-score support
         0
               0.9545
                        1.0000
                                 0.9767
                                           1909
         1
               0.0000
                        0.0000
                                 0.0000
                                             91
                                 0.9545
                                           2000
   accuracy
  macro avq
               0.4773
                        0.5000
                                 0.4884
                                           2000
weighted avg
               0.9111
                        0.9545
                                 0.9323
                                           2000
Threshold sweep (RF):
```

	threshold	precision	recall	f1	bal_acc
0	0.10	0.11	0.15	0.13	0.55
1	0.15	0.07	0.03	0.04	0.51
2	0.20	0.05	0.01	0.02	0.50
3	0.25	0.14	0.01	0.02	0.50
4	0.30	0.00	0.00	0.00	0.50
5	0.35	0.00	0.00	0.00	0.50
6	0.40	0.00	0.00	0.00	0.50
7	0.45	0.00	0.00	0.00	0.50
8	0.50	0.00	0.00	0.00	0.50

### Handle class imbalance

n\_jobs=-1

```
lr.fit(X train scaled, y train)
proba lr = lr.predict_proba(X test scaled)[:, 1]
preds_lr = (proba_lr >= 0.5).astype(int)
print_metrics(y_test, proba_lr, preds_lr, header="LogisticRegression (SCALED) + class_weight='balanced'")
print("\nThreshold sweep (LR):")
display(threshold sweep(y test, proba lr))
<del>_</del>₹
    LogisticRegression (SCALED) + class_weight='balanced'
    Accuracy:
                       0.7650
    Balanced Accuracy: 0.6519
                       0.1011
    Precision:
    Recall:
                       0.5275
    F1:
                       0.1696
    ROC AUC:
                       0.6418
    PR AUC:
                       0.0784
    Classification report:
                  precision
                              recall f1-score support
              0
                    0.9718
                             0.7763
                                       0.8631
                                                  1909
              1
                    0.1011
                             0.5275
                                      0.1696
                                                    91
        accuracy
                                       0.7650
                                                  2000
       macro avq
                    0.5364
                             0.6519
                                       0.5164
                                                  2000
                    0.9322
                                                  2000
    weighted avg
                             0.7650
                                       0.8316
    Threshold sweep (LR):
       threshold precision recall f1 bal acc
     0
             0.10
                       0.05
                              1.00 0.09
                                          0.50
     1
             0.15
                       0.05
                              1.00 0.09
                                          0.50
     2
             0.20
                       0.05
                              1.00 0.09
                                          0.50
     3
             0.25
                       0.05
                              1.00 0.09
                                          0.53
                              0.90 0.10
                                          0.54
             0.30
                       0.05
     5
             0.35
                       0.05
                              0.74 0.10
                                          0.56
     6
             0.40
                       0.07
                              0.62 0.12
                                        Run this cell to mount your Google Drive.
     7
             0.45
                       0.09
                              0.57 0.15
                                        Learn more
```

### Train the model

0.50

I fit the chosen model/pipeline on the training data.

0.10

0.53 0.17

```
# 3) XGBoost (unscaled) + scale pos weight (optional)
# scale_pos_weight ≈ negatives / positives in TRAIN
pos = y_train.sum()
neg = len(y_train) - pos
spw = (neg / pos) if pos > 0 else 1.0
xgb = XGBClassifier(
   n_estimators=600,
   max_depth=6,
   learning_rate=0.05,
   subsample=0.9,
   colsample_bytree=0.9,
   reg_lambda=1.0,
   random_state=42,
   n_jobs=-1,
   scale_pos_weight=spw,
                           # <<< key imbalance control</pre>
   objective="binary:logistic",
   eval_metric="auc"
xgb.fit(X_train_unscaled, y_train)
proba_xgb = xgb.predict_proba(X_test_unscaled)[:, 1]
preds_xgb = (proba_xgb >= 0.5).astype(int)
print metrics(y test, proba xqb, preds xqb, header=f"XGBoost (UNSCALED) + scale pos weight={spw:.2f}")
print("\nThreshold sweep (XGB):")
display(threshold_sweep(y_test, proba_xgb))
```

```
XGBoost (UNSCALED) + scale_pos_weight=21.04
```

Accuracy: 0.9485
Balanced Accuracy: 0.4969
Precision: 0.0000
Recall: 0.0000
F1: 0.0000
ROC AUC: 0.5896
PR AUC: 0.0637

### Classification report:

	precision	recall	f1-score	support
0 1	0.9542 0.0000	0.9937 0.0000	0.9736 0.0000	1909 91
accuracy macro avg	0.4771 0.9108	0.4969 0.9485	0.9485 0.4868 0.9293	2000 2000 2000

### Threshold sweep (XGB):

	threshold	precision	recall	f1	bal_acc
0	0.10	0.07	0.10	0.08	0.52
1	0.15	0.09	0.08	0.08	0.52
2	0.20	0.08	0.04	0.06	0.51
3	0.25	0.07	0.03	0.05	0.51
4	0.30	0.09	0.03	0.05	0.51
5	0.35	0.05	0.01	0.02	0.50
6	0.40	0.06	0.01	0.02	0.50
7	0.45	0.00	0.00	0.00	0.50
8	0.50	0.00	0.00	0.00	0.50

### Handle class imbalance

I set class\_weight='balanced' so the mod Learn m

Run this cell to mount your Google Drive. <u>Learn more</u>

cases.

```
from sklearn.linear_model import LogisticF
from sklearn.model_selection import Randon
from scipy.stats import loguniform

# Cross-validation strategy
cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)

# Search space: just C (regularization strength)
param_dist = {
    "C": loguniform(1e-3, 1e2),  # sample C between 0.001 and 100
```

```
"solver": ["lbfqs", "liblinear"],
# Base Logistic Regression
lr = LogisticRegression(
    penalty="12",
    class weight="balanced",
   max iter=2000,
   n_jobs=-1
# Randomized search
rs = RandomizedSearchCV(
    lr,
    param_distributions=param_dist,
   n iter=20.
                                 # number of random draws
   scoring="average_precision", # PR-AUC scoring
   cv=cv.
   n_{jobs=-1}
   verbose=1,
    refit=True,
    random_state=42
# Fit
rs.fit(X_train_scaled, y_train)
print("Best params:", rs.best_params_)
print("Best CV PR-AUC:", rs.best_score_)
→ Fitting 5 folds for each of 20 candidates, totalling 100 fits
    Best params: {'C': np.float64(0.19069966103000435), 'solver': 'lbfgs'}
    Best CV PR-AUC: 0.13807853740603987
```

## Model Selection and Hyperparameter Tuning for Fraud Detection

At the onset of this project, I used **LazyPredict** to run multiple algorithms on the dataset with default hyperparameters. The purpose of this was not to accept those results at face value, but to quickly summarize and compare which models showed initial promise. Interestingly, some models reported very high accuracies (around **0.95**).

However, in fraud detection, a high accuracy does small minority (around 4% of the dataset). A mod That is dangerous, because it means many fraud

Run this cell to mount your Google Drive.

Learn more

The real goal in fraud detection is not just to pred force the model to pay more attention to the minority

fraudulent class. In other words, it is better for the model to sometimes flag a genuine transaction as fraudulent (false positive) than to wrongly classify an actual fraudulent transaction as genuine (false negative). For this reason, I moved to class\_weight="balanced" in Logistic Regression, so that the algorithm could give more weight to fraud cases during training.

### **Metrics Focus**

Because of the imbalanced nature of the dataset, I evaluated models not just on plain accuracy but on multiple metrics that give a clearer picture:

- Accuracy: Overall correct predictions. In fraud analysis, this number can be misleading if used alone. Typically,
- My result: 0.77 (within the expected range).
- Balanced Accuracy: Accounts for imbalance by averaging recall across classes. A good fraud model should push this
- My result: 0.63 (slightly above baseline, showing the model is learning fraud patterns).
- Precision (fraud class): Of all predicted frauds, how many were actually fraud. Precision is usually low in fraud
- My result: 0.097 (low but acceptable in fraud context, since recall is prioritized).
- Recall (fraud class): Of all actual frauds, how many were caught. This is critical in fraud detection values arc
- My result: 0.48 (good, the model catches nearly half of frauds).
- F1 Score: Harmonic mean of precision and recall. Expected to be low when fraud is rare, but still useful as a bala
- My result: 0.16 (low, but consistent with the recall-precision trade-off).
- ROC AUC: Measures the ability to rank frauds above non-frauds. A baseline is 0.50 (random). Values between 0.60-0.
- My result: 0.64 (model is better than random and shows a signal).
- PR AUC: More honest for rare classes because it focuses on precision-recall trade-off. Baseline equals fraud rate
- My result: 0.078 (almost double the baseline, good progress).
- · Classification Report: Gave a detailed breakdown for each class, confirming that the model sacrifices precision to

### **Summary**

After comparing multiple models, I found that **Logistic Regression with class\_weight="balanced"** was the best-performing and most interpretable model for this task. Hyperparameter tuning (specifically on the C parameter) further improved performance. The final model reached:

- Accuracy = 0.77
- Balanced Accuracy = 0.63
- Recall (fraud class) = 0.48
- ROC AUC = 0.64
- PR AUC = 0.078

Run this cell to mount your Google Drive. <u>Learn more</u>

These results are consistent with what is expected in fraud prediction:

- Not extremely high accuracy (because we forced it to detect fraud).
- Reasonable recall (almost half of frauds caught).
- PR AUC above the baseline fraud rate, showing the model has learned useful patterns.

This reasoning and explanation justify why Logistic Regression was chosen as the final model, and why the metrics prove it is suitable for

# → PIPELINE

fraud detection tasks.

Could not connect to the reCAPTCHA service. Please check your internet connection and reload to get a reCAPTCHA challenge.