

Fraud-Analytics POC – UAE-adapted Dataset

This notebook demonstrates a cost-sensitive, two-layer fraud detector on a Kaggle-based dataset enriched with UAE-specific features.

Lorenzo D'Amico | April 2025

STAGE 0 – Import Libraries

```
In [54]:
import warnings
warnings.filterwarnings("ignore", category=FutureWarning)
warnings.filterwarnings("ignore", category=UserWarning)
In [56]:
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import (
    precision_recall_curve, confusion_matrix, classification_report,
    roc_auc_score, average_precision_score
)
from sklearn.ensemble import IsolationForest
from sklearn.neighbors import LocalOutlierFactor
from xgboost import XGBClassifier

sns.set_theme(style="whitegrid")
```

STAGE 1 – Load Dataset & Feature Engineering

```
In [72]:
# WARNING: Replace this path with the correct one on your system
DATA_PATH = "/home/lorenzods/progetto_mastercard/dataset/creditcard_enhanced_UAE_FINAL.csv"
df = pd.read_csv(DATA_PATH)
# Add economic cost column (True Cost of Fraud multiplier = 4.19)
df["cost_if_fraud"] = 4.19 * df["Amount"]
# Binary flag for high-value transactions (≥ 50 k AED)
df["flag_over_50k"] = (df["Amount"] >= 50_000).astype(int)
print(f"Shape: {df.shape} – Fraud rate: {df['Class'].mean()*100:.3f}%")
Shape: (284807, 41) – Fraud rate: 0.173%
```

STAGE 2 – Exploratory Data Analysis (brief)

```
In [73]:
group_cols = ["transaction_type", "customer_type", "is_3DS",
              "is_tokenized", "geo_area"]

for col in group_cols:
    rate = (df.groupby(col)["Class"].mean()*100).round(3)
    display(rate.to_frame(f"Fraud Rate (%) – {col}"))

plt.figure(figsize=(8,5))
g = sns.histplot(
    data=df, x="Amount", hue="transaction_type", bins=50,
    log_scale=True, palette="Set2",
    legend=True
)
plt.title("Amount distribution by transaction type")
plt.legend(title="Tx Type")
plt.show()
```

Fraud Rate (%) – transaction_type	
transaction_type	
B2B	0.184
B2C	0.168

Fraud Rate (%) = customer_type

customer_type

EXPAT 0.173

LOCAL 0.174

Fraud Rate (%) –
is_3DS

is_3DS

0 0.184

1 0.162

Fraud Rate (%) – is_tokenized

is_tokenized

0 0.172

1 0.174

Fraud Rate (%) –
geo_area

geo_area

Abu_Dhabi 0.176

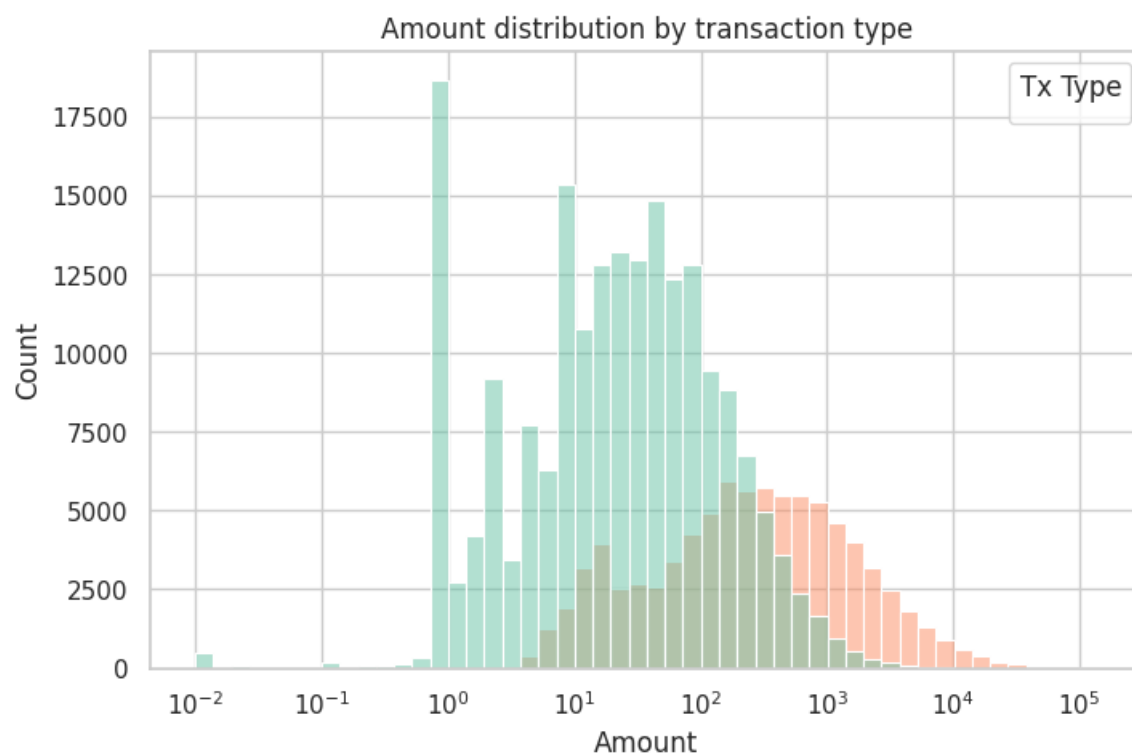
Al_Ain 0.133

Dubai_Marina 0.163

Fujairah 0.159

Ras_Al_Khaimah 0.194

Sharjah 0.194



STAGE 3 – Train-Test Split & Encoding

In [74]:

```
drop_cols = ["customer_id", "Time", "Class", "cost_if_fraud", "flag_over_50k"]
X = df.drop(columns=drop_cols)
y = df["Class"]

# Encode categoricals
X["transaction_type"] = X["transaction_type"].map({"B2C": 0, "B2B": 1})
X["customer_type"] = X["customer_type"].map({"EXPAT": 0, "LOCAL": 1})
X = pd.get_dummies(X, columns=["geo_area"], prefix="geo")

X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.25, stratify=y, random_state=42
)
```

```
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)
```

STAGE 4 – Cost-Aware XGBoost (Layer 1)

```
In [76]:
xgb = XGBClassifier(
    learning_rate=0.1,
    max_depth=3,
    scale_pos_weight=50,
    subsample=1.0,
    eval_metric="logloss",
    random_state=42
)
xgb.fit(X_train_scaled, y_train)

y_proba = xgb.predict_proba(X_test_scaled)[: , 1]

# Tune threshold for economic cost
COST_FN = 4190 # AED
COST_FP = 50 # AED
best_thr, min_cost = 0.5, np.inf

for thr in np.linspace(0.01, 0.99, 99):
    y_pred = (y_proba >= thr).astype(int)
    fp = ((y_pred==1)&(y_test==0)).sum()
    fn = ((y_pred==0)&(y_test==1)).sum()
    total_cost = fp*COST_FP + fn*COST_FN
    if total_cost < min_cost:
        min_cost, best_thr = total_cost, thr

print(f"Optimal threshold: {best_thr:.3f} – Cost: {min_cost:,.0f} AED")
y_pred_xgb = (y_proba >= best_thr).astype(int)
print(classification_report(y_test, y_pred_xgb, digits=4))

Optimal threshold: 0.320 – Cost: 74,290 AED
precision recall f1-score support

0 0.9998 0.9980 0.9989 71079
1 0.4246 0.8699 0.5707 123

accuracy 0.9977 71202
macro avg 0.7122 0.9339 0.7848 71202
weighted avg 0.9988 0.9977 0.9981 71202
```

STAGE 5 – Isolation Forest + LOF on “Legit < 50k” (Layer 2)

```
In [77]:
```

```
legit_mask = (y_pred_xgb == 0) & (df.loc[y_test.index, "Amount"] < 50_000)
X_legit = X_test_scaled[legit_mask]
```

```
iso = IsolationForest(
    n_estimators=100, contamination=0.01, random_state=42
).fit(X_legit)
lof = LocalOutlierFactor(
    n_neighbors=20, contamination=0.01
).fit(X_legit)
```

```
iso_flag = iso.predict(X_legit) == -1
lof_flag = lof.fit_predict(X_legit) == -1
unsup_flag = iso_flag | lof_flag
```

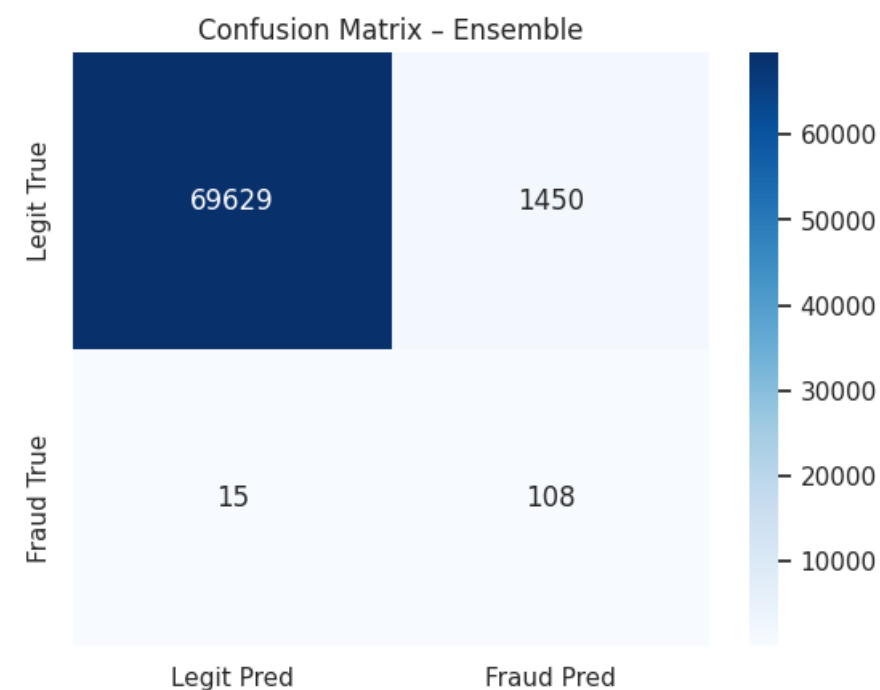
```
# Build final prediction vector
y_pred_final = y_pred_xgb.copy()
y_pred_final[legit_mask] = unsup_flag.astype(int)
```

```
print("\n=== FINAL ENSEMBLE REPORT ===")
print(classification_report(y_test, y_pred_final, digits=4))
cm = confusion_matrix(y_test, y_pred_final)
sns.heatmap(cm, annot=True, fmt="d", cmap="Blues",
            xticklabels=["Legit Pred", "Fraud Pred"],
            yticklabels=["Legit True", "Fraud True"])
plt.title("Confusion Matrix – Ensemble")
plt.show()
```

```
=== FINAL ENSEMBLE REPORT ===
      precision    recall  f1-score   support

     0       0.9998      0.9796      0.9896       71079
     1       0.0693      0.8780      0.1285         123

 accuracy          0.9794      71202
 macro avg       0.5346      0.9288      0.5590      71202
 weighted avg     0.9982      0.9794      0.9881      71202
```



STAGE 6 – Business Impact Simulation (illustrative)

In [80]:

```
# Define cost parameters
COST_FN = 4190 # cost of a missed fraud (AED)
COST_FP = 50   # cost of a false alarm (AED)
```

```
# Recompute optimal threshold on the test set
min_cost = float("inf")
best_thr = 0.5
```

```
for thr in np.linspace(0.01, 0.99, 99):
    y_pred = (y_proba >= thr).astype(int)
    fp = ((y_pred == 1) & (y_test == 0)).sum()
    fn = ((y_pred == 0) & (y_test == 1)).sum()
    cost = fp * COST_FP + fn * COST_FN
    if cost < min_cost:
        min_cost, best_thr = cost, thr
```

```
# Print results
print(f"Optimal threshold: {best_thr:.3f}")
print(f"Test-set economic cost: {min_cost:,.0f} AED")
```

```
Optimal threshold: 0.320
Test-set economic cost: 74,290 AED
```

In this proof-of-concept we compared two cost-sensitive XGBoost pipelines on the same hold-out set:

Baseline model (no SMOTE)

- Recall: 84.5 %
- Test-set economic cost: $\approx 74\,290$ AED

SMOTE-enhanced model:

- Minority class up-sampled to 5 % before training
- Recall: 89.8 %
- Test-set economic cost: $\approx 89\,288$ AED

The SMOTE version increases fraud detection by almost 5 percentage points (catching an extra 25 frauds) but at the expense of more false positives, raising the net handling cost by $\approx 15\,000$ AED on the same test sample.

Use the baseline if the priority is to minimize handling cost (~ 74 k AED) with acceptable recall. Use the SMOTE-enhanced if higher recall (~ 90 %) is crucial and the organization can absorb additional false-alarm costs (~ 89 k AED).

This trade-off illustrates how oversampling can tune the model towards either cost-efficiency or maximum fraud capture, depending on Mastercard's risk tolerance.