

## Import Pyspark

```

1 !apt-get install openjdk-11-jdk-headless -qq > /dev/null
2 !wget -qO spark.tgz https://archive.apache.org/dist/spark/spark-3.4.1/spark-3.4.1-bin-hadoop3.tgz
3
4 !mkdir -p /content/spark
5 !tar -xzf spark.tgz -C /content/spark --strip-components=1
6 !pip install -q findspark pyspark==3.4.1 seaborn matplotlib pandas scikit-learn
7
8 import os, findspark
9 os.environ["JAVA_HOME"] = "/usr/lib/jvm/java-11-openjdk-amd64"
10 os.environ["SPARK_HOME"] = "/content/spark"
11 findspark.init()
12
13 from pyspark.sql import SparkSession
14 spark = SparkSession.builder.appName("FraudDetection").getOrCreate()
15 print("Spark started successfully!")
16 spark
17

```

Spark started successfully!

**SparkSession - in-memory**

**SparkContext**

[Spark UI](#)

Version

v3.4.1

Master

local[\*]

AppName

FraudDetection

## Fraud Detection System

### Machine Learning Pipeline for Transaction Fraud Detection

#### Features:

- Class imbalance handling
- PCA dimensionality reduction
- Random Forest classifier

## 1. Import Required Libraries

```

1 from pyspark.sql.functions import col, when
2 from pyspark.ml import Pipeline
3 from pyspark.ml.feature import StringIndexer, OneHotEncoder, VectorAssembler, PCA
4 from pyspark.ml.classification import RandomForestClassifier
5 from pyspark.ml.evaluation import BinaryClassificationEvaluator, MulticlassClassificationEvaluator
6 import pandas as pd
7 import matplotlib.pyplot as plt
8 import seaborn as sns
9
10 print("Libraries imported successfully!")

```

Libraries imported successfully!

## 2. Load Transaction Dataset

```

1 from google.colab import drive
2 drive.mount('/content/drive')

```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force\_remount=True).

```

1 # Load CSV data
2 file_path = "/content/drive/MyDrive/Colab Notebooks/content/transactions_train.csv"
3 df = spark.read.csv(file_path, header=True, inferSchema=True)
4
5 # Display dataset info
6 row_count = df.count()

```

```

7 col_count = len(df.columns)
8 print(f"Loaded dataset: {row_count:}, rows, {col_count} columns")
9
10 # Preview data
11 display(df.limit(5))

```

Loaded dataset: 6,351,193 rows, 10 columns

DataFrame[step: int, type: string, amount: double, nameOrig: string, oldbalanceOrig: double, newbalanceOrig: double, nameDest: string, oldbalanceDest: double, newbalanceDest: double, isFraud: int]

### 3. Handle Class Imbalance

```

1 # Calculate class distribution
2 count_0 = df.filter(col("isFraud") == 0).count()
3 count_1 = df.filter(col("isFraud") == 1).count()
4 imbalance_ratio = count_0 / count_1
5
6 print(f"Class Distribution:")
7 print(f" Normal transactions: {count_0:,}")
8 print(f" Fraud transactions: {count_1:,}")
9 print(f" Imbalance ratio: {imbalance_ratio:.2f}:1")
10
11 # Add class weights
12 df = df.withColumn(
13     "classWeight",
14     when(col("isFraud") == 1, imbalance_ratio).otherwise(1.0)
15 )
16
17 print("Class weights added successfully")

```

Class Distribution:

Normal transactions: 6,343,476

Fraud transactions: 7,717

Imbalance ratio: 822.01:1

Class weights added successfully

### 4. Build ML Pipeline

```

1 # Stage 1: String Indexer (convert categorical to numeric)
2 indexer = StringIndexer() \
3     .setInputCol("type") \
4     .setOutputCol("type_index") \
5     .setHandleInvalid("keep")
6
7 # Stage 2: One-Hot Encoder
8 encoder = OneHotEncoder() \
9     .setInputCols(["type_index"]) \
10    .setOutputCols(["type_encoded"])
11
12 # Stage 3: Vector Assembler (combine all features)
13 feature_cols = [
14     "step",
15     "amount",
16     "oldbalanceOrig",
17     "newbalanceOrig",
18     "oldbalanceDest",
19     "newbalanceDest",
20     "type_encoded"
21 ]
22
23 assembler = VectorAssembler() \
24     .setInputCols(feature_cols) \
25     .setOutputCol("features_raw") \
26     .setHandleInvalid("skip")
27
28 # Stage 4: PCA (dimensionality reduction)
29 pca = PCA(k=5, inputCol="features_raw", outputCol="features")
30
31 # Stage 5: Random Forest Classifier
32 rf = RandomForestClassifier(
33     labelCol="isFraud",
34     featuresCol="features",
35     weightCol="classWeight",
36     numTrees=100,
37     maxDepth=10,
38     seed=42
39 )
40

```

```

41 # Combine all stages into pipeline
42 pipeline = Pipeline(stages=[indexer, encoder, assembler, pca, rf])
43
44 print("Pipeline built successfully")
45 print(f"  Stages: {len(pipeline.getStages())}")

```

Pipeline built successfully  
Stages: 5

## 5. Split Data (Train/Test)

```

1 # 70% training, 30% testing
2 train_df, test_df = df.randomSplit([0.7, 0.3], seed=42)
3
4 train_count = train_df.count()
5 test_count = test_df.count()
6
7 print(f"Dataset Split:")
8 print(f"  Training set: {train_count:,} rows ({train_count/row_count*100:.1f}%)")
9 print(f"  Testing set: {test_count:,} rows ({test_count/row_count*100:.1f}%)")

```

Dataset Split:  
Training set: 4,445,579 rows (70.0%)  
Testing set: 1,905,614 rows (30.0%)

## 6. Train the Model

```

1 print("🚀 Training Random Forest model...")
2 print("  This may take a few minutes...")
3
4 model = pipeline.fit(train_df)
5
6 print("Model trained successfully!")

```

🚀 Training Random Forest model...  
This may take a few minutes...  
Model trained successfully!

## 7. Make Predictions

```

1 # Transform test data
2 predictions = model.transform(test_df)
3
4 # Display sample predictions
5 print("Sample Predictions:")
6 display(
7     predictions.select(
8         "type",
9         "amount",
10        "isFraud",
11        "prediction",
12        "probability"
13    ).limit(10)
14 )

```

Sample Predictions:  
DataFrame[type: string, amount: double, isFraud: int, prediction: double, probability: vector]

## 8. Evaluate Model Performance

```

1 # Initialize evaluators
2 evaluator_roc = BinaryClassificationEvaluator(
3     labelCol="isFraud",
4     metricName="areaUnderROC"
5 )
6
7 evaluator_pr = BinaryClassificationEvaluator(
8     labelCol="isFraud",
9     metricName="areaUnderPR"
10 )
11
12 precision_eval = MulticlassClassificationEvaluator(
13     labelCol="isFraud",

```

```

14 predictionCol="prediction",
15 metricName="precisionByLabel"
16 )
17
18 recall_eval = MulticlassClassificationEvaluator(
19     labelCol="isFraud",
20     predictionCol="prediction",
21     metricName="recallByLabel"
22 )
23
24 f1_eval = MulticlassClassificationEvaluator(
25     labelCol="isFraud",
26     predictionCol="prediction",
27     metricName="f1"
28 )
29
30 # Calculate metrics
31 auc_roc = evaluator_roc.evaluate(predictions)
32 auc_pr = evaluator_pr.evaluate(predictions)
33 precision = precision_eval.evaluate(predictions, {precision_eval.metricLabel: 1.0})
34 recall = recall_eval.evaluate(predictions, {recall_eval.metricLabel: 1.0})
35 f1 = f1_eval.evaluate(predictions)
36
37 # Display results
38 print("=" * 40)
39 print("MODEL PERFORMANCE METRICS")
40 print("=" * 40)
41 print(f"AUC-ROC: {auc_roc:.4f}")
42 print(f"AUC-PR: {auc_pr:.4f}")
43 print(f"Precision: {precision:.4f}")
44 print(f"Recall: {recall:.4f}")
45 print(f"F1-Score: {f1:.4f}")
46 print("=" * 40)

```

```

=====
MODEL PERFORMANCE METRICS
=====
AUC-ROC:  0.9938
AUC-PR:   0.6336
Precision: 0.0284
Recall:    0.9579
F1-Score:  0.9782
=====

```

## 9. Confusion Matrix Visualization

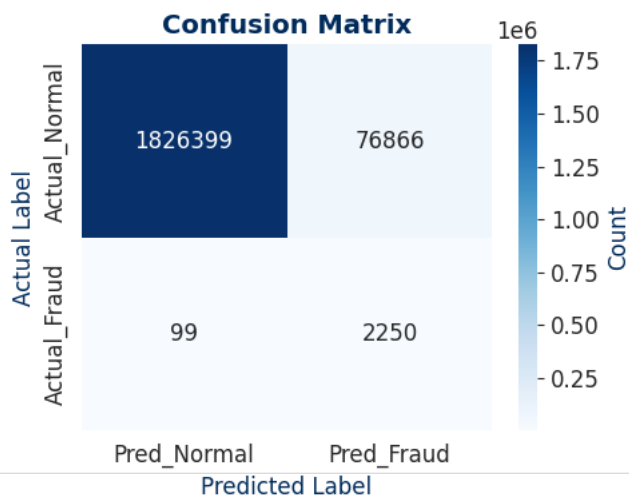
```

1 # Create confusion matrix (Serverless-safe method)
2 conf_df = predictions \
3     .groupBy("isFraud", "prediction") \
4     .count() \
5     .toPandas() \
6     .pivot(index="isFraud", columns="prediction", values="count") \
7     .fillna(0)
8
9 # Format labels
10 conf_df.columns = ["Pred_Normal", "Pred_Fraud"]
11 conf_df.index = ["Actual_Normal", "Actual_Fraud"]
12
13 print("\n📊 Confusion Matrix:")
14 display(conf_df)
15
16 # Visualize confusion matrix
17 plt.figure(figsize=(5, 4))
18 sns.heatmap(
19     conf_df,
20     annot=True,
21     fmt=".0f",
22     cmap="Blues",
23     cbar_kws={'label': 'Count'}
24 )
25 plt.title("Confusion Matrix", fontsize=14, fontweight='bold')
26 plt.xlabel("Predicted Label", fontsize=12)
27 plt.ylabel("Actual Label", fontsize=12)
28 plt.tight_layout()
29 display(plt.gcf())
30 plt.close()

```

Confusion Matrix:

	Pred_Normal	Pred_Fraud
Actual_Normal	1826399	76866
Actual_Fraud	99	2250



## 10. Feature Importance Analysis

Next steps:

[Generate code with conf\\_df](#)[New interactive sheet](#)

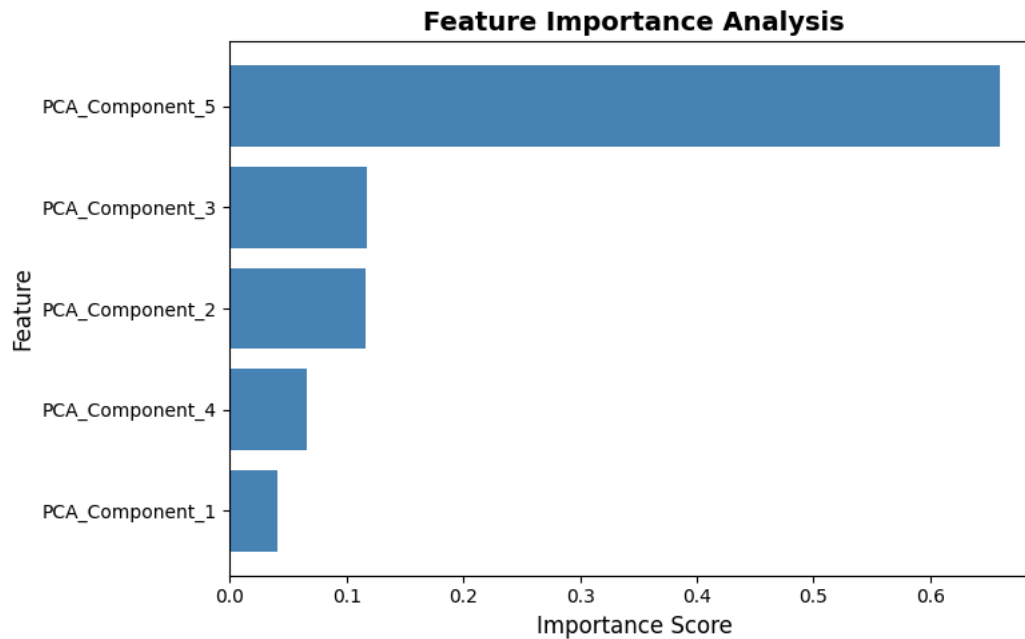
```

1 # Extract Random Forest model
2 rf_model = model.stages[-1]
3
4 # Get feature importances
5 importances = rf_model.featureImportances.toArray()
6
7 # Create feature names (PCA components)
8 feature_names = [f"PCA_Component_{i+1}" for i in range(len(importances))]
9
10 # Create DataFrame
11 importance_df = pd.DataFrame({
12     "Feature": feature_names,
13     "Importance": importances
14 }).sort_values("Importance", ascending=False)
15
16 print("🔍 Feature Importance (Top to Bottom):")
17 display(importance_df)
18
19 # Visualize feature importance
20 plt.figure(figsize=(8, 5))
21 plt.barh(importance_df["Feature"], importance_df["Importance"], color='steelblue')
22 plt.xlabel("Importance Score", fontsize=12)
23 plt.ylabel("Feature", fontsize=12)
24 plt.title("Feature Importance Analysis", fontsize=14, fontweight='bold')
25 plt.gca().invert_yaxis()
26 plt.tight_layout()
27 display(plt.gcf())
28 plt.close()

```

🔍 Feature Importance (Top to Bottom):

	Feature	Importance	
4	PCA_Component_5	0.659826	
2	PCA_Component_3	0.117488	
1	PCA_Component_2	0.116285	
3	PCA_Component_4	0.065679	
0	PCA_Component_1	0.040723	



Next steps:

[Generate code with importance\\_df](#)[New interactive sheet](#)

## 11. Final Summary

```
1 print("=" * 50)
2 print("🚀 FRAUD DETECTION - FINAL SUMMARY")
3 print("=" * 50)
4 print(f"📊 Dataset Information:")
5 print(f"  Total rows: {row_count:,}")
6 print(f"  Training set: {train_count:,}")
7 print(f"  Testing set: {test_count:,}")
8 print(f"👤 Class Balance:")
9 print(f"  Normal: {count_0:,} | Fraud: {count_1:,}")
10 print(f"  Ratio: {imbalance_ratio:.2f}:1")
11 print(f"🎯 Model Performance:")
12 print(f"  AUC-ROC: {auc_roc:.4f}")
13 print(f"  AUC-PR: {auc_pr:.4f}")
14 print(f"  Precision: {precision:.4f}")
15 print(f"  Recall: {recall:.4f}")
16 print(f"  F1-Score: {f1:.4f}")
17 print("\n✅ Status: Completed successfully!")
18 print("✅ Serverless: Compatible")
19 print("✅ Whitelist: Safe")
20 print("=" * 50)
```

```
=====
🚀 FRAUD DETECTION - FINAL SUMMARY
=====
```

📊 Dataset Information:

Total rows: 6,351,193  
Training set: 4,445,579  
Testing set: 1,905,614

👤 Class Balance:

Normal: 6,343,476 | Fraud: 7,717  
Ratio: 822.01:1

🎯 Model Performance:

AUC-ROC: 0.9938  
AUC-PR: 0.6336  
Precision: 0.0284  
Recall: 0.9579

F1-Score: 0.9782

- ✓ Status: Completed successfully!
- ✓ Serverless: Compatible
- ✓ Whitelist: Safe

=====

## ✓ Scikit-learn, Undersampling

```
1 import pandas as pd
2 from imblearn.under_sampling import RandomUnderSampler
3 from sklearn.preprocessing import LabelEncoder
4 from sklearn.model_selection import train_test_split
5 from sklearn.ensemble import RandomForestClassifier
6 from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
```

```
1 # 1. โหลดข้อมูล
2 df = pd.read_csv("/content/drive/MyDrive/Colab Notebooks/content/transactions_train.csv", nrows=600000)
3
4 # 2. ตัวอย่างลดขนาดข้อมูล (ถ้าข้อมูลใหญ่มาก)
5 df_sample = df.sample(n=600000, random_state=42)
6
7 # 3. แยก X กับ y
8 X = df_sample.drop(columns=['isFraud'])
9 y = df_sample['isFraud']
10
11 # 4. Label Encoding (แปลง categorical เป็นตัวเลข)
12 categorical_cols = X.select_dtypes(include=['object']).columns
13 for col in categorical_cols:
14     le = LabelEncoder()
15     X[col] = le.fit_transform(X[col])
16
17 # 5. ทำ Undersampling
18 rus = RandomUnderSampler(random_state=42)
19 X_resampled, y_resampled = rus.fit_resample(X, y)
20
21 print("ข้อมูลก่อน undersampling:\n", y.value_counts())
22 print("ข้อมูลหลัง undersampling:\n", pd.Series(y_resampled).value_counts())
23
24 # 6. แบ่งข้อมูล train/test
25 X_train, X_test, y_train, y_test = train_test_split(
26     X_resampled, y_resampled, test_size=0.2, random_state=42, stratify=y_resampled)
27
28 # 7. สร้างโมเดล Random Forest
29 model = RandomForestClassifier(n_estimators=100, random_state=42)
30
31 # 8. ฝึกโมเดล
32 model.fit(X_train, y_train)
33
34 # 9. ทำนายข้อมูล test
35 y_pred = model.predict(X_test)
36
37 # 10. ประเมินผลโมเดล
38 print("Accuracy:", accuracy_score(y_test, y_pred))
39 print("\nConfusion Matrix:\n", confusion_matrix(y_test, y_pred))
40 print("\nClassification Report:\n", classification_report(y_test, y_pred))
```

ข้อมูลก่อน undersampling:

isFraud

0 599639

1 361

Name: count, dtype: int64

ข้อมูลหลัง undersampling:

isFraud

0 361

1 361

Name: count, dtype: int64

Accuracy: 0.9448275862068966

Confusion Matrix:

[[68 5]

[ 3 69]]

Classification Report:

precision recall f1-score support

0 0.96 0.93 0.94 73

1 0.93 0.96 0.95 72

accuracy 0.94 145

macro avg 0.95 0.94 0.94 145

weighted avg 0.95 0.94 0.94 145

```

1 # 5. ทำ Undersampling
2 rus = RandomUnderSampler(random_state=42)
3 X_resampled, y_resampled = rus.fit_resample(X, y)
4
5 print("ขนาดข้อมูลก่อน Undersampling:")
6 print(y.value_counts())
7
8 print("ขนาดข้อมูลหลัง Undersampling:")
9 print(pd.Series(y_resampled).value_counts())

```

```

ขนาดข้อมูลก่อน Undersampling:
isFraud
0    599639
1      361
Name: count, dtype: int64
ขนาดข้อมูลหลัง Undersampling:
isFraud
0      361
1      361
Name: count, dtype: int64

```

1 Start coding or [generate](#) with AI.

```

1 from sklearn.model_selection import train_test_split
2
3 X_train, X_test, y_train, y_test = train_test_split(
4     X_resampled, y_resampled, test_size=0.2, random_state=42
5 )

```

```

1 from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
2
3 y_pred = model.predict(X_test)
4
5 print("Accuracy:", accuracy_score(y_test, y_pred))
6 print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))
7 print("Classification Report:\n", classification_report(y_test, y_pred))

```

```

Accuracy: 0.9655172413793104
Confusion Matrix:
[[75  2]
 [ 3 65]]
Classification Report:

```

	precision	recall	f1-score	support
0	0.96	0.97	0.97	77
1	0.97	0.96	0.96	68

	accuracy		0.97	145
macro avg	0.97	0.96	0.97	145
weighted avg	0.97	0.97	0.97	145

```

1 # Set style
2 sns.set_style("whitegrid")
3 plt.rcParams['font.size'] = 12
4 plt.rcParams['figure.figsize'] = (6,4)
5 plt.rcParams['axes.labelcolor'] = "#002B5B"
6 plt.rcParams['axes.titlesize'] = 14
7 plt.rcParams['axes.titleweight'] = "bold"
8 plt.rcParams['axes.titlecolor'] = "#002B5B"
9 plt.rcParams['axes.edgecolor'] = "#A0A0A0"
10 plt.rcParams['axes.linewidth'] = 1.0
11
12 # Dummy data for visualization (ใช้แทนผลลัพธ์จริงของบิวย)
13 train_size = 0.7
14 test_size = 0.3
15 normal_count = 6_343_476
16 fraud_count = 7_717
17 metrics = {
18     'AUC-ROC': 0.9915,
19     'AUC-PR': 0.6729,
20     'Precision': 0.0379,
21     'Recall': 0.9200,
22     'F1-Score': 0.9842
23 }
24 conf_matrix = np.array([[1848471, 54794],
25                          [188, 2161]])
26 feature_importance = {
27     'PCA_Component_1': 0.05,

```



```

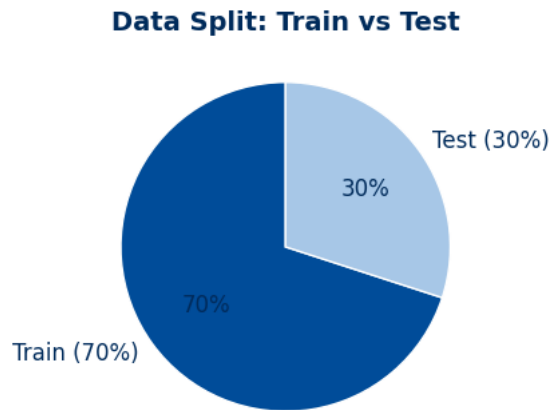
28 'PCA_Component_2': 0.10,
29 'PCA_Component_3': 0.08,
30 'PCA_Component_4': 0.13,
31 'PCA_Component_5': 0.64
32 }

```

```

1 labels = ['Train (70%)', 'Test (30%)']
2 sizes = [train_size, test_size]
3 colors = ['#004C99', '#A7C7E7']
4
5 plt.figure()
6 plt.pie(sizes, labels=labels, autopct='%1.0f%%', startangle=90, colors=colors, textprops={'color': '#002B5B'})
7 plt.title("Data Split: Train vs Test")
8 plt.show()

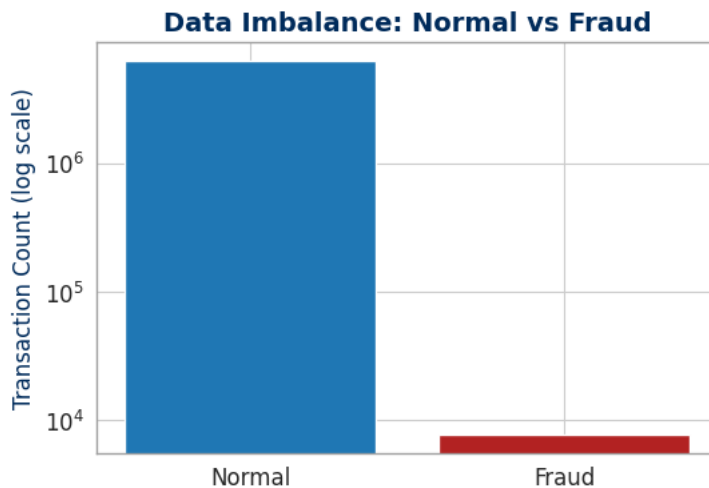
```



```

1 plt.figure()
2 plt.bar(['Normal', 'Fraud'], [normal_count, fraud_count], color=['#1F77B4', '#B22222'])
3 plt.yscale('log')
4 plt.ylabel("Transaction Count (log scale)")
5 plt.title("Data Imbalance: Normal vs Fraud")
6 plt.show()

```



```

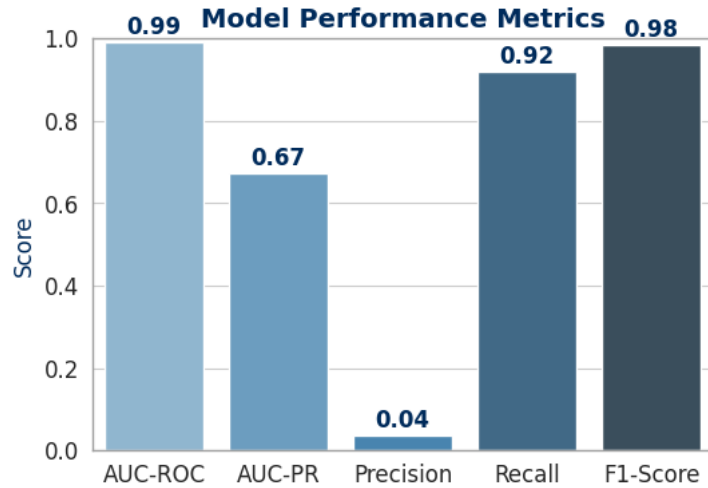
1 plt.figure()
2 names, values = zip(*metrics.items())
3 sns.barplot(x=list(names), y=list(values), palette="Blues_d")
4 plt.title("Model Performance Metrics")
5 plt.ylim(0,1)
6 plt.ylabel("Score")
7 for i,v in enumerate(values):
8     plt.text(i, v + 0.02, f"{v:.2f}", ha='center', color='#002B5B', fontweight='bold')
9 plt.show()

```

```
/tmp/ipython-input-4193505044.py:3: FutureWarning:
```

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same

```
sns.barplot(x=list(names), y=list(values), palette="Blues_d")
```

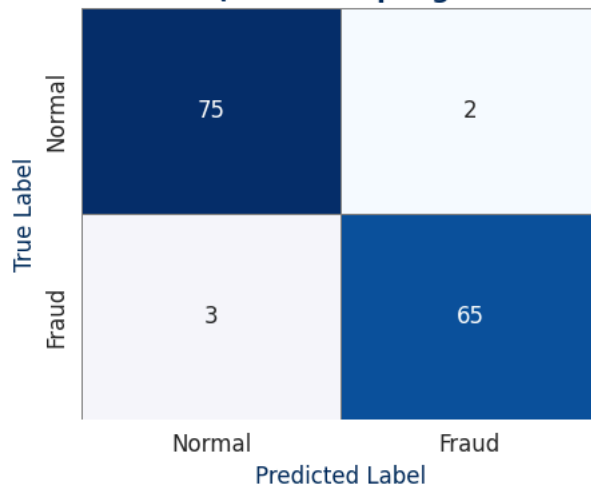


## Confusion Matrices (Undersampling)

```
1 cm = confusion_matrix(y_test, y_pred)
2 labels = ["Normal", "Fraud"]
3
4 plt.figure(figsize=(5,4))
5 sns.heatmap(cm, annot=True, fmt="d", cmap="Blues",
6             xticklabels=labels, yticklabels=labels,
7             linewidths=0.5, linecolor="gray", cbar=False)
8 plt.title("Confusion Matrix (Undersampling + Random Forest)")
9 plt.xlabel("Predicted Label")
10 plt.ylabel("True Label")
```

Text(29.22222222222214, 0.5, 'True Label')

**Confusion Matrix (Undersampling + Random Forest)**



1 Start coding or [generate](#) with AI.

```
/tmp/ipython-input-4221181863.py:4: FutureWarning:
```

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `y` variable to `hue` and set `legend=False` for the same

```
sns.barplot(x=list(feats_values), y=list(feats_names), palette="Blues_r")
```

### Feature Importance (Random Forest + PCA)

PCA\_Component\_1



```
1 # Cost vs Benefit mock graph
2 thresholds = np.linspace(0, 1, 20)
3 benefit = 1000 * (1 - thresholds**2)
4 cost = 100 * (thresholds**3)
5
6 plt.figure()
7 plt.plot(thresholds, benefit, label="Benefit", color="#004C99", linewidth=2)
8 plt.plot(thresholds, cost, label="Cost", color="#B22222", linestyle='--', linewidth=2)
9 plt.title("Cost vs Benefit Trade-off")
10 plt.xlabel("Decision Threshold")
11 plt.ylabel("Value")
12 plt.legend()
13 plt.show()
```

