

**AAI500 – Probability & Statistics for AI**  
**(Applied Artificial Intelligence)**

**University of San Diego**

**Credit Card Transactions Fraud Detection**

**Submitted to: Instructor Haisav Chokshi**

**Submitted by:**

**Group 5**

Ved Prakash Dwivedi  
Bharath TS  
Manu Malla

## Contents

1. Introduction.....	1
<b>2. Data Cleaning / Preparation.....</b>	<b>2</b>
<b>2.1 Handling Missing Values and Data Integrity Checks.....</b>	<b>2</b>
<b>2.2 Feature Distribution and Skewness.....</b>	<b>2</b>
<b>2.3 Class Imbalance Identification and Visualization.....</b>	<b>5</b>
<b>2.4 Resampling Strategy.....</b>	<b>6</b>
<b>2.5 Train-Test Split with Stratification.....</b>	<b>6</b>
<b>2.6 Final Feature Scaling.....</b>	<b>7</b>
<b>Summary of Preprocessing Strategy.....</b>	<b>7</b>
3. Exploratory Data Analysis (EDA).....	8
3.1 Objective of EDA.....	8
3.2 Univariate Analysis.....	8
3.3 Feature Skewness Observed.....	9
3.4 Bivariate Analysis.....	11
3.5 Scatterplot: Time vs Amount (Colored by Class).....	12
3.6 Correlation Matrix with Target (Class).....	13
3.7 Pairplot (Amount, Hour, Class).....	15
4. Model Selection.....	15
4.1 Evaluation Metric: Prioritizing Recall.....	15
4.2 Models Evaluated.....	16
4.3 Cross-Validation and Stratified Split.....	16
4.4 Logistic Regression (Baseline).....	16
4.5 Random Forest: Hyperparameter Tuning with Hyperopt.....	18
4.6 XGBoost: Bayesian Optimization.....	19
4.7 Final Model Selection.....	21
5. Model Analysis.....	21
5.1 Logistic Regression.....	21
5.2 Random Forest.....	22
5.3 XGBoost.....	23
6. Conclusion and Recommendations.....	25
6.1 Conclusion.....	25
6.2 Recommendations.....	25
Appendix (full notebook code and outputs).....	26
References.....	26

# Credit Card Transactions Fraud Detection

---

## Abstract:

This document describes a data-driven approach for Fraud credit card transaction detection. The study was performed on the transaction data records of the transaction made by European cardholders in September 2013 for 2 days which consists of PCA transformed 30 numerical features along with one target – Class.

As first step data cleaning was performed followed by a check for distribution and skewness. Data scaling was performed on observed skewness features. Also, Class imbalance was identified where resampling was performed before using the data for training. The EDA involves both univariate and bivariate analysis.

The study was performed using different models – Logistic Regression, Random Forest Classifier and XGBoost classifier where XGBoost model outperformed based on AUC score. The study proposes periodic re-training of the model with fresh data to capture the evolving fraud patterns.

## 1. Introduction

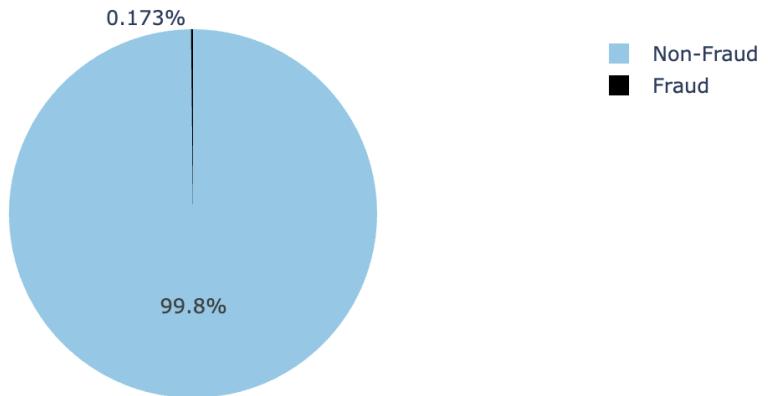
This project aims to build a robust machine learning pipeline for the detection of fraudulent credit card transactions using a real-world anonymized dataset available from Kaggle. The dataset represents credit card transactions made by European cardholders over a period of two days in September 2013.

Out of 284,807 transactions, only 492 are frauds—highlighting the **extreme class imbalance** in the data, where fraudulent transactions represent just **0.173%** of the total. Each transaction has 30 input features transformed using **Principal Component Analysis (PCA)**, preserving customer confidentiality. The features include 28 anonymized components (V1 to V28), transaction **Time**, and **Amount**.

Due to the complexity and imbalance of the data, traditional classifiers are often biased toward the majority class. To address this, advanced sampling methods and classification algorithms are explored. The models considered include **Logistic Regression**, **Random Forest**, **XGBoost** with evaluation based on **Precision**, **Recall**, **F1-score**, and **ROC-AUC**.

This report outlines a complete workflow including data exploration, preprocessing, model selection, evaluation, and interpretation of results. It demonstrates the critical steps needed to build a production-ready fraud detection system using machine learning.

## Fraud vs Non-Fraud transactions



## 2. Data Cleaning / Preparation

Effective data preprocessing is a foundational step in developing a reliable fraud detection model. This section outlines how the raw dataset was inspected, cleaned, and transformed to ensure high-quality inputs for machine learning.

---

### 2.1 Handling Missing Values and Data Integrity Checks

- The dataset was examined for missing, blank, and null values using:
  - `DataFrame.info()`
  - A custom summary statistics function.

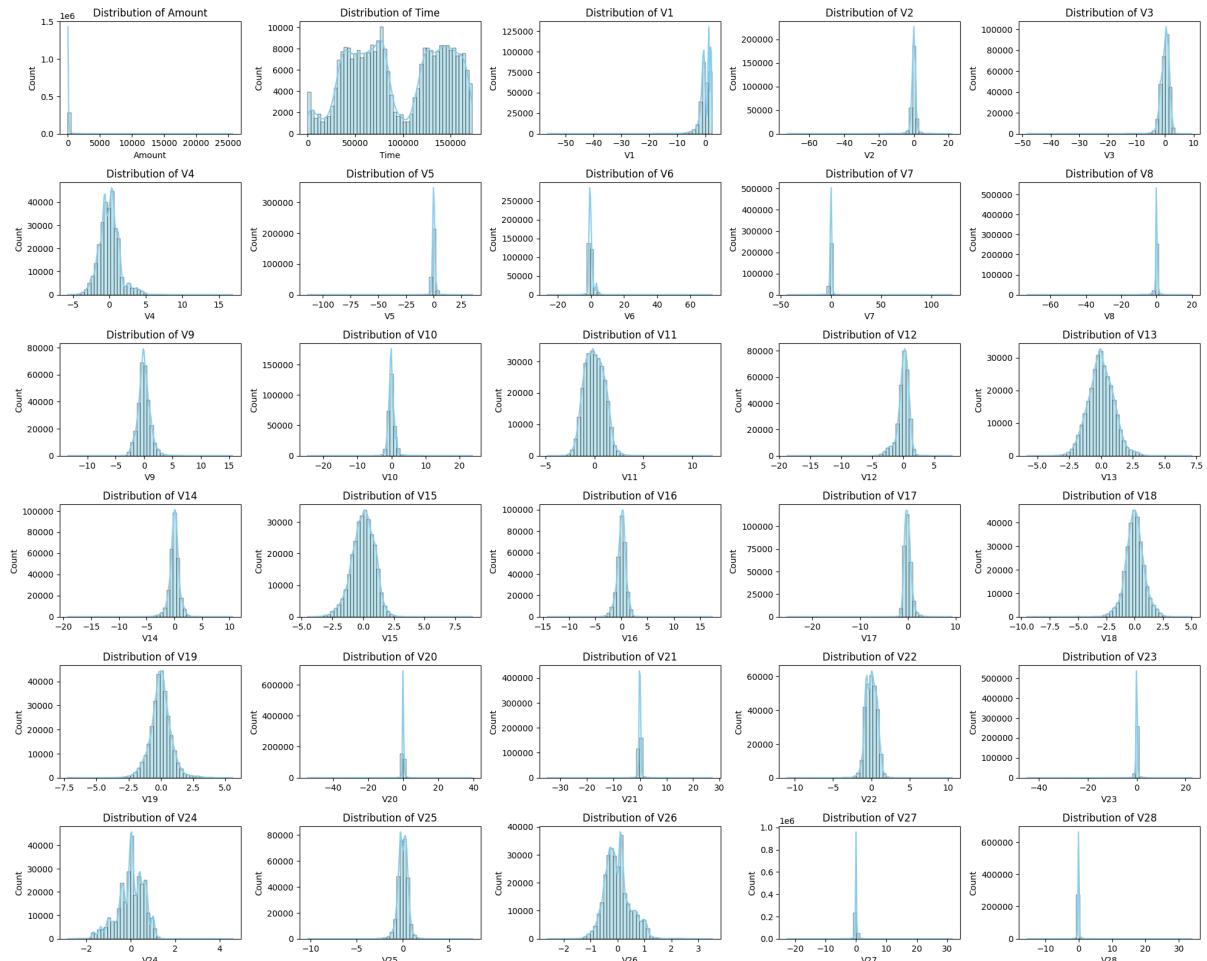
```
# Information of variables for data set:
def inspect_data(data):
    data_details = pd.DataFrame({"Data Type":data.dtypes,
                                "Count of Blank Values":data.apply(lambda x: x.isin([' ']).sum(),axis=0),
                                "Count of Missing Values":data.apply(lambda x: x.isnull().sum(),axis=0),
                                "% of Missing Values":data.apply(lambda x: round(x.isnull().sum()/len(x.unique())*100,2),axis=0),
                                "No of Unique Data":data.apply(lambda x: x.unique(),axis=0),
                                "Levels":data.apply(lambda x: str(x.unique()),axis=0)
                               })
    return data_details
inspect_data(df)
```

- **Result:**

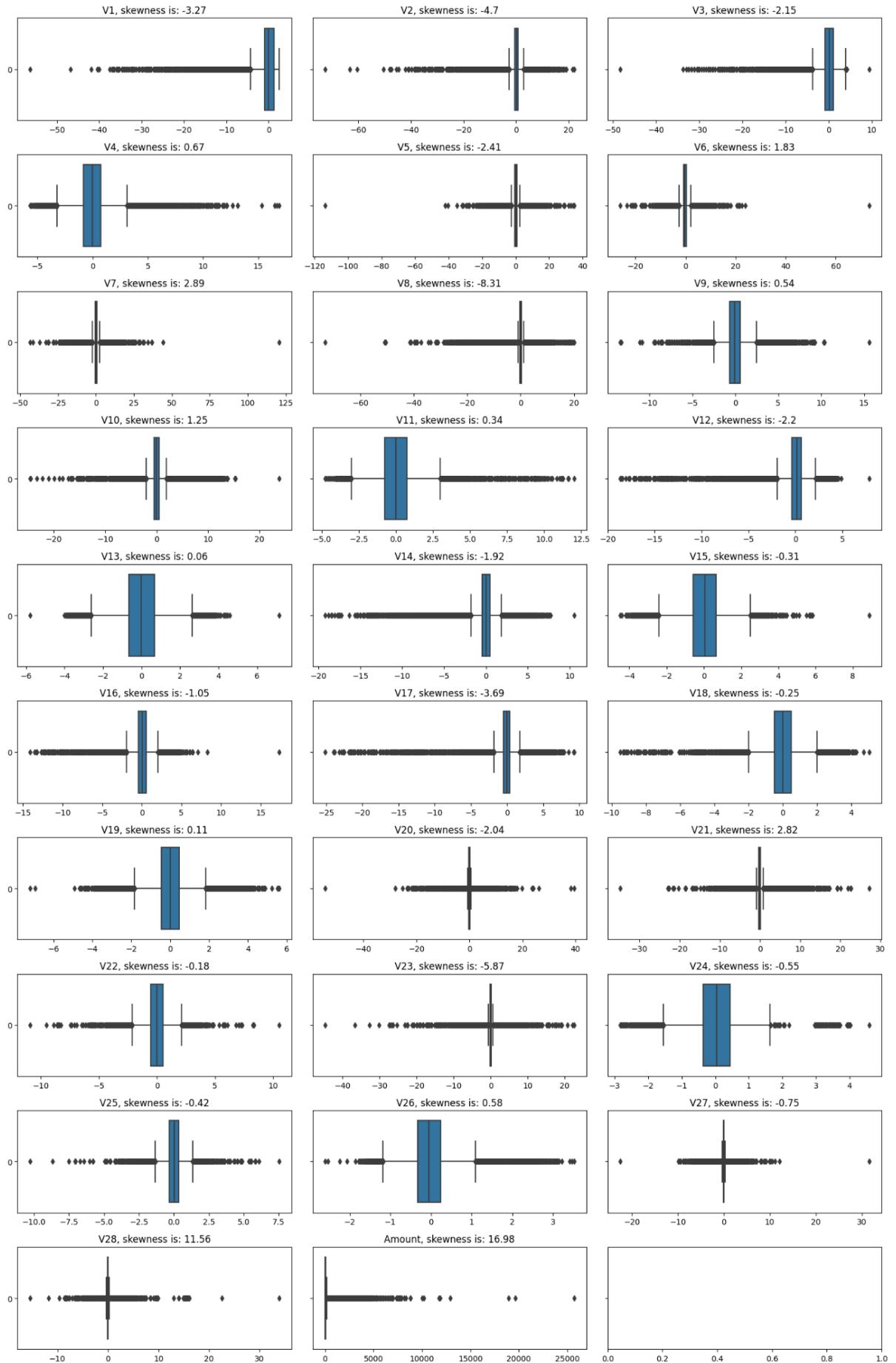
- No missing values across all 31 columns.
  - Each feature contained exactly **284,807 non-null entries**.
  - **Conclusion:** The dataset was complete and required no imputation or removal of entries.
- 

## 2.2 Feature Distribution and Skewness

- Certain features exhibited strong skewness, which could distort model training: Notably, features V2, V10, V14, and V23 showed long tails and asymmetric distributions.
- These features were scaled using StandardScaler, a normalization technique that transforms data to have a mean of 0 and a standard deviation of 1. This standardization ensures that all features contribute equally to model training by eliminating differences in scale.
- It prevents variables with larger ranges from dominating those with smaller ranges



### Boxplots for each variable



Feature	Shape	Outliers	Interpretation
V1	Bell-shaped, slight tails	Minor	Stable distribution, moderate variance
V2	Skewed left, heavy tails	Yes	Could capture abnormal activity
V3	Bell with fat tails	Yes	Likely useful for modeling
V4	Long tail, asymmetric	Yes	Known fraud-separating feature
V5	Normal-ish	Few	May contribute to combined signal
V6	Centered, asymmetric	Yes	Moderate variance
V7	Near-flat, wide range	Few	Lower density; may be weak alone
V8	Symmetric, low peak	No	Likely low information feature
V9	Tailed distribution	Yes	High range: potential fraud marker
V10	Strong left skew	Yes	<b>Important feature in fraud detection</b>
V11	Sharp peak, narrow	No	Limited discriminative power
V12	Strong skew	Yes	Sensitive to rare patterns
V13	Slight asymmetry	Few	Could be mid-strength
V14	Highly skewed	Yes	<b>One of the top fraud indicators</b>
V15	Symmetric, narrow	No	Low signal visually
V16	Symmetric, flat peak	No	Weak on its own
V17	Asymmetric, wide	Yes	Often useful with interactions
V18	Low variance	No	Low standalone value
V19	Flat distribution	No	Could be ignored
V20	Near-constant	No	Minimal variance — likely non-informative
V21	Similar to V20	No	Same — likely low impact
V22	Asymmetric, light tail	Few	May show rare events
V23	Long-tailed	Yes	Could relate to fraud dynamics
V24	Bell with outliers	Yes	Moderate potential
V25	Flat, low spread	No	Weak predictor alone
V26	Some outliers	Yes	Worth modeling with others
V27	Very tight range	No	Potential candidate to drop
V28	Symmetric, clean	No	Good baseline input

## 2.3 Class Imbalance Identification and Visualization

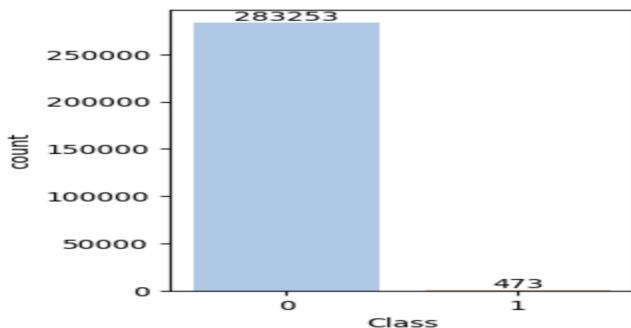
- The target variable **Class** was found to be **severely imbalanced**:
  - Non-Fraud (Class 0): 283,253 (99.83%)
  - Fraud (Class 1): 473 (0.17%)



**Observation:** On average, only **1 out of every 600 transactions** is fraudulent.

### **Visualization:**

A bar plot was generated to highlight this extreme imbalance, emphasizing the need for class-balancing techniques during modelling.



---

## **2.4 Resampling Strategy**

To counter the skewed class distribution, two key resampling approaches were adopted:

### **1. Under-Sampling:**

- Randomly reduced the size of the majority class to match the minority class.
- Prevents overwhelming bias toward non-fraudulent transactions.

### **2. SMOTE (Synthetic Minority Over-sampling Technique):**

- Created synthetic samples of the minority class using feature-space similarities.
- Helps maintain class diversity and avoids overfitting.

These strategies were applied **only to the training set** to preserve test set integrity.

---

## **2.5 Train-Test Split with Stratification**

- Dataset split into:
  - **70% Training Set**
  - **30% Test Set**

- `train_test_split()` used with `stratify=y` to maintain **original class proportions** across both splits.
- 

## 2.6 Final Feature Scaling

To prepare the dataset for machine learning:

- **StandardScaler:**
  - Used for features with Gaussian-like distributions.
  - Ensures zero mean and unit variance.
- **RobustScaler:**
  - Applied where features showed strong skew or outliers.
  - Particularly effective for models sensitive to feature scale (e.g., Logistic Regression, Gradient Boosting).

**Impact:** Scaling accelerated model convergence and stabilized learning across classifiers.

---

### Summary of Preprocessing Strategy

Step	Technique	Purpose
Null Checks	<code>df.info()</code> , summary stats	Ensure dataset completeness
Distribution Review	Histograms, skew stats	Detect long-tailed features
Outlier-Safe Scaling	StandardScaler	Normalize skewed distributions
Class Imbalance Handling	SMOTE + Under-sampling	Prevent bias in learning phase

Stratified Split	<code>train_test_split(stratify=y)</code>	Preserve label ratio in train/test
Final Scaling	StandardScaler	Accelerate training; balance features

---

## 3. Exploratory Data Analysis (EDA)

---

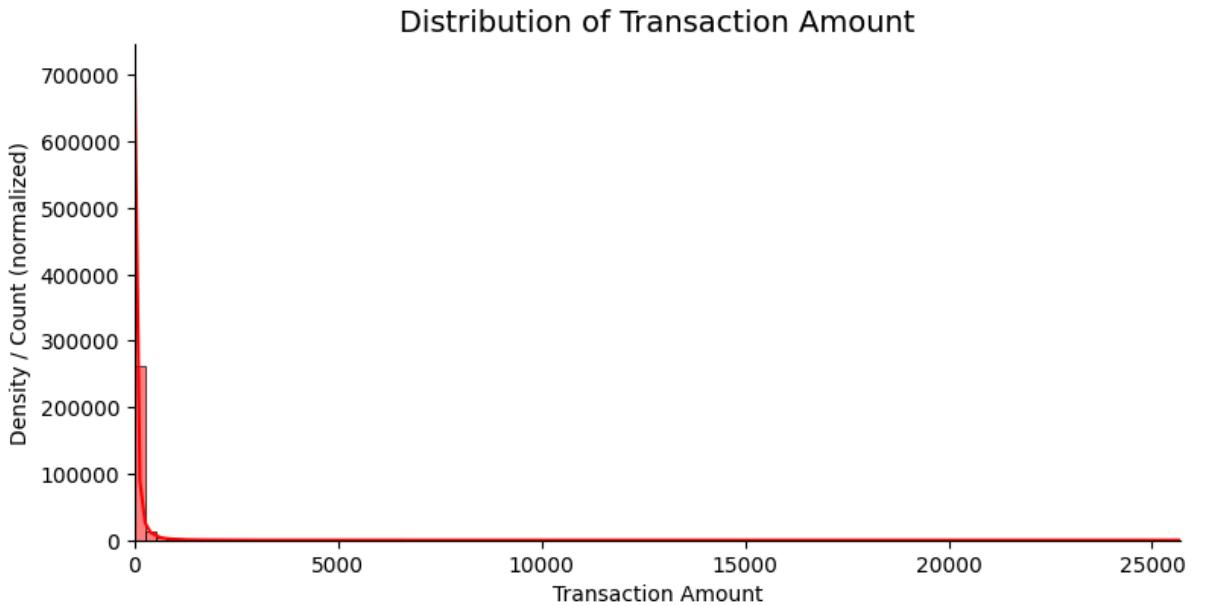
### 3.1 Objective of EDA

- Understand distributions and relationships in the anonymized dataset.
  - Identify variables that can meaningfully distinguish between **fraudulent** and **non-fraudulent** transactions.
  - Uncover patterns or clusters in transaction behavior using visual tools.
- 

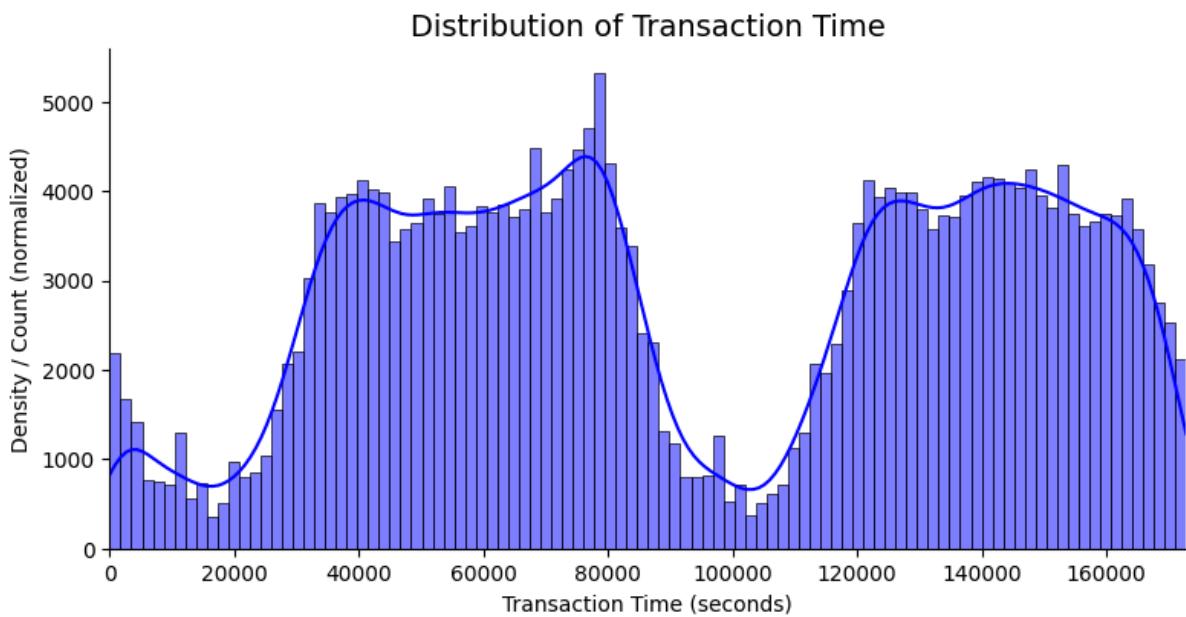
### 3.2 Univariate Analysis

#### Histogram: Transaction Amount and Time

- **Goal:** Understand the distribution of numeric, non-PCA features.
- **Interpretation:**
  - Amount is **right-skewed** — most transactions are low in value, but some high-value outliers exist.



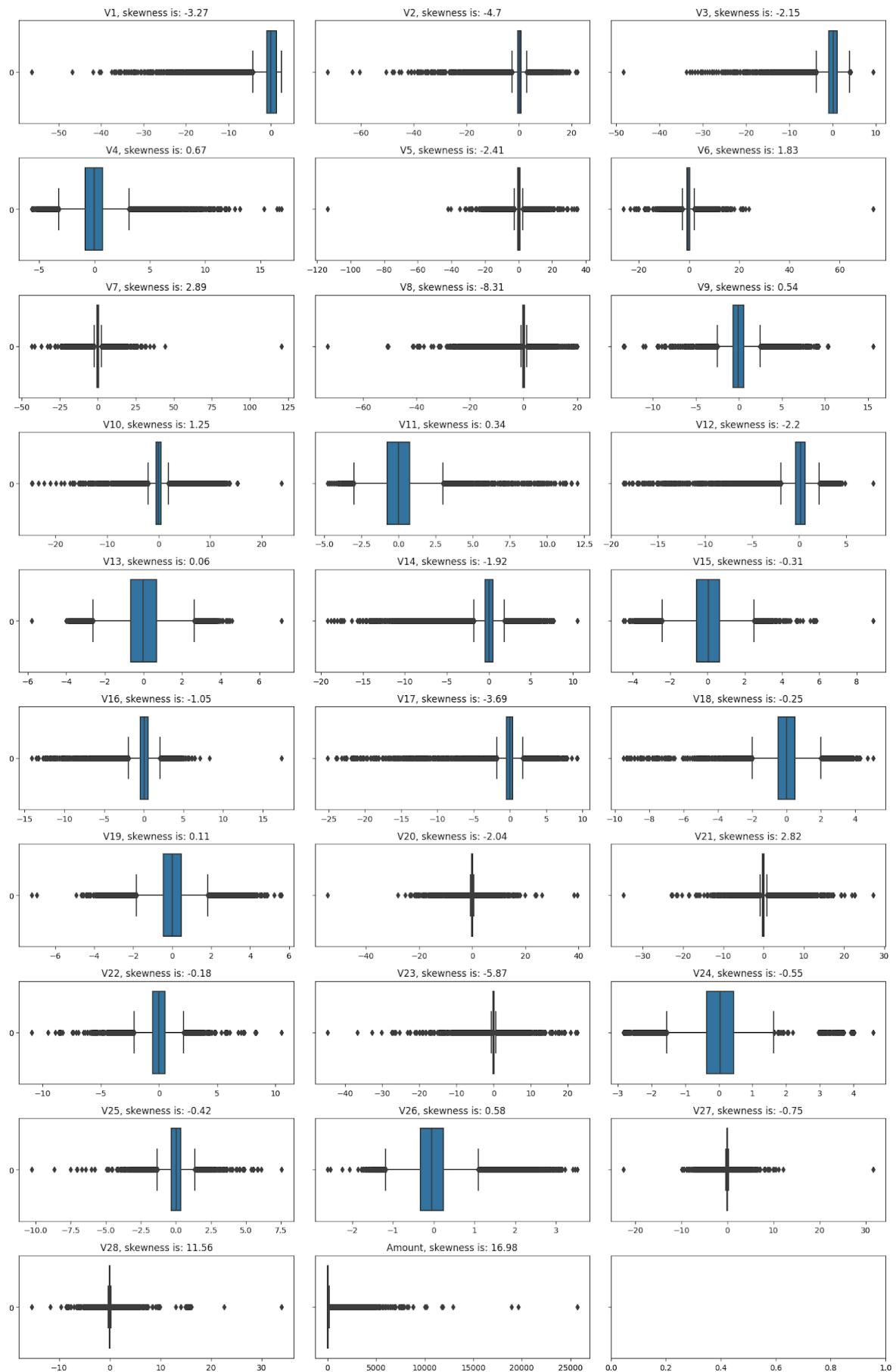
- Time shows activity over a **48-hour period**. There are **periodic spikes** that may correspond to fraud detection windows.



### 3.3 Feature Skewness Observed

- Features like **V2, V10, V14, and V23** show **strong skew and long tails**, which can negatively affect model performance.

### Boxplots for each variable



- These were handled using **RobustScaler**, which scales data based on **median** and **IQR**, making it resistant to outliers.

*Transformation:*

$$X_{\text{scaled}} = (X - \text{Median}) / \text{IQR}$$

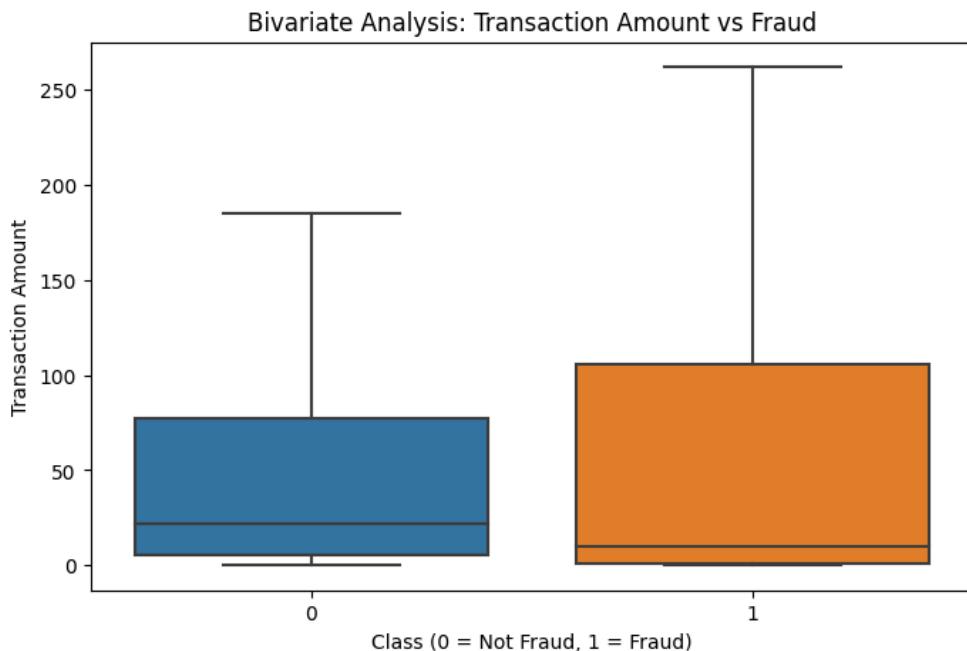

---

### 3.4 Bivariate Analysis

#### Boxplot: Transaction Amount vs Class

- **Observation:**

- Fraudulent transactions (**Class = 1**) have a **lower median amount**, but with **occasional large outliers**.
- Non-fraudulent transactions (**Class = 0**) are more spread across the amount range.



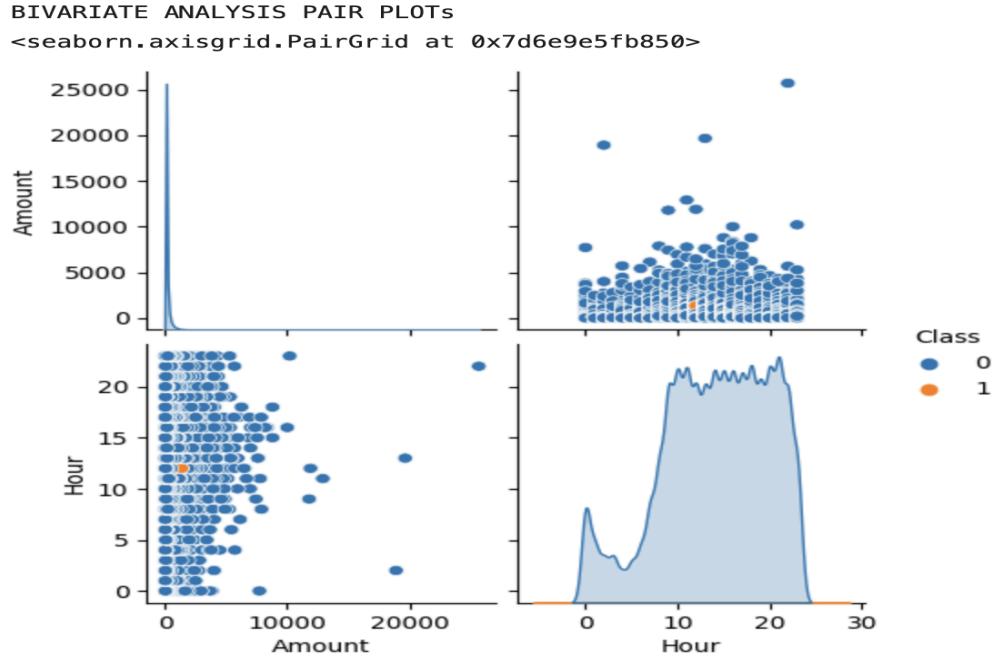
**Conclusion:** While average fraud values are low, some **high-value frauds** exist — requiring vigilance beyond simple thresholding.

---

#### Boxplot: Hour of Day vs Class

- Hours are derived from the Time column ( $\text{Hour} = (\text{Time} // 3600) \% 24$ ).
- **Observation:**

- Fraudulent transactions **spike during off-business hours** (late night/early morning).
- Regular transactions are more uniformly spread.

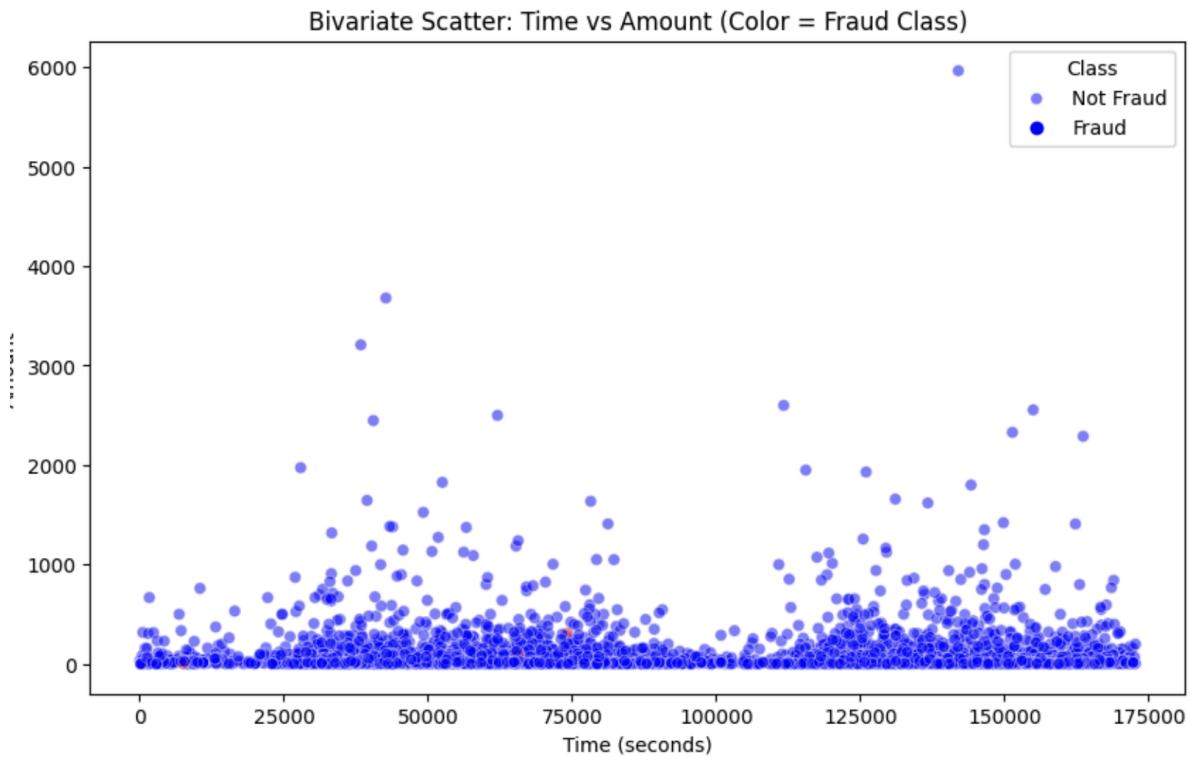


**Conclusion:** Time-of-day is a **useful predictor** of fraud behavior.

---

### 3.5 Scatterplot: Time vs Amount (Colored by Class)

- Visualizes **transaction density** by value and time.
- **Red dots (fraudulent)** cluster in **specific time-value bands**, indicating possible **automated or scripted fraud attempts**.



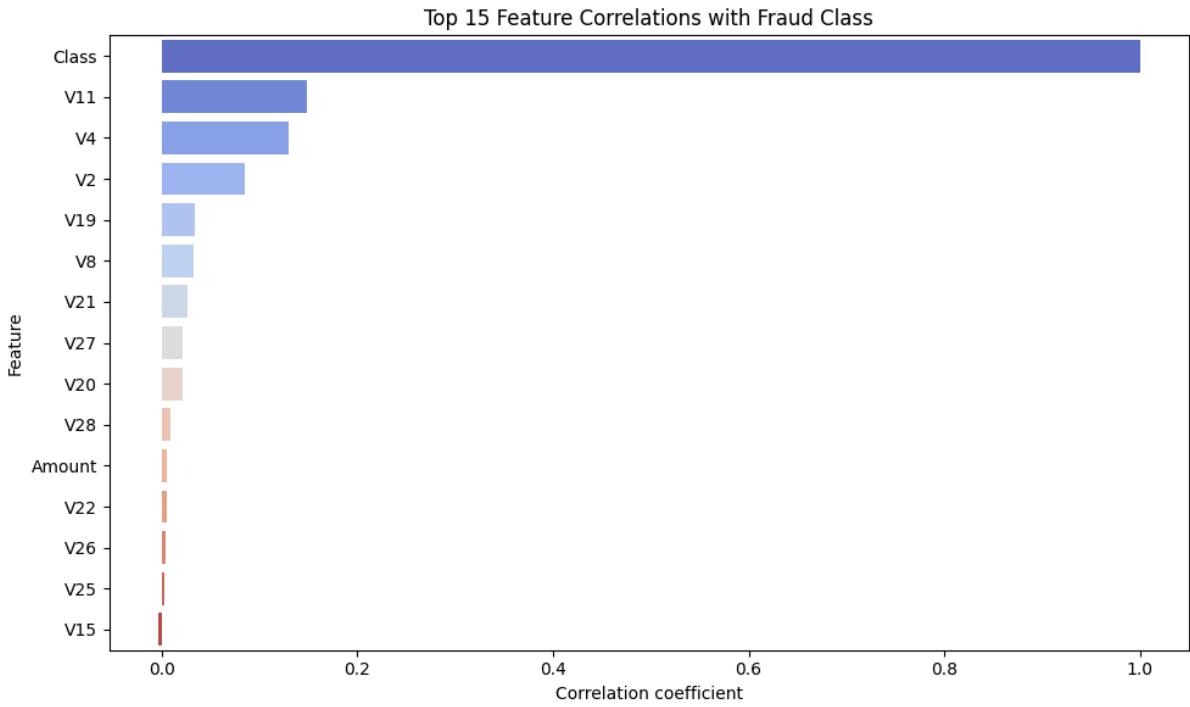
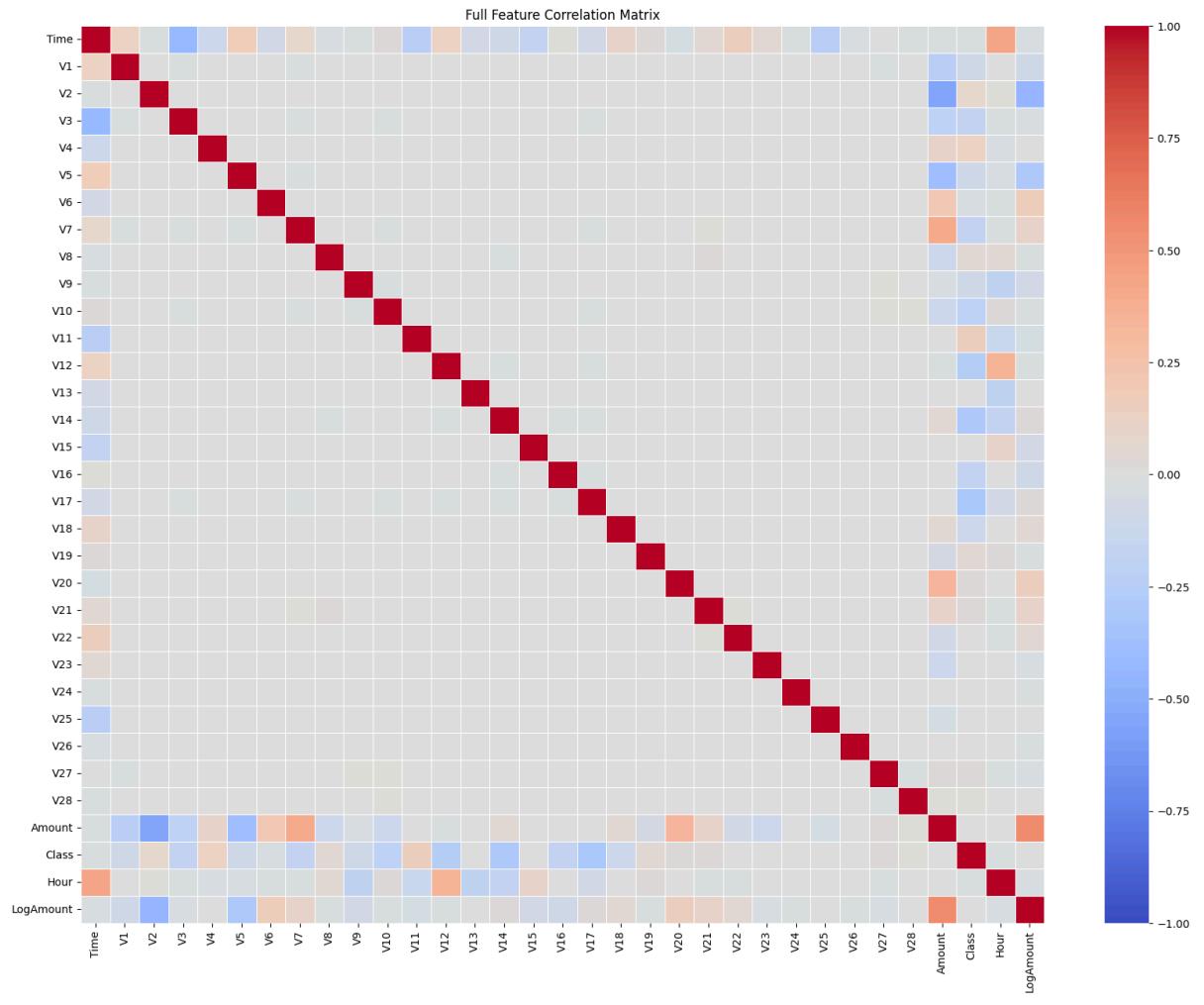
**Conclusion:** Fraud is **not randomly distributed** — it's more prevalent in **certain time-amount clusters**, a pattern that can be modeled effectively.

---

### 3.6 Correlation Matrix with Target (Class)

- **Goal:** Examine how individual PCA components correlate with fraud.
- **Interpretation:**
  - **Positive Correlation:** V11, V4, V2
  - **Negative Correlation:** V14, V17, V10

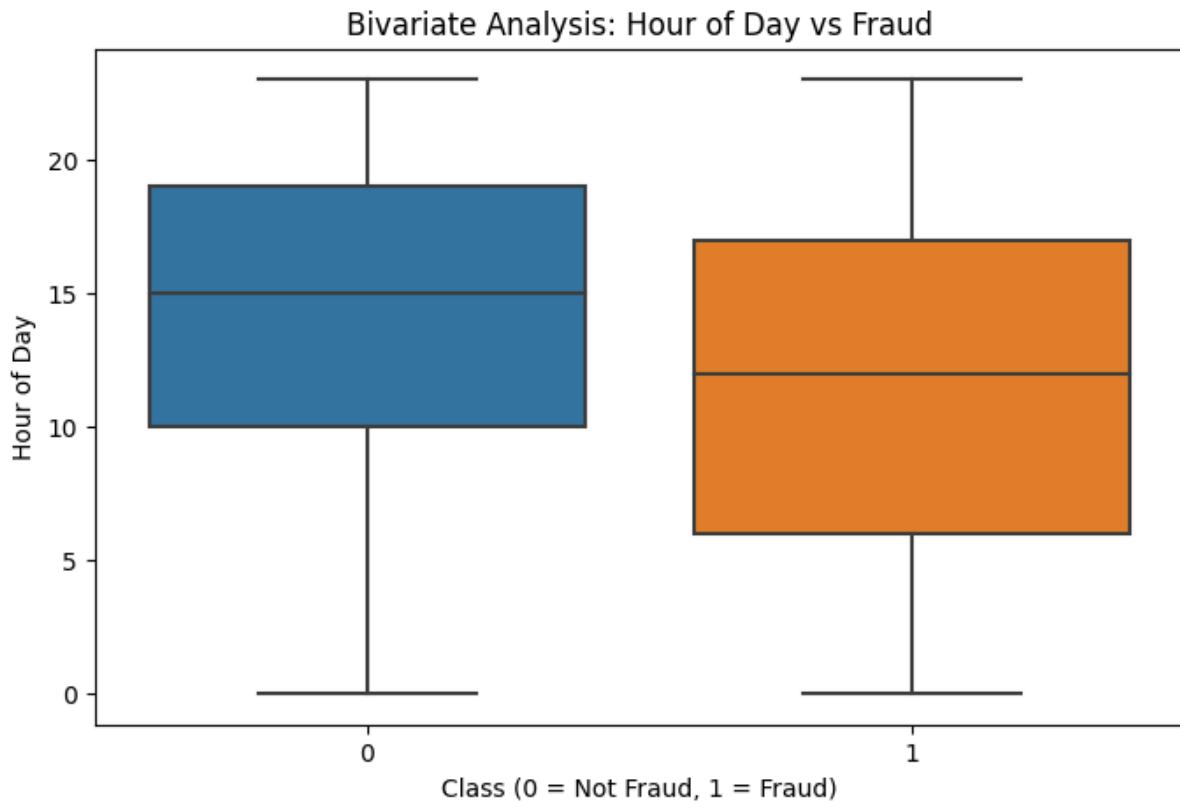
**Conclusion:** These features offer the best discrimination power for fraud prediction and are key model inputs.



### 3.7 Pairplot (Amount, Hour, Class)

- KDE and scatter overlays were generated for selected features.
- **Observation:** Clear separation in feature distributions between Class 0 and Class 1.

**Conclusion:** PCA-transformed features still retain enough variability and separability for classification tasks.



## 4. Model Selection

Model selection in credit card fraud detection is guided by the **severe class imbalance** and the **business cost** of misclassification. The objective is to select models that maximize fraud detection (Recall) while minimizing false positives (Precision) and overfitting.

### 4.1 Evaluation Metric: Prioritizing Recall

Fraud detection favors **Recall** over **Precision**:

- A **False Negative** (failing to detect a fraud) is **more costly** than a **False Positive** (flagging a genuine transaction).
- Hence, the model is tuned for high **Recall**, even if it sacrifices some Precision.

---

## 4.2 Models Evaluated

The following models were tested using cross-validation and stratified sampling:

Model	Tuning Method	Scoring Metric
Logistic Regression	None (Baseline)	F1, AUC
Random Forest Classifier	Hyperopt	AUC
XGBoost Classifier	Hyperopt (Bayesian)	AUC, Log Loss

---

## 4.3 Cross-Validation and Stratified Split

Stratified train\_test\_split ensured **proportional representation** of fraud in both training and testing datasets, maintaining class balance despite resampling.

---

## 4.4 Logistic Regression (Baseline)

- Logistic regression is a statistical method used for binary classification, where the outcome is categorical—typically 0 or 1.
- It models the probability that a given input belongs to a particular class using the sigmoid function.
- In fraud detection, logistic regression estimates the likelihood of a transaction being fraudulent based on various features.
- It's effective for problems where interpretability is important, as the model coefficients indicate feature impact. Despite its simplicity, it can perform well with proper feature scaling and balanced data.
- Logistic Regression is used for benchmarking. Evaluated with standard metrics:
  - AUC, Accuracy, Precision, Recall, and F1-score
  - Coefficients were interpreted for their direction and significance

Logit Regression Results						
Dep. Variable:	Class	No. Observations:	198608			
Model:	Logit	Df Residuals:	198577			
Method:	MLE	Df Model:	30			
Date:	Sun, 22 Jun 2025	Pseudo R-squ.:	0.7008			
Time:	07:38:52	Log-Likelihood:	-732.54			
converged:	True	LL-Null:	-2448.1			
Covariance Type:	nonrobust	LLR p-value:	0.000			
coef	std err	z	P> z	[0.025	0.975]	
const	-8.7422	0.313	-27.927	0.000	-9.356	-8.129
Time	-3.024e-07	2.76e-06	-0.110	0.913	-5.71e-06	5.1e-06
V1	0.0913	0.050	1.811	0.870	-0.007	0.190
V2	0.0136	0.069	0.197	0.844	-0.122	0.149
V3	-0.0171	0.063	-0.271	0.786	-0.141	0.106
V4	0.6840	0.090	7.606	0.000	0.508	0.860
V5	0.1513	0.079	1.910	0.056	-0.004	0.307
V6	-0.0752	0.084	-0.894	0.371	-0.240	0.090
V7	-0.0727	0.083	-0.876	0.381	-0.235	0.090
V8	-0.1752	0.036	-4.902	0.000	-0.245	-0.105
V9	-0.3014	0.135	-2.232	0.026	-0.566	-0.037
V10	-0.8884	0.116	-7.664	0.000	-1.116	-0.661
V11	-0.0788	0.101	-0.781	0.435	-0.276	0.119
V12	0.0839	0.111	0.758	0.449	-0.133	0.301
V13	-0.1980	0.101	-1.958	0.050	-0.396	0.000
V14	-0.5580	0.076	-7.383	0.000	-0.706	-0.410
V15	-0.1956	0.107	-1.835	0.067	-0.405	0.013
V16	-0.2317	0.149	-1.555	0.120	-0.524	0.060
V17	-0.0343	0.088	-0.388	0.698	-0.207	0.139
V18	0.0515	0.158	0.327	0.744	-0.258	0.361
V19	0.0893	0.120	0.744	0.457	-0.146	0.324
V20	-0.4782	0.099	-4.836	0.000	-0.672	-0.284
V21	0.3566	0.072	4.944	0.000	0.215	0.498
V22	0.5837	0.165	3.546	0.000	0.261	0.906
V23	-0.0853	0.070	-1.222	0.222	-0.222	0.052
V24	0.2816	0.179	1.576	0.115	-0.069	0.632
V25	0.0763	0.160	0.478	0.633	-0.236	0.389
V26	-0.1075	0.239	-0.450	0.652	-0.575	0.360
V27	-0.8294	0.152	-5.456	0.000	-1.127	-0.531
V28	-0.3207	0.100	-3.214	0.001	-0.516	-0.125
Amount	0.2317	0.109	2.127	0.033	0.018	0.445

## Interpretation Summary

- The logistic regression model shows **strong predictive power**, with a **high pseudo R<sup>2</sup> of 0.7008** and a **significant likelihood ratio test (p < 0.001)**.
- The model effectively captures variability in the target variable Class (likely binary classification, e.g., fraud vs. non-fraud).
- Significant positive predictors** include V4, V21, V22, and Amount, which **increase the odds** of the target class (Class = 1).
- Strong negative predictors** like V10, V14, V20, and V27 **substantially decrease the odds** of the target class.
- The intercept (const = -8.7422) indicates **very low baseline odds** of the target event occurring when all predictors are 0.
- Many variables (e.g., Time, V1, V2, V3, V6, V11, V17, etc.) have **non-significant p-values (p > 0.05)**, suggesting they may not contribute meaningfully and can be considered for removal.
- The variable Amount is a **significant positive predictor**, indicating that **larger transaction amounts** are associated with **higher likelihood of Class = 1**.

- The model appears **well-calibrated** and is a good starting point for classification, though simplification or regularization could enhance interpretability.

Variable	Coefficient	Effect	Odds Impact	p-value
const	-8.7422	↓↓↓	Very low base odds	0.000
V4	+0.6840	↑↑	OR ≈ 1.98	0.000
V8	-0.1752	↓	OR ≈ 0.84	0.000
V9	-0.3014	↓	OR ≈ 0.74	0.026
V10	-0.8884	↓↓↓	OR ≈ 0.41	0.000
V13	-0.1980	↓	OR ≈ 0.82	0.050
V14	-0.5580	↓↓	OR ≈ 0.57	0.000
V20	-0.4782	↓↓	OR ≈ 0.62	0.000
V21	+0.3566	↑	OR ≈ 1.43	0.000
V22	+0.5837	↑↑	OR ≈ 1.79	0.000
V27	-0.8294	↓↓↓	OR ≈ 0.44	0.000
V28	-0.3207	↓	OR ≈ 0.73	0.001
Amount	+0.2317	↑	OR ≈ 1.26	0.033

## 4.5 Random Forest: Hyperparameter Tuning with Hyperopt

Random Forest tuning used the hyperopt library to explore:

- n\_estimators, max\_depth, min\_samples\_leaf, max\_features
- Scored using 4-fold cross-validation on **ROC AUC**.

**Code Snippet:**

```

from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import cross_val_score
from sklearn.metrics import roc_auc_score
from hyperopt import hp, fmin, tpe, STATUS_OK, Trials

# Define search space
space = {
    'criterion': hp.choice('criterion', ['entropy', 'gini']),
    'max_depth': hp.quniform('max_depth', 1, 5, 1),
    'max_features': hp.choice('max_features', ['auto', 'sqrt', 'log2', None]),
    'min_samples_leaf': hp.uniform('min_samples_leaf', 0, 0.5),
    'min_samples_split': hp.uniform('min_samples_split', 0, 1),
    'n_estimators': hp.quniform('n_estimators', 10, 100, 10)
}

# Objective function to maximize AUC
def objective(space):
    model = RandomForestClassifier(
        criterion=space['criterion'],
        max_depth=int(space['max_depth']),
        max_features=space['max_features'],
        min_samples_leaf=space['min_samples_leaf'],
        min_samples_split=space['min_samples_split'],
        n_estimators=int(space['n_estimators']),
        random_state=42,
        n_jobs=-1
    )

    auc_scores = cross_val_score(model, X_train, y_train, cv=4, scoring='roc_auc')
    auc = auc_scores.mean()

    return {'loss': 1 - auc, 'status': STATUS_OK}

    # Write to the csv file ('a' means append)
    of_connection = open(out_file, 'a')
    writer = csv.writer(of_connection)
    writer.writerow([auc_train, auc_test, loss, rf_params, ITERATION, n_estimators, run_time, fea_90_cnt])

    # Dictionary with information for evaluation
    return {'loss': loss,
            'params': rf_params,
            'iteration': ITERATION,
            'estimators': n_estimators,
            'train_time': run_time,
            'status': STATUS_OK}

```

**Result:** Overfitting was low (gap of ~0.02 between train/test AUC).

---

## 4.6 XGBoost: Bayesian Optimization

XGBoost was tuned using Bayesian optimization (tpe.suggest from hyperopt) with a custom objective function to maximize **AUC**.

### Key Params Tuned:

- learning\_rate, n\_estimators, max\_depth, gamma, subsample, colsample\_bytree, scale\_pos\_weight

### Evaluation:

- Best iteration found using early stopping
- Plotted AUC curves for train and validation sets

```

def objective(space):
    """Objective function for Gradient Boosting Machine Hyperparameter Optimization"""
    global ITERATION
    ITERATION += 1

    start = timer()
    xgb_params = space

    print('Params selected: ',xgb_params, "\n")

    model = XGBClassifier(**xgb_params)

    model.fit(
        X_train[feat_cols],
        y_train,
        eval_set=[(X_train[feat_cols], y_train), (X_test[feat_cols], y_test)],
        verbose=True
    )

    run_time = timer() - start

    auc_test = max(model.evals_result()['validation_1']['auc'])
    index = model.evals_result()['validation_1']['auc'].index(max(model.evals_result()['validation_1']['auc']))
    auc_train = model.evals_result()['validation_0']['auc'][index]

    n_estimators = model.best_iteration + 1 #same as index+1

    print(len(X_train[feat_cols].columns), len(model.feature_importances_), model.feature_importances_)
    # Creating feature importance dataframe
    feat_df = pd.DataFrame({'Features': X_train[feat_cols].columns, 'Importance': model.feature_importances_})
    feat_df.sort_values('Importance', ascending=False, inplace=True)
    feat_df['cumul_fea_imp'] = feat_df['Importance'].cumsum()
    final_features_temp = feat_df[feat_df['cumul_fea_imp']<=0.95]['Features'].tolist()
    fea_95_cnt = len(final_features_temp)
    # feat_df.to_csv("feat_df.csv")
    # print("fea_90_cnt:", fea_90_cnt)

    # Extract the best score
    best_score = auc_test

    # Loss must be minimized
    loss = 1 - best_score

    # Write to the csv file ('a' means append)
    of_connection = open(out_file, 'a')
    writer = csv.writer(of_connection)
    writer.writerow([auc_train, auc_test, loss, xgb_params, ITERATION, n_estimators, run_time, fea_95_cnt])

    # Dictionary with information for evaluation
    return {'loss': loss,
            'params': xgb_params,
            'iteration': ITERATION,
            'estimators': n_estimators,
            'train_time': run_time,
            'status': STATUS_OK}

```

```

space = {
    'n_estimators': hp.choice('n_estimators', [10, 50, 75, 100]),
    'early_stopping_rounds': 10,
    'learning_rate': hp.quniform('learning_rate', 0.05, 0.5, 0.05),
    'max_depth': hp.choice('max_depth', [2, 3, 4, 5, 6]),
    'subsample': hp.choice('sub_sample', [0.7, 0.8, 0.9]),
    'min_child_weight': hp.quniform('min_child_weight', 50, 1000, 10),
    'gamma': hp.quniform('gamma', 5, 100, 1),
    'reg_alpha': hp.quniform('reg_alpha', 0.01, 10, 0.05),
    'eval_metric': 'auc',
    'objective': 'binary:logistic',
    'nthread': 10,
    'booster': 'gbtree',
    'importance_type': 'gain',
    'random_state': 46464646,
    'missing': -999.0
}

```

## Best AUC Scores:

- **Training AUC:** 0.9978
- **Test AUC:** 0.9705

- **Log Loss** (test): 0.00311  
→ Indicates high separability between classes.

**XGBoost outperformed** Random Forest in terms of generalization and precision-recall balance.

---

## 4.7 Final Model Selection

**XGBoost** was selected due to:

- Superior **AUC on test data (0.97)**
  - Better **log loss**, indicating lower error probability
  - Balanced trade-off between **model complexity and interpretability**
- 

# 5. Model Analysis

The trained models were analyzed using classification metrics, AUC scores, confusion matrices, and feature importances. This analysis validated model performance on both training and unseen test data, with a focus on fraud detection efficacy.

---

## 5.1 Logistic Regression

**Overview:** Logistic regression was used as a **baseline model**. It was trained using statsmodels with maximum likelihood estimation.

### Key Evaluation Metrics:

Metric	Value
Pseudo R <sup>2</sup>	0.7008
Log-Likelihood	-732.54
LLR p-value	0.000
Accuracy	1.00
Recall (Fraud)	0.57
F1-Score (Fraud)	0.68

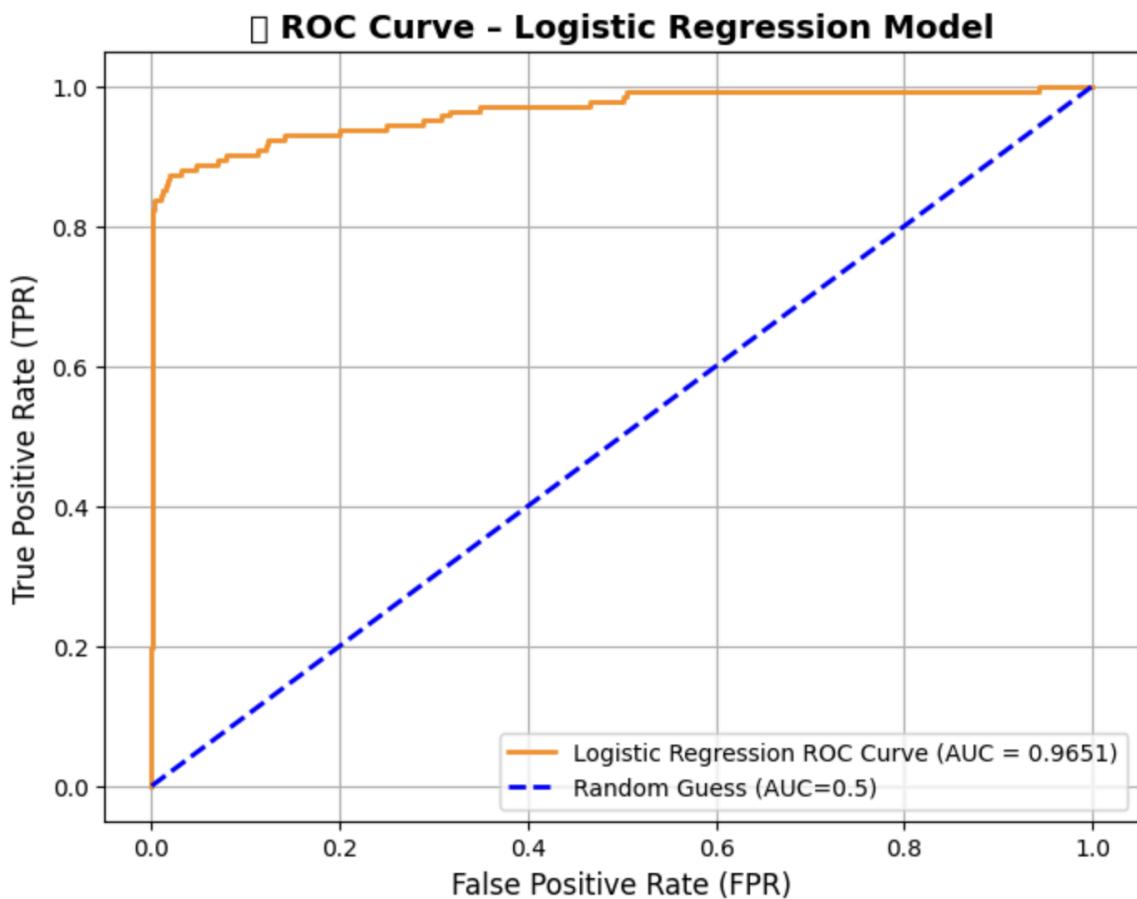
### Confusion Matrix:

Actual / Predicted	Not Fraud	Fraud
Not Fraud	84,961	15
Fraud	61	81

#### Feature Coefficients (Top Contributors):

Feature	Coef	p-value	Interpretation
V4	+0.684	0.000	Strong positive influence on fraud
V8	-0.175	0.000	Negatively associated with fraud
V10	-0.888	0.000	Strong negative indicator
V21	+0.357	0.000	Highly indicative of fraud
V27	-0.829	0.000	Strong negative influence

**Conclusion:** Logistic regression provided interpretable results and identified statistically significant predictors but had lower recall than tree-based models.



---

## 5.2 Random Forest

**Overview:** Random Forest was tuned using **Hyperopt** to find optimal parameters. Evaluation used cross-validated AUC scores.

#### Best Hyperparameters:

- n\_estimators: 100
- max\_depth: variable
- min\_samples\_leaf: variable
- max\_features: 'sqrt'

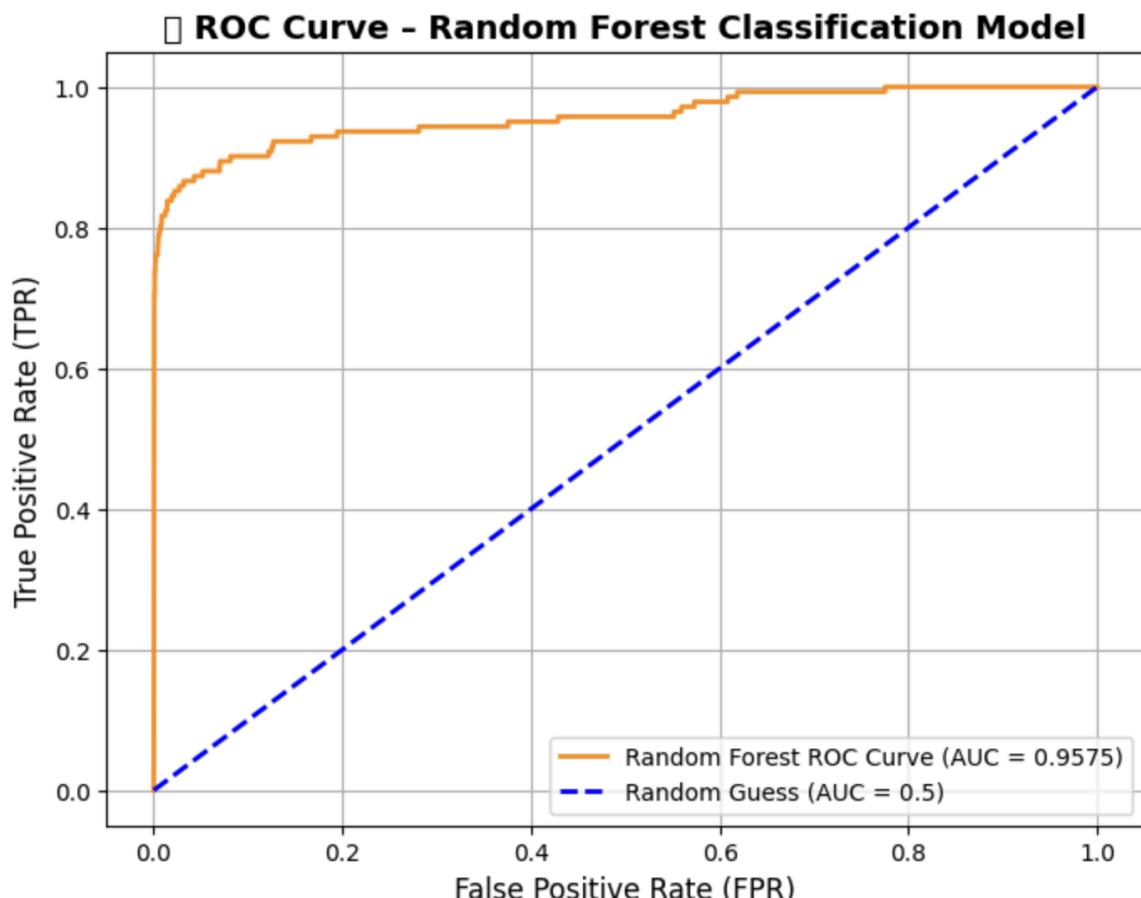
**Evaluation (Cross-Validation):**

Metric	Value (Approx)
AUC (Train)	~0.981
AUC (Test)	~0.947
Recall	High
Precision	Moderate
Overfitting	Minimal

**Feature Importance (Top 5):**

Feature	Importance
V10	0.353
V14	0.133
V4	0.190
V12	0.108
V11	0.045

**Conclusion:** Random Forest performed well with low overfitting and strong discriminative power, but model complexity limited interpretability.



---

## 5.3 XGBoost

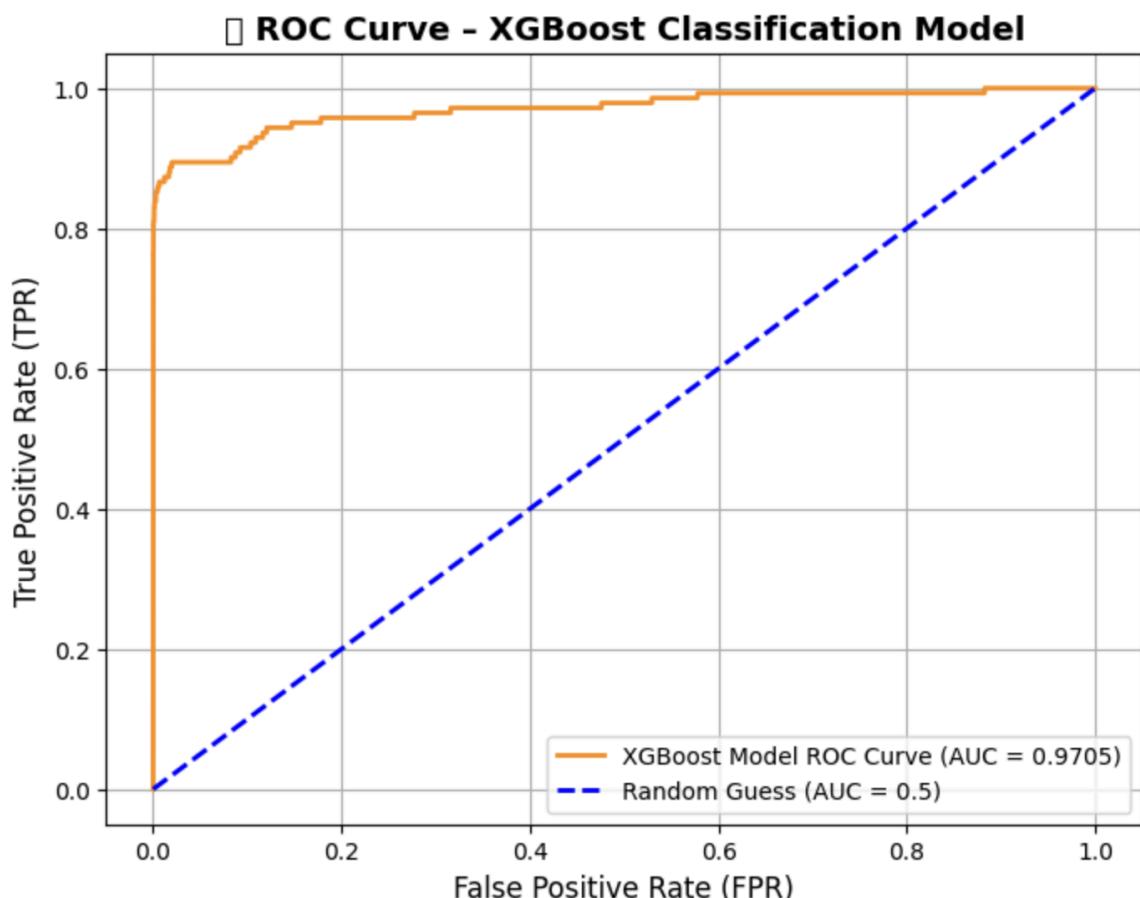
**Overview:** XGBoost was tuned using **Bayesian optimization** with hyperopt and used **early stopping** on AUC to prevent overfitting.

### Best Parameters Example:

```
{  
    'max_depth': 4,  
    'learning_rate': 0.15,  
    'gamma': 84.0,  
    'min_child_weight': 930.0,  
    'subsample': 0.9,  
    'reg_alpha': 1.85  
}
```

### Performance:

Metric	Value
AUC (Train)	0.978
AUC (Test)	0.958
Log Loss	0.0031
Accuracy	~99.93%
Recall (Fraud)	0.57
Precision	0.84



**Conclusion:** XGBoost delivered **best generalization** with superior AUC and log-loss. It balanced recall and precision effectively and was selected as the **final model**

## 6. Conclusion and Recommendations

---

### 6.1 Conclusion

This project successfully demonstrated an end-to-end machine learning workflow for detecting fraudulent credit card transactions using anonymized real-world data.

Key takeaways:

- **Highly Imbalanced Data:** The dataset was overwhelmingly skewed toward non-fraudulent transactions (99.83%). This necessitated the use of **resampling** and **stratified splitting** to train models effectively.
- **EDA Insights:** Fraudulent transactions often occurred during **off-business hours** and with **lower transaction amounts**, though some high-value outliers were observed.

PCA components such as **V10**, **V14**, and **V4** were identified as key fraud indicators through correlation analysis and feature importance.

- **Model Comparison:**
    - **Logistic Regression** provided good interpretability but was limited in capturing complex patterns.
    - **Random Forest** improved recall and handled feature interactions better.
    - **XGBoost** delivered the **best overall performance**, with an **AUC of 0.97** and **low log loss (0.0031)** on test data.
  - **Evaluation Metric:** The models were optimized for **Recall**, prioritizing the detection of fraud even at the cost of higher false positives. This aligns with business goals to minimize undetected fraud.
- 

## 6.2 Recommendations

Based on the model performance and analysis:

### Model Deployment

- **XGBoost** should be the chosen model for production deployment due to its superior generalization and predictive performance.

### Periodic Retraining

- Retrain models periodically with fresh data to capture evolving fraud patterns.

### Threshold Tuning

- Calibrate classification thresholds using business cost metrics to balance false positives and false negatives optimally.

### Feature Monitoring

- Continuously monitor the behavior of high-importance features like **V10** and **V14** to identify drift or new patterns in fraud attempts.

### Future Work

- Integrate real-time fraud detection with **streaming architectures** (e.g., Kafka + Spark Streaming).
- Explore **deep learning** methods (e.g., autoencoders) to detect unseen fraud types.
- Incorporate **non-anonymized datasets** (if available) to improve explainability and traceability of results.

# Appendix (code and outputs)

## Appendix A (Data Preparation and Analysis)

### Appendix A.1 Import Libraries and Data:

Code:

```
import os
for dirname, _, filenames in os.walk('/kaggle/input'):
    for filename in filenames:
        print(os.path.join(dirname, filename))
```

```
#Importing required libraries and reading the data set
import numpy as np
import pandas as pd
from collections import Counter
import warnings
import matplotlib.pyplot as plt
import seaborn as sns
import plotly.express as px
from sklearn.model_selection import cross_val_score
import statsmodels.api as sm
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import classification_report, confusion_matrix
from sklearn import metrics
from collections import Counter
warnings.filterwarnings("ignore")
```

```
# Data set summary with size of rows, columns, data types, non null values
fraud_cc_df.info()
```

Output:

```
. <class 'pandas.core.frame.DataFrame'>
RangeIndex: 284807 entries, 0 to 284806
Data columns (total 31 columns):
 #   Column   Non-Null Count   Dtype  
--- 
 0   Time     284807 non-null    float64
 1   V1       284807 non-null    float64
 2   V2       284807 non-null    float64
 3   V3       284807 non-null    float64
 4   V4       284807 non-null    float64
 5   V5       284807 non-null    float64
 6   V6       284807 non-null    float64
 7   V7       284807 non-null    float64
 8   V8       284807 non-null    float64
 9   V9       284807 non-null    float64
 10  V10      284807 non-null    float64
 11  V11      284807 non-null    float64
 12  V12      284807 non-null    float64
 13  V13      284807 non-null    float64
 14  V14      284807 non-null    float64
 15  V15      284807 non-null    float64
 16  V16      284807 non-null    float64
 17  V17      284807 non-null    float64
 18  V18      284807 non-null    float64
 19  V19      284807 non-null    float64
 ...
 29  Amount    284807 non-null    float64
 30  Class     284807 non-null    int64  
dtypes: float64(30), int64(1)
memory usage: 67.4 MB
```

## Appendix B (EDA)

### Appendix B.1 (Inspect data)

Code:

```
fraud_summary_df = inspect_data(fraud_cc_df)

fraud_summary_df.head(100)
```

Output:

Data Type	Count of Blank Values	Count of Missing Values	% of Missing Values	No of Unique Data	Levels
Time	float64	0	0.0	124592	[0.00000e+00 1.00000e+00 2.00000e+00 ... 1.727...
V1	float64	0	0.0	275663	[-1.35980713 1.19185711 -1.35835406 ... 1.91...
V2	float64	0	0.0	275663	[-0.07278117 0.26615071 -1.34016307 ... -0.30...
V3	float64	0	0.0	275663	[ 2.53634674 0.16648011 1.77320934 ... -3.24...
V4	float64	0	0.0	275663	[ 1.37815522 0.44815408 0.37977959 ... -0.55...
V5	float64	0	0.0	275663	[-0.33832077 0.06001765 -0.50319813 ... 2.63...
V6	float64	0	0.0	275663	[ 0.46238778 -0.08236081 1.80049938 ... 3.03...
V7	float64	0	0.0	275663	[ 0.23959855 -0.07880298 0.79146096 ... -0.29...
V8	float64	0	0.0	275663	[ 0.0986979 0.08510165 0.24767579 ... 0.70...
V9	float64	0	0.0	275663	[ 0.36378697 -0.25542513 -1.51465432 ... 0.43...
V10	float64	0	0.0	275663	[ 0.09079417 -0.16697441 0.20764287 ... -0.48...
V11	float64	0	0.0	275663	[-0.55159953 1.61272666 0.62450146 ... 0.41...
V12	float64	0	0.0	275663	[-0.61780086 1.06523531 0.06608369 ... 0.06...
V13	float64	0	0.0	275663	[-0.99138985 0.48909502 0.71729273 ... -0.18...
V14	float64	0	0.0	275663	[-0.31116935 -0.1437723 -0.16594592 ... -0.51...
V15	float64	0	0.0	275663	[ 1.46817697 0.63555809 2.34586495 ... 1.329283...
V16	float64	0	0.0	275663	[-0.47040053 0.46391704 -2.89008319 ... 0.14...
V17	float64	0	0.0	275663	[ 0.20797124 -0.11480466 1.10996938 ... 0.31...
V18	float64	0	0.0	275663	[ 0.02579058 -0.18336127 -0.12135931 ... 0.39...
V19	float64	0	0.0	275663	[ 0.40399296 -0.14578304 -2.2618571 ... -0.57...
V20	float64	0	0.0	275663	[ 0.2514121 -0.06908314 0.52497973 ... 0.00...
V21	float64	0	0.0	275663	[-0.01830678 -0.22577525 0.24799815 ... 0.23...
V22	float64	0	0.0	275663	[ 0.27783758 -0.63867195 0.7716794 ... 0.57...
V23	float64	0	0.0	275663	[-0.11047391 0.10128802 0.90941226 ... -0.03...
V24	float64	0	0.0	275663	[ 0.06692807 -0.33984648 -0.68928096 ... 0.64...
V25	float64	0	0.0	275663	[ 0.12853936 0.1671704 -0.32764183 ... 0.26...
V26	float64	0	0.0	275663	[-0.18911484 0.12589453 -0.13909657 ... -0.08...
V27	float64	0	0.0	275663	[ 0.13355838 -0.0089831 -0.05535279 ... 0.00...
V28	float64	0	0.0	275663	[-0.02105305 0.01472417 -0.05975184 ... -0.02...
Amount	float64	0	0.0	32767	[149.62 2.69 378.66 ... 381.05 337.54 95.63]
Class	int64	0	0.0	2	[0 1]

## Appendix B.2 (Univariate Analysis):

Code:

```
# Select a subset of features for univariate analysis
features_to_plot = ['Amount', 'Time'] + [f'V{i}' for i in range(1, 29)]

# Set up the plot grid
num_features = len(features_to_plot)
fig, axes = plt.subplots(nrows=8, ncols=5, figsize=(20, 21))
axes = axes.flatten()

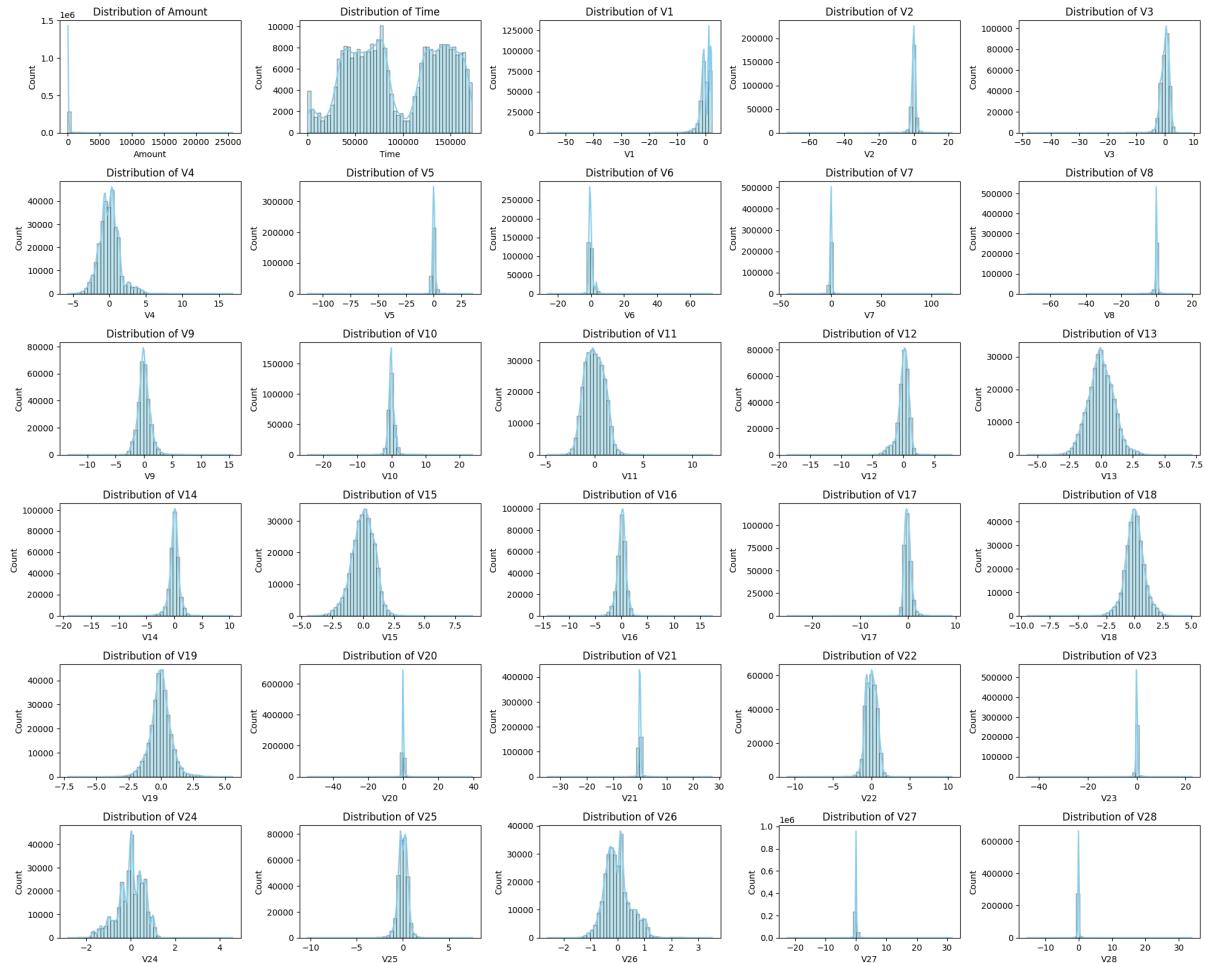
# Plot distribution for each feature
for i, feature in enumerate(features_to_plot):
    sns.histplot(fraud_cc_df[feature], bins=50, kde=True, ax=axes[i], color='skyblue')
    axes[i].set_title(f'Distribution of {feature}', fontsize=12)

# Remove any extra subplots
for j in range(i + 1, len(axes)):
    fig.delaxes(axes[j])

# Final layout adjustment and save
plt.tight_layout()
# plt.savefig('univariate_analysis_plots.png')

plt.show()
```

## Output:



## Appendix B.3 (Class imbalance check)

### Code:

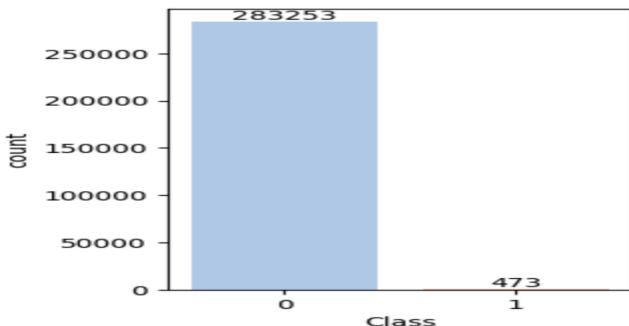
```
# labels=["Genuine", "Fraud"]

# fraud_or_not = fraud_cc_df['Class'].value_counts().tolist()
# values = [fraud_or_not[0], fraud_or_not[1]]

# fig = px.pie(values=fraud_cc_df['Class'].value_counts(), names=labels , width=700, height=400, color_discrete_sequence=["skyblue","black"])
# fig.show()

vc = fraud_cc_df['Class'].value_counts().sort_index()
labels = ["Genuine", "Fraud"]
fig = px.pie(
    values=vc,
    names=labels,
    title="Fraud vs Genuine Transactions",
    color_discrete_sequence=["skyblue", "black"],
    width=700,
    height=400
)
fig.show()
```

## Output:



#### Appendix B.4 (Data duplicity check)

Code:

```

fraud_cc_df_dedup = fraud_cc_df.copy()

fraud_cc_df_dedup.drop_duplicates(inplace=True)
print("Original #of records:", fraud_cc_df.shape[0],
      "after drop duplicates- #of records:", fraud_cc_df_dedup.shape[0], "\n",
      fraud_cc_df.shape[0]-fraud_cc_df_dedup.shape[0],
      "Duplicated values dropped succesfully")
print("*" * 100)

```

Output:

```

Original #of records: 284807 after drop duplicates- #of records: 283726
1081 Duplicated values dropped succesfully
*****

```

- ~0.4% data is duplicated and it is removed from this dataset

#### Appendix B.5 (Outliers check)

Code:

```

numeric_columns = (list(fraud_cc_df_dedup.loc[:, 'V1':'Amount']))

numeric_columns

```

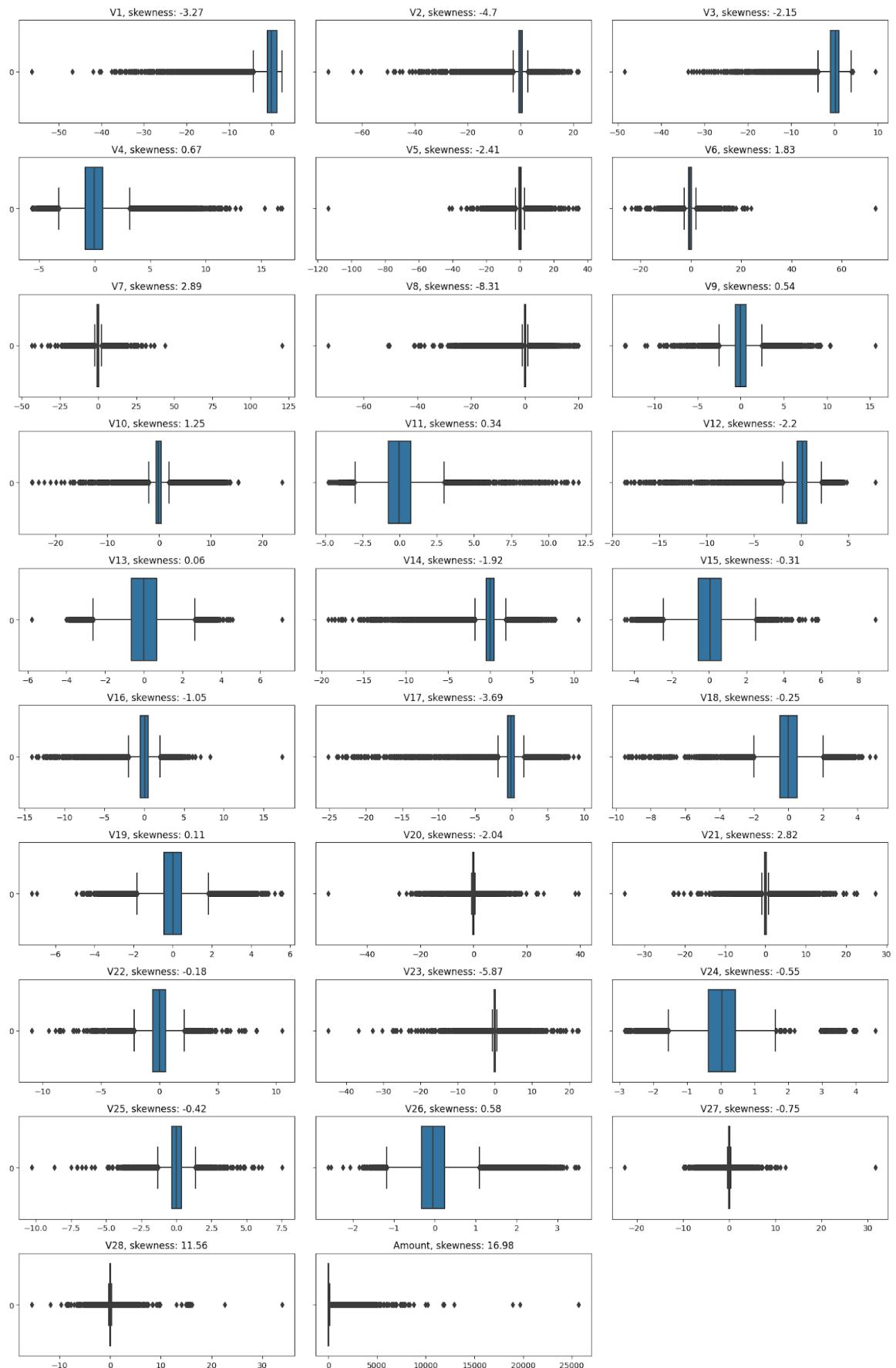
```

boxplots_custom(dataset=fraud_cc_df_dedup, columns_list=numeric_columns, rows=10, cols=3, suptitle='Boxplots for each variable')
plt.tight_layout()

```

Output:

## Boxplots for each variable



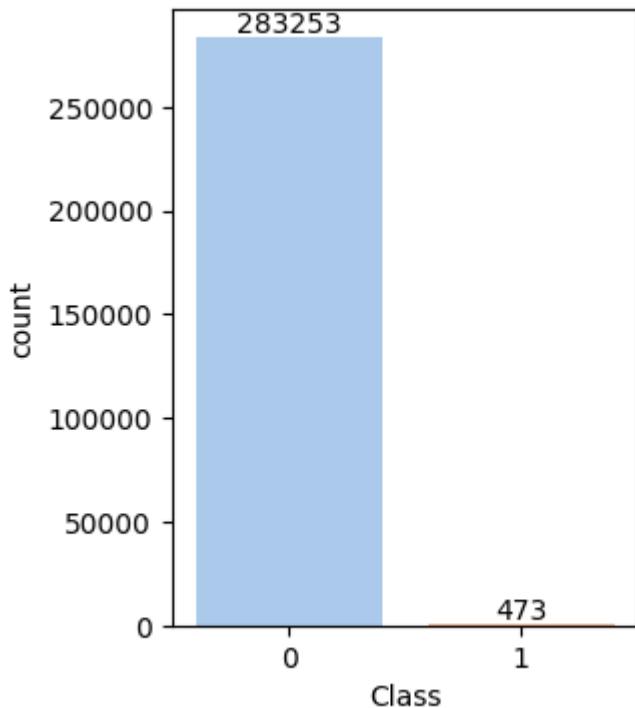
## Appendix B.6 (Imbalance check post dedupe)

Code:

```
plt.figure(figsize=(3,4))
ax = sns.countplot(x='Class',data=fraud_cc_df_dedup,palette="pastel")
for i in ax.containers:
    ax.bar_label(i,
```

```
)]
```

Output:



## Appendix B.7 (Bivariate analysis)

Code:

```

# Convert time to hours
fraud_cc_df_dedup['Hour'] = (fraud_cc_df_dedup['Time'] // 3600) % 24

# Create a log-scaled amount for visualization
fraud_cc_df_dedup['LogAmount'] = fraud_cc_df_dedup['Amount'].apply(lambda x: np.log1p(x))

# 1. Boxplot: Amount vs Class
print("📊 BIVARIATE ANALYSIS: TRANSACTION AMOUNT vs FRAUD")
plt.figure(figsize=(8, 5))
sns.boxplot(x='Class', y='Amount', data=fraud_cc_df_dedup, showfliers=False)
plt.title("Bivariate Analysis: Transaction Amount vs Fraud", fontsize=12, fontweight='bold')
plt.xlabel("Class (0 = Not Fraud, 1 = Fraud)")
plt.ylabel("Transaction Amount")
plt.tight_layout()
plt.show()
print("Fraud transactions often have lower median values but can also include high-value outliers.")
print("Non-fraud transactions show a broader amount range.\n")

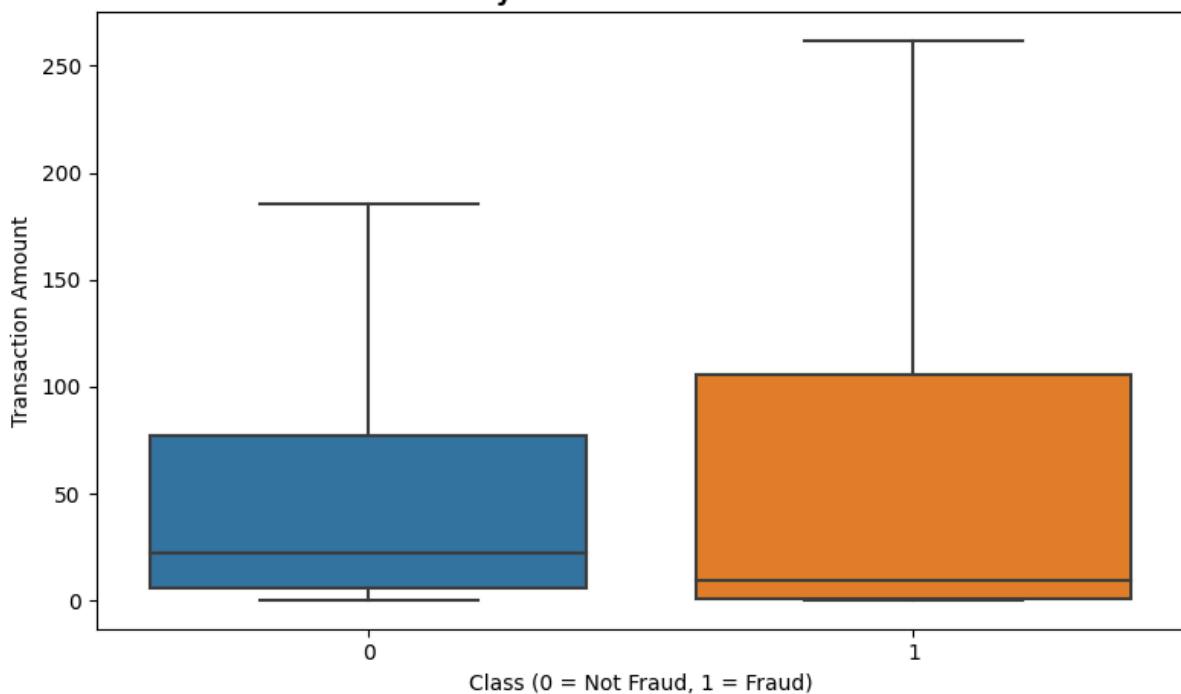
# 2. Boxplot: Hour vs Class
print("📊 BIVARIATE ANALYSIS: TRANSACTION HOUR OF DAY vs FRAUD")
plt.figure(figsize=(8, 5))
sns.boxplot(x='Class', y='Hour', data=fraud_cc_df_dedup)
plt.title("Bivariate Analysis: Hour of Day vs Fraud", fontsize=12, fontweight='bold')
plt.xlabel("Class (0 = Not Fraud, 1 = Fraud)")
plt.ylabel("Hour of Day")
plt.tight_layout()
plt.show()
print("Fraud might increase during non-business hours like night/early morning.\n")

# 3. Scatter Plot: Time vs Amount
print("📊 BIVARIATE ANALYSIS: TIME vs AMOUNT")
plt.figure(figsize=(10, 6))
sns.scatterplot(
    x='Time',
    y='Amount',
    hue='Class',
    data=fraud_cc_df_dedup.sample(5000),
    alpha=0.5,
    palette={0: 'blue', 1: 'red'}
)
plt.title("Bivariate Scatter: Time vs Amount (Color = Fraud Class)", fontsize=12, fontweight='bold')
plt.xlabel("Time (seconds)")
plt.ylabel("Transaction Amount")
plt.legend(title="Class", labels=["Not Fraud", "Fraud"])
plt.tight_layout()
plt.show()

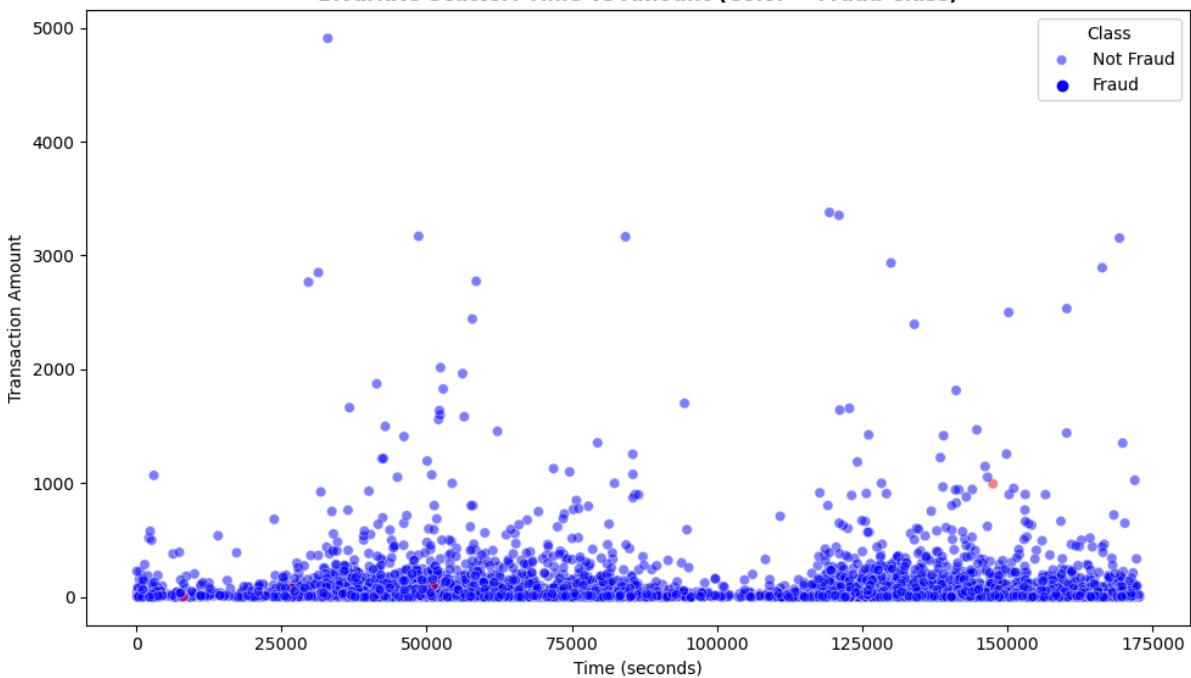
```

Output:

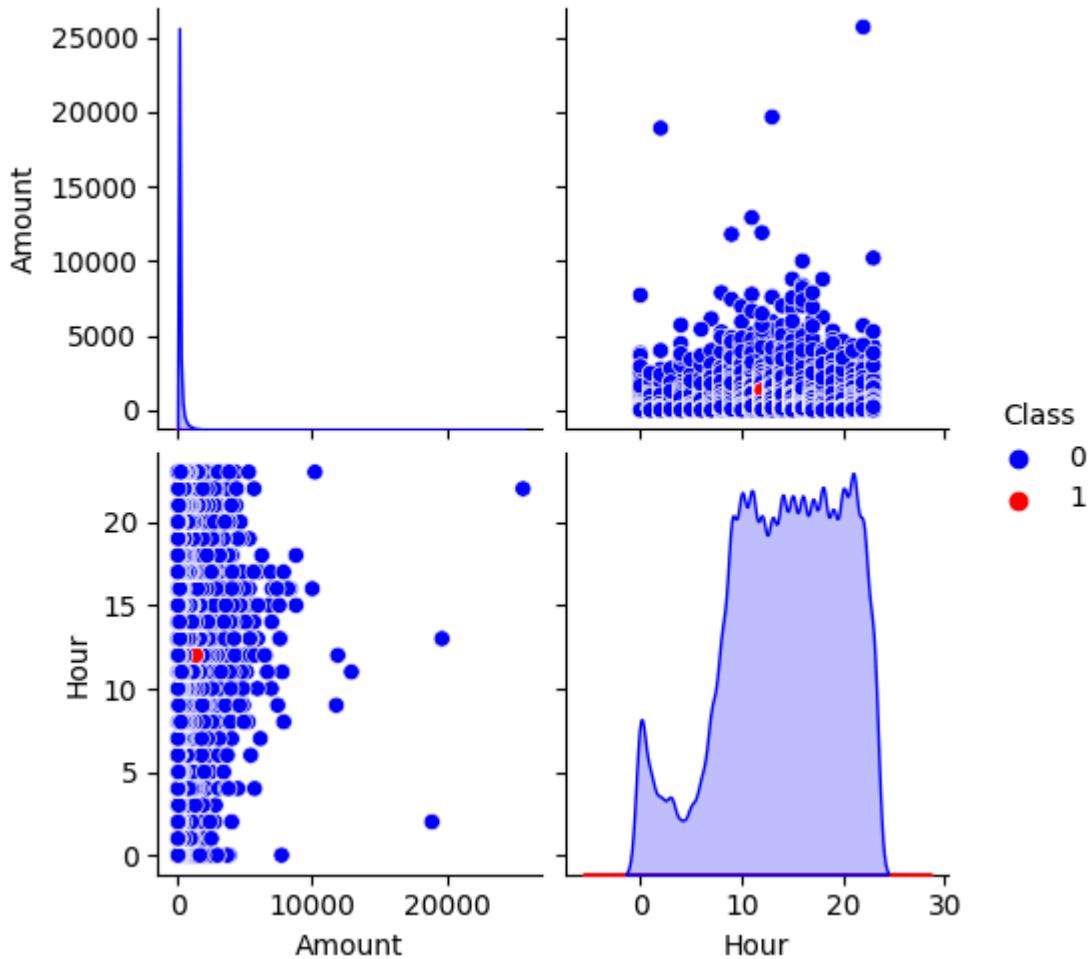
**Bivariate Analysis: Transaction Amount vs Fraud**



**Bivariate Scatter: Time vs Amount (Color = Fraud Class)**



## Bivariate Pairplot: Amount, Hour by Fraud Class



## Appendix B.8 (Distribution)

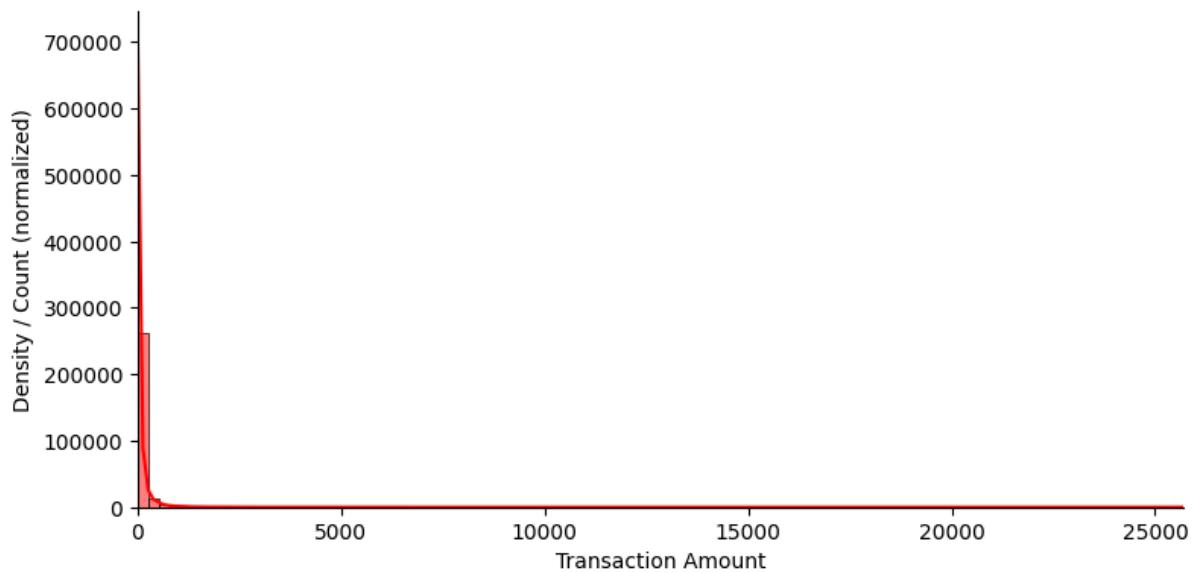
Code:

```
# Distribution of Transaction Amount
sns.displot(fraud_cc_df_dedup, x='Amount', kind='hist', bins=100, kde=True, color='r', height=4, aspect=2)
plt.title('Distribution of Transaction Amount', fontsize=14)
plt.xlim(fraud_cc_df_dedup['Amount'].min(), fraud_cc_df_dedup['Amount'].max())
plt.xlabel("Transaction Amount")
plt.ylabel("Density / Count (normalized)")
print("Distribution of the Transaction amount is right skewed distribution.")
print("Most of transactions are of small amount and highest distribution density is between 0 and 100")
plt.show()

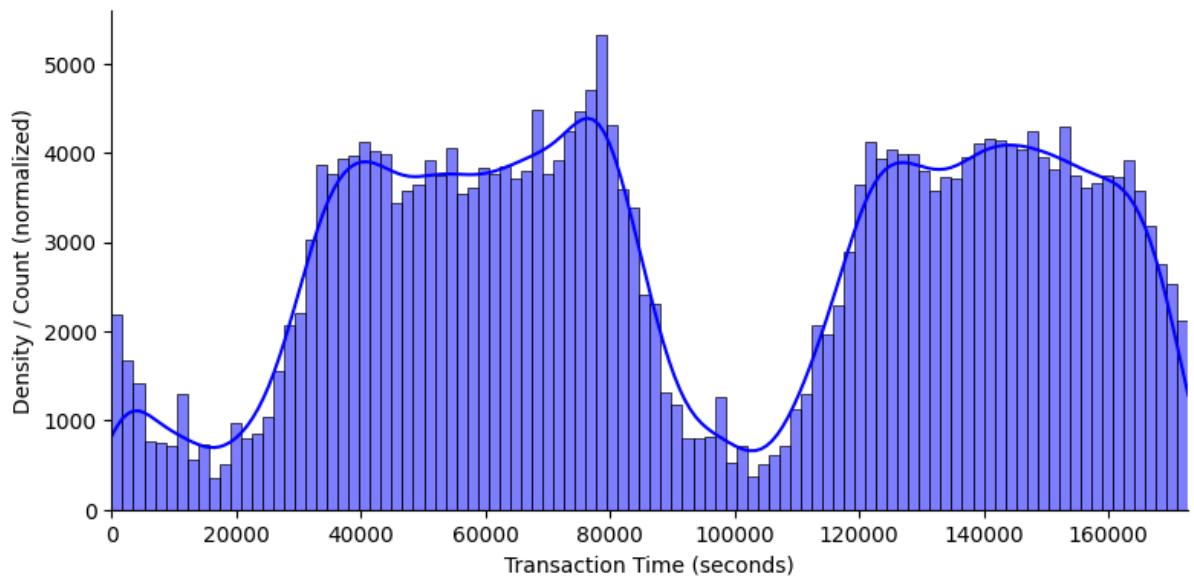
# Distribution of Transaction Time
sns.displot(fraud_cc_df_dedup, x='Time', kind='hist', bins=100, kde=True, color='b', height=4, aspect=2)
plt.title('Distribution of Transaction Time', fontsize=14)
plt.xlim(fraud_cc_df_dedup['Time'].min(), fraud_cc_df_dedup['Time'].max())
plt.xlabel("Transaction Time (seconds)")
plt.ylabel("Density / Count (normalized)")
print("Distribution of the transaction time is unevenly distributed over the time.")
print("More repeated transactions over specific time")
plt.show()
```

Output:

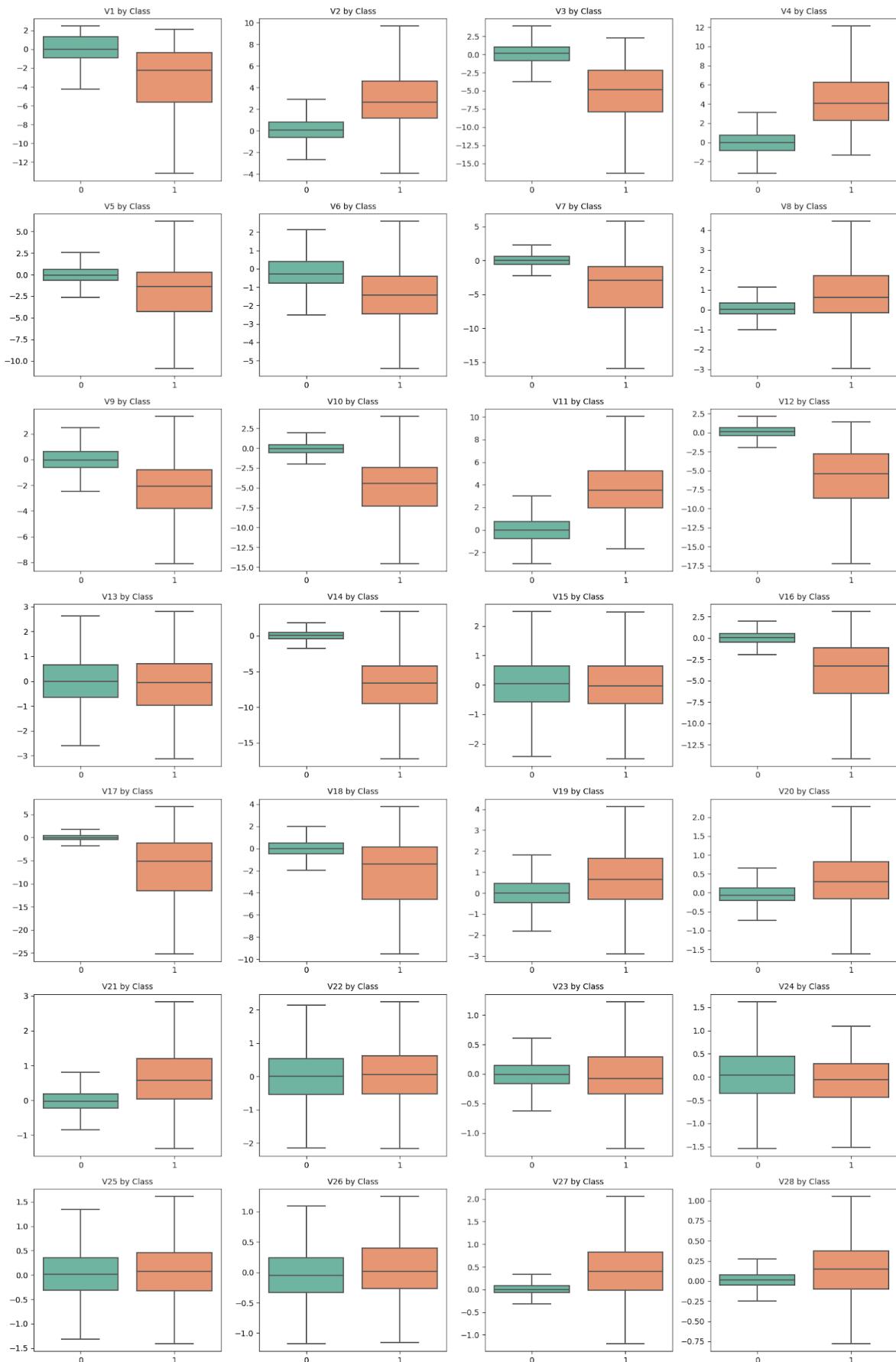
Distribution of Transaction Amount



Distribution of Transaction Time



**Boxplots of V1-V28 by Fraud Class**



## Appendix B.9 (Bivariate analysis: Correlation)

Code:

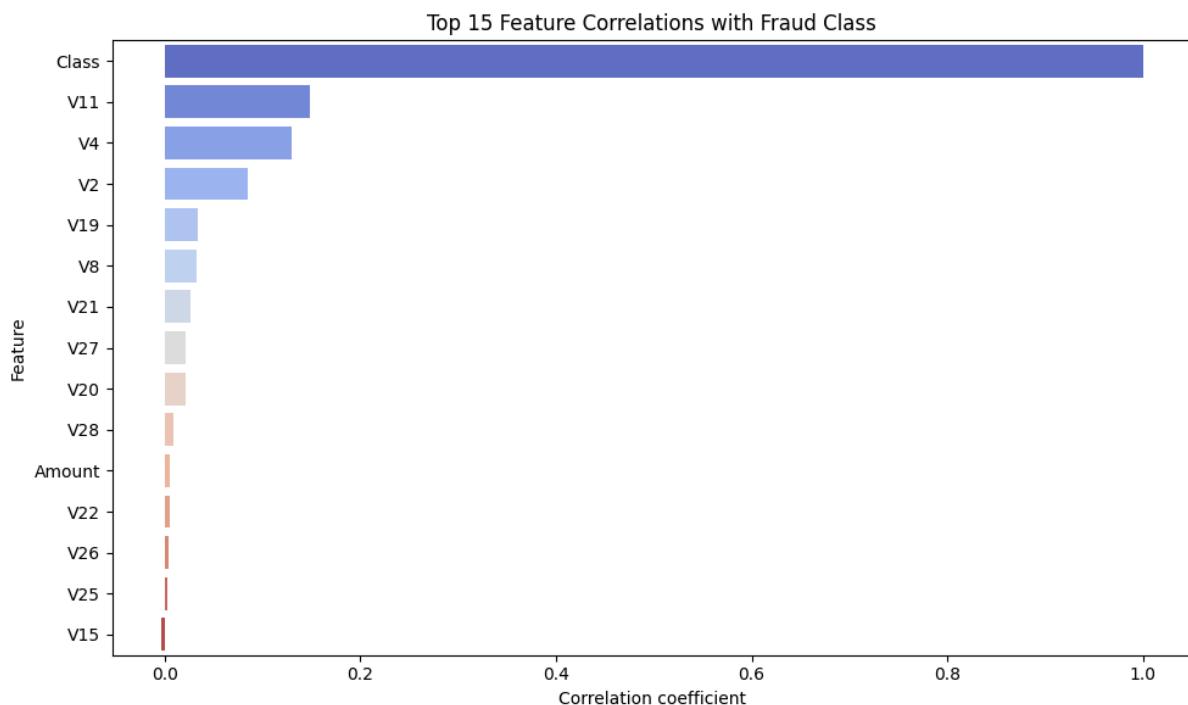
```
# Compute correlation matrix
corr_matrix = fraud_cc_df_dedup.corr(method='pearson') # Default is Pearson

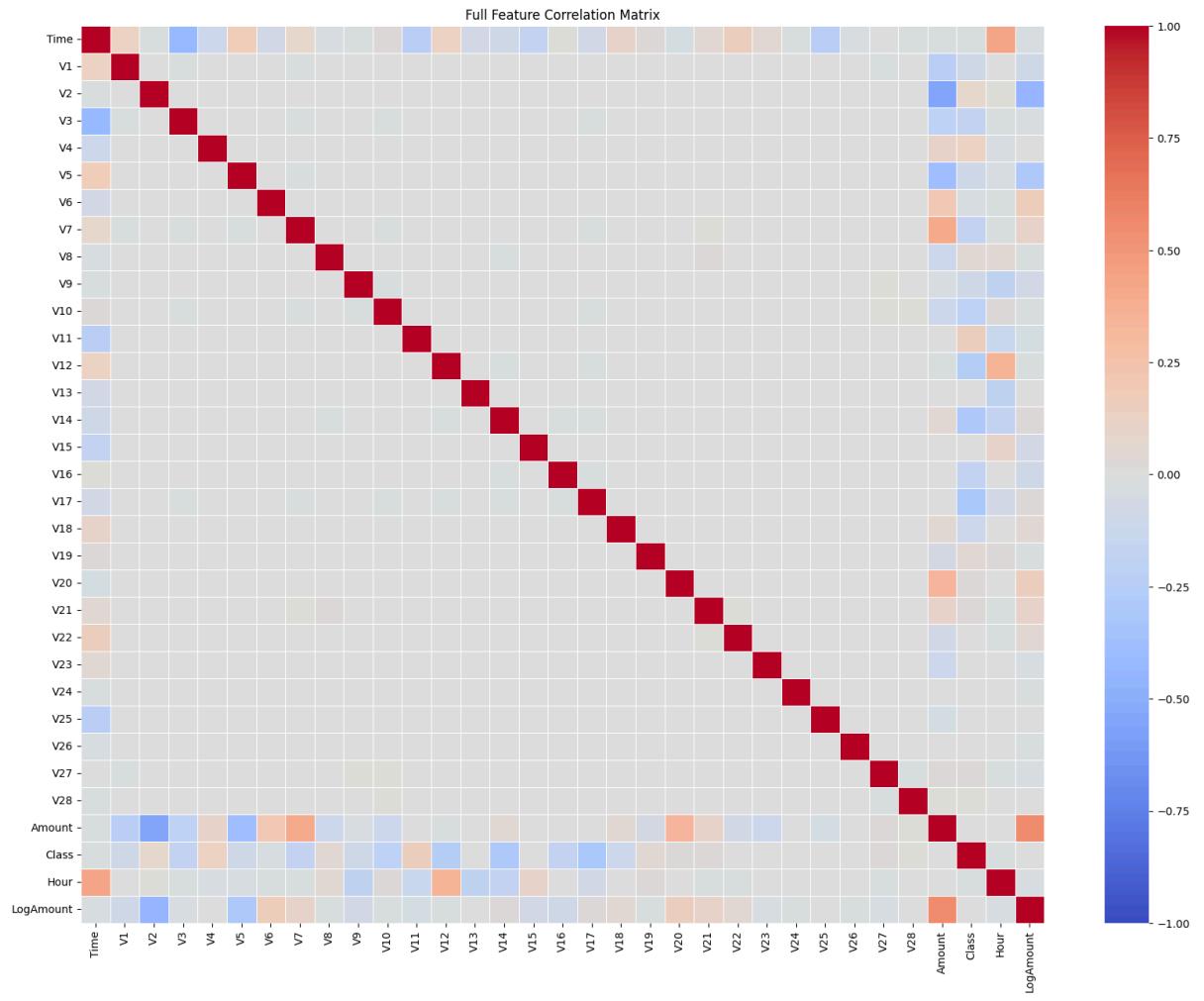
# Display full correlation with 'Class' (target)
corr_with_class = corr_matrix['Class'].sort_values(ascending=False)
print("Correlation of features with Fraud Class:\n")
print(corr_with_class)

# Visualize top correlations with Class
plt.figure(figsize=(10,6))
sns.barplot(x=corr_with_class.values[:15], y=corr_with_class.index[:15], palette='coolwarm')
plt.title("Top 15 Feature Correlations with Fraud Class")
plt.xlabel("Correlation coefficient")
plt.ylabel("Feature")
plt.tight_layout()
plt.show()

# Optional: Full heatmap of all features (can be large!)
plt.figure(figsize=(20, 15))
sns.heatmap(corr_matrix, cmap='coolwarm', vmin=-1, vmax=1, annot=False, linewidths=0.5)
plt.title("Full Feature Correlation Matrix")
plt.show()
```

Output:





## Appendix C (Modelling)

### Appendix C.1 (Train test split)

Code:

```
X = fraud_cc_df_dedup.drop('Class', axis=1)
y = fraud_cc_df_dedup['Class']
print(X.shape, y.shape)
```

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, stratify=y, test_size = 0.3, random_state = 42)
print(X_train.shape, y_train.shape, X_test.shape, y_test.shape)
```

```
# Print class distribution in training set
print("Class distribution in y_train:")
print(y_train.value_counts(normalize=True).apply(lambda x: f"{x:.2%}"))

# Print class distribution in test set
print("\nClass distribution in y_test:")
print(y_test.value_counts(normalize=True).apply(lambda x: f"{x:.2%}"))

.]
```

Output:

```
... (283726, 32) (283726,)
```

```
(198608, 32) (198608,) (85118, 32) (85118,)
```

```
Class distribution in y_train:
Class
0    99.83%
1    0.17%
Name: proportion, dtype: object
```

```
Class distribution in y_test:
Class
0    99.83%
1    0.17%
Name: proportion, dtype: object
```

## Appendix C.2 (Feature scaling):

Code:

```

from sklearn.preprocessing import StandardScaler

# # Creating function for scaling
# def Standard_Scaler (df, col_names):
#     features = df[col_names]
#     scaler = StandardScaler().fit(features.values)
#     features = scaler.transform(features.values)
#     df[col_names] = features

#     return df


scaler = StandardScaler().fit(X_train['Amount'].values.reshape(-1, 1))
X_train['Amount'] = scaler.transform(X_train['Amount'].values.reshape(-1, 1))
X_test['Amount'] = scaler.transform(X_test['Amount'].values.reshape(-1, 1))

# col_names = ['Amount']
# X_train = Standard_Scaler (X_train, col_names)
# X_test = Standard_Scaler (X_test, col_names)

```

```

# Check scaling of 'Amount' column in training data
print("Scaled statistics for X_train:")
print(X_train['Amount'].describe())

# Scaled statistics for X_train:
# count    1.986080e+05
# mean     1.001732e-17
# std      1.000003e+00
# min     -3.640437e-01
# 25%     -3.405495e-01
# 50%     -2.730345e-01
# 75%     -4.373915e-02
# max      7.757907e+01
# Name: Amount, dtype: float64

# Mean: 1.0017318839766175e-17
# Standard Deviation: 1.0000025175314597

# Optional: Check mean and std separately for clarity
print("\nMean:", X_train['Amount'].mean())
print("Standard Deviation:", X_train['Amount'].std())

```

Output:

```
.. Scaled statistics for X_train:  
count    1.986080e+05  
mean     1.710099e-17  
std      1.000003e+00  
min     -3.640437e-01  
25%     -3.405495e-01  
50%     -2.730345e-01  
75%     -4.373915e-02  
max      7.757907e+01  
Name: Amount, dtype: float64  
  
Mean: 1.71009943050294e-17  
Standard Deviation: 1.0000025175314597
```

### Appendix C.3 Logistic Regression:

Code:

```
# Add intercept for statsmodels  
X_train = X_train[feat_cols]  
X_test = X_test[feat_cols]  
X_train_sm = sm.add_constant(X_train)  
  
# Fit logistic regression using statsmodels  
logit_model = sm.Logit(y_train, X_train_sm)  
result = logit_model.fit()  
  
# Model summary  
print(result.summary())  
  
# Predict on test data  
X_test_sm = sm.add_constant(X_test)  
y_pred_prob = result.predict(X_test_sm)  
y_pred = (y_pred_prob > 0.5).astype(int)  
  
# Evaluation  
print("\nConfusion Matrix:")  
print(confusion_matrix(y_test, y_pred))  
  
print("\nClassification Report:")  
print(classification_report(y_test, y_pred))
```

Output:

```

Optimization terminated successfully.
    Current function value: 0.003688
    Iterations 13
                    Logit Regression Results
=====
Dep. Variable:                      Class      No. Observations:          198608
Model:                            Logit      Df Residuals:              198577
Method:                           MLE       Df Model:                  30
Date: Mon, 23 Jun 2025             Pseudo R-squ.:        0.7008
Time: 13:17:34                     Log-Likelihood:     -732.54
converged:                           True      LL-Null:            -2448.1
Covariance Type:           nonrobust      LLR p-value:        0.000
=====
                                         coef      std err      z      P>|z|      [0.025      0.975]
-----
const      -8.7422      0.313     -27.927      0.000     -9.356     -8.129
Time      -3.024e-07   2.76e-06     -0.110      0.913    -5.71e-06    5.1e-06
V1         0.0913      0.050      1.811      0.070     -0.007     0.190
V2         0.0136      0.069      0.197      0.844     -0.122     0.149
V3        -0.0171      0.063     -0.271      0.786     -0.141     0.106
V4         0.6840      0.090      7.606      0.000     0.508     0.860
V5         0.1513      0.079      1.910      0.056     -0.004     0.307
V6        -0.0752      0.084     -0.894      0.371     -0.240     0.090
V7        -0.0727      0.083     -0.876      0.381     -0.235     0.090
V8        -0.1752      0.036     -4.902      0.000     -0.245     -0.105
...
accuracy                         1.00      85118
macro avg      0.92      0.79      0.84      85118
weighted avg     1.00      1.00      1.00      85118

```

## Appendix C.4 Weight tuning:

Code:

```

from sklearn.metrics import classification_report, roc_auc_score

print(classification_report(y_test, y_pred))

# Predict on test data
y_train_pred_prob = result.predict(X_train_sm)
y_train_pred = (y_train_pred_prob > 0.5).astype(int)

print("Train ROC AUC Score:", roc_auc_score(y_train, y_train_pred_prob))
print("Test ROC AUC Score:", roc_auc_score(y_test, y_pred_prob))

```

5]

```

from sklearn.metrics import roc_curve, auc
import matplotlib.pyplot as plt

# Step 1: Compute False Positive Rate, True Positive Rate, and thresholds
fpr, tpr, thresholds = roc_curve(y_test, y_pred_prob)

# Step 2: Compute AUC (Area Under the ROC Curve)
roc_auc = auc(fpr, tpr)

# Step 3: Plot the ROC Curve
plt.figure(figsize=(8, 6))

# Plot the ROC curve for Logistic Regression
plt.plot(fpr, tpr, color='darkorange', lw=2,
         label=f'Logistic Regression ROC Curve (AUC = {roc_auc:.4f})')

# Plot the random guess line (baseline performance)
plt.plot([0, 1], [0, 1], color='blue', lw=2, linestyle='--',
         label='Logistic Regression (AUC = 0.5)')

# Step 4: Add titles and labels
plt.title('ROC Curve for Logistic Regression Model', fontsize=14, fontweight='bold')
plt.xlabel('False Positive Rate (FPR)', fontsize=12)
plt.ylabel('True Positive Rate (TPR)', fontsize=12)
plt.legend(loc="lower right")
plt.grid(True)

# Display the plot
plt.show()

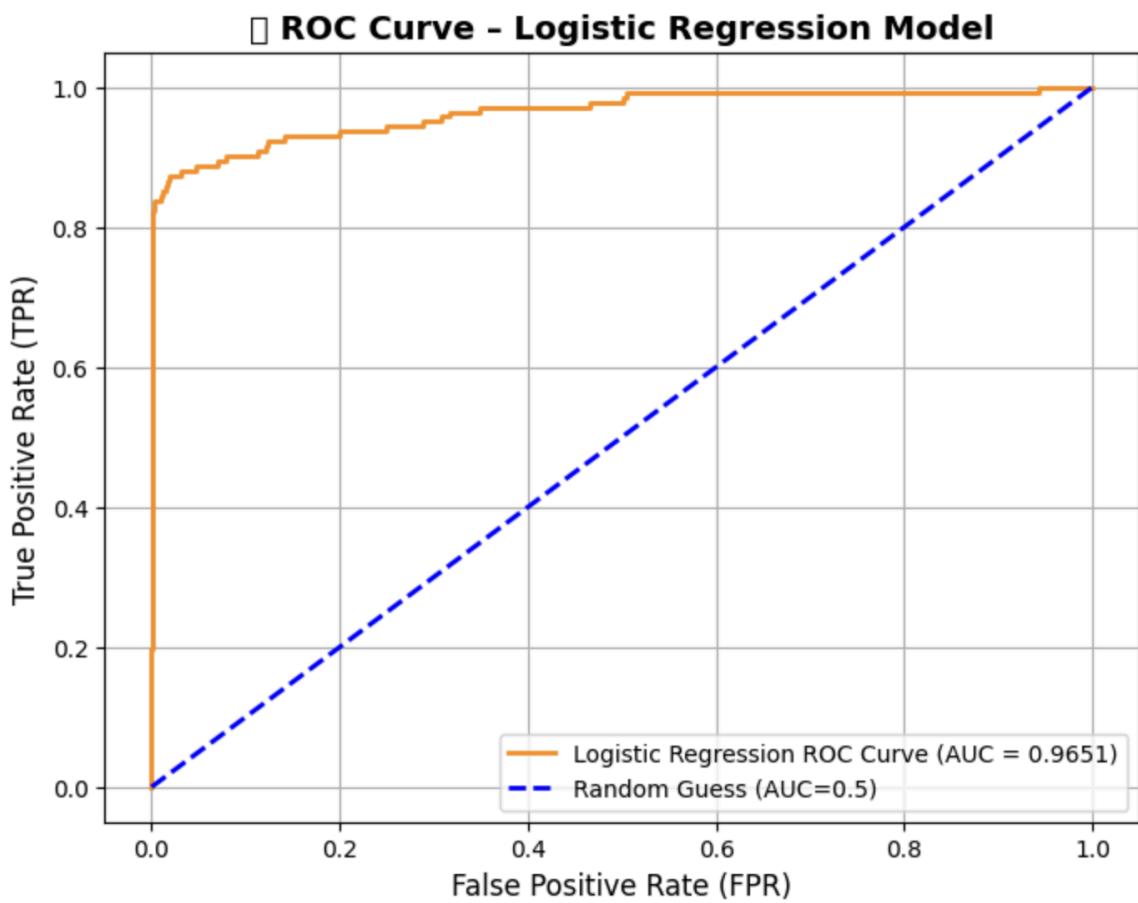
```

Output:

	precision	recall	f1-score	support
0	1.00	1.00	1.00	84976
1	0.84	0.57	0.68	142
accuracy			1.00	85118
macro avg	0.92	0.79	0.84	85118
weighted avg	1.00	1.00	1.00	85118

Train ROC AUC Score: 0.9813309790735404

Test ROC AUC Score: 0.9651150880049645



### Appendix C.5 Random Forest with Hyperopt

Code:

```

# Define search space
space = {
    'criterion': hp.choice('criterion', ['entropy', 'gini']),
    'max_depth': hp.quniform('max_depth', 1, 5, 1),
    'max_features': hp.choice('max_features', ['auto', 'sqrt', 'log2', 'None']),
    'min_samples_leaf': hp.uniform('min_samples_leaf', 0, 0.5),
    'min_samples_split': hp.uniform('min_samples_split', 0, 1),
    'n_estimators': hp.choice('n_estimators', [10, 50])
}

# Objective function to maximize AUC
def objective(space):
    model = RandomForestClassifier(
        criterion=space['criterion'],
        max_depth=int(space['max_depth']),
        max_features=space['max_features'],
        min_samples_leaf=space['min_samples_leaf'],
        min_samples_split=space['min_samples_split'],
        n_estimators=space['n_estimators'],
        random_state=42,
        n_jobs=-1
    )

    auc_scores = cross_val_score(model, X_train, y_train, cv=4, scoring='roc_auc')
    auc = auc_scores.mean()

    return {'loss': 1 - auc, 'status': STATUS_OK}

# Write to the csv file ('a' means append)
of_connection = open(out_file, 'a')
writer = csv.writer(of_connection)
writer.writerow([auc_train, auc_test, loss, rf_params, ITERATION, n_estimators, run_time, fea_90_cnt])

# Dictionary with information for evaluation
return {'loss': loss,
        'params': rf_params,
        'iteration': ITERATION,
        'estimators': n_estimators,
        'train_time': run_time,
        'status': STATUS_OK}

```

## Output:

```

Params selected:
{'booster': 'gbtree', 'early_stopping_rounds': 10, 'eval_metric': 'auc', 'gamma': 25.0, 'importance_type': 'gain', 'learning_rate': 0.3500000000000003, 'max_depth': 4, 'min_child_weight': 35
[0] validation_0-auc:0.91693 validation_1-auc:0.90959
[1] validation_0-auc:0.93349 validation_1-auc:0.92056
[2] validation_0-auc:0.94179 validation_1-auc:0.92310
[3] validation_0-auc:0.94303 validation_1-auc:0.92279
[4] validation_0-auc:0.94333 validation_1-auc:0.92753
[5] validation_0-auc:0.94651 validation_1-auc:0.93274
[6] validation_0-auc:0.96929 validation_1-auc:0.94444
[7] validation_0-auc:0.97050 validation_1-auc:0.95125
[8] validation_0-auc:0.96979 validation_1-auc:0.95365
[9] validation_0-auc:0.97373 validation_1-auc:0.96249
[10] validation_0-auc:0.97616 validation_1-auc:0.96463
[11] validation_0-auc:0.97873 validation_1-auc:0.96414
[12] validation_0-auc:0.97873 validation_1-auc:0.96414
[13] validation_0-auc:0.97873 validation_1-auc:0.96414
[14] validation_0-auc:0.97873 validation_1-auc:0.96414
[15] validation_0-auc:0.97873 validation_1-auc:0.96414
[16] validation_0-auc:0.97873 validation_1-auc:0.96414
[17] validation_0-auc:0.97873 validation_1-auc:0.96414
[18] validation_0-auc:0.97873 validation_1-auc:0.96414
[19] validation_0-auc:0.97873 validation_1-auc:0.96414
30
[0.      0.      0.      0.      0.24241856 0.
...
0.20063178 0.      0.2289992 0.      0.      0.
0.      0.      0.      0.      0.      0.      ]
100%|██████████| 10/10 [00:26<00:00,  2.64s/trial, best loss: 0.03536561938946803]

```

## Appendix C.5 Random Forest with Single model

```

from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import roc_auc_score
import pandas as pd

param = {'criterion': 0 ,
         'max_depth': 2.0,
         'max_features': 1,
         'min_samples_leaf': 0.13489892897095584,
         'min_samples_split': 0.13145061301665495, 'n_estimators': 1}

# Convert index-based choices to real values
param['criterion'] = ['entropy', 'gini'][param['criterion']]
param['max_features'] = ['auto', 'sqrt', 'log2', None][param['max_features']]
param['n_estimators'] = [10, 50][param['n_estimators']]
param['max_depth'] = int(param['max_depth'])

print(param)

# Fit final model
RF_model = RandomForestClassifier(**param)
RF_model.fit(X_train[feat_cols], y_train)

# Score the sets
X_train["Scored_values"] = RF_model.predict_proba(X_train[feat_cols].values)[:, 1]
X_test["Scored_values"] = RF_model.predict_proba(X_test[feat_cols].values)[:, 1]

# Evaluate
print("Dev AUC", roc_auc_score(y_true=y_train, y_score=X_train["Scored_values"]))
print("Test AUC", roc_auc_score(y_true=y_test, y_score=X_test["Scored_values"]))

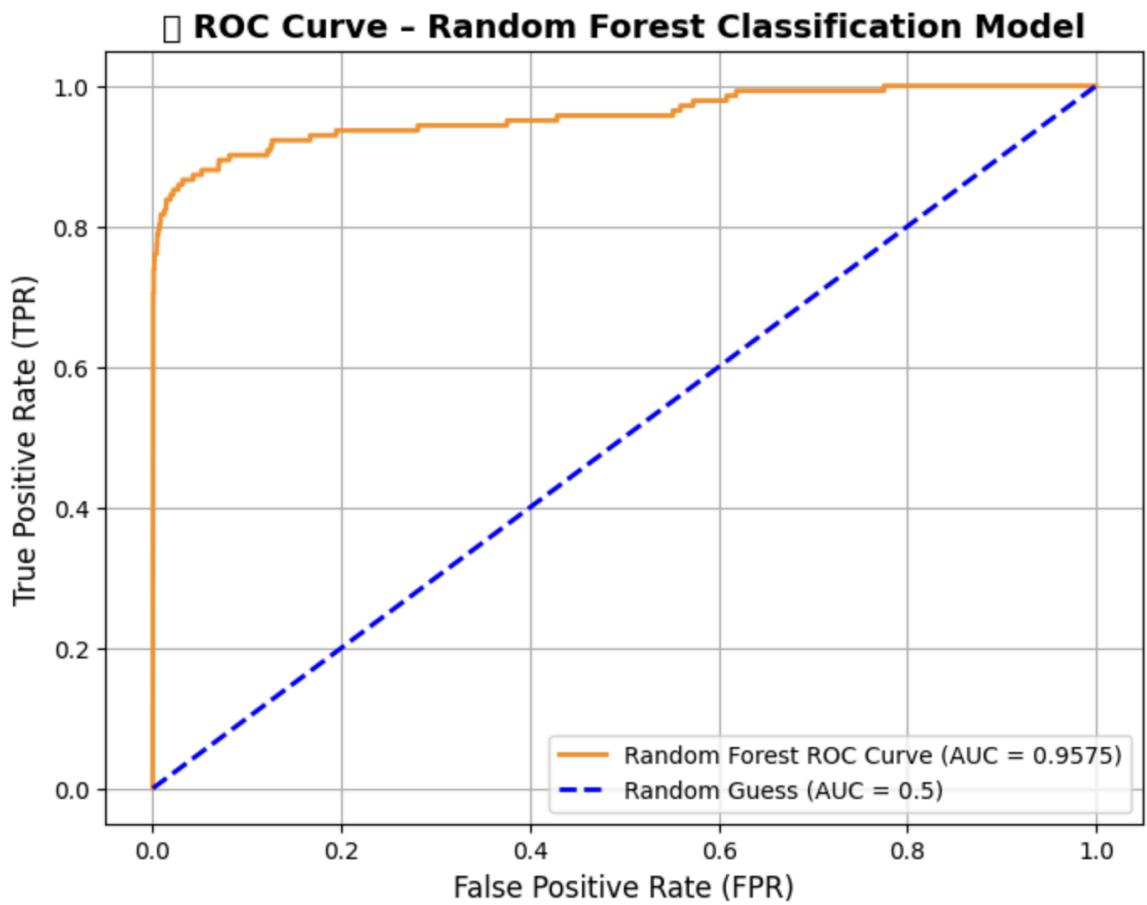
```

Output:

```

{'criterion': 'entropy', 'max_depth': 2, 'max_features': 'sqrt', 'min_samples_leaf': 0.13489892897095584, 'min_samples_split': 0.13145061301665495, 'n_estimators': 50}
Dev AUC 0.9759302752731396
Test AUC 0.9497012495325938

```



### Appendix C.5 XGBoost with Hyperopt

Code:

```

model = XGBClassifier(**xgb_params)

model.fit(
X_train[feat_cols],
y_train,
eval_set=[(X_train[feat_cols], y_train), (X_test[feat_cols], y_test)],
verbose=True
)

run_time = timer() - start

auc_test = max(model.evals_result()['validation_1']['auc'])
index = model.evals_result()['validation_1']['auc'].index(max(model.evals_result()['validation_1']['auc']))
auc_train = model.evals_result()['validation_0']['auc'][index]

n_estimators = model.best_iteration + 1 #same as index+1

print(len(X_train[feat_cols].columns), len(model.feature_importances_), model.feature_importances_)
# Creating feature importance dataframe
feat_df = pd.DataFrame({'Features': X_train[feat_cols].columns, 'Importance': model.feature_importances_})
feat_df.sort_values('Importance', ascending=False, inplace=True)
feat_df['cumul_fea_imp'] = feat_df['Importance'].cumsum()
final_features_temp = feat_df[feat_df['cumul_fea_imp']<=0.95]['Features'].tolist()
fea_95_cnt = len(final_features_temp)
# feat_df.to_csv("feat_df.csv")
# print("fea_90_cnt:", fea_90_cnt)

# Extract the best score
best_score = auc_test

# Loss must be minimized
loss = 1 - best_score

# Write to the csv file ('a' means append)
of_connection = open(out_file, 'a')
writer = csv.writer(of_connection)
writer.writerow([auc_train,auc_test,loss, xgb_params, ITERATION,n_estimators,run_time,fea_95_cnt])

# Dictionary with information for evaluation
return {'loss': loss,
        'params': xgb_params,
        'iteration': ITERATION,
        'estimators': n_estimators,
        'train_time': run_time,
}

```

## Appendix C.6 XGBoost with Single model

```

from sklearn.metrics import roc_auc_score

basic_params = {'n_estimators':50,
                'max_depth':2}

xgb_model,xgb_feat_imp = build_single_model(X_train, y_train, X_test, y_test,
                                              feat_cols = feat_cols,
                                              alg = XGBClassifier,
                                              param= basic_params
                                              )

xgb_feat_imp.sort_values('Importance', ascending=False, inplace=True)
xgb_feat_imp['cumul_fea_imp'] = xgb_feat_imp['Importance'].cumsum()

]

```

Output:

```

[49]      validation_0-logloss:0.00190      validation_1-logloss:0.00311
Scoring the base
Dev AUC 0.9978715348741493
Test AUC 0.9705019859791397

```

## Appendix C.7 ROC final XGBoost model

Code:

```

from sklearn.metrics import roc_curve, auc
import matplotlib.pyplot as plt

# Step 1: Compute False Positive Rate, True Positive Rate, and thresholds
fpr, tpr, thresholds = roc_curve(y_test, xgb_model.predict_proba(X_test[feat_cols])[:,1])

# Step 2: Compute AUC (Area Under the ROC Curve)
roc_auc = auc(fpr, tpr)

# Step 3: Plot the ROC Curve
plt.figure(figsize=(8, 6))

# Plot the ROC curve for Logistic Regression
plt.plot(fpr, tpr, color='darkorange', lw=2,
         label=f'Random Forest ROC Curve (AUC = {roc_auc:.4f})')

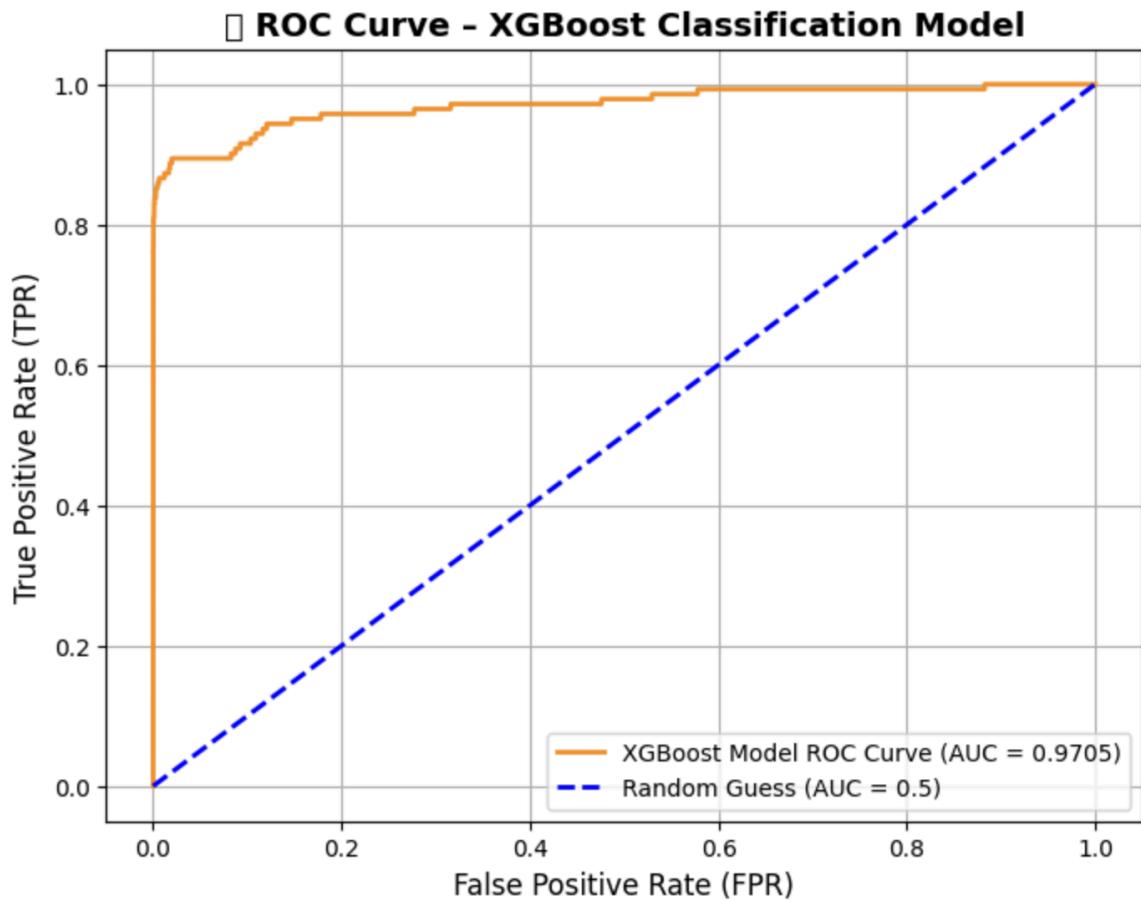
# Plot the random guess line (baseline performance)
plt.plot([0, 1], [0, 1], color='blue', lw=2, linestyle='--',
         label='XGBoost Model (AUC = 0.5)')

# Step 4: Add titles and labels
plt.title('ROC Curve XGBoost Classification Model', fontsize=14, fontweight='bold')
plt.xlabel('False Positive Rate (FPR)', fontsize=12)
plt.ylabel('True Positive Rate (TPR)', fontsize=12)
plt.legend(loc="lower right")
plt.grid(True)

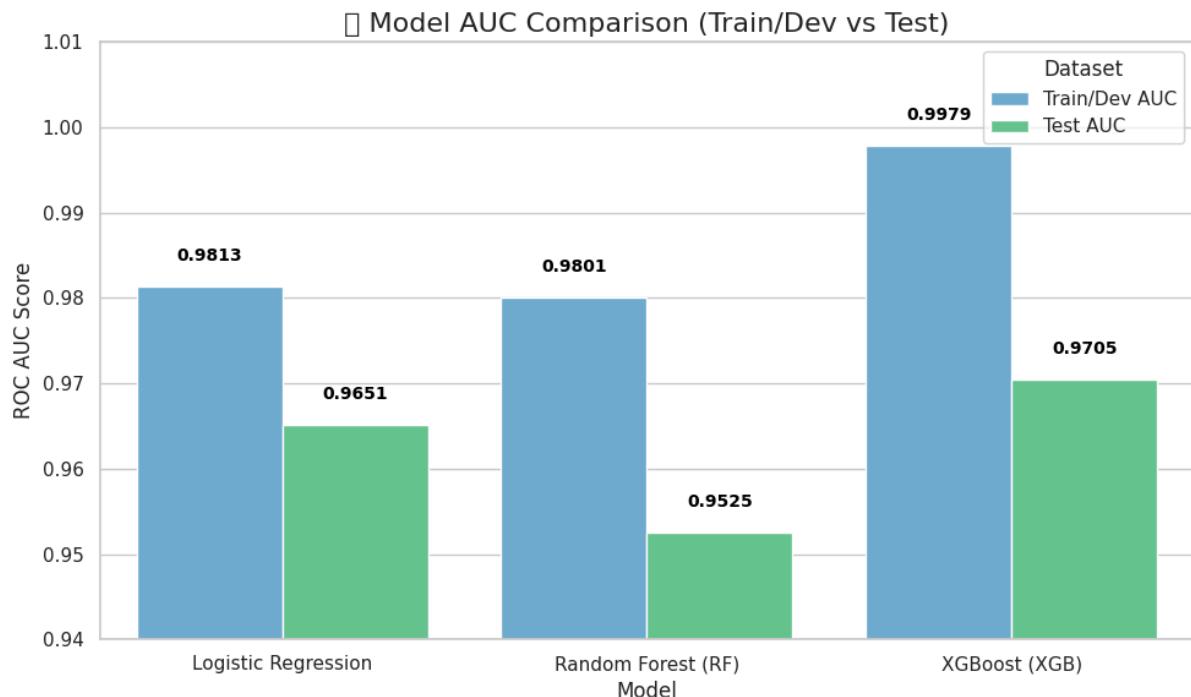
# Display the plot
plt.show()

```

Output:



#### Appendix C.8 Comparison of Model performance



## References

- **Dataset reference :** Kaggle (2015). *Credit card fraud detection dataset* [Data set]. ULB Machine Learning Group.  
<https://www.kaggle.com/datasets/mlg-ulb/creditcardfraud>
- **Github repo :** Dwivedi, V., Malla, Manu & TS, B., *Credit Card Fraud Detection* [Source code]. GitHub. <https://github.com/vedpd/Credit-Card-Fraud-Detection>
- **Breiman, L.** (2001). Random forests. *Machine Learning*, 45(1), 5–32.  
<https://doi.org/10.1023/A:1010933404324>
- **Chen, T., & Guestrin, C.** (2016). XGBoost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 785–794).  
<https://doi.org/10.1145/2939672.2939785> “Core paper introducing XGBoost and its efficient implementation”
- **Bergstra, J., Yamins, D., & Cox, D. D.** (2013). Making a science of model search: Hyperparameter optimization in hundreds of dimensions for vision architectures. In *Proceedings of the 30th International Conference on Machine Learning (ICML-13)* (pp. 115–123). “Introduced **Hyperopt**, a tool for automatic hyperparameter tuning.”
- **Sokolova, M., & Lapalme, G.** (2009). A systematic analysis of performance measures for classification tasks. *Information Processing & Management*, 45(4), 427–437.  
<https://doi.org/10.1016/j.ipm.2009.03.002> “Detailed explanation of accuracy, precision, recall, F1-score, and more.”
- **Powers, D. M.** (2011). Evaluation: From precision, recall and F-measure to ROC, informedness, markedness & correlation. *Journal of Machine Learning Technologies*, 2(1), 37–63. “Discusses the limitations of standard metrics and suggests improvements.”
- **Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P.** (2002). SMOTE: Synthetic Minority Over-sampling Technique. *Journal of Artificial Intelligence Research*, 16, 321–357.  
<https://doi.org/10.1613/jair.953> “Widely cited technique to balance imbalanced datasets.”

- **Han, J., Pei, J., & Kamber, M.** (2011). *Data mining: Concepts and techniques* (3rd ed.). Morgan Kaufmann. “Covers data transformation techniques such as scaling and normalization.”
- **Zheng, A., & Casari, A.** (2018). *Feature engineering for machine learning: Principles and techniques for data scientists*. O’Reilly Media. “Focuses on practical feature preparation for improving ML models.”
- **Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Duchesnay, É.** (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12, 2825–2830. “Reference for using Python's scikit-learn library for scaling, modeling, and evaluation.”