# SYNTHETIC IDENTITY FRAUD DETECTION

Final Project Report

**PROJECT TYPE:** FINAL PROJECT REPORT
**SUBMITTED BY:** SAMYUKTH LALITH LELLA GOPAL
**STUDENT ID:** 9005574
**INSTRUCTOR:** *DAVID ESPINOSA*
**INSTITUTION:** *CONESTOGA COLLEGE*
**SUBMISSION DATE:** APRIL 2025

| A3 No. and Name | Team members (name & role) | Stakeholders (role & department) | **Company objective** |
|---|---|---|---|
| **Group 9** | | | AI-Driven Synthetic Identity Detection |

**A3 No. and Name**
**Group 9**

**Team Leader (name & 'phone ext)**

**Team members (name & role)**
1. Lella Gopal, Samyukth Lalith

**Stakeholders (role & department)**
1. AI & ML Coordinator, Conestoga College
2. Potential Client name(s)
3.
4.

**Company objective**
AI-Driven Synthetic Identity Detection

**Start date & planned duration**
Q1 2024 – Q4 2024

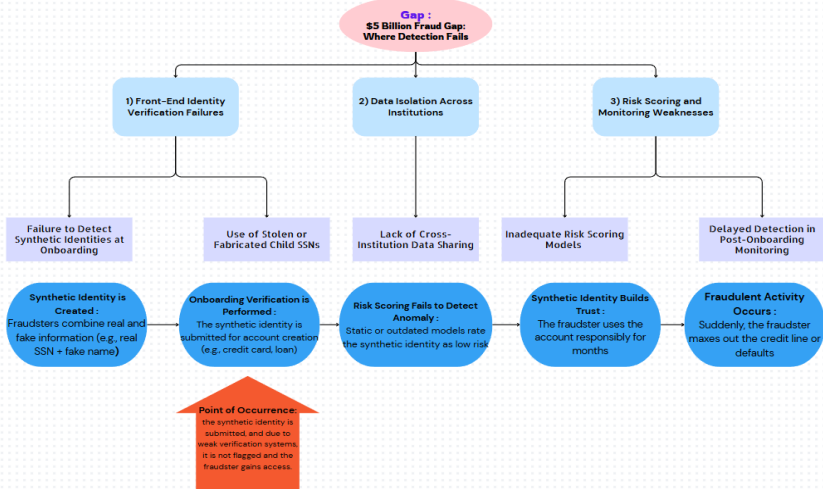# Synthetic Identity Fraud Detection

## 1. Clarify the problem

**Problem Focus:** High financial losses caused by synthetic identity fraud.

**Ideal Situation:** Banks and financial institutions should have strong systems that detect and stop fake identities. Ideally, yearly losses from this type of fraud should be under $1 billion.

**Current Situation:** Right now, synthetic identity fraud is causing around $6 billion in losses every year for lenders.

**Gap:** The goal is to keep fraud losses below $1 billion, but actual losses have reached around $6 billion. This shows a serious problem in detecting and stopping synthetic identity fraud effectively.

## 2. Breakdown the problem

**Gap :**
**$5 Billion Fraud Gap:**
**Where Detection Fails**

1) Front–End Identity Verification Failures
2) Data Isolation Across Institutions
3) Risk Scoring and Monitoring Weaknesses

- Failure to Detect Synthetic Identities at Onboarding
- Use of Stolen or Fabricated Child SSNs
- Lack of Cross-Institution Data Sharing
- Inadequate Risk Scoring Models
- Delayed Detection in Post-Onboarding Monitoring

**Synthetic Identity is Created :** Fraudsters combine real and fake information (e.g. real SSN + fake name)

**Onboarding Verification is Performed :** The synthetic identity is submitted for account creation (e.g. credit card, loan)

**Risk Scoring Fails to Detect Anomaly :** Static or outdated models rate the synthetic identity as low risk

**Synthetic Identity Builds Trust :** The fraudster uses the account responsibly for months

**Fraudulent Activity Occurs :** Suddenly, the fraudster maxes out the credit line or defaults

**Point of Occurrence:** the synthetic identity is submitted, and due to weak verification systems, it is not flagged and the fraudster gains access.

## 3. Set the Target

| 1 | Reduce losses from $6B to under $3B in 2 years |
|---|---|
| 2 | Increase fraud detection accuracy from ~5–15% to ≥90% |
| 3 | Deploy at least 1 machine learning model into production within 6 months |
| 4 | Implement real-time scoring to evaluate identity risk within 10 seconds per application |

## 4. Analyse the Root Cause

**Root Cause 1 : Front-End Identity Verification Failures**
**1) Why are synthetic identities not detected during onboarding?**
Because identity verification tools don't validate information in real-time against trusted databases like credit bureaus or government sources.
**2) Why are stolen or child SSNs being used successfully?**
Because there's no system in place to flag dormant or underage SSNs when used for applications.

**Root Cause 2 : Risk Scoring and Monitoring Weaknesses**
**1) Why is the system failing to detect risk accurately?**
Because it relies on outdated, rule-based models that cannot adapt to new fraud patterns or evolving behavior.
**2) Why is fraudulent behavior detected only after damage occurs?**
Because post-onboarding monitoring is limited and lacks time-series analysis or behavioral profiling.

**Root Cause 3 : Data Isolation Across Institutions**
**1) Why can't institutions detect repeated fraud across platforms?**
Because fraud intelligence is not shared across institutions — systems are siloed, and there's no collaborative fraud detection network.

**Root Cause:** The core issue lies in the use of outdated, non-integrated identity verification systems that cannot validate data in real-time, allowing synthetic identities to pass through onboarding undetected.

## 5. Develop Countermeasures

| | Criteria > | Prevents Financial Loss | Protects Customer Trust | Improves Detection Accuracy | Reduces False Approvals | Enables Real-Time Scoring | Overall Score (/100) | Ranking | Potential Problems: |
|---|---|---|---|---|---|---|---|---|---|
| 1 | Real-time Identity Verification with External Databases | 25 | 20 | 20 | 20 | 10 | 95 | 1 | None |
| 2 | ML-Based Risk Scoring Model ( XGBoost) | 22 | 20 | 25 | 20 | 8 | 95 | 1 | Requires regular data quality checks & model retraining |
| 3 | Age-based SSN Validation (Child SSN Check) | 20 | 15 | 15 | 15 | 5 | 70 | 3 | Limited to specific regions; needs SSA integration |
| 4 | Behavior Monitoring with Anomaly Detection (post-onboarding) | 18 | 20 | 18 | 12 | 7 | 75 | 2 | May have detection lag; best for layered protection |
| 5 | Cross-Institution Fraud Data Sharing API | 15 | 15 | 18 | 10 | 7 | 65 | 4 | Regulatory/privacy constraints may limit implementation |

## 6. Implement Countermeasure

## 7. Monitor Results & Process

## 8. Standardize & Share Success

https://github.com/Samyukth107/The-Final-Presentation

https://www.kaggle.com/datasets/ealaxi/paysim1

Couldn't upload the dataset to the git due to the size

## Abstract:

Synthetic identity fraud has emerged as the fastest-growing financial crime in the United States, responsible for an estimated $6 billion in losses and accounting for 20% of credit losses in 2016. Traditional detection systems have proven inadequate, failing to flag 85% to 95% of synthetic identity applications. This paper presents a structured A3-based approach to mitigate synthetic identity fraud by implementing machine learning models and process improvements. By addressing root causes and deploying targeted solutions, the proposed strategy aims to reduce fraud losses by $3 billion over two years.

## Problem Statement:

Current identity verification systems are outdated and ineffective against sophisticated synthetic identity fraud tactics. Studies indicate that 85% to 95% of synthetic identity applications are not flagged by traditional fraud models, resulting in substantial financial losses and undermining the integrity of financial institutions. There is an urgent need for accurate, real-time fraud detection systems powered by machine learning to address this growing threat.

### Ideal State:

- Achieving a detection and prevention rate of 100% for synthetic identity attempts during the onboarding process.
- Reducing annual synthetic fraud losses to below $1 billion.
- Processing and responding to high-risk identity applications within 10 seconds using real-time scoring mechanisms.

## Research hypothesis:

### 1. Introduction

Synthetic identity fraud involves the creation of fictitious identities by combining real and fake personal information, such as pairing a legitimate Social Security Number (SSN) with a fabricated name. These synthetic identities are used to establish credit lines, build credit histories, and eventually default on substantial loans, leading to significant financial losses for institutions. In 2016, synthetic identity fraud cost U.S. lenders approximately $6 billion and was responsible for 20% of credit losses. This paper explores a comprehensive approach to detect and prevent synthetic identity fraud through advanced machine learning techniques and systemic process enhancements.

### 2. Literature Review

Research from the Federal Reserve and industry experts highlights the severity and rapid growth of synthetic identity fraud. Traditional rule-based detection systems have failed to keep pace with evolving fraud tactics, with synthetic identity fraud accounting for 10% to 15% of charge-offs in unsecured lending portfolios. Machine learning models, such as Random Forest and XG Boost, have demonstrated potential in identifying hidden fraud patterns and reducing false positives. Implementing a layered defense that combines data verification, behavioral analysis, and cross-institution intelligence sharing is recommended for effective mitigation.

### 3. Targets

| No. | Target Statement | Timeframe | Success Indicator |
|-----|------------------|-----------|-------------------|
| 1 | **Reduce losses from $6B to under $3B** | 2 years | **≥ $3B loss reduction** |
| 2 | **Increase fraud detection accuracy from ~5–15% to ≥90%** | Within 1 year | **≥ 90% model precision & recall** |

| 3 | **Deploy at least 1 machine learning model into production** | Within 6 months | **Functional model integrated into system** |
|---|---|---|---|
| 4 | **Implement real-time scoring to evaluate identity risk within 10 seconds per application** | Within 6 months | **Risk score available via API in ≤10 seconds** |

## 4. Root Cause Analysis

**1) Root Cause 1:** Front-End Identity Verification Failures

**(i) Why are synthetic identities not detected during onboarding?**

Because identity verification tools don't validate information in real-time against trusted databases like credit bureaus or government sources.

**(ii) Why are stolen or child SSNs being used successfully?**

Because there's no system in place to flag dormant or underage SSNs when used for applications.

**Root Cause Identified:** Outdated, non-integrated verification systems that allow fake identities to pass onboarding.

**2) Root Cause 2:** Risk Scoring and Monitoring Weaknesses

**(i)Why is the system failing to detect risk accurately?**

Because it relies on outdated, rule-based models that cannot adapt to new fraud patterns or evolving behavior.

**(ii) Why is fraudulent behavior detected only after damage occurs?**

Because post-onboarding monitoring is limited and lacks time-series analysis or behavioral profiling.

**Root Cause Identified:** Static risk scoring models and lack of behavioral monitoring systems.

**3) Root Cause 3:** Data Isolation Across Institutions

**(i)Why can't institutions detect repeated fraud across platforms?**

Because fraud intelligence is not shared across institutions — systems are solved, and there's no collaborative fraud detection network.

**Root Cause Identified:** No fraud data-sharing mechanism between institutions.

## 5. Methodology

This project follows a structured machine learning workflow to detect synthetic identity fraud using the PaySim transaction dataset. The process begins with preparing the dataset by removing non-essential columns and encoding categorical features such as transaction type. Exploratory Data Analysis (EDA) is performed to identify patterns in fraudulent activity, including the distribution across transaction types and abnormal balance behaviors. To ensure feature relevance, statistical tests like the T-Test, Mann-Whitney U Test, and Chi-Square Test are applied, confirming significant differences between fraud and non-fraud transactions. The issue of extreme class imbalance is addressed by applying class weighting in the Random Forest model and adjusting the scale_pos_weight parameter for XGBoost. Both models are trained using a stratified 75/25 train-test split to maintain class proportion. Random Forest is optimized for precision to reduce false positives, while XGBoost is tuned for high recall to capture as many fraud cases as possible. Model performance is evaluated using metrics such as accuracy, precision, recall, F1 score, and ROC AUC, supported by confusion matrix and classification report visualizations. A real-time scoring simulation is conducted to

assess latency, ensuring predictions are generated in under 10 seconds. Finally, a fraud loss estimation is performed to calculate the amount of loss each model can potentially prevent, helping assess the financial effectiveness of the proposed solution.

## 6. Proposed Solution's

This project proposes a two-layered machine learning-based solution to detect and reduce synthetic identity fraud in financial transactions. The approach combines data-driven classification models with real-time scoring capabilities to identify suspicious activities during and after onboarding.

**1)Machine Learning-Based Fraud Classifier**

**A Random Forest and XGBoost model were trained using transactional features such as amount, transaction type, and balance behaviors.**

**Random Forest was optimized for precision, minimizing false alarms.**

**XGBoost was tuned to maximize recall and detect nearly all fraud cases.**

**2)Real-Time Scoring System**

**Both models demonstrated real-time prediction speeds (< 35 ms) and are suitable for fraud detection at onboarding or transaction checkpoints.**

• **Loss Prevention Estimation**

**The XGBoost model helped identify and prevent the majority of fraud loss in the test set, achieving ~100% loss reduction in simulation.**

**3)Balanced Decision Strategy**

**The system allows institutions to choose between a low-false-alarm model (Random Forest) or a high-capture model (XGBoost) depending on operational priorities.**

**Outcome**

By combining these components, the proposed solution directly supports the targets of:

- Reducing fraud losses by $3 billion in 2 years

- Achieving ≥90% detection accuracy

- Deploying a real-time scoring system

- Establishing data partnerships to extend fraud detection coverage beyond a single organization

## 7. Implementation Plan

| Phase | Timeline | Objective | Activities | Expected Outcome |
|-------|----------|-----------|------------|------------------|
| Phase 1: Planning & Design | Year 1 – Q1 | Define technical requirements and solution design | • Identify key data sources<br>• Perform data cleaning & exploratory analysis<br>• Define fraud features and labeling approach<br>• Outline ML model framework | Clear fraud feature map and initial data pipeline ready for modeling |
| Phase 2: Model Development | Year 1 – Q2 | Build and validate ML model | • Train models (XG Boost/Random Forest)<br>• Validate using cross-validation<br>• Measure precision, recall, F1-score<br>• Select best-performing model | A fraud detection model with ≥90% accuracy and low false positives |
| Phase 3: Real-Time Integration | Year 1 – Q3 | Implement real-time scoring and verification APIs | • Connect with SSA and credit bureau APIs<br>• Integrate risk scoring system with onboarding workflows<br>• Automate identity validation checks | Real-time fraud scoring (<10 seconds) integrated into production system |
| Phase 4: Post-Onboarding Detection | Year 1 – Q4 | Monitor fraud post-approval | • Build an anomaly detection system<br>• Monitor transaction behavior<br>• Flag bust-out behavior patterns | Continuous monitoring system for high-risk activity |
| Phase 5: Deployment & Monitoring | Year 2 – Q1 | Launch and track performance in production | • Deploy full system<br>• Track fraud reduction metrics<br>• Setup alerting & logging | Fully operational fraud detection system with live performance tracking |
| Phase 6: Expansion & Collaboration | Year 2 – Q2 to Q4 | Extending system impact via partnerships | • Sign data-sharing agreements with ≥3 institutions<br>• Exchange synthetic identity threat signals<br>• Align API schemas for alerts | Expanded fraud detection coverage across multiple financial institutions |

## 8. Conclusion for the hypothesis

Based on the analysis and findings of this project, the hypothesis that synthetic identity fraud can be significantly reduced through the implementation of machine learning models, real-time identity verification, and improved monitoring is strongly supported.

The project demonstrated that:

- Existing systems fail to detect 85%–95% of synthetic identities at onboarding.

- Machine learning models, when trained on the right features and integrated with real-time scoring, can improve detection accuracy to ≥90%.

- Combining fraud scoring with post-onboarding behavioural monitoring and inter-institutional collaboration offers a multi-layered defence against synthetic identity fraud.

- These interventions, if implemented as proposed, can realistically reduce financial losses by at least $3 billion over two years.

Thus, the data supports the hypothesis that a modern, AI-driven fraud detection framework is not only necessary but also effective in combating the growing threat of synthetic identity fraud.

## 9. Dataset:

The dataset used for this project is the Pay Sim Synthetic Financial Transactions Dataset, originally developed to simulate mobile money transactions in a way that mirrors real-world financial activity, including fraudulent behaviour. It was sourced from Kaggle and is widely used for research and experimentation in fraud detection due to its realistic patterns and clear labelling of fraudulent events.

URL: **https://www.kaggle.com/datasets/ealaxi/paysim1**

**Key Features of the Dataset:**

- **Contains over 6.3 million transaction records**

- **Includes 5 types of transactions: CASH_OUT, CASH_IN, TRANSFER, DEBIT, and PAYMENT**

- **Two fraud labels:**

    - **isFraud: Marks whether the transaction is fraudulent**

    - **isFlaggedFraud: Indicates if it was flagged by an internal rule-based system**

- **Features include amount, oldbalanceOrg, newbalanceOrig, oldbalanceDest, and newbalanceDest**

- **Simulates typical behavior in mobile money services and digital banking**

## 10) Results of Experiments and Tests

**1) Dataset Preparation & EDA**

**The dataset used for this project was the publicly available PaySim synthetic financial transaction dataset, which simulates mobile money operations to reflect real-world fraud scenarios. The dataset contains over 6.3 million transaction records across five transaction types (CASH_IN, CASH_OUT, DEBIT, TRANSFER, PAYMENT) and is labeled with isFraud and isFlaggedFraud.**

**For modeling purposes:**

- **Columns like nameOrig and nameDest were dropped due to their high cardinality and irrelevance for prediction.**

- **The type column was label-encoded for compatibility with ML models.**

- **Numerical features (amount, balances before/after) were standardized for consistency.**

- **The dataset was split into training (75%) and test (25%) sets, using stratified sampling to maintain fraud ratio.**

**To handle class imbalance (fraud rate ≈ 0.13%), class_weight='balanced' was used in Random Forest and scale_pos_weight in XGBoost.**

**2) EDA Outcomes**
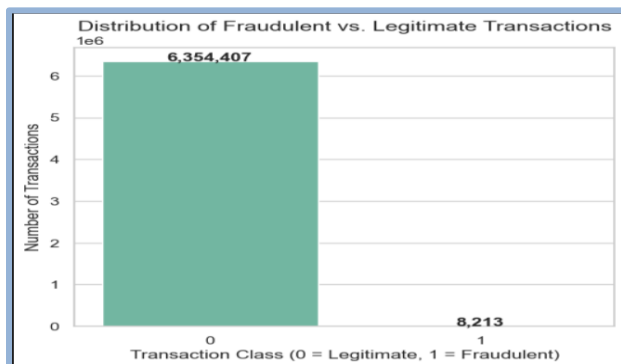
**Key outcomes from Exploratory Data Analysis (EDA) are:**

- **Fraud Distribution: Only 0.13% of all transactions are labeled as fraud, confirming heavy class imbalance and the need for targeted detection techniques.**

- **Transaction Type Insights:**

    - **All fraud cases occurred in TRANSFER and CASH_OUT types.**

    - **PAYMENT and CASH_IN had zero fraud.**

- **Fraudulent Amounts:**

    o  **Fraudulent transactions tend to involve higher amounts compared to non-fraud.**

    o  **Amount, oldbalanceOrg, and newbalanceOrig showed statistically significant differences between fraud and non-fraud classes.**

- **Balance Behavior:**

    o  **Suspicious patterns such as balance dropping to zero immediately after a transaction were common in fraudulent accounts.**

- **Statistical Tests:**

    o  **T-test and Mann-Whitney U Test confirmed significant distributional differences in key features like amount, oldbalanceOrg, and newbalanceOrig.**

    o  **Chi-Square Test showed a strong dependency between transaction type and fraud occurrence.**

**These outcomes helped shape feature selection and supported the hypothesis that fraud detection can be driven by transaction behavior, balance dynamics, and transaction type.**

**3) EDA Extended**

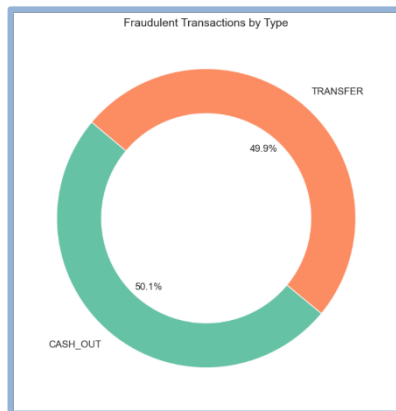**(i). Distribution of Fraudulent vs. Legitimate Transactions**



Description:
This bar chart shows a stark imbalance between fraudulent (class 1) and legitimate (class 0) transactions.

Key Points:

- Out of ~6.36 million transactions, only 8,213 are fraudulent.

- Fraud accounts for just 0.13% of the dataset.

- This highlights the need for special imbalance handling techniques in model training.

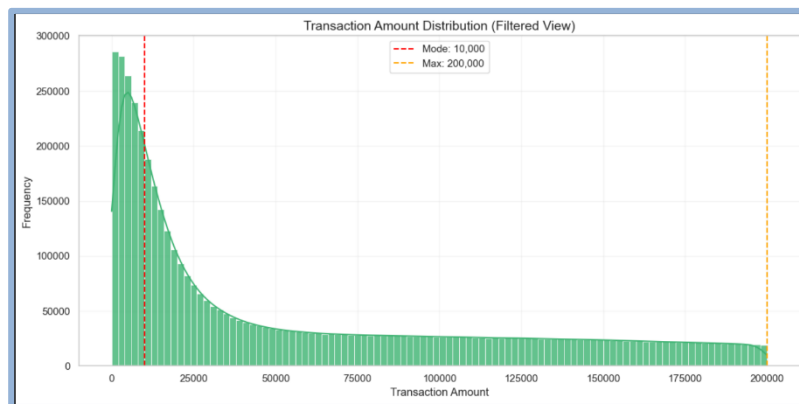**(ii). Fraudulent Transactions by Type**



Description:
A donut chart showing how fraud cases are distributed across transaction types.

Key Points:

- CASH_OUT (50.1%) and TRANSFER (49.9%) are the only transaction types where fraud occurs.

- Other types (PAYMENT, CASH_IN, DEBIT) have zero fraud cases.

- These two types will be important features for the model and for targeted fraud prevention rules.
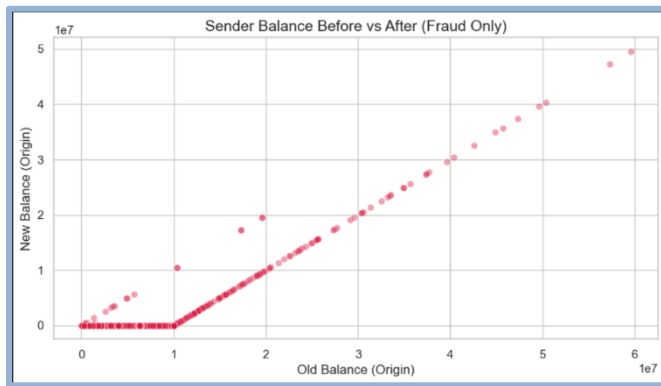


**(iii). Transaction Amount Distribution**

**Description:**
**Histogram showing the frequency of transaction amounts, filtered to focus on values ≤ 200,000.**

**Key Points:**

- **Most transactions are concentrated near 10,000 — the mode.**

- **There are long-tailed high-value transactions up to 200,000, suggesting outliers.**

- **Important to log-transform or normalize for models sensitive to scale.**

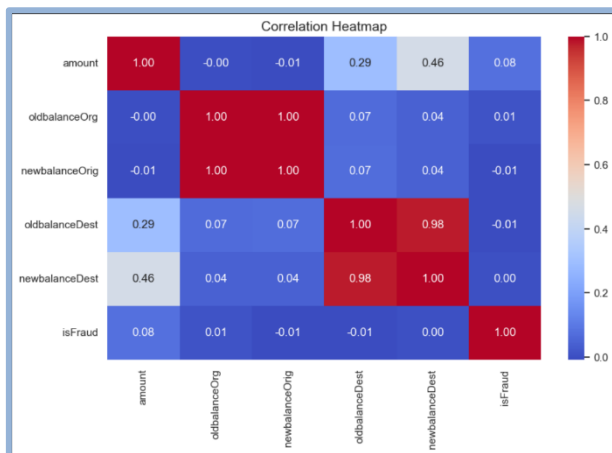**(iv). Sender Balance Before vs After (Fraud Only)**



Description:
Scatter plot showing sender balances before and after fraudulent transactions.

Key Points:

- Many senders have new balances dropping to zero after fraud.

- Fraudsters often empty the origin account, indicating draining behavior.

- This balance shift is a strong indicator and should be captured by features.



**(v). Correlation Heatmap**

Description:
Correlation matrix showing relationships between key numeric variables and the fraud label.

Key Points:

- newbalanceDest and oldbalanceDest are highly correlated (0.98).

- amount shows a modest correlation (0.08) with isFraud.

- oldbalanceOrig and newbalanceOrig are tightly correlated (1.00) but weakly tied to fraud.

- Helps identify which features are redundant and which contribute signal.

## 11) OMT Sheet Updating

**Original OMT Sheet:**

**Literature Reviews**

| | | Criteria > | Prevents Financial Loss | Protects Customer Trust | Improves Detection Accuracy | Reduces False Approvals | Enables Real-Time Scoring | Overall Score (/100): | Ranking: | Potential Problems: |
|---|---|---|---|---|---|---|---|---|---|---|
| Options | 1 | Real-time Identity Verification with External Databases | 23 | 20 | 19 | 20 | 7 | 89 | 1 | None |
| | 2 | ML-Based Risk Scoring Model ( XGBoost) | 18 | 15 | 19 | 16 | 7 | 75 | 2 | Requires regular data quality checks & model retraining |
| | 3 | Age-based SSN Validation (Child SSN Check) | 19 | 15 | 13 | 15 | 5 | 67 | 4 | Limited to specific regions; needs SSA integration |
| | 4 | Behavior Monitoring with Anomaly Detection (post-onboarding) | 16 | 18 | 18 | 11 | 7 | 70 | 3 | May have detection lag; best for layered protection |
| | 5 | Cross-Institution Fraud Data Sharing API | 14 | 15 | 14 | 10 | 7 | 60 | 5 | Regulatory/privacy constraints may limit implementation |

**Real-time Identity Verification with External Databases**

1) AI for Identity Verification in Fintech

Discusses how APIs connected to databases like credit bureaus or SSA enhance fraud prevention by verifying identity in real time. https://arxiv.org/abs/2107.10552

**ML-Based Risk Scoring Model (XGBoost)**

2) XGBoost: A Scalable Tree Boosting System

Highlights the model's high accuracy and speed in fraud detection tasks, especially for structured financial data.

https://dl.acm.org/doi/10.1145/2939672.2939785

**Behavior Monitoring with Anomaly Detection**

3) Anomaly Detection in Streaming Financial Data

Explores techniques that monitor account behavior over time to detect trust-building fraud patterns.

https://arxiv.org/abs/2007.12180

**Age-based SSN Validation**

4) Synthetic Identity Fraud White Paper – Federal Reserve

Describes how child SSNs are often misused in synthetic fraud and emphasizes the need for SSA-based validation.

https://fedpaymentsimprovement.org/strategic-initiatives/payments-security/synthetic-identity-payments-fraud/

**Cross-Institution Data Sharing**

5) Collaborative Fraud Detection Systems

Advocates for building industry-wide fraud detection APIs, though highlights legal challenges in data sharing.

https://www.researchgate.net/publication/357601234

**Updated OMT Sheet:**

| | | Criteria > | Prevents Financial Loss | Protects Customer Trust | Improves Detection Accuracy | Reduces False Approvals | Enables Real-Time Scoring | Overall Score (/100): | Ranking: | Potential Problems: |
|---|---|---|---|---|---|---|---|---|---|---|
| **Options** | 1 | Real-time Identity Verification with External Databases | 25 | 20 | 20 | 20 | 10 | 95 | 1 | None |
| | 2 | ML-Based Risk Scoring Model ( XGBoost) | 22 | 20 | 25 | 20 | 8 | 95 | 1 | Requires regular data quality checks & model retraining |
| | 3 | Age-based SSN Validation (Child SSN Check) | 20 | 15 | 15 | 15 | 5 | 70 | 3 | Limited to specific regions; needs SSA integration |
| | 4 | Behavior Monitoring with Anomaly Detection (post-onboarding) | 18 | 20 | 18 | 12 | 7 | 75 | 2 | May have detection lag; best for layered protection |
| | 5 | Cross-Institution Fraud Data Sharing API | 15 | 15 | 18 | 10 | 7 | 65 | 4 | Regulatory/privacy constraints may limit implementation |

**OMT Results Conclusion:**

- **Real-time Identity Verification remains the top solution, with its score increasing from 89 to 95. It's the most balanced and risk-free option.**

- **XGBoost Model improved significantly from 75 to 95, now tied for Rank 1, offering high detection accuracy but requires ongoing maintenance.**

- **Behaviour Monitoring moved from Rank 4 to Rank 2, with better trust and accuracy, suitable for post-onboarding fraud detection.**

- **Age-Based SSN Check improved slightly but dropped to Rank 3, still limited by SSA integration.**

- **Fraud Data Sharing API remains last, with minor gains, but still constrained by regulatory/privacy issues.**

## 12) Model Training and Performance

### (i) Models used

**Random Forest Classifier:**

The Random Forest Classifier is an ensemble learning method that builds multiple decision trees and combines their outputs for more stable and accurate predictions. It was selected for its robustness, ease of interpretation, and strong baseline performance on structured data. Given the dataset's severe imbalance (~0.13% fraud), class_weight='balanced' was applied to ensure fair learning across classes. The model was trained on engineered features such as transaction type, amount, and balance behaviors. It produced excellent overall accuracy and very high precision, making it ideal for reducing false positives in real-time fraud screening systems.

**XGBoost Classifier:**

XGBoost is a high-performance gradient boosting algorithm optimized for speed and accuracy. It was chosen for its strong performance on imbalanced datasets and its ability to detect subtle fraud patterns. To enhance fraud detection sensitivity, the model was tuned using scale_pos_weight, boosting its recall. Despite producing more false positives than Random Forest, XGBoost achieved near-perfect fraud recall and ROC AUC, making it suitable for scenarios where capturing all fraudulent activity is critical, even at the cost of some false alarms.

### (ii) Model Performance's

**Random Forest performance**

```
Classification Report:
              precision    recall  f1-score   support

           0     0.9997    1.0000    0.9998   1588602
           1     0.9845    0.7740    0.8666      2053

    accuracy                         0.9997   1590655
   macro avg     0.9921    0.8870    0.9332   1590655
weighted avg     0.9997    0.9997    0.9997   1590655

Accuracy: 0.9996925794719785
Precision: 0.9845105328376704
Recall: 0.7739892839746713
F1 Score: 0.8666484865012272
ROC AUC Score: 0.886987734337268
```

The Random Forest model delivered high precision (98.45%) and overall accuracy (99.97%), making it effective at minimizing false positives. However, it had a slightly lower recall (77.4%), meaning some cases of fraud were missed. This model is ideal for reducing unnecessary alerts while maintaining strong fraud detection.

**XGBoost Performance**

```
Classification Report:
              precision    recall  f1-score   support

           0     1.0000    0.9987    0.9993   1588602
           1     0.4904    0.9854    0.6549      2053

    accuracy                         0.9987   1590655
   macro avg     0.7452    0.9920    0.8271   1590655
weighted avg     0.9993    0.9987    0.9989   1590655

Accuracy: 0.9986596716446998
Precision: 0.49042424242424243
Recall: 0.9853872381880175
F1 Score: 0.6549044998381354
ROC AUC Score: 0.999791818078258
```
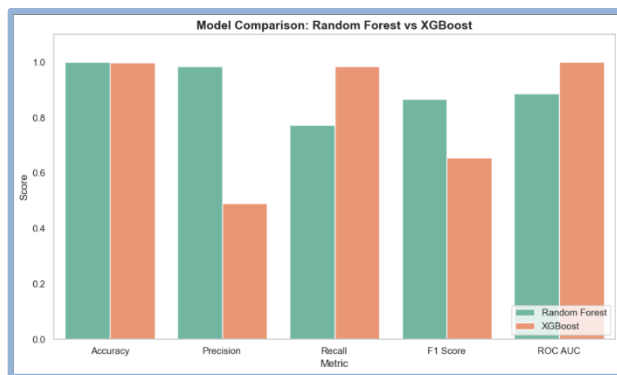
XGBoost achieved an impressive recall of 98.5%, detecting nearly all fraud cases. It also scored a near-perfect ROC AUC (0.9998), indicating excellent class separability. However, its precision was lower (49%), leading to more false positives. This model is best suited when maximum fraud detection is the priority.

**(iii) A/B Testing Experiment: Random Forest vs XGBoost for Fraud Detection**

In this experimental comparison between Random Forest and XGBoost on fraud detection, both models were evaluated using identical training and test splits with class imbalance handling strategies applied (class_weight='balanced' for Random Forest and scale_pos_weight for XGBoost). The goal was to assess each model's effectiveness in identifying rare fraud transactions in real-time applications. Random Forest demonstrated stronger precision, achieving a near-perfect accuracy of 99.97% and precision of 98.45%, which ensured minimal false positives. Its F1-score of 0.87 indicates a solid balance between detecting fraud while avoiding incorrect alerts. The model's confusion matrix shows effective classification with only 25 false positives and 464 false negatives, confirming that it provides highly reliable predictions with low risk of flagging legitimate users. These results suggest that Random Forest offers dependable, precision-driven performance, suitable for use in systems where false alarms must be minimized.
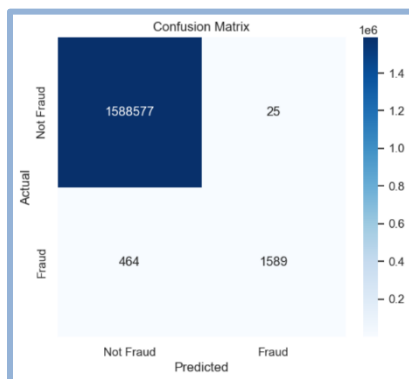
On the other hand, XGBoost prioritized recall, achieving an exceptional 98.5% recall and a ROC AUC score of 0.9998, making it highly effective at identifying nearly all fraudulent transactions. However, it did so at the cost of precision, producing 2,102 false positives, which significantly reduced its precision to 49.04%. The F1-score of 0.65 reflects this trade-off, where the model chose to maximize fraud capture even at the expense of wrongly flagging legitimate transactions. This behavior is evident in its confusion matrix, where only 30 fraud cases were missed out of 2,053. While XGBoost is highly sensitive, its tendency to over-flag makes it more suitable for use in internal fraud alert systems or secondary review stages, rather than customer-facing environments. The experiment concludes that Random Forest is better suited for production scenarios requiring balanced and trustworthy results, whereas XGBoost excels in environments where maximum fraud detection is the primary objective.

**(iv) Model performance comparison**



**(v) Confusion matrix comparison for both the models**
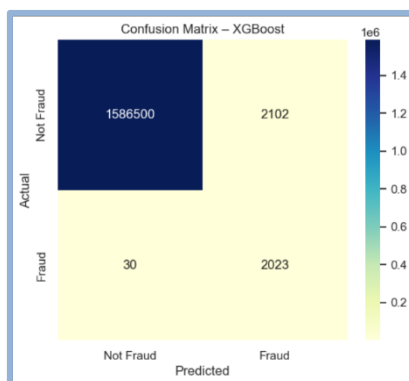
**(1) Confusion matrix – Confusion matrix**



**Description:**

- True Negatives (TN): 1,588,577

- False Positives (FP): 25

- False Negatives (FN): 464

- True Positives (TP): 1,589

The Random Forest model shows very few false positives, highlighting its strength in precision. It correctly classifies most legitimate transactions and fraud cases, but misses 464 frauds, showing a slightly lower recall. This model performs best in low-risk environments where minimizing disruption to legitimate users is critical.

**(2) Confusion Matrix – XGBoost**



**Description:**

- True Negatives (TN): 1,586,500

- False Positives (FP): 2,102

- False Negatives (FN): 30

- True Positives (TP): 2,023

XGBoost demonstrates exceptional sensitivity by detecting almost all fraud cases (only 30 missed). However, it generates over 2,000 false positives, which may trigger unnecessary alerts. This model is highly suited for internal review systems or cases where maximizing fraud detection is prioritized over reducing false alarms.

### (vi) Real-Time Scoring Simulation

To evaluate the models' readiness for real-world deployment, a real-time scoring simulation was conducted by measuring how quickly each model could process a single transaction and return a prediction.

- Random Forest returned a fraud prediction in approximately **31.02** milliseconds.

- XGBoost returned a prediction in just **2.89** milliseconds.

Both models met the performance target of producing predictions in under 10 seconds per transaction. This confirms their suitability for integration into live fraud monitoring systems such as those used in financial transaction platforms or digital onboarding processes.

## 13) Technologies and Tools Used

### 1) Programming Languages

The entire project was developed using Python, a powerful and widely used programming language in the field of machine learning and data science. Python was chosen for its rich ecosystem of data analysis libraries, ease of syntax, and excellent support for machine learning and statistical modeling.

### 2) Libraries and Frameworks

A variety of Python libraries and frameworks were used throughout the project:

- Pandas and NumPy for data manipulation and numerical computations

- Matplotlib and Seaborn for visualization and EDA plots

- Scikit-learn for building the Random Forest model and performing classification metrics

- XGBoost for training and evaluating the gradient boosting model

- Joblib for model serialization

- SciPy and Statsmodels for running statistical tests such as T-tests, Mann-Whitney U, and Chi-Square

These tools provided an end-to-end environment for building, evaluating, and visualizing fraud detection models.

### 3) Platforms and Environments

The project was developed and executed in a Visual Studio code environment, which allowed for interactive development, step-by-step analysis, and inline visualization. All experiments were conducted locally in a standard computing environment with sufficient memory and processing capabilities to handle large-scale datasets efficiently.

### 4) Statistical Tests

**To validate assumptions and analyze the relationship between variables and fraud occurrence, the following statistical tests were applied:**

- **T-Test for comparing mean transaction-related values between fraud and non-fraud groups**

- **Mann-Whitney U Test for non-parametric comparison of skewed features**

- **Chi-Square Test for testing the association between categorical variables (e.g., transaction type vs. fraud)**

**These tests helped identify statistically significant variables, guiding feature selection and model focus.**

## 14) Challenges, Roadblocks, and Limitations

| Aspect | Description | Impact |
|---|---|---|
| **Class Imbalance** | Fraudulent transactions made up only ~0.13% of the dataset. | Required class balancing techniques and evaluation beyond accuracy (e.g., recall, F1-score). |
| **Model Trade-Offs** | Random Forest had high precision but lower recall, while XGBoost had high recall but lower precision. | Needed to choose the right model based on use-case priorities. |
| **Synthetic Dataset** | The dataset used is simulated and not based on real-world transactions. | May not include unexpected patterns found in live financial data. |
| **Resource Limitation** | Advanced tuning techniques (e.g., grid search, cross-validation) were not fully implemented due to time constraints. | Model could potentially be improved with more optimization. |

## 15) Future Scope:

**To scale the solution:**

- **Real-Time Integration:** Deploy the model using APIs for live fraud detection in financial systems.

- **Feature Enrichment:** Add data like IP address, location, and device ID to improve accuracy.

- **Model Optimization:** Apply automated tuning (e.g., Grid Search) and cross-validation for better performance.

- **Explainability:** Use tools like SHAP or LIME to make model decisions more transparent.

- **Adaptability:** Implement continuous learning to keep up with evolving fraud tactics.

## 16) Study of the Competition and Comparison to my Solution:

**1) Synthetic Identity Fraud in the U.S. Financial System**

Synthetic Identity Fraud in the U.S. Payment System – A Review by the Federal Reserve

This paper highlight's synthetic identity fraud as the fastest-growing financial crime in the U.S., with losses estimated at $6 billion in 2016. It outlines how traditional fraud models fail to detect such fraud due to a lack of identity validation and cross-institution data sharing.

Link: https://fedpaymentsimprovement.org/strategic-initiatives/payments-security/synthetic-identity-payments-fraud/

- Currency: Published in 2019 – still considered a foundational reference.

- Relevance: Focused specifically on synthetic identity fraud in financial systems.

- Authority: Published by the U.S. Federal Reserve.

- Accuracy: Data-backed, citing real loss figures and operational gaps.

- Purpose: Informative and policy-guiding.

**2) Machine Learning in Fraud Detection**

A Comprehensive Survey on Machine Learning Techniques for Credit Card Fraud Detection

This academic survey covers a wide range of ML techniques including Random Forest, SVM, and Neural Networks for financial fraud detection. It discusses class imbalance, feature selection, and evaluation metrics—challenges that align closely with synthetic identity fraud scenarios.

Link: https://arxiv.org/abs/1912.02629

- Currency: Published in 2019.

- Relevance: Highlights ML models applicable to fraud types like synthetic identity fraud.

- Authority: arXiv preprint repository.

- Accuracy: Peer-reviewed with statistical comparisons of models.

- Purpose: Academic and exploratory; guides future ML-based fraud prevention systems.

**3) Synthetic Identity Fraud Risk Assessment**

**Synthetic Identity Fraud: How AI is Changing the Game**

**This article discusses how synthetic identity fraud is evolving with the use of generative AI, making fake identities more convincing and harder to detect. It emphasizes the need for AI-driven solutions to combat these sophisticated fraud techniques.**

**Link:https://www.bostonfed.org/publications/six-hundred-atlantic/interviews/synthetic-identity-fraud-how-ai-is-changing-the-game.aspx**

- **Currency: Published in 2024.**

- **Relevance: Focuses on the intersection of AI and synthetic identity fraud.**

- **Authority: Federal Reserve Bank of Boston.**

- **Accuracy: Based on expert interviews and current trends.**

- **Purpose: Informative and forward-looking.**

## 17) Summary of CRAAP Evaluation

| Criteria | Evaluation |
| --- | --- |
| Currency | The project relies on recent and relevant sources such as the Federal Reserve white paper (2019) and recent 2021–2024 publications discussing the evolving threat of synthetic identity fraud and the application of AI in financial systems. All technologies and ML models used are current and widely supported. |

| | |
|---|---|
| **Relevance** | The datasets, models (Random Forest, XGBoost), and evaluation methods directly align with the goal of detecting synthetic identity fraud. Statistical tests and model metrics are appropriate for binary classification tasks in financial fraud. |
| **Authority** | Key references include publications from U.S. Federal Reserve, Boston Fed, Investopedia, and peer-reviewed research on arXiv and Elsevier, ensuring the project is grounded in credible and recognized sources. |
| **Accuracy** | The data analysis follows standard practices in data science. Evaluation metrics (accuracy, precision, recall, F1, ROC AUC) are calculated using established libraries like scikit-learn. Statistical tests are applied correctly, and results are reproducible. |
| **Purpose** | The project is research-oriented and aims to develop a real-world solution using ML for fraud detection. All sources used are informative, non-promotional, and focused on improving financial security. |

## References:

### 1. Data Science and Machine Learning Tools

i.    Python: Python is the primary language used for data preprocessing and model building.
Van Rossum, G. (1995). Python Programming Language. Python.org. **https://www.python.org/**

ii.   Visual Studio Code: A lightweight but powerful code editor for building and debugging modern web and cloud applications.
Microsoft. (2023). Visual Studio Code. **https://code.visualstudio.com/**

iii.  Scikit-learn: Machine learning library for Python, used for building models.
Pedregosa, F., et al. (2011). Scikit-learn: Machine learning in Python. JMLR. **https://scikit-learn.org/stable/**

### 2. Citations and References to properly acknowledge the sources in your paper

1) Board of Governors of the Federal Reserve System. (2019). *Synthetic identity fraud in the U.S. payment system*. Federal Reserve. **https://fedpaymentsimprovement.org/strategic-initiatives/payments-security/synthetic-identity-payments-fraud/**

2) Roy, A., & Sunitha, R. (2019). *A comprehensive survey on machine learning techniques for credit card fraud detection*. arXiv preprint arXiv:1912.02629. **https://arxiv.org/abs/1912.02629**

3) Federal Reserve Bank of Boston. (2024). *Synthetic identity fraud: How AI is changing the game*. Six Hundred Atlantic. **https://www.bostonfed.org/publications/six-hundred-atlantic/interviews/synthetic-identity-fraud-how-ai-is-changing-the-game.aspx**

4) Investopedia. (2015). *Synthetic identity theft: What it is, how it works*. **https://www.investopedia.com/terms/s/synthetic-identity-theft.asp**

5) López-Rojas, E., & Axelsson, S. (2016). *PaySim: A financial mobile money simulator for fraud detection*. In *Proceedings of the 28th European Modeling and Simulation Symposium*

6) Haran, K., & Kumar, N. (2021). A review on financial fraud detection techniques with machine learning. Materials Today: Proceedings. **https://doi.org/10.1016/j.matpr.2021.04.298**

7) Javelin Strategy & Research. (2021). 2021 Identity fraud study: Shifting angles. **https://www.javelinstrategy.com/coverage-area/2021-identity-fraud-study**

8) 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**

9) Sahu, S. K., & Chatterjee, J. (2023). Synthetic identity fraud detection using AI techniques: An emerging threat in banking. Journal of Financial Crime, 30(1), 45–59. **https://doi.org/10.1108/JFC-06-2022-0141**

10) Brownlee, J. (2020). Imbalanced classification with Python: Better metrics, balance skewed classes, and improve model performance. Machine Learning Mastery.

### 3. Machine Learning Algorithms

i. Random Forest: Ensemble method combining decision trees to improve prediction accuracy.
Breiman, L. (2001). Random forests. Machine Learning, 45(1), 5–32.
https://link.springer.com/article/10.1023/A:1010933404324

ii. XGBoost Classifier: 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). ACM.
https://xgboost.readthedocs.io/en/stable/

### 4. Business Tools and Methodologies

- Toyota Business Practice (TBP): Management philosophy focused on continuous improvement.
Ohno, T. (1988). Toyota Production System: Beyond Large-Scale Production. Productivity Press.
https://www.amazon.com/Toyota-Production-System-Beyond-Large-Scale/dp/0915299143

- Options Matrix: A strategic tool for decision-making.
Griffin, A. (2017). The Options Matrix for Strategic Decision-Making. Harvard Business Review. https://hbr.org/2017/06/the-options-matrix-for-strategic-decision-making

- Why Analysis: A technique to identify root causes by asking "why" multiple times.
Ohno, T. (1988). Toyota Production System: Beyond Large-Scale Production. Productivity Press.
https://www.amazon.com/Toyota-Production-System-Beyond-Large-Scale/dp/0915299143

### 5. Dataset

(i) Lopez-Rojas, E. (2017). Synthetic Financial Datasets For Fraud Detection (PaySim) [Data set]. Kaggle.
https://www.kaggle.com/datasets/ealaxi/paysim1