

IEEE-CIS Fraud Detection Challenge

Comparative Study of SVM and Decision Tree Binary Classification

Khashayar Zardou

Dept. Computer Science & Software Engineering

Concordia University

Montreal, Canada

khashayar.zardoui@mail.concordia.ca

ID: 40052568

Paolo Junior Angeloni

Dept. Computer Science & Software Engineering

Concordia University

Montreal, Canada

p_ange@live.concordia.ca

ID: 25976944

Abstract—The objective of this project was to develop a machine learning pipeline capable of identifying fraudulent credit card transactions within the IEEE-CIS Fraud Detection dataset. The primary challenge was the extreme class imbalance (approx. 3.5% fraud vs. 96.5% legitimate), requiring models that prioritize Precision (minimizing customer friction via false alarms) while maintaining high Recall (capturing actual fraud).

I. THE FRAUD-DETECTION PIPELINE

We implemented a robust data processing pipeline involving missing value imputation, label encoding for categorical variables, and standard scaling. To address the computational challenges of the large dataset (590k+ entries), we utilized Principal Component Analysis (PCA) for dimensionality reduction on linear models. We evaluated three distinct classification architectures:

- 1) Linear Support Vector Machine (LinearSVC): A baseline linear classifier utilizing the Primal solver and PCA
- 2) Decision Tree: A non-linear, rules-based classifier tuned for depth and split criteria
- 3) XGBoost: An advanced gradient-boosting ensemble method utilizing histogram-based optimization
- 1) Data loading and exploration (EDA)
- 2) Removal, Imputation, Label Encoding and Scaling
- 3) Model training and performance assessment

II. EXPLORATORY DATA ANALYSIS

A. Data Structure Inspection

- 1) Before any data transformation, we observed the `train` and `test` datasets had a mixture of `float64`, `int64` and `object` types
- 2) missing values
description goes here...
- 3) target balance
description goes here...

B. Statistical Summary & Visualizations

Fig. 1. some image here

TABLE I
SOME STATS...

Metric	Value
one	...
two (%)	...
three	...
four	...

C. Findings & Hypotheses

...

...

III. DATA PRE-PROCESSING & CLEANING

A. Imputation & Removal

- 1) ...
- 2) ...
- 3) ...

B. Normalize & Scale Features

- 1) ...
- 2) ...
- 3) ...

C. Encoding Categorical Features

- 1) ...
- 2) ...
- 3) ...

IV. MODELS

intro to the models used

A. Support Vector Machine (SVM) Classifier

Due to the size of the dataset (590,000+ samples), a standard SVM with a non-linear kernel ($O(n^3)$) was computationally infeasible. We opted for a LinearSVC ($O(n)$) to utilize the entire training set. To satisfy the hyperparameter tuning requirement2, we tuned the Regularization parameter (C) and the Loss function (Hinge vs. Squared Hinge) instead of the kernel. Additionally, we applied Principal Component Analysis (PCA) to the SVM input to reduce dimensionality, which resolved convergence issues and significantly improved training speed. experiment hyperparameters (C, gamma, kernel etc)

cross-validation and validation splits to evaluate performance results using different hyperparameters training and test metrics: confusion matrix, precision, recall, F1-score, and accuracy

B. Decision Tree Classifier

experiment hyperparameters (max depth, min samples split, cross-validation and validation splits to evaluate performance results using different hyperparameters training and test metrics: confusion matrix, precision, recall, F1-score, and accuracy

The Decision Tree model demonstrated signs of mild overfitting. While it achieved a high F1-score of 0.68 on the training set, this dropped to 0.57 on the validation set. Specifically, the Precision for fraud detection fell from 92% (training) to 77% (validation), indicating that some of the decision rules learned were specific to the training noise and did not generalize well. Furthermore, the Recall remained low in both sets (0.53 training vs. 0.45 validation), suggesting that a single decision tree lacks the complexity required to capture the full variety of fraudulent patterns in this dataset.

C. XGBoost Classifier

The XGBoost classifier proved to be the superior model. While the Decision Tree suffered from overfitting (high variance), XGBoost demonstrated robust generalization.

On the Validation set, XGBoost achieved a Precision of 0.93, meaning it generated very few false positives (false alarms), which is critical for maintaining user trust. Simultaneously, it achieved a Recall of 0.56, capturing the majority of fraud instances. The F1-Score of 0.70 (Validation) significantly outperforms the Decision Tree (0.57) and indicates that the Gradient Boosting method successfully captured complex, non-linear relationships that the simpler models missed.

V. MODEL COMPARISON

TABLE II
SOME STATS...

Metric	SVM	Decision Tree
one
two (%)
three
four

a) discuss similarities & differences. use table: The LinearSVC provided a baseline AUC of 0.815, but struggled with convergence times and lacked the complexity to model non-linear fraud patterns.

The Decision Tree achieved a higher AUC of 0.848, but exhibited signs of overfitting (high training accuracy vs. lower validation precision), confirming that a single tree has high variance.

XGBoost emerged as the vastly superior model, achieving an AUC-ROC of 0.963 and an F1-Score of 0.70. It successfully balanced a high Precision (93%) with a Recall of 56%,

significantly outperforming the other models in identifying fraud without disrupting legitimate users.

ACKNOWLEDGMENT

We would like to thank Professor Arash Azarfar and Firat Oncel for their guidance and support throughout this project. Large Language Models, like Google's Gemini were used in an educational context to further understand the resources for this research.

REFERENCES

- [1] Numpy, “NumPy API Documentation,” [Online]. Available: <https://numpy.org/doc/stable/>. [Accessed: Nov. 25, 2025].
- [2] matplotlib, “Matplotlib API Documentation,” [Online]. Available: <https://matplotlib.org/stable/index.html>. [Accessed: Nov. 25, 2025].
- [3] pandas, “pandas API Documentation,” [Online]. Available: <https://pandas.pydata.org/docs/>. [Accessed: Nov. 26, 2025].
- [4] seaborn, “seaborn API Documentation,” [Online]. Available: <https://seaborn.pydata.org/api.html>. [Accessed: Nov. 26, 2025].
- [5] data preprocessing, “sklearn API Documentation,” [Online]. Available: <https://scikit-learn.org/stable/modules/preprocessing.html>. [Accessed: Nov. 28, 2025].
- [6] data imputation, “sklearn API Documentation,” [Online]. Available: <https://scikit-learn.org/stable/modules/impute.html>. [Accessed: Nov. 28, 2025].
- [7] label encoding, “sklearn API Documentation,” [Online]. Available: https://scikit-learn.org/stable/modules/preprocessing_targets.html#label-encoding. [Accessed: Nov. 28, 2025].
- [8] hyperparameter tuning using GridSearchCV, “sklearn API Documentation,” [Online]. Available: https://scikit-learn.org/stable/modules/grid_search.html#grid-search. [Accessed: Nov. 30, 2025].
- [9] Support Vector Machines, “sklearn API Documentation,” [Online]. Available: <https://scikit-learn.org/stable/modules/svm.html>. [Accessed: Nov. 30, 2025].
- [10] Principal Component Analysis (PCA), “sklearn API Documentation,” [Online]. Available: <https://scikit-learn.org/stable/modules/decomposition.html#pca>. [Accessed: Dec. 02, 2025].
- [11] Decision Trees , “sklearn API Documentation,” [Online]. Available: <https://scikit-learn.org/stable/modules/tree.html>. [Accessed: Nov. 30, 2025].
- [12] XGBoostClassifier, “XGBoost API Documentation,” [Online]. Available: <https://xgboost.readthedocs.io/en/stable/>. [Accessed: Nov. 30, 2025].
- [13] metrics, “sklearn API Documentation,” [Online]. Available: <https://scikit-learn.org/stable/api/sklearn.metrics.html>. [Accessed: Dec. 02, 2025].