

MACHINE LEARNING BASED TOLL FRAUD DETECTION WITH HIGHWAY DASHBOARD

¹Surya. R, ²Mr. R. Balamurugan

¹M.Sc. Student, ²Assistant Professor

^{1,2}Department of Data Science and Business Analytics,

Vels Institute of Science, Technology and Advanced Studies (VISTAS), Chennai, India

suryaraja2403@gmail.com

Abstract — Toll fraud on national highways is a growing financial challenge for transportation authorities worldwide. This paper presents a machine learning based system for detecting fraudulent toll transactions, integrated with a real-time Highway Dashboard for operational monitoring. The proposed system applies an ensemble of three supervised classifiers, Random Forest, Support Vector Machine, and XGBoost, to classify toll transactions as legitimate or fraudulent using vehicle metadata, transaction timing, lane usage patterns, and payment behavior features. A web-based Highway Dashboard built on Flask and React provides administrators with live transaction monitoring, automated fraud alerts, vehicle blacklisting, and historical analytics. Experimental evaluation on a simulated dataset of 250,000 transactions demonstrates that the ensemble approach achieves an accuracy of 97.4% with an AUC-ROC of 0.992. The system is designed to be scalable, interpretable, and deployable within existing Electronic Toll Collection infrastructure. This work contributes an end-to-end practical solution to the problem of toll evasion and manipulation in modern highway networks.

Keywords — toll fraud detection; machine learning; random forest; XGBoost; anomaly detection; highway dashboard; electronic toll collection; classification; real-time monitoring.

I. INTRODUCTION

Transportation infrastructure, particularly highway toll systems, plays a vital role in a country's economic framework. In India and many other nations, tolls collected at highway plazas serve as a primary source of funding for road construction, maintenance, and expansion projects. The introduction of Electronic Toll Collection systems such as FASTag has brought greater efficiency and reduced congestion at toll booths. However, alongside these technological advancements, new forms of fraud have also emerged, threatening the financial integrity of toll operations.

Toll fraud covers a wide range of activities including the use of duplicate or cloned RFID tags, unauthorized lane switching, tampered vehicle registration details, repeated evasion by the same vehicles, and collusion between toll operators and vehicle owners. Traditional rule-based detection systems fall short when dealing with the dynamic and increasingly sophisticated nature of these fraudulent patterns. There is a clear need for intelligent, data-driven approaches that can adapt to evolving fraud strategies.

Machine learning offers a promising avenue for addressing these challenges. By training models on historical transaction data, it becomes possible to learn the subtle patterns that distinguish fraudulent activity from legitimate toll usage. Unlike rigid rule systems, machine learning models can generalize from past examples and flag previously unseen anomalies with high accuracy.

This paper presents a fraud detection framework that combines multiple machine learning classifiers with an interactive Highway Dashboard. The dashboard serves as the operational interface for transport administrators, providing real-time visibility into transaction flows, immediate fraud alerts, and analytical summaries that inform policy decisions.

The main contributions of this work are: (i) an ensemble machine learning pipeline combining Random Forest, SVM, and XGBoost for toll fraud classification, (ii) a feature engineering framework using transaction-level and behavioral features, (iii) a real-time Highway Dashboard for live fraud monitoring and alert management, (iv) a vehicle blacklist system for preemptive intervention, and (v) experimental evaluation on a 250,000-transaction simulated dataset demonstrating 97.4% accuracy.

II. LITERATURE SURVEY

Research on fraud detection using machine learning has grown significantly over the past decade, spanning domains such as credit card fraud, insurance fraud, telecommunications, and healthcare billing. Several studies have addressed parts of this problem but none have integrated a complete toll fraud detection pipeline with an operational highway dashboard.

Study 1: Bhattacharyya et al. (2011) demonstrated that ensemble methods, particularly Random Forest, outperform individual classifiers in detecting credit card fraud, establishing a benchmark that has influenced subsequent research. Their work showed that combining multiple weak learners produces a more robust fraud detector than any single model alone.

Study 2: Ozbayoglu et al. (2020) examined machine learning approaches for detecting anomalous behavior in transit systems and found that gradient-boosted trees were particularly effective when combined with temporal features derived from trip sequences. Their work highlighted the importance of contextual features such as time-of-day patterns, concepts directly applicable to highway toll fraud.

Study 3: Kumar et al. (2019) explored the vulnerabilities in FASTag-based toll systems and proposed a rule-based anomaly flagging

mechanism. However, their approach relied on static thresholds that were prone to both false positives and missed detections when fraud patterns shifted over time.

Study 4: Liu et al. (2008) proposed Isolation Forest for unsupervised anomaly detection where labeled fraud data is scarce. The method has since been adapted for transaction-level fraud detection in domains including transportation, providing a useful baseline for comparison with supervised approaches.

Study 5: Zhang et al. (2021) combined deep autoencoders with supervised classifiers in a hybrid architecture to improve fraud recall while controlling false alarm rates. The present study builds on these threads to develop an end-to-end system with a live operational dashboard tailored for highway toll administration.

III. PROPOSED WORK

The core idea behind the proposed system is straightforward. Existing toll monitoring systems either rely on manual inspection that is too slow for real-time intervention, or on rigid rule-based filters that cannot adapt to evolving fraud patterns. The proposed system addresses both limitations by applying ensemble machine learning to classify every incoming toll transaction automatically, while presenting results through an interactive dashboard that keeps the administrator in control of final intervention decisions.

The system starts by ingesting raw toll transaction records from Electronic Toll Collection plaza sensors and backend systems. Each record contains vehicle ID, RFID tag number, lane identifier, timestamp, vehicle class, payment method, transaction amount, and toll plaza geolocation. These raw attributes are enriched through feature engineering to produce 28 input features for the detection engine.

The detection engine applies three classifiers in a majority-voting ensemble. Random Forest captures feature interactions and handles noisy labels well. Support Vector Machine provides strong generalization in high-dimensional spaces. XGBoost manages class imbalance effectively through iterative boosting that emphasizes difficult-to-classify samples. The final fraud prediction is the majority vote across all three models.

The Highway Dashboard receives fraud predictions via a REST API and presents them to administrators in real time. Transactions flagged as high-risk trigger automated alerts to the responsible toll plaza. Vehicles repeatedly flagged are added to a blacklist for preemptive monitoring. This human-in-the-loop design is the core principle separating the proposed system from a fully automated blackbox fraud engine.

IV. DESIGN SYSTEM

The system design follows a layered architecture starting from raw data ingestion and ending at the administrator taking a decision on the dashboard. Each layer has a specific responsibility and passes its output to the next without creating tight coupling between components.

The Data Layer ingests raw transaction records from toll plaza ETC sensors and preprocesses them through cleaning, normalization, and feature extraction. The Feature Engineering Layer constructs 28 features per transaction including derived behavioral attributes such as transaction frequency per vehicle per day, lane mismatch score, and a 7-day rolling anomaly score. The Detection Engine applies the three-classifier ensemble and produces a fraud probability score and binary label. The API Layer exposes these predictions through versioned REST endpoints. The Dashboard Layer presents live transaction feeds, fraud alerts, blacklist management, and historical analytics to the administrator.

The architecture ensures that the administrator is always the final decision maker. The system detects and flags but never independently escalates or blocks a vehicle without a human reviewing the alert first. This deliberate design choice maintains accountability and reduces the operational risk of automated false positives.

V. PROJECT PHASES / METHODOLOGY

The system was built across four phases. Each phase targeted a specific layer and produced a working, testable component before the next phase began.

PHASE 1: DATA COLLECTION AND PREPROCESSING

The first phase focused on establishing the data foundation. A simulated dataset of 250,000 toll transactions was constructed spanning six months of operations across multiple highway plazas. Fraudulent transactions were seeded to represent four major fraud categories: cloned RFID tag usage, payment amount manipulation, unauthorized vehicle class declaration, and repeated tag reuse across non-adjacent plazas within implausible time windows. The dataset contained 8,000 fraudulent records representing 3.2% of the total.

Data preprocessing included missing value imputation, outlier capping at the 99th percentile, and min-max normalization of numerical features. Categorical variables were one-hot encoded. SMOTE was applied to the training set to balance the class distribution before model training.

PHASE 2: DETECTION ENGINE DEVELOPMENT

Three classifiers were trained independently using 5-fold cross-validation. Random Forest used 200 estimators, maximum tree depth of 15, and minimum samples per leaf of 5. The SVM model used an RBF kernel with $C=10$ and $\gamma=0.01$. XGBoost was configured with 150 boosting rounds, learning rate of 0.05, and maximum tree depth of 6. The final ensemble used majority voting across all three model predictions.

Feature importance analysis identified the most discriminative features as the time gap between consecutive transactions for the same RFID tag, the deviation from expected toll amount for the declared vehicle class, and the 7-day rolling anomaly score.

PHASE 3: DASHBOARD DEVELOPMENT

The dashboard was developed as a responsive web application using Flask backend, PostgreSQL database, and React frontend. The dashboard communicates with the detection engine via REST API, receiving fraud probability scores and classification labels for each incoming transaction.

The key features of the dashboard are:

- Real-Time Transaction Monitor showing a live feed with color-coded fraud risk scores refreshing every 5 seconds
- Alert Management Panel generating automated notifications to plaza administrators for high-risk transactions
- Vehicle Blacklist Management automatically flagging repeat offenders for preemptive monitoring
- Analytics and Reporting Module generating daily, weekly, and monthly fraud summaries with trend charts
- Performance Metrics Panel showing fraud detection rate, false positive rate, and revenue protected estimates

PHASE 4: TESTING AND DEPLOYMENT

All modules were integrated, tested, and deployed. The dashboard handled 500 concurrent transactions per minute with an average API response latency of 210 milliseconds. End-to-end latency from transaction arrival to dashboard alert was under 2 seconds in 95% of test cases.

The system was tested under the following scenarios:

- Normal transaction sessions to validate classification accuracy
- Simulated cloned RFID tag attacks to verify fraud recall
- High-volume load tests to confirm pipeline stability
- Extended sessions to confirm no memory leaks or service degradation

VI. IMPLEMENTATION

System Integration: The process of system integration brought all modules developed across each phase into a unified deployable application. Every module plays a critical role in the end-to-end fraud detection pipeline. Integration testing confirmed that data flows correctly between each layer without interruption or data loss.

The various components of the integrated system include:

- Transaction Ingestion Connector – pulls raw ETC records from toll plaza sensors and backend systems
- Feature Engineering Module – constructs 28 per-transaction features from raw attributes and behavioral history
- Random Forest Classifier – tree-based ensemble model trained on balanced transaction data
- SVM Classifier – RBF kernel model providing high-dimensional generalization
- XGBoost Classifier – gradient-boosted model with boosting iterations emphasizing hard fraud samples
- Majority Voting Ensemble – aggregates three model predictions into final fraud label and probability
- Flask Backend – serves all detection outputs through versioned REST endpoints
- PostgreSQL Database – stores transaction records, fraud labels, and alert history
- React Frontend – renders live transaction feed, alert panels, blacklist manager, and analytics charts

The Flask backend continuously serves updated prediction data through versioned endpoints. The React frontend renders all panels including the live transaction feed and the alert management interface. The administrator views the dashboard, checks flagged transactions, reviews supporting evidence, and decides whether to dispatch a field intervention manually.

Testing: The test results showed the system efficiently processes transactions and generates correct fraud classifications without delays. Classification latency averaged under 45 milliseconds per transaction on standard server hardware.

Deployment: The system was deployed on a server running Flask and Gunicorn behind an Nginx reverse proxy. The pipeline ran continuously without requiring a manual restart. All REST endpoints remained available throughout and the dashboard rendered correctly at different screen resolutions.

VII. RESULTS AND DISCUSSION

All three classifiers were evaluated on the held-out test set of 37,500 transactions. Table 1 presents the dataset summary statistics and Table 2 presents the performance metrics. Given the class imbalance, F1-Score and AUC-ROC are the primary evaluation metrics.

Table 1: Dataset Summary Statistics

Parameter	Value	Description
Total Transactions	2,50,000	Full dataset size
Fraudulent Transactions	8,000 (3.2%)	Labeled positive class
Legitimate Transactions	2,42,000 (96.8%)	Labeled negative class
Features per Record	28	After feature engineering
Training Set	1,75,000	70% stratified split
Validation Set	37,500	15% stratified split
Test Set	37,500	15% stratified split
Coverage Period	6 months	Jan–Jun 2024

Table 2: Model Performance Comparison on Test Set

Model	Accuracy	Precision	Recall	F1-Score	AUC-ROC
Random Forest	96.1%	94.3%	91.7%	92.98%	0.981
SVM (RBF)	94.8%	92.1%	89.4%	90.73%	0.967
XGBoost	96.7%	95.1%	93.2%	94.14%	0.988
Ensemble (Voting)	97.4%	96.2%	94.8%	95.50%	0.992

The ensemble voting approach achieved the highest performance across all metrics, with an accuracy of 97.4% and an AUC-ROC of 0.992. The high recall of 94.8% is particularly significant in the fraud detection context, as it reflects the system’s ability to correctly identify most actual fraud cases, minimizing the risk of undetected fraudulent transactions.

Feature importance analysis from the Random Forest model revealed that the most discriminative features were the time gap between consecutive transactions for the same RFID tag, the deviation from expected toll amount for the declared vehicle class, and the 7-day rolling anomaly score.

Error analysis on misclassified samples showed that the majority of false negatives involved sophisticated fraud attempts where individual transaction attributes fell within normal ranges but the cumulative behavioral pattern over time was anomalous. This suggests that sequence-based features could further improve recall in future iterations.

VIII. CONCLUSION

The proposed machine learning based toll fraud detection system was built to solve a specific and growing problem in highway toll management. The proposed system applies an ensemble of Random Forest, SVM, and XGBoost classifiers to transaction-level data, achieving an accuracy of 97.4% and an AUC-ROC of 0.992 on the held-out test set.

The Highway Dashboard complements the detection engine by providing transport administrators with actionable, real-time insight into fraud activity, reducing the time from detection to intervention. The integration of alert management, vehicle blacklisting, and analytical reporting into a single interface represents a meaningful improvement over the fragmented, manual monitoring processes currently used in many toll systems.

Going forward the plan is to extend the system to handle streaming transaction data at national scale using Apache Kafka and Apache Spark for distributed processing. Sequence modelling approaches such as LSTM networks will be investigated to better capture long-range behavioral patterns. Incorporating ANPR-based vehicle image verification is also a promising direction that could further strengthen detection capabilities.

IX. REFERENCES

JOURNAL REFERENCES

- [1] Bhattacharyya, S., Jha, S., Tharakunnel, K., and Westland, J. C., “Data Mining for Credit Card Fraud: A Comparative Study,” *Decision Support Systems*, vol. 50, no. 3, pp. 602–613, 2011.
- [2] Ozbayoglu, A. M., Gudelek, M. U., and Sezer, O. B., “Deep Learning for Financial Applications: A Survey,” *Applied Soft Computing*, vol. 93, 106384, 2020.
- [3] Kumar, R., Singh, A., and Gupta, N., “FASTag Vulnerabilities and Anomaly Detection in Electronic Toll Collection Systems,” *International Journal of Transportation Engineering*, vol. 7, no. 2, pp. 145–158, 2019.
- [4] Liu, F. T., Ting, K. M., and Zhou, Z. H., “Isolation Forest,” in *Proceedings of the 8th IEEE International Conference on Data Mining*, pp. 413–422, 2008.
- [5] Zhang, L., Wang, Q., and Chen, H., “Hybrid Deep Learning Architecture for Real-Time Fraud Detection in Toll Systems,” *IEEE*

Transactions on Intelligent Transportation Systems, vol. 22, no. 9, pp. 5531–5542, 2021.

- [6] Chawla, N. V., Bowyer, K. W., Hall, L. O., and Kegelmeyer, W. P., “SMOTE: Synthetic Minority Over-Sampling Technique,” *Journal of Artificial Intelligence Research*, vol. 16, pp. 321–357, 2002.
- [7] Breiman, L., “Random Forests,” *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [8] Chen, T. and Guestrin, C., “XGBoost: A Scalable Tree Boosting System,” in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 785–794, 2016.

WEB REFERENCES

- [1] NHAI, “FASTag Transaction Statistics and Toll Plaza Reports.” Available: <https://www.nhai.gov.in/>
- [2] Ministry of Road Transport and Highways. Available: <https://morth.nic.in/>
- [3] scikit-learn Documentation. Available: <https://scikit-learn.org/>
- [4] XGBoost Documentation. Available: <https://xgboost.readthedocs.io/>
- [5] Flask Documentation. Available: <https://flask.palletsprojects.com/>
- [6] React Documentation. Available: <https://react.dev/>
- [7] PostgreSQL Documentation. Available: <https://www.postgresql.org/docs/>
- [8] Python Software Foundation. Available: <https://docs.python.org/3/>

Copyright & License:

© Authors retain the copyright of this article. This work is published under the Creative Commons Attribution 4.0 International License (CC BY 4.0), permitting unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.