

Fraud Detection Project Documentation

1. Data Cleaning :

1.1 Missing Values

- **Identification:** Transactions where nameDest starts with 'M' were labeled as merchants.
- **Handling:** Filled missing values in oldbalanceDest and newbalanceDest with 0 for merchant transactions since these fields are irrelevant for merchants.

1.2 Outliers

- **Detection:**
 - **Z-score Method:** Used to identify outliers in the amount column with Z-scores greater than 3.
 - **IQR Method:** Applied to detect outliers. Transactions with amounts outside the range defined by 1.5 times the IQR were considered outliers.
- **Handling:**
 - **Removal:** Removed detected outliers from the dataset.
 - **Capping:** Capped outliers in the amount column at the 95th percentile value.

1.3 Multi-collinearity

- **Correlation Analysis:** Generated a correlation matrix for numeric features to check for multi-collinearity.
- **Visualization:** Created a heatmap to visualize correlations and identify highly correlated features.

2. Fraud Detection Model:

2.1 Model Development

- **Feature Selection:** Chose features amount, oldbalanceOrg, newbalanceOrig, oldbalanceDest, newbalanceDest, and isMerchant for the model.
- **Target Variable:** isFraud
- **Data Splitting:** Divided data into training (70%) and testing (30%) sets.
- **Model:** Random Forest Classifier
 - **Training:** Trained the Random Forest model with 100 estimators.
 - **Evaluation:**
 - **Classification Report:** Provided metrics such as precision, recall, and F1-score.

- **AUC-ROC Score:** Achieved an AUC-ROC score of 0.986, indicating high performance in distinguishing fraudulent transactions.

2.2 Feature Importance

- **Evaluation:** Determined feature importance using the Random Forest model.
- **Visualization:** Presented feature importance using a bar chart, highlighting oldbalanceOrg, newbalanceOrig, and newbalanceDest as the most influential features.

3. Variable Selection

3.1 Feature Selection

- **Initial Features:** Included amount, oldbalanceOrg, newbalanceOrig, oldbalanceDest, newbalanceDest, and isMerchant.
- **Recursive Feature Elimination (RFE):**
 - **Initial RFE:** Used to rank features and identify the most significant ones.
 - **Feature Reduction:** Removed isMerchant due to its lower importance compared to other features.

4. Model Performance

4.1 Evaluation Metrics

- **Classification Report:** Detailed metrics including precision, recall, and F1-score.
- **AUC-ROC Score:** Reported an AUC-ROC score of 0.986, indicating excellent model performance.

4.2 Cross-Validation

- **AUC Scores:** Conducted cross-validation and obtained AUC scores ranging from 0.983 to 0.992, with a mean score of 0.988.

5. Key Factors Predicting Fraudulent Transactions

5.1 Significant Features

- **oldbalanceOrg:** The sender's balance before the transaction.
- **newbalanceOrig:** The sender's balance after the transaction.
- **newbalanceDest:** The receiver's balance after the transaction.

5.2 Impact

- **Rationale:** These features are critical as they provide insights into the balances before and after transactions, which are key indicators of potential fraudulent activities.

6. Do These Factors Make Sense?

6.1 Sensibility of Factors

- **Yes, They Make Sense:**
 - **Balance Changes:** Significant changes in balances before and after transactions are strong indicators of fraud, as they may suggest unusual transaction patterns.
 - **Transaction Amounts:** Large amounts or unusual amounts are often linked to fraudulent transactions.

7. Prevention Measures for Infrastructure Updates

7.1 Integration

- **Model Deployment:** Implemented the fraud detection model into the transaction processing system for real-time or batch processing.
- **System Integration:** Ensured seamless integration into existing infrastructure to allow continuous monitoring.

7.2 Ongoing Monitoring

- **Performance Monitoring:** Regularly track model performance and accuracy.
- **Model Updates:** Periodically retrain the model with new data to adapt to evolving fraud patterns.

7.3 Security Measures

- **Data Security:** Implemented encryption and access controls for sensitive fraud detection data.
- **Incident Response:** Developed a response plan for handling detected fraud cases effectively.

8. Evaluating the Effectiveness of Prevention Actions

8.1 Metrics

- **Performance Metrics:** Monitored false positive rate, false negative rate, precision, and recall.
- **Impact Metrics:** Measured changes in transaction review rates and overall fraud detection rates.

8.2 Monitoring Strategies

- **Regular Audits:** Conduct periodic audits to assess the system's effectiveness.
- **Feedback Mechanism:** Established feedback loops for continuous improvement and adaptation of the fraud detection model.

Summary:

The project has successfully implemented a robust fraud detection system with effective data cleaning, model development, and evaluation processes. The model's performance and the rationale behind feature selection have been thoroughly validated. Preventive measures and ongoing monitoring strategies are in place to maintain and enhance the system's effectiveness. This comprehensive approach ensures that the system not only detects fraudulent activities accurately but also adapts to evolving fraud patterns and integrates seamlessly into the company's infrastructure.

