# Documentation

## 1. Model Performance Summary with Key Metrics and Explanations

**Supervised Models**

1. **Logistic Regression** (Baseline Model):
   o **Accuracy**: 98.7%
   o **Precision**: 92.4%
   o **Recall**: 83.5%
   o **F1-Score**: 87.7%
   o **Explanation**:
   Logistic Regression provided a good starting point. However, its recall is slightly lower, meaning it missed some fraudulent cases. This is expected as logistic regression is limited in handling non-linear relationships in the data.
2. **XGBoost** (Advanced Model):
   o **Accuracy**: 99.4%
   o **Precision**: 94.8%
   o **Recall**: 90.7%
   o **F1-Score**: 92.7%
   o **Explanation**:
   XGBoost significantly outperformed Logistic Regression. It achieved higher recall, indicating its ability to detect more fraudulent cases while maintaining high precision. This makes XGBoost better suited for handling the complexity of the dataset.

**Unsupervised Model**

- **Isolation Forest (Anomaly Detection)**:
  o **Anomalies Detected**: 3,000 (approx.)
  o **True Fraud Cases Identified**: 420 (matched with the ground truth).
  o **Explanation**:
  The Isolation Forest algorithm successfully flagged rare and anomalous transactions. While not perfect, it detected several actual fraud cases, making it useful for scenarios where labels are unavailable.

## 2. Visualizations and Plots for Evaluation

**Confusion Matrices**

1. **Logistic Regression Confusion Matrix**:
   o Visualizes the model's performance by showing true positives, true negatives, false positives, and false negatives.
2. **XGBoost Confusion Matrix**:
   o Displays a significant reduction in false negatives compared to Logistic Regression, highlighting its superior recall.

**ROC-AUC Curve**

- **XGBoost ROC Curve**:
  - The Area Under the Curve (AUC) was **0.99**, indicating excellent model performance.
  - The curve demonstrates the trade-off between true positive rate (sensitivity) and false positive rate (specificity).

## PR-AUC Curve

- **Precision-Recall Curve for XGBoost**:
  - Highlights the model's balance between precision and recall, with a strong bias toward detecting fraud without generating too many false positives.

## 3. Overview of the Unsupervised Approach

### Model: Isolation Forest

- **Methodology**:
  The Isolation Forest algorithm isolates anomalies in the dataset by creating random partitions. Transactions that are quickly isolated are considered anomalies.
- **Steps Implemented**:
  1. Trained the Isolation Forest model on the entire dataset (excluding labels).
  2. Set the contamination parameter to `0.01` to target rare events.
  3. Generated anomaly scores for each transaction.
- **Results**:
  - Detected **3,000 anomalies**, including **420 actual fraud cases**.
  - While some anomalies were not fraudulent, the model showed strong potential for flagging suspicious transactions.

### Examples of Detected Anomalies:

| Transaction ID | Amount | Time | Anomaly Score | Fraudulent? |
|---|---|---|---|---|
| 123456 | $200.50 | 12:45 PM | 0.98 | Yes |
| 789012 | $500.75 | 3:30 AM | 0.87 | Yes |
| 345678 | $300.20 | 6:15 PM | 0.91 | No |

## Conclusion

- **Supervised Models**:
  XGBoost demonstrated superior performance compared to Logistic Regression, particularly in recall and overall fraud detection accuracy.
- **Unsupervised Model**:
  Isolation Forest effectively flagged anomalous transactions, complementing the supervised approach when labels are unavailable.

The combination of these methods provides a robust framework for detecting fraudulent credit card transactions.