# Healthcare Provider Fraud Detection - Technical Report

**German International University of Applied Sciences**
**Informatics and Computer Science**
**Machine Learning, Winter 2025**
**Project 2: DataOrbit - HealthCare Provider Fraud Detection**

## Table of Contents

## 1. Project Overview

### 1.1 Objective

This project aims to develop an intelligent fraud detection system for the Centers for Medicare & Medicaid Services (CMS) to identify potentially fraudulent healthcare providers. The system must:

- Detect fraudulent providers from multi-table claims data
- Handle severe class imbalance (approximately 9.35% of providers are labeled fraudulent)
- Provide explainable predictions for investigators and regulators
- Demonstrate business value by prioritizing high-risk providers effectively

### 1.2 Problem Statement

Healthcare fraud costs the U.S. healthcare system over $68 billion annually. CMS can currently only investigate a small fraction of suspicious cases, allowing many fraudulent activities to go undetected. Existing systems rely on basic rule-based methods that capture obvious patterns but fail to identify more sophisticated fraud schemes.

**Types of Healthcare Fraud:**

- Billing for services never rendered
- Upcoding - billing for higher-cost procedures than those performed
- Unbundling - billing separately for procedures that should be combined
- Submitting claims for deceased patients
- Prescribing unnecessary treatments for financial gain
- Engaging in kickback or referral schemes

## 1.3 Dataset Description

The project uses the Healthcare Provider Fraud Detection dataset containing anonymized Medicare data:

**Files Included:**

- `Train_Beneficiarydata.csv` - Demographics, coverage, and chronic conditions for each patient (BeneID)
- `Train_Inpatientdata.csv` - Hospital admission claims with financial, procedural, and physician details
- `Train_Outpatientdata.csv` - Outpatient claim data (visits, tests, minor procedures)
- `Train_labels.csv` - Provider-level fraud labels (Yes or No)

**Key Identifiers:**

- `BeneID` links patients to claims
- `Provider` links claims to the fraud label

**Dataset Statistics:**

- Total Providers: 5,410
- Fraudulent Providers: 506 (9.35%)
- Non-Fraudulent Providers: 4,904 (90.65%)
- Total Beneficiaries: 138,556
- Inpatient Claims: 40,474
- Outpatient Claims: 517,737

---

# 2. Data Understanding & Exploration

## 2.1 Data Quality Assessment

**Missing Values Analysis:**

- **Beneficiary Data:** Only `DOD` (Date of Death) has missing values (137,135 missing), which is expected as most patients are alive
- **Inpatient Data:** Missing values in physician fields (`AttendingPhysician`, `OperatingPhysician`, `OtherPhysician`) and diagnosis codes, with `OtherPhysician` having the most (35,784 missing)
- **Outpatient Data:** Similar pattern with more missing values in diagnosis codes and physician fields, particularly `OperatingPhysician` (427,120 missing)

**Data Completeness:**

- All providers in labels have corresponding claims data
- 100% coverage of beneficiaries in claims data
- No orphaned records detected

## 2.2 Exploratory Data Analysis

**2.2.1 Target Variable Distribution**

The dataset exhibits significant class imbalance:

- **Fraud Rate:** 9.35% (506 providers)
- **Non-Fraud Rate:** 90.65% (4,904 providers)
- **Imbalance Ratio:** ~9.69:1

This imbalance required special handling strategies throughout the modeling process.

**2.2.2 Beneficiary Demographics**

**Age Distribution:**

- Mean Age: 73.7 years
- Median Age: 74.3 years
- Range: 26.1 to 101.0 years

**Gender Distribution:**

- Male: 42.9% (59,450)
- Female: 57.1% (79,106)

**Chronic Conditions (Top 5):**

1. Ischemic Heart Disease: 67.6%
2. Diabetes: 60.2%
3. Heart Failure: 49.4%
4. Depression: 35.6%
5. Alzheimer's: 33.2%

**2.2.3 Claim Amount Analysis**

**Inpatient Claims:**

- Total Claims: 40,474
- Total Amount: $408,297,020
- Average Claim: $10,087.88
- Median Claim: $7,000.00
- Max Claim: $125,000.00

**Outpatient Claims:**

- Total Claims: 517,737
- Total Amount: $148,246,120
- Average Claim: $286.33
- Median Claim: $80.00
- Max Claim: $102,500.00

**Key Insight:** Inpatient claims are significantly higher in value but lower in frequency compared to outpatient claims.

**2.2.4 Fraud vs Non-Fraud Comparison**

**Inpatient Claims:**

- Fraud providers: Mean $10,310.59, Median $7,000.00
- Non-Fraud providers: Mean $9,782.60, Median $7,000.00
- Difference: 5.4% higher average for fraud providers

**Outpatient Claims:**

- Fraud providers: Mean $287.19, Median $80.00
- Non-Fraud providers: Mean $285.84, Median $80.00
- Difference: 0.5% higher average for fraud providers

**Provider-Level Totals:**

- Fraud providers average total inpatient amount: $548,382.98
- Non-Fraud providers average total inpatient amount: $101,094.74
- **Critical Finding:** Fraud providers have approximately 5.4x higher total inpatient claim amounts

### 2.2.5 Geographic Patterns

Analysis revealed significant geographic variation in fraud rates:

- States analyzed: 52
- States with fraud rate > overall average: 43
- Highest fraud rate: State 21.0 (42.9%)
- Several states showed fraud rates above 30%, indicating potential regional patterns

### 2.2.6 Temporal Patterns

**Date Range:** November 27, 2008 to December 31, 2009

**Monthly Trends:**

- Monthly fraud rate relatively stable: 37.75% - 38.50%
- Month with highest fraud rate: August (38.50%)
- Month with lowest fraud rate: July (37.75%)

**Quarterly Distribution:**

- Q1: 89,078 non-fraud, 54,929 fraud claims
- Q2: 89,423 non-fraud, 55,320 fraud claims
- Q3: 86,872 non-fraud, 53,586 fraud claims
- Q4: 80,042 non-fraud, 48,961 fraud claims

## 2.3 Data Relationships

**Join Keys Verified:**

- `BeneID` successfully links beneficiaries to claims
- `Provider` successfully links claims to fraud labels
- All providers in labels have corresponding claims

- 100% beneficiary coverage in claims data

**Granularity Levels:**

- **Beneficiary Level:** Individual patient records
- **Claim Level:** Individual inpatient/outpatient claims
- **Provider Level:** Aggregated features per provider (modeling unit)

---

# 3. Data Preprocessing & Feature Engineering

## 3.1 Date Preprocessing

**Beneficiary Dates:**

- Converted DOB (Date of Birth) to datetime
- Calculated Age as of December 31, 2009
- Handled missing DOD (Date of Death) values

**Inpatient Dates:**

- Converted `ClaimStartDt`, `ClaimEndDt`, `AdmissionDt`, `DischargeDt` to datetime
- Calculated `LengthOfStay` = `DischargeDt` - `AdmissionDt`

**Outpatient Dates:**

- Converted `ClaimStartDt`, `ClaimEndDt` to datetime

## 3.2 Feature Engineering Strategy

Features were aggregated at the **provider level** to create a single record per provider, as this is the unit of prediction (fraud label is at provider level).

### 3.2.1 Inpatient Claim Features

**Financial Aggregations:**

- `IP_TotalClaimAmt`: Sum of all inpatient claim amounts
- `IP_AvgClaimAmt`: Mean claim amount
- `IP_MedianClaimAmt`: Median claim amount
- `IP_StdClaimAmt`: Standard deviation of claim amounts
- `IP_ClaimCount`: Total number of inpatient claims
- `IP_TotalDeductible`: Sum of deductibles paid
- `IP_AvgDeductible`: Mean deductible amount

**Length of Stay Features:**

- `IP_AvgLOS`: Average length of stay
- `IP_MedianLOS`: Median length of stay
- `IP_TotalLOS`: Total length of stay days

**Physician Features:**

- `IP_UniqueAttendingPhys`: Number of unique attending physicians
- `IP_UniqueOperatingPhys`: Number of unique operating physicians
- `IP_UniqueOtherPhys`: Number of unique other physicians

### 3.2.2 Outpatient Claim Features

**Financial Aggregations:**

- `OP_TotalClaimAmt`: Sum of all outpatient claim amounts
- `OP_AvgClaimAmt`: Mean claim amount
- `IP_MedianClaimAmt`: Median claim amount
- `OP_StdClaimAmt`: Standard deviation of claim amounts
- `OP_ClaimCount`: Total number of outpatient claims
- `OP_TotalDeductible`: Sum of deductibles paid
- `OP_AvgDeductible`: Mean deductible amount

**Physician Features:**

- `OP_UniqueAttendingPhys`: Number of unique attending physicians
- `OP_UniqueOperatingPhys`: Number of unique operating physicians
- `OP_UniqueOtherPhys`: Number of unique other physicians

### 3.2.3 Combined Features

**Total Aggregations:**

- `TotalClaimAmt`: Sum of IP and OP claim amounts
- `TotalClaimCount`: Sum of IP and OP claim counts
- `AvgClaimAmt`: Overall average claim amount

**Ratio Features:**

- `IP_OP_Ratio_Claims`: Ratio of inpatient to outpatient claim counts
- `IP_OP_Ratio_Amount`: Ratio of inpatient to outpatient claim amounts

### 3.2.4 Beneficiary Features (Aggregated by Provider)

**Patient Demographics:**

- `UniquePatients`: Number of unique patients per provider
- `AvgPatientAge`: Average patient age
- `MedianPatientAge`: Median patient age
- `StdPatientAge`: Standard deviation of patient age
- `PctMale`: Percentage of male patients

**Geographic Features:**

- `UniqueStates`: Number of unique states
- `UniqueCounties`: Number of unique counties

**Chronic Condition Percentages:**

- `PctAlzheimer`: Percentage of patients with Alzheimer's
- `PctHeartfailure`: Percentage of patients with heart failure
- `PctKidneyDisease`: Percentage of patients with kidney disease
- `PctCancer`: Percentage of patients with cancer
- `PctDiabetes`: Percentage of patients with diabetes

**Reimbursement Features:**

- `AvgIPAnnualReimb`: Average annual inpatient reimbursement per patient
- `AvgOPAnnualReimb`: Average annual outpatient reimbursement per patient

## 3.3 Missing Value Handling

**Strategy:**

- Count/Sum/Total features: Filled with 0 (provider has no claims of that type)
- Average/Mean features: Filled with median value
- Ratio features: Calculated with 0 replacement for missing denominators

## 3.4 Feature Selection

**Selection Criteria:**

- Features with significant correlation (p-value < 0.05) with target
- Features with correlation > 0.1 with target
- Removed highly correlated duplicate features (correlation > 0.8)

**Final Feature Set:** 43 features selected for modeling

**Top 10 Most Important Features (by correlation):**

1. `IP_TotalClaimAmt` (0.533)
2. `IP_TotalLOS` (0.526)
3. `IP_TotalDeductible` (0.525)
4. `IP_ClaimCount` (0.525)
5. `UniquePatients` (0.394)
6. `UniqueCounties` (0.373)
7. `OP_TotalClaimAmt` (0.338)
8. `OP_ClaimCount` (0.336)
9. `IP_UniqueOperatingPhys` (0.334)
10. `OP_TotalDeductible` (0.329)

**Key Insight:** Inpatient claim volume and total amounts are the strongest predictors of fraud, suggesting fraudulent providers tend to submit more and higher-value inpatient claims.

---

# 4. Class Imbalance Strategy

## 4.1 Imbalance Analysis

The dataset exhibits significant class imbalance:

- **Imbalance Ratio:** 9.69:1 (non-fraud:fraud)
- **Fraud Rate:** 9.35%

This imbalance poses challenges:

- Models may bias toward predicting the majority class
- Standard accuracy metrics become misleading
- Need for specialized metrics (Precision, Recall, F1, PR-AUC)

## 4.2 Strategies Implemented

### 4.2.1 Class Weighting

Applied class weights inversely proportional to class frequency:

- Non-Fraud weight: 1.0
- Fraud weight: ~9.69

**Models using class weighting:**

- Logistic Regression (Weighted)
- Random Forest (Weighted)
- Gradient Boosting (Weighted)

### 4.2.2 SMOTE (Synthetic Minority Oversampling Technique)

Applied SMOTE to training data to generate synthetic fraud samples:

- Oversamples minority class to balance distribution
- Creates synthetic examples using k-nearest neighbors

**Models using SMOTE:**

- Random Forest - SMOTE

### 4.2.3 Baseline Models (No Imbalance Handling)

Trained models without explicit imbalance handling to establish baseline performance:

- Decision Tree - Baseline
- Random Forest - Baseline
- Gradient Boosting - Baseline
- Logistic Regression - Baseline
- SVM - Baseline

## 4.3 Metric Selection for Imbalanced Data

**Primary Metrics:**

- **Precision:** Ability to correctly identify fraud (minimize false positives)
- **Recall:** Ability to catch all fraud cases (minimize false negatives)
- **F1-Score:** Harmonic mean of precision and recall

- **ROC-AUC:** Overall discriminatory ability
- **PR-AUC:** Performance in imbalanced setting (more informative than ROC-AUC)

**Rationale:**

- Accuracy is misleading with 90%+ majority class
- PR-AUC focuses on positive class performance
- Precision/Recall trade-off critical for business decisions

---

# 5. Algorithm Selection & Modeling

## 5.1 Algorithm Evaluation

Multiple algorithms were evaluated to balance interpretability, computational feasibility, robustness to imbalance, and suitability for mixed data types.

### 5.1.1 Decision Tree - Baseline

**Rationale:**

- Fast training and prediction
- Highly interpretable
- Baseline for comparison
- Handles non-linear relationships

**Configuration:**

- `max_depth=6`
- `min_samples_split=50`

### 5.1.2 Random Forest - Baseline & Variants

**Rationale:**

- Ensemble method improves stability
- Handles non-linear relationships
- Robust to overfitting
- Can handle mixed data types

**Variants Tested:**

- **Baseline:** `n_estimators=20`, `max_depth=6`
- **Robust:** Tuned to reduce overfitting
- **Weighted:** Class weights applied
- **SMOTE:** Trained on SMOTE-resampled data

### 5.1.3 Gradient Boosting - Baseline & Weighted

**Rationale:**

- Strong predictive power
- Handles complex patterns
- Sequential learning from errors

**Configuration:**

- `n_estimators=20`
- `learning_rate=0.15`
- `max_depth=3`

### 5.1.4 Logistic Regression - Baseline & Weighted

**Rationale:**

- Highly interpretable (coefficients)
- Fast training and prediction
- Linear relationships
- Good baseline for comparison

**Configuration:**

- `learning_rate=0.02`
- `max_iter=300`
- Requires feature scaling

### 5.1.5 SVM - Baseline

**Rationale:**

- Effective in high-dimensional spaces
- Handles non-linear relationships with kernels
- Requires scaled features

**Configuration:**

- `C=1.0`
- `gamma='scale'`

## 5.2 Model Training Process

**Data Splitting:**

- Training Set: 80% (4,328 providers)
- Test Set: 20% (1,082 providers)
- Stratified sampling to maintain class distribution

**Feature Scaling:**

- StandardScaler applied for algorithms requiring scaled features (Logistic Regression, SVM)
- Scaling fit on training data, applied to test data

**Validation Strategy:**

- Train/test split with stratification
- Cross-validation for hyperparameter tuning where applicable

## 5.3 Hyperparameter Tuning

**Approach:**

- Manual tuning based on computational constraints
- Focus on preventing overfitting
- Balance between model complexity and performance

**Key Parameters Tuned:**

- Tree depth (max_depth)
- Minimum samples per split
- Number of estimators (for ensemble methods)
- Learning rate (for gradient boosting)
- Regularization parameters

---

# 6. Model Evaluation

## 6.1 Evaluation Metrics

All models were evaluated using comprehensive metrics appropriate for imbalanced data:

- **Precision:** TP / (TP + FP) - Minimize false positives
- **Recall:** TP / (TP + FN) - Minimize false negatives
- **F1-Score:** 2 × (Precision × Recall) / (Precision + Recall)
- **ROC-AUC:** Area under ROC curve
- **PR-AUC:** Area under Precision-Recall curve

## 6.2 Model Comparison Results

| Model | Precision | Recall | F1-Score | ROC-AUC | PR-AUC |
|---|---|---|---|---|---|
| **Logistic Regression - Baseline** | **0.7191** | **0.6337** | **0.6737** | **0.9556** | **0.7502** |
| Random Forest - Robust | 0.7808 | 0.5644 | 0.6552 | 0.9584 | 0.7464 |
| Logistic Regression - Balanced & Interpretable | 0.4973 | 0.9109 | 0.6434 | 0.9620 | 0.7638 |
| Random Forest (Weighted) | 0.4712 | 0.8911 | 0.6164 | 0.9563 | 0.6834 |
| Random Forest - SMOTE | 0.4611 | 0.8812 | 0.6054 | 0.9464 | 0.6183 |
| Random Forest - Baseline | 0.7656 | 0.4851 | 0.5939 | 0.9239 | 0.7023 |
| Decision Tree - Baseline | 0.7206 | 0.4851 | 0.5799 | 0.7329 | 0.3977 |
| Logistic Regression (Weighted) | 0.4099 | 0.9010 | 0.5635 | 0.9557 | 0.7416 |

| Model | Precision | Recall | F1-Score | ROC-AUC | PR-AUC |
|-------|-----------|--------|----------|---------|--------|
| Gradient Boosting (Weighted) | 0.8125 | 0.2574 | 0.3910 | 0.9575 | 0.6763 |
| SVM - Baseline | 0.9545 | 0.2079 | 0.3415 | 0.7926 | 0.5097 |
| Gradient Boosting - Baseline | 0.0000 | 0.0000 | 0.0000 | 0.7638 | 0.5095 |

## 6.3 Best Model Selection

**Selected Model: Logistic Regression - Baseline**

**Justification:**

1. **Highest F1-Score (0.6737):** Best balance between precision and recall
2. **Strong ROC-AUC (0.9556):** Excellent discriminatory ability
3. **Good PR-AUC (0.7502):** Strong performance on imbalanced data
4. **Interpretability:** Coefficients provide clear feature importance
5. **Computational Efficiency:** Fast training and prediction
6. **Robust Performance:** Consistent across metrics

**Trade-offs:**

- **Precision (0.7191):** Good - correctly identifies 71.91% of flagged providers as fraudulent
- **Recall (0.6337):** Moderate - catches 63.37% of all fraudulent providers
- **Interpretability:** High - linear model with clear coefficients
- **Speed:** Very fast - suitable for real-time deployment

## 6.4 Alternative Model Analysis

**Random Forest - Robust:**

- Highest precision (0.7808) but lower recall (0.5644)
- Better for minimizing false positives
- Less interpretable than logistic regression

**Logistic Regression - Balanced & Interpretable:**

- Highest recall (0.9109) but lower precision (0.4973)
- Catches most fraud but flags many legitimate providers
- Useful when missing fraud is more costly than false positives

## 6.5 Confusion Matrix Analysis

**Logistic Regression - Baseline:**

```
                Predicted
 Actual      Non-Fraud   Fraud
 Non-Fraud      759        26
 Fraud           43        38
```

**Interpretation:**

- **True Positives:** 38 (correctly identified fraud)
- **True Negatives:** 759 (correctly identified non-fraud)
- **False Positives:** 26 (legitimate providers flagged)
- **False Negatives:** 43 (fraudulent providers missed)

## 6.6 Cost-Based Analysis

**Assumptions:**

- Cost of False Positive (FP): $100 (investigation cost for legitimate provider)
- Cost of False Negative (FN): $500 (missed fraud cost)

**Logistic Regression - Baseline Cost:**

- FP Cost: 26 × $100 = $2,600
- FN Cost: 43 × $500 = $21,500
- **Total Cost: $24,100**

**Random Forest - Robust Cost:**

- FP Cost: 1 × $100 = $100
- FN Cost: 23 × $500 = $11,500
- **Total Cost: $11,600** (Lower cost, but lower recall)

**Trade-off:** Random Forest - Robust has lower total cost but misses more fraud cases, which may not be acceptable from a regulatory perspective.

## 6.7 ROC and Precision-Recall Curves

**ROC Curve Analysis:**

- All models show strong discriminatory ability (ROC-AUC > 0.75)
- Logistic Regression - Baseline: ROC-AUC = 0.9556
- Random Forest - Robust: ROC-AUC = 0.9584 (slightly higher)

**Precision-Recall Curve Analysis:**

- More informative for imbalanced data
- Logistic Regression - Baseline: PR-AUC = 0.7502
- Logistic Regression - Balanced: PR-AUC = 0.7638 (higher recall focus)

# 7. Error Analysis

## 7.1 Error Distribution

Using the **Random Forest - Robust** model (best on validation set):

- **False Positives:** 1 sample (legitimate provider flagged as fraud)

- **False Negatives:** 23 samples (fraudulent providers missed)

## 7.2 False Positive Analysis

**Characteristics of False Positives:**

- Legitimate providers incorrectly flagged as fraudulent
- Likely patterns:
    - High total claim amounts (similar to fraud patterns)
    - High number of unique patients
    - High claim counts
    - Multiple physicians involved

**Case Study - False Positive Provider:**

- **IP_TotalClaimAmt:** $918,000 (very high)
- **IP_ClaimCount:** 91 claims
- **UniquePatients:** 203 patients
- **UniqueCounties:** 14 counties
- **IP_TotalLOS:** 552 days

**Analysis:** This provider exhibits patterns similar to fraudulent providers (high claim amounts, many patients, high LOS). However, this could be a legitimate large hospital or medical center. The model correctly identifies high-risk patterns but may need additional context (provider type, size) to distinguish legitimate large providers from fraudulent ones.

**Business Impact:**

- Unnecessary investigation cost
- Potential reputational damage to legitimate provider
- Resource waste

**Mitigation Strategies:**

1. Add provider type/size features
2. Use ensemble with multiple models
3. Implement threshold tuning based on provider characteristics
4. Add manual review step for high-value cases

## 7.3 False Negative Analysis

**Characteristics of False Negatives:**

- Fraudulent providers missed by the model
- Likely patterns:
    - Lower total claim amounts (below typical fraud threshold)
    - Moderate patient counts
    - Patterns that mimic legitimate providers

**Case Study - False Negative Provider 1:**

- **IP_TotalClaimAmt:** $57,080 (moderate)
- **IP_ClaimCount:** 9 claims
- **UniquePatients:** 96 patients
- **IP_TotalLOS:** 53 days
- **OP_TotalClaimAmt:** $41,880

**Analysis:** This provider has moderate claim amounts and patient counts, making it less obvious than high-volume fraud cases. The fraud may be more subtle (e.g., upcoding, unnecessary procedures) rather than high-volume billing.

**Case Study - False Negative Provider 2:**

- **IP_TotalClaimAmt:** $53,000 (moderate)
- **IP_ClaimCount:** 8 claims
- **UniquePatients:** 80 patients
- **OP_TotalClaimAmt:** $0 (no outpatient claims)

**Analysis:** This provider has very low activity (only 8 inpatient claims), which may not trigger fraud detection patterns. The fraud may involve quality issues (unnecessary procedures) rather than volume.

**Business Impact:**

- Financial loss from undetected fraud
- Continued fraudulent activity
- System vulnerability

**Mitigation Strategies:**

1. Feature engineering for subtle fraud patterns (ratios, anomalies)
2. Lower prediction threshold to increase recall
3. Ensemble methods to catch diverse fraud patterns
4. Time-series analysis for unusual patterns over time
5. Anomaly detection for outliers in feature space

## 7.4 Error Pattern Insights

**Common Characteristics of Errors:**

**False Positives:**

- High-volume, high-value providers
- Multiple geographic locations
- Diverse patient populations
- Legitimate large medical facilities

**False Negatives:**

- Moderate-volume providers
- Lower claim amounts
- Subtle fraud patterns
- Providers that mimic legitimate behavior

## 7.5 Refinement Recommendations

1. **Feature Engineering:**

   - Add provider type/size indicators
   - Create anomaly scores (deviation from provider norms)
   - Time-based features (trends, seasonality)
   - Patient-to-claim ratios
   - Geographic concentration metrics

2. **Model Improvements:**

   - Ensemble of multiple models (voting/stacking)
   - Threshold tuning based on business costs
   - Cost-sensitive learning with FP/FN cost weights
   - Two-stage model (high-risk screening + detailed analysis)

3. **Data Enhancements:**

   - Provider metadata (type, size, specialty)
   - Historical fraud patterns
   - Peer comparison features
   - Regulatory action history

4. **Deployment Strategy:**

   - Risk scoring with manual review for borderline cases
   - Continuous monitoring and model retraining
   - Feedback loop from investigations
   - A/B testing of different thresholds

---

# 8. Trials and Experiments

## 8.1 Experimentation Log

**Experiment 1: Baseline Models Without Imbalance Handling**

**Objective:** Establish baseline performance **Results:**

- Low recall on fraud class (0.21-0.49)
- High precision but missed most fraud cases
- **Insight:** Imbalance handling is critical

**Experiment 2: SMOTE Oversampling**

**Objective:** Balance classes through synthetic samples **Results:**

- Recall improved significantly (0.88-0.91)
- Precision decreased (0.46-0.50)
- **Insight:** SMOTE helps catch more fraud but increases false positives

**Experiment 3: Class Weighting**

**Objective:** Penalize misclassifying minority class **Results:**

- Similar to SMOTE: high recall, lower precision
- Faster than SMOTE (no data augmentation)
- **Insight:** Class weighting is efficient alternative to SMOTE

**Experiment 4: Hyperparameter Tuning**

**Objective:** Optimize model parameters **Results:**

- F1-score improved by 2-3% with tuning
- Reduced overfitting
- **Insight:** Careful tuning important but not sufficient alone

**Experiment 5: Feature Selection**

**Objective:** Reduce complexity and multicollinearity **Results:**

- Reduced from 50 to 43 features
- Minimal accuracy loss
- Faster training
- **Insight:** Feature selection improves efficiency without significant performance loss

**Experiment 6: Model Comparison**

**Objective:** Find best algorithm for this problem **Results:**

- Logistic Regression best balance of metrics
- Random Forest best precision
- Ensemble methods not significantly better
- **Insight:** Simpler models can outperform complex ones with proper tuning

## 8.2 Key Learnings

1. **Class imbalance is the primary challenge** - Required specialized handling
2. **Feature engineering is critical** - Provider-level aggregations essential
3. **Interpretability matters** - Logistic regression coefficients provide insights
4. **Precision-Recall trade-off** - Business context determines optimal balance
5. **Simple models can excel** - Logistic regression competitive with ensemble methods
6. **Error analysis reveals patterns** - Guides future improvements

---

# 9. Conclusion

## 9.1 Summary of Achievements

This project successfully developed a fraud detection system for healthcare providers with the following achievements:

1. **Comprehensive Data Analysis:** Thorough exploration of multi-table Medicare data, identifying key patterns and relationships
2. **Robust Feature Engineering:** Created 50 provider-level features capturing financial, operational, and demographic patterns
3. **Effective Imbalance Handling:** Implemented multiple strategies (SMOTE, class weighting) to address 9.35% fraud rate
4. **Model Comparison:** Evaluated 11 different model configurations across 5 algorithm types
5. **Best Model Selection:** Identified Logistic Regression - Baseline as optimal with F1-Score of 0.6737
6. **Error Analysis:** Conducted detailed analysis of false positives and false negatives with case studies

## 9.2 Model Performance

**Best Model: Logistic Regression - Baseline**

- **Precision:** 0.7191 (71.91% of flagged providers are fraudulent)
- **Recall:** 0.6337 (63.37% of fraudulent providers detected)
- **F1-Score:** 0.6737 (balanced performance)
- **ROC-AUC:** 0.9556 (excellent discriminatory ability)
- **PR-AUC:** 0.7502 (strong performance on imbalanced data)

## 9.3 Business Impact

**Potential Benefits:**

- Identifies 63.37% of fraudulent providers
- Reduces investigation workload by 71.91% precision
- Provides interpretable predictions for regulators
- Can be deployed in production with fast inference

**Limitations:**

- Misses 36.63% of fraudulent providers (false negatives)
- Flags 2.4% of legitimate providers (false positives)
- Requires continuous monitoring and retraining
- May need domain-specific refinements

## 9.4 Recommendations for Deployment

1. **Immediate Deployment:**

   - Use Logistic Regression - Baseline for initial deployment
   - Implement risk scoring with manual review for borderline cases
   - Set up monitoring dashboard for model performance

2. **Short-term Improvements:**

   - Add provider metadata features
   - Implement ensemble voting
   - Tune threshold based on investigation capacity
   - Create feedback loop from investigations

3. **Long-term Enhancements:**

   ○ Continuous model retraining with new data
   ○ Deep learning for complex pattern detection
   ○ Integration with external data sources
   ○ Real-time fraud detection pipeline

## 9.5 Final Thoughts

This project demonstrates that machine learning can effectively identify fraudulent healthcare providers, even with severe class imbalance. The selected Logistic Regression model provides a good balance between performance and interpretability, making it suitable for regulatory and investigative use. However, fraud detection is an ongoing challenge that requires continuous improvement, monitoring, and adaptation to evolving fraud patterns.

The error analysis reveals that the model struggles with:

- Legitimate high-volume providers (false positives)
- Subtle, moderate-volume fraud (false negatives)

Future work should focus on feature engineering to distinguish these cases and potentially implement a two-stage detection system for different fraud types.

---

# Appendix A: Feature List

**Final 43 Features Selected for Modeling:**

1. IP_TotalClaimAmt
2. IP_AvgClaimAmt
3. IP_MedianClaimAmt
4. IP_StdClaimAmt
5. IP_ClaimCount
6. IP_TotalDeductible
7. IP_AvgDeductible
8. IP_AvgLOS
9. IP_MedianLOS
10. IP_TotalLOS
11. IP_UniqueAttendingPhys
12. IP_UniqueOperatingPhys
13. IP_UniqueOtherPhys
14. OP_TotalClaimAmt
15. OP_AvgClaimAmt
16. OP_MedianClaimAmt
17. OP_StdClaimAmt
18. OP_ClaimCount
19. OP_TotalDeductible
20. OP_AvgDeductible
21. OP_UniqueAttendingPhys

22. OP_UniqueOperatingPhys
23. OP_UniqueOtherPhys
24. UniquePatients
25. AvgPatientAge
26. MedianPatientAge
27. StdPatientAge
28. PctMale
29. UniqueStates
30. UniqueCounties
31. PctAlzheimer
32. PctHeartfailure
33. PctKidneyDisease
34. PctCancer
35. PctDiabetes
36. AvgIPAnnualReimb
37. AvgOPAnnualReimb
38. TotalClaimAmt
39. TotalClaimCount
40. AvgClaimAmt
41. IP_OP_Ratio_Claims
42. IP_OP_Ratio_Amount
43. [Additional derived features]

# Appendix B: Model Configurations

**Logistic Regression - Baseline:**

- Algorithm: Logistic Regression (from scratch implementation)
- Learning Rate: 0.02
- Max Iterations: 300
- Feature Scaling: StandardScaler
- Class Weights: None (balanced by default handling)

**Random Forest - Robust:**

- Algorithm: Random Forest (from scratch implementation)
- N Estimators: 20
- Max Depth: 6
- Min Samples Split: 50
- Feature Scaling: None required

**End of Technical Report**