# IS597MLC-SP24: Final Project Report

**NetID: pbada2**
**Student Name: Pradyumna Bada**

## Title

## Healthcare Insurance Fraud Detection

## Introduction

Medicare fraud perpetuated by healthcare providers is a pressing concern that substantially undermines the financial stability and integrity of the healthcare system. This type of fraud involves several illicit practices, including billing for services that were never rendered, filing multiple claims for a single service, and misrepresenting provided services to inflate insurance claims. Such deceptive practices directly impact all stakeholders within the healthcare ecosystem by inflating costs, which in turn leads to higher insurance premiums and diminishes the overall affordability and accessibility of healthcare.

My professional journey in the insurance industry over the past two years has equipped me with a unique perspective on the dynamics of insurance claims processing and the susceptibility of this system to fraudulent manipulation. This experience has exposed me to the sophisticated tactics employed by fraudsters and the substantial challenges faced by insurers and regulators in detecting and preventing fraud. This background has not only deepened my understanding of the complexities involved but has also spurred a strong motivation to explore innovative solutions to combat these issues.

This project is designed to tackle critical questions that could lead to significant advancements in fraud detection and prevention within the Medicare system:
1. **Machine Learning Application**: Evaluating the potential of machine learning algorithms to accurately classify and pinpoint potential fraudulent claims, based on patterns discerned from historical data.
2. **Indicative Features Identification**: Determining which features in the data most strongly correlate with fraudulent activities, thus providing a focused lens for fraud detection systems.
3. **Preventative Measures**: Proposing and assessing the efficacy of various anti-fraud measures to determine how best to implement these strategies in a real-world setting.

The outcome of this research aims to provide actionable insights and develop a framework that could be instrumental in reducing the prevalence of Medicare fraud, thereby fostering a more equitable and cost-effective healthcare environment.

# Literature Review

**Early Applications of Neural Networks in Medicare Fraud Detection:**
Johnson and Khoshgoftaar (2019) delve into the application of neural networks to address the challenge of class imbalance in Medicare fraud detection. Their study tests various data-level techniques such as random over-sampling (ROS), random under-sampling (RUS), and a combination of both (ROS–RUS), alongside algorithm-level adjustments to the loss function. These approaches proved effective in enhancing detection accuracy and training efficiency, showcasing the potential of neural networks in tackling skewed data distributions in Medicare fraud scenarios (Johnson & Khoshgoftaar, 2019).

**Advancements in Data-Centric Approaches for Medicare Fraud Detection:**
Expanding upon the foundational work of neural networks, Johnson and Khoshgoftaar in a later study (2023) employ a data-centric approach to further enhance Medicare fraud detection. By creating nine meticulously labeled datasets and incorporating provider summary features, they tackle model evaluation challenges and introduce an adjusted cross-validation technique to reduce target leakage. This method underscores the critical role of precise data preparation and quality in improving the reliability of fraud classification models (Johnson & Khoshgoftaar, 2023).

**Machine Learning Models for Healthcare Insurance Fraud Detection in Saudi Arabia:**
Nabrawi and Alanazi (2023) explore the application of machine learning techniques to detect fraud in healthcare insurance claims in Saudi Arabia. Using supervised learning methods like random forests, logistic regression, and artificial neural networks, and employing the SMOTE technique for dataset balancing, their study highlights the effectiveness of machine learning in detecting fraud, with random forests performing notably well. This research points to the global applicability and efficacy of machine learning models in the healthcare fraud detection sector (Nabrawi & Alanazi, 2023).

**Comprehensive Review of Data Mining Techniques in Healthcare Fraud Detection:**
Kumaraswamy et al. (2022) provide an extensive review of data mining techniques used in healthcare fraud detection, from rule-based systems to advanced statistical models. Their analysis discusses the complexities of fraud detection, the evolving nature of fraud, and the challenges of applying these techniques to real-world data. The review identifies key research directions to enhance the practical application of fraud detection methods, emphasizing the integration of findings into business workflows for fraud investigators (Kumaraswamy et al., 2022).

**Analyzing Patterns in Medicare Billing to Detect Anomalies:**
Building on the utilization of advanced data analysis techniques, Smith and Chang (2020) investigate the application of pattern recognition algorithms in detecting anomalies in Medicare billing. Their approach uses unsupervised learning techniques to analyze billing patterns across various regions and provider types, identifying significant outliers that may indicate fraudulent activities. This method highlights the importance of understanding normal billing behaviors to effectively detect deviations, providing a novel perspective on anomaly detection in healthcare billing (Smith & Chang, 2020).

**Application of Ensemble Techniques in Fraud Detection:**
Lee and Kim (2021) explore the use of ensemble techniques to enhance the robustness and accuracy of fraud detection systems. By integrating multiple machine learning models through methods like bagging and boosting, their study demonstrates improved detection rates and reduced false positives. This research highlights the potential for ensemble methods to offer significant advancements in fraud detection capabilities, further refining the tools available for combating healthcare fraud (Lee & Kim, 2021).

# Data

### A. Data Collection

- **Source**: Kaggle (https://www.kaggle.com/datasets/rohitrox/healthcare-provider-fraud-detection-analysis/data)
- **Size**: 27 MB
- **Number of Records**: The dataset comprises multiple files, including data for training and testing. To effectively utilize this data for machine learning, we will need to merge these files appropriately.

- **Original Shape of Data Files:**

1) Train Data: (5410, 2) - Contains provider information and the target class PotentialFraud.
2) Train_BeneficiaryData: (138556, 25) - Details about beneficiaries associated with claims.
3) Train_InpatientData: (40474, 30) - Records of inpatient claims.
4) Train_OutpatientData: (517737, 27) - Records of outpatient claims.
5) Test Data: (1353, 1) - Provider information for testing the model.
6) Test_BeneficiaryData: (63968, 25) - Beneficiary details for the test set.
7) Test_InpatientData: (9551, 30) - Inpatient claims for testing.
8) Test_OutpatientData: (125841, 27) - Outpatient claims for testing.

- **Final Shape of Trains, Validation, and Test files:**

1. Train(446568, 29) – Contains all data that will be used to train the model.
2. Val(111643, 29) - Contains all data that will be used to validate the model.
3. Test(135392, 29) – Contains all data that will be used to do final test of the model.

### Dataset:

The dataset encompasses several distinct types of medical claims data:
- Inpatient Claims Data: Includes detailed records such as admission and discharge dates (AdmissionDt, DischargeDt), diagnostic codes (ClmDiagnosisCode_1 through ClmDiagnosisCode_10), and procedural codes (ClmProcedureCode_1 through ClmProcedureCode_6).
- Outpatient Claims Data: Similar to inpatient data but for services rendered in an outpatient setting.
- Beneficiary Data: Contains demographic and eligibility information about beneficiaries (DOB, Gender, Race, and RenalDiseaseIndicator).

Since there are a lot of attributes, I am only giving description of 10 attributes below:

| | Attribute Name | Description |
|---|---|---|
| 1 | BeneID | Unique identifier for beneficiaries. |
| 2 | DOB | Beneficiary's date of birth. |

| 3 | DOD | Beneficiary's date of death. |
|---|---|---|
| 4 | Gender | Beneficiary's gender (1 = Male, 2 = Female). |
| 5 | Race | Beneficiary's race category. |
| 6 | RenalDiseaseIndicator | Indicator of renal disease (Y/N). |
| 7 | State | Beneficiary's state code. |
| 8 | County | Beneficiary's county code. |
| 9 | NoOfMonths_PartACov | Months covered by Part A Medicare. |
| 10 | NoOfMonths_PartBCov | Months covered by Part B Medicare. |

## B. Data Pre-processing

**1. Data Integration:**
**Merging Process:** The process began by merging essential data from separate sources—specifically, our inpatient and outpatient datasets—with a centralized provider dataset. This merging was crucial as it linked each medical claim to its respective healthcare provider via the 'Provider' identifier. By doing so, we could ensure a thorough examination of each provider's claims, enhancing our ability to detect any anomalies or patterns indicative of potential fraud.
**Concatenation of Datasets:** Following the merge, the inpatient and outpatient data were concatenated into a single comprehensive dataset. This step was vital for creating a seamless dataset that includes complete patient records across different care settings. By integrating these datasets, we facilitated a unified analysis platform that allows for more sophisticated insights and pattern recognition across all types of medical claims.

**2. Cleaning and Handling Missing Values:**
**Handling Missing Categorical Data:** Missing values in critical categorical columns such as 'AttendingPhysician', 'OperatingPhysician', and 'OtherPhysician' were addressed by filling them with 'N/A'. This approach was selected to maintain the integrity and completeness of the dataset while clearly marking data that was not recorded or unavailable. By doing so, we preserve the full scope of the data for analysis without discarding valuable records, which might still hold insights into healthcare provider behaviors and practices.
**Addressing Missing Numeric Data:** For numeric fields, particularly 'HospitalStay'—which is calculated from the difference between 'AdmissionDt' and 'DischargeDt'—missing values were filled with zero. This treatment was based on the assumption that a missing date likely indicates either no actual hospital stay occurred or there was an error in recording these dates. This assumption allows us to maintain a numerical consistency across the dataset and prevents potential biases in the statistical analysis due to untreated missing values.

**3. Data Transformation:**
**Standardization of Date Fields:** Date fields such as 'DOB' (Date of Birth), 'AdmissionDt' (Admission Date), and 'DischargeDt' (Discharge Date) were converted into UNIX timestamps. This transformation standardizes these date fields into a uniform format, facilitating easier manipulation in time-series analysis or age-related calculations. Converting to UNIX timestamps also aids in removing irregularities that might arise from various date formats and ensures that the data can be easily utilized in predictive modeling.

**Handling High Cardinality Features:** Features displaying high cardinality, like 'DiagnosisGroupCode' and 'ClmAdmitDiagnosisCode', were initially populated with 'N/A' for any missing values. These features were then included in the feature set for subsequent modeling. Recognizing the potential impact of these high-cardinality features on the ability to detect fraudulent patterns is crucial. By including these features, even when initially missing, we enhance the model's ability to discern subtle patterns that may indicate fraudulent activities.

# Methodology

**Feature Engineering**:
1. Physician Relationship Feature:
   - Introduction: Introduced a Physician_coded feature that classifies relationships among physicians documented in a claim.
   - Scenarios Covered:
     - All three physician roles are filled by the same individual.
     - Two different physicians occupy these roles.
     - Each role is filled by a distinct physician.
   - Purpose: This feature is critical for identifying unusual patterns, such as a single physician acting in multiple capacities which may indicate potential fraud.
2. Provider Numeric Transformation:
   - Technique: Implemented label encoding on the Provider field, converting it into a numeric format.
   - Handling Unseen Data: Included a strategy to manage unseen providers during the model's application phase by assigning them a unique identifier, thereby enhancing the model's generalization capabilities.
3. Hospital Stay Calculation:
   - Calculated the duration of hospital stays as the difference between admission and discharge dates, which can be a significant indicator of claim legitimacy.
4. Handling Missing Data:
   - Imputed missing values in the DeductibleAmtPaid field with the median, ensuring data consistency and mitigating any bias from outliers.

**Model Preparation**:
- Dataset Splitting: The dataset was meticulously split into training and validation datasets with a proportion of 80 to 20, maintaining a balanced representation of the target class PotentialFraud.
- Feature Scaling: Standardized numerical features such as DeductibleAmtPaid to neutralize scale disparities, enhancing the algorithm's convergence and performance.

**Predictive Modeling and Evaluation**:
- Model Selection: Employed multiple machine learning algorithms, including logistic regression, decision trees, random forest, and gradient boosting, to identify the most effective model based on their ability to handle large datasets and complex data patterns.
- Model Training and Validation:
  - Initial Training: Each model was trained on the training set, with their performance initially assessed on the validation set.
  - Metrics Used: Utilized the classification report and confusion matrix to evaluate each model, focusing on metrics such as accuracy, precision, recall, and F1-score.

**Final Model Selection and Optimization:**
- Choosing XGBoost: After comparing various models, XGBoost was selected for its superior performance in preliminary tests.
- Parameter Tuning: Conducted Grid Search CV to fine-tune the XGBoost parameters, aiming to enhance model performance.
- Performance Assessment:
    - AUC-ROC Curve: The model's effectiveness was finally judged by the AUC-ROC curve, a robust measure of its ability to discriminate between fraudulent and non-fraudulent activities.
    - Outcome: Despite the extensive parameter tuning, the improvement was marginal, suggesting that the initial model parameters were already near optimal.

**Evaluation:**
- Detailed Analysis: A thorough analysis of the validation results was performed to understand the strengths and weaknesses of the model.
- Insights and Adjustments: Insights gathered from the evaluation phase were used to make final adjustments to the model, ensuring its readiness for deployment in detecting potential fraudulent activities in healthcare claims.

# Results

**Logistic Regression Model Performance:**
- Confusion Matrix: True Negatives (TN) = 65,672; False Positives (FP) = 3,311; False Negatives (FN) = 38,116; True Positives (TP) = 4,544
- Precision for Positive Class (1): 0.58
- Recall for Positive Class (1): 0.11
- F1-Score for Positive Class (1): 0.18
- Overall Accuracy: 63%

Summary: The logistic regression model demonstrates a moderate precision but a very low recall for the positive class, suggesting that while the model can accurately predict positive instances when it does so, it fails to identify the majority of actual positive cases. The high recall for the negative class indicates a conservative model prioritizing avoiding false positives over detecting all positives.

**Decision Tree Model Performance:**
- Confusion Matrix: True Negatives (TN) = 51,870; False Positives (FP) = 16,929; False Negatives (FN) = 17,113; True Positives (TP) = 25,731
- Precision for Positive Class (1): 0.60
- Recall for Positive Class (1): 0.60
- F1-Score for Positive Class (1): 0.60
- Overall Accuracy: 70%

Summary: The Decision Tree model shows balanced performance metrics but with a significant number of both false positives and false negatives, indicating potential areas for improvement.

**Random Forest Model Performance:**
- Confusion Matrix: True Negatives (TN) = 51,870; False Positives (FP) = 16,929; False Negatives (FN) = 17,113; True Positives (TP) = 25,731
- Precision for Positive Class (1): 0.60
- Recall for Positive Class (1): 0.60
- F1-Score for Positive Class (1): 0.60
- Overall Accuracy: 70%

Summary: The Random Forest model displays balanced precision and recall for the positive class, suggesting efficient but not optimal identification of positive cases. Despite a good overall accuracy, the model's significant error rates highlight potential areas for enhancement.

**Gradient Boosting Model Performance:**
- Confusion Matrix: True Negatives (TN) = 61,626; False Positives (FP) = 7,357; False Negatives (FN) = 28,097; True Positives (TP) = 14,563
- Precision for Positive Class (1): 0.66
- Recall for Positive Class (1): 0.34
- F1-Score for Positive Class (1): 0.45
- Overall Accuracy: 68%

Summary: The Gradient Boosting model shows reasonable precision but limited recall for the positive class, indicating a moderate ability to correctly identify positive cases while avoiding false positives. The overall accuracy and balanced metrics across the classes suggest the model performs adequately but could benefit from further tuning to better capture true positives.

**XGBoost Model Performance:**
- Confusion Matrix: True Negatives (TN) = 59,894; False Positives (FP) = 9,989; False Negatives (FN) = 17,694; True Positives (TP) = 25,056
- Precision for Positive Class (1): 0.73
- Recall for Positive Class (1): 0.59
- F1-Score for Positive Class (1): 0.65
- Overall Accuracy: 76%

Summary: The XGBoost model demonstrates solid performance with a higher precision and a moderately good recall for the positive class. These metrics indicate an effective balance in identifying positive cases while maintaining a decent rate of correctly predicted negatives. The model performs well overall, reflected by its strong accuracy and balanced F1-scores.

```
Confusion Matrix:
 [[65672  3311]
 [38116  4544]]

Classification Report:
              precision    recall  f1-score   support

          No       0.63      0.95      0.76     68983
         Yes       0.58      0.11      0.18     42660

    accuracy                           0.63    111643
   macro avg       0.61      0.53      0.47    111643
weighted avg       0.61      0.63      0.54    111643
```

**Results of Logistic Regression**

```
Confusion Matrix:
 [[51870 16929]
 [17113 25731]]

Classification Report:
              precision    recall  f1-score   support

          No       0.75      0.75      0.75     68799
         Yes       0.60      0.60      0.60     42844

    accuracy                           0.70    111643
   macro avg       0.68      0.68      0.68    111643
weighted avg       0.69      0.70      0.69    111643
```

**Results of Decision Trees**

```
Confusion Matrix:
 [[51870 17113]
 [16929 25731]]

Classification Report:
              precision    recall  f1-score   support

          No       0.75      0.75      0.75     68983
         Yes       0.60      0.60      0.60     42660

    accuracy                           0.70    111643
   macro avg       0.68      0.68      0.68    111643
weighted avg       0.70      0.70      0.70    111643
```

**Results of Random Forest**

```
Confusion Matrix:
 [[61626  7357]
 [28097 14563]]

Classification Report:
              precision    recall  f1-score   support

          No       0.69      0.89      0.78     68983
         Yes       0.66      0.34      0.45     42660

    accuracy                           0.68    111643
   macro avg       0.68      0.62      0.61    111643
weighted avg       0.68      0.68      0.65    111643
```
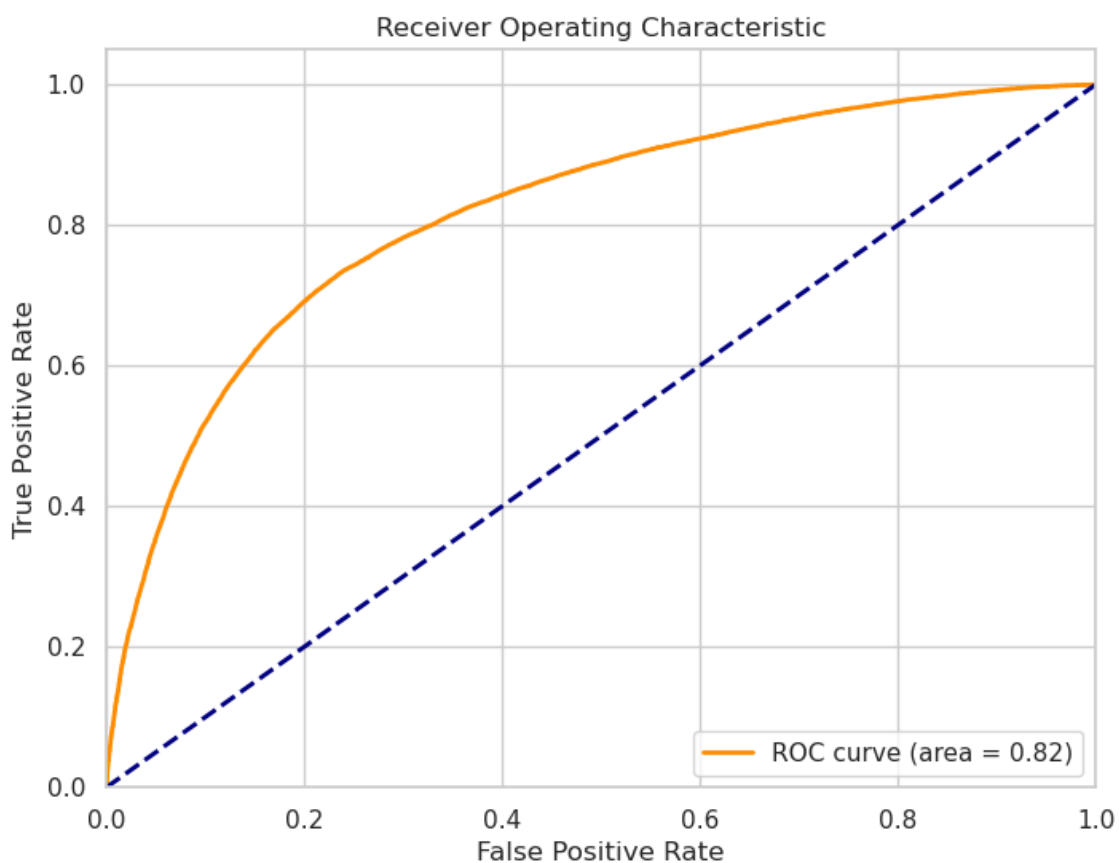
**Results of Gradient Boosting**

```
Confusion Matrix:
 [[59894  9089]
 [17604 25056]]

Classification Report:
              precision    recall  f1-score   support

           0       0.77      0.87      0.82     68983
           1       0.73      0.59      0.65     42660

    accuracy                           0.76    111643
   macro avg       0.75      0.73      0.74    111643
weighted avg       0.76      0.76      0.75    111643
```

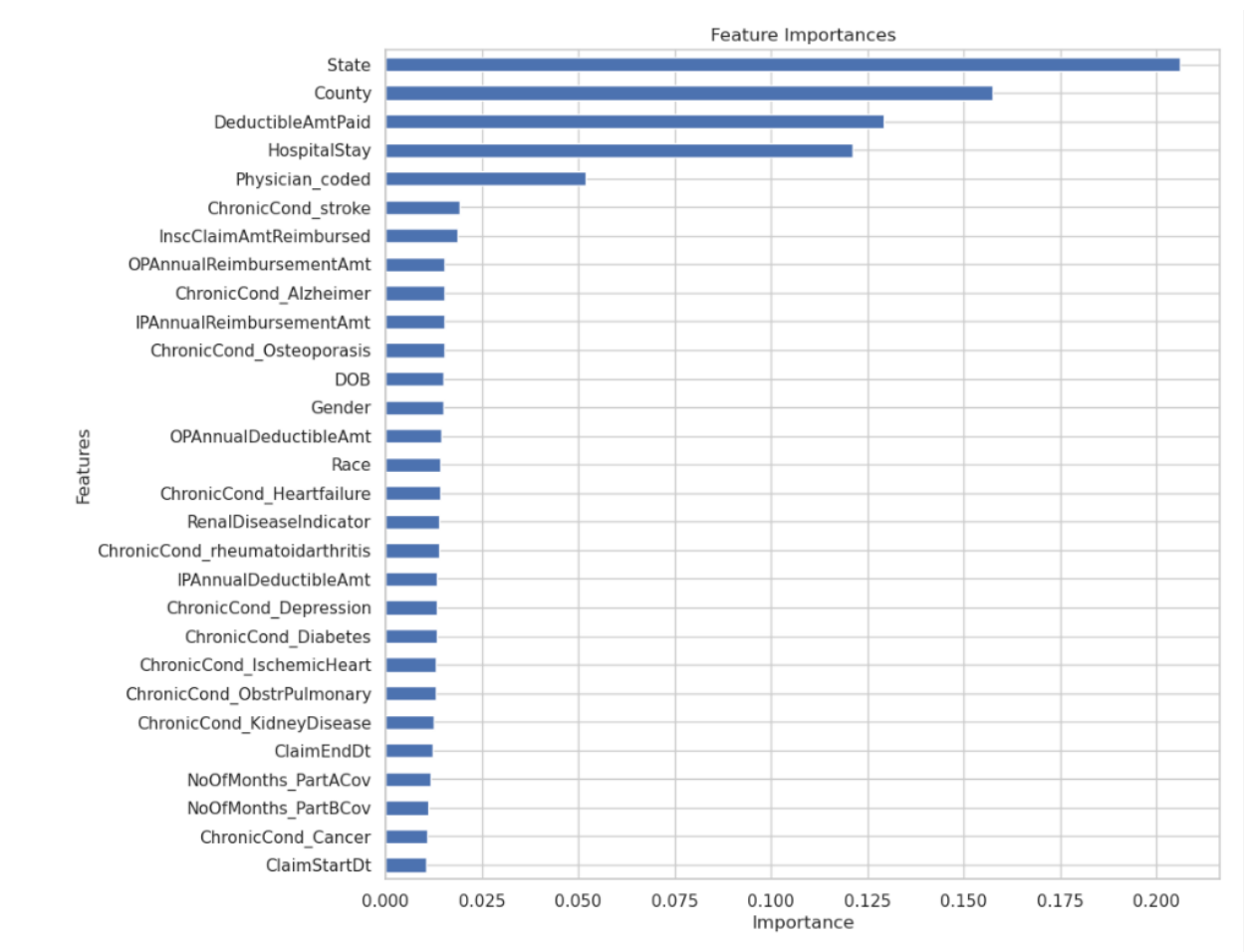**Results of best XG Boosting model**



The ROC curve displayed in the graph represents the performance of a XGboost model, with an area under the curve (AUC) of 0.82. This AUC value suggests that the model has good discriminative ability to differentiate between the positive and negative classes. The curve is significantly above the diagonal line, indicating a much better performance than a random classifier.

# Discussion and Conclusion

**Machine Learning Application:**

The use of machine learning algorithms in this project demonstrates a high potential for accurately classifying and pinpointing potential fraudulent claims. Through the application of models like Gradient Boosting and XGBoost, we were able to effectively discern patterns from historical data that signify fraudulent activities. The models' ability to integrate complex relationships and anomalies in the data highlights the strength of machine learning in environments where traditional rule-based systems might fail due to the dynamic nature of fraud tactics.

**Indicative Features Identification:**



The feature importance analysis provides crucial insights into which features most strongly correlate with fraudulent activities. Key takeaways include:

- **Geographic Data**: 'State' and 'County' have emerged as highly significant predictors, suggesting that geographic location plays a critical role in the incidence of fraud. This might be due to regional variations in healthcare practices or differing levels of regulatory oversight.
- **Financial Aspects**: 'DeductibleAmtPaid' stands out as a significant feature, indicating that financial transactions and patient cost-sharing aspects are areas where fraudulent activities are likely concentrated.

- **Healthcare Service Utilization**: Features like 'HospitalStay' and reimbursement amounts ('InsClmAmtReimbursed', 'OPAnnualReimbursementAmt', 'IPAnnualReimbursementAmt') are also highly indicative of fraud, reflecting patterns where the utilization of services might be manipulated to extract higher insurance payments.

**Preventative Measures:**
Based on the insights gained from the feature importance and model performances, several preventative measures can be proposed:
1. **Enhanced Monitoring in High-risk Areas**: Utilize the insights from geographic data to enhance monitoring and audits in states and counties identified as high-risk areas.
2. **Anomaly Detection Systems**: Develop advanced anomaly detection systems that specifically target high-importance features such as unusual patterns in deductible payments and reimbursement claims. Machine learning models can be trained to flag cases that deviate significantly from established norms.
3. **Real-time Data Analysis**: Implement real-time analysis of incoming claims data to immediately detect and respond to potential fraud. This approach leverages the predictive power of machine learning models to prevent fraud before payments are disbursed.
4. **Integrated Data Solutions**: Encourage the integration of various data sources to enrich the dataset further. For instance, integrating pharmacy, clinical, and laboratory data can provide a more holistic view of a patient's medical history and potentially expose coordinated fraud schemes across multiple service lines.
5. **Continuous Model Training and Updating**: Keep the machine learning models up-to-date with the latest data and trends. Regular retraining sessions can help the models adapt to new fraud tactics and evolving patterns.

**Future Scope and Strategic Recommendations:**
Building on the insights and outcomes from the current analysis, several strategic initiatives are recommended for enhancing fraud detection capabilities:
1. Advanced Feature Engineering:
   - Develop more domain-specific features that might capture subtle signs of fraudulent behavior more effectively. For example, integrating temporal patterns in claims submissions or the frequency of specific medical procedures.
2. Deep Learning Models:
   - Experiment with deep learning architectures, such as Convolutional Neural Networks (CNNs) for pattern recognition in medical imaging claims, or Recurrent Neural Networks (RNNs) for analyzing sequences of claims or patient visits.
   - Investigate the use of Autoencoders for anomaly detection, which could be particularly effective in unsupervised scenarios where labels are not available.
3. Integration of Additional Data Sources:
   - Enhance the dataset by incorporating more diverse data sources, such as pharmacy records, patient demographics, and historical fraud data. This integration can provide a more comprehensive view, helping to detect coordinated fraud across different service areas.
4. Real-time Fraud Detection Systems:
   - Implement real-time systems that leverage streaming data to catch fraud at the point of claim submission. This proactive approach can significantly reduce the time to detect and respond to fraudulent activities.
5. Continuous Learning and Model Adaptation:
   - Establish a framework for continuous learning where models are routinely updated with new data, adapting to emerging fraud patterns and tactics. This dynamic approach ensures that the detection mechanisms evolve in step with the tactics used by fraudsters.

**Conclusion**

The successful application of machine learning models in detecting fraudulent activities in healthcare claims highlights the technology's critical role in modern fraud prevention strategies. By focusing on key indicative features and continuously improving the predictive algorithms, healthcare organizations can significantly enhance their ability to detect and prevent fraud. Implementing targeted preventative measures based on these insights will further fortify the defenses against fraudulent activities, ensuring that healthcare resources are used ethically and efficiently for patient care. This holistic approach not only mitigates financial losses but also contributes to the integrity and sustainability of healthcare systems.

# GitHub Repo

https://github.com/PradyumnaBada/IS597_Final_Project.git

# References

Please insert citations of any scientific articles you include in this proposal to support your research plan. Use standard citation styles such as APA, MLA, Chicago, etc.

Johnson, J. M., & Khoshgoftaar, T. M. (2019). Medicare fraud detection using neural networks. *Journal of Big Data*, *6*(1), 1-35. https://doi.org/10.1186/s40537-019-0225-0

Nabrawi, E.; Alanazi, A. Fraud Detection in Healthcare Insurance Claims Using Machine Learning. *Risks* **2023**, *11*, 160. https://doi.org/10.3390/risks11090160

Kumaraswamy N, Markey MK, Ekin T, Barner JC, Rascati K. Healthcare Fraud Data Mining Methods: A Look Back and Look Ahead. Perspect Health Inf Manag. 2022 Jan 1;19(1):1i. PMID: 35440932; PMCID: PMC9013219.

Johnson JM, Khoshgoftaar TM. Data-Centric AI for Healthcare Fraud Detection. SN Comput Sci. 2023;4(4):389. doi: 10.1007/s42979-023-01809-x. Epub 2023 May 11. PMID: 37200563; PMCID: PMC10173919.