

WCD Machine Learning Project

...

# Healthcare Provider Fraud Detection Analysis

Cheryl Chien

05.13.2023

- Intro
- ML in Healthcare & Types of Health Insurance Fraud
- ML workflow
- Target Variable
- Dataset EDA
- Features Used for Modeling
- Data Preprocessing
- Model Selection
- Summary

# Intro

## Why this topic?

- Want to work on healthcare IT and analysis
- Financial factors play an important role in healthcare decisions.
- Preventing health insurance fraud can ensure that resources are used fairly
- Prevent insurance premium increase

# ML in Healthcare & Types of Health Insurance Fraud

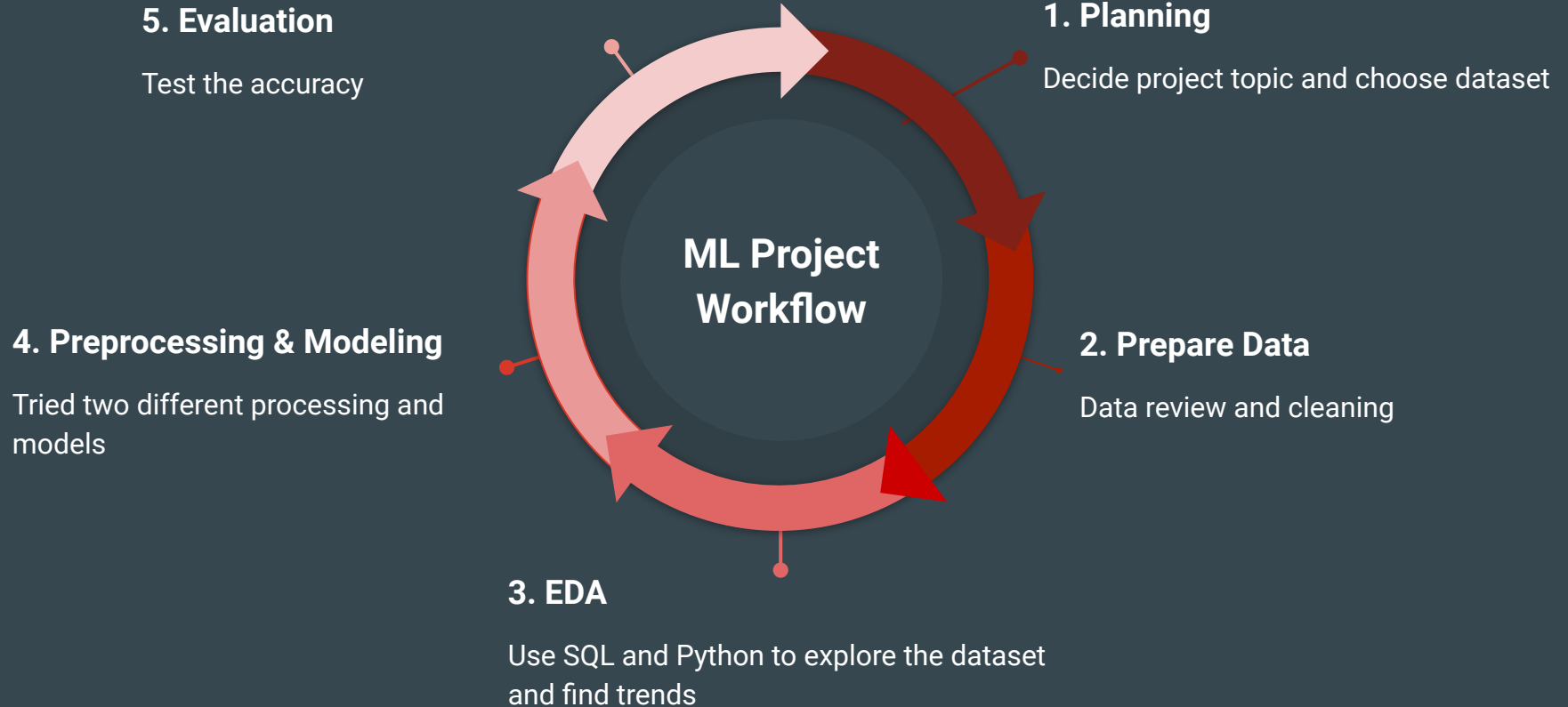
## ML in Healthcare

1. Fraud prediction
2. Risk assessment
3. Dx and Rx evaluation
4. Healthcare workforce automation

## Health Insurance Fraud

1. By provider, beneficiary, prescription
2. Billing for services that were not provided.
3. Duplicate submission of a claim for the same service.
4. Misrepresenting the service provided.
5. Charging for a more complex or expensive service than was actually provided.

# ML Project Workflow



# Target Variable

1. **Predict the potentially fraudulent providers-- YES/NO**
2. Discover patterns of potentially fraudulent providers
3. Health insurance beneficiaries who are at greater risk of abuse, ex. Age, Health condition...

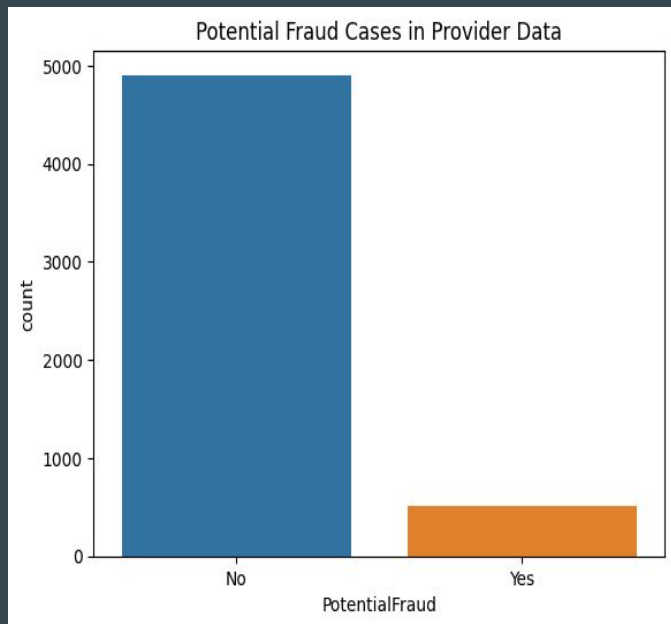
# About Dataset

- **Data source:** Kaggle [HEALTHCARE PROVIDER FRAUD DETECTION ANALYSIS](#)
- 4 Tables: Provider data, Beneficiary Details Data, Inpatient Data, Outpatient Data

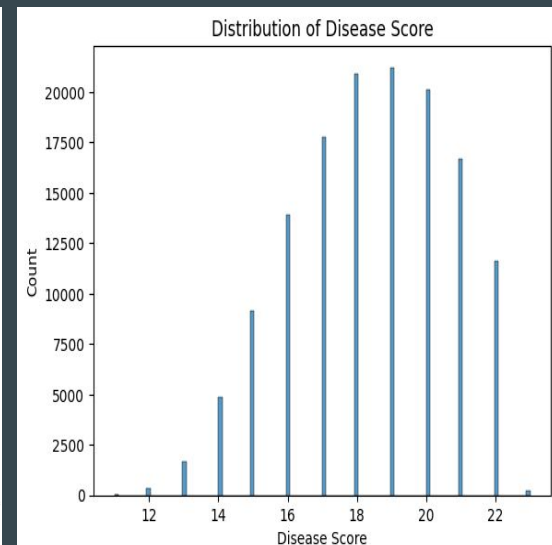
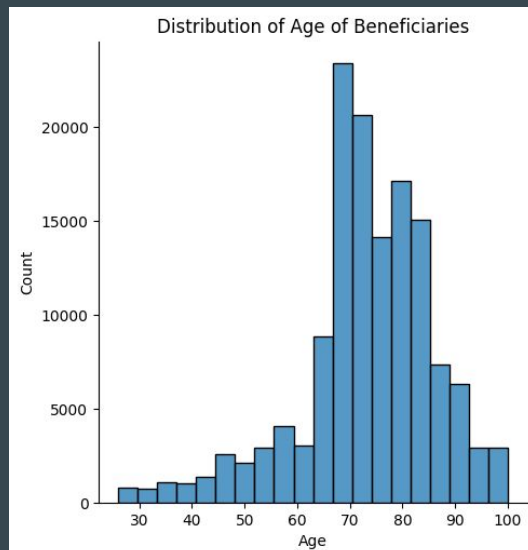
01	Provider Data	<ul style="list-style-type: none"><li>• Shape: 5410, 2</li><li>• 5410 unique providers</li><li>• Potential Fraud Yes/No</li></ul>
02	Beneficiary Details Data	<ul style="list-style-type: none"><li>• Shape: 138556, 25</li><li>• 138556 unique beneficiaries</li><li>• KYC like demographic, health conditions</li></ul>
03	Inpatient Data	<ul style="list-style-type: none"><li>• Shape: 40474, 30</li><li>• 40474 unique claimID</li><li>• Dates, provider, diagnosis, \$ Reimburse</li></ul>
04	Outpatient Data	<ul style="list-style-type: none"><li>• Shape: 517737, 27</li><li>• 517737 unique claimID</li><li>• Dates, provider, diagnosis, \$ Reimburse</li></ul>

# Exploratory Data Analysis--1

## Provider Data



## Beneficiary Data

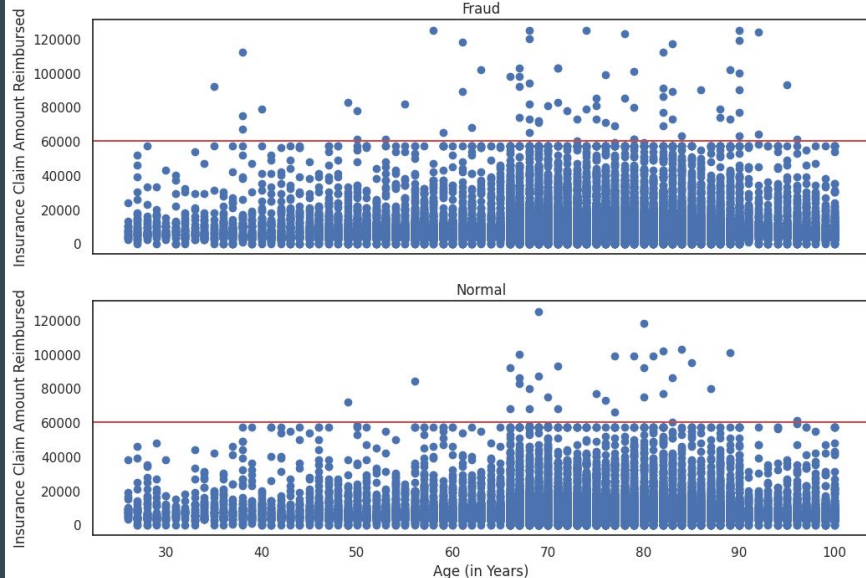




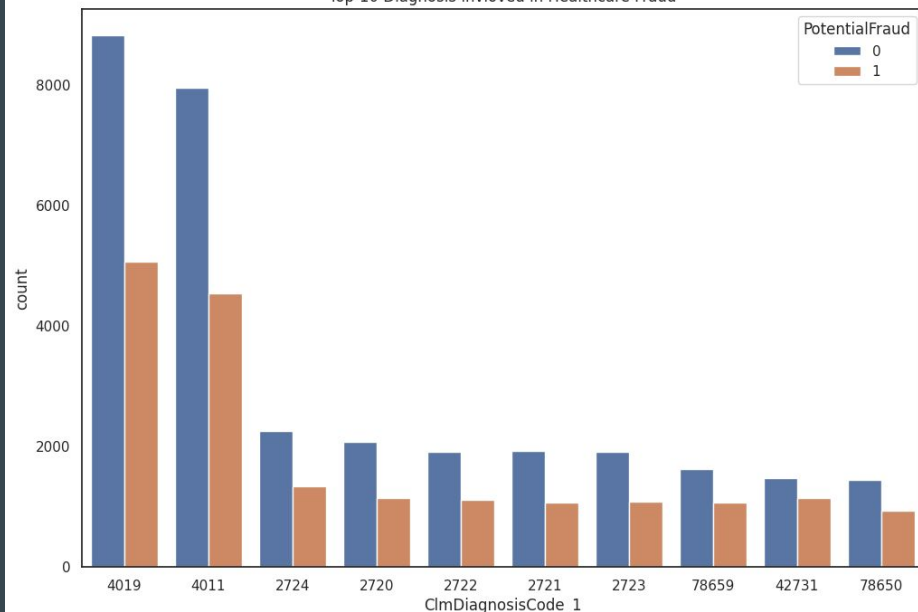
# Exploratory Data Analysis--2

## Joining all table: Provider, Beneficiary, IP, OP Data

Insurance Claim Amount Reimbursed Vs Age



Top-10 Diagnosis involved in Healthcare Fraud



# Features Used for Modeling

Original Merged dataset shape: 558211 rows, 61 features  
dtypes: bool(1); datetime64[ns](4), float64(8), int64(27), object(21)

## >> Drop features

- **10 Unrelated features:** 'BeneID', 'ClaimID', 'ClaimStartDt', 'ClaimEndDt', 'has\_claimed\_op', 'DOB', 'DOD', 'NoOfMonths\_PartACov', 'NoOfMonths\_PartBCov', 'has\_claimed\_ip'
- **20 More than 50% null:** 'AttendingPhysician', 'OperatingPhysician', 'OtherPhysician', 'ClmDiagnosisCode\_4', 'ClmDiagnosisCode\_5', 'ClmDiagnosisCode\_6', 'ClmDiagnosisCode\_7', 'ClmDiagnosisCode\_8', 'ClmDiagnosisCode\_9', 'ClmDiagnosisCode\_10', 'ClmProcedureCode\_1', 'ClmProcedureCode\_2', 'ClmProcedureCode\_3', 'ClmProcedureCode\_4', 'ClmProcedureCode\_5', 'ClmProcedureCode\_6', 'ClmAdmitDiagnosisCode', 'AdmissionDt', 'DischargeDt', 'DiagnosisGroupCode'

# Dataset used for Training

After dropping the features, the dataset will be used for training has 558211 rows and 31 features

>> Separate X and y data

X = Drop Provider and PotentialFraud columns

y = PotentialFraud

>>> Split the data into train and test sets

X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=0.2, random\_state=42)

# Data Preprocessing 1

1. Replace PotentialFraud value from YES/NO to 1 and 0
2. Replace RenalDiseaseIndicator value from YES/ 0 to 1 and 0

3. **TargetEncoder:**

Transform Object ClmDiagnosisCode\_1, ClmDiagnosisCode\_2, and ClmDiagnosisCode\_3

4. **Filling missing value:**

ClmDiagnosisCode\_1, ClmDiagnosisCode\_2, and ClmDiagnosisCode\_3, and DeductibleAmtPaid

# Model Selection

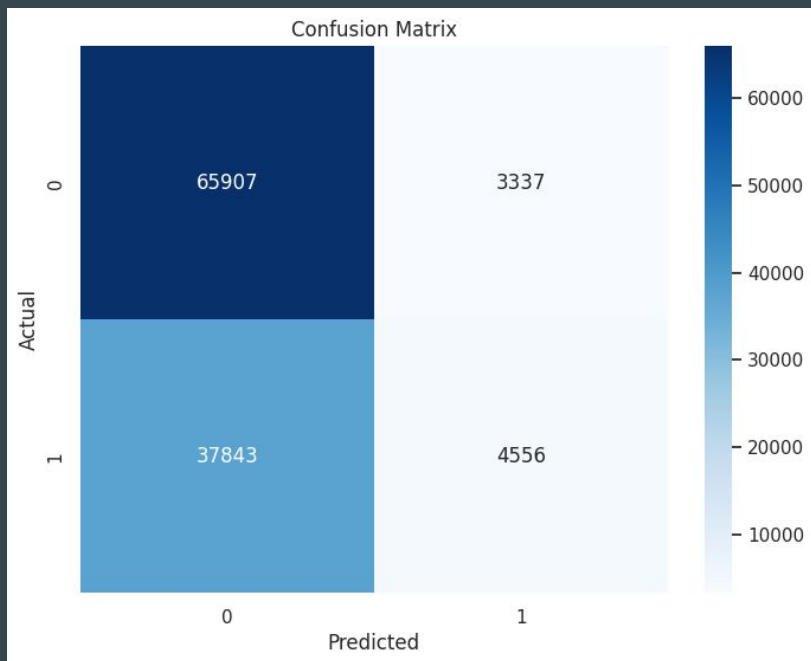
**Target variable:** Predict the potentially fraudulent providers-- **YES/NO**

**Features are labeled>> Supervised machine learning**

**Tried Logistic Regression and Random Forest Classifier**

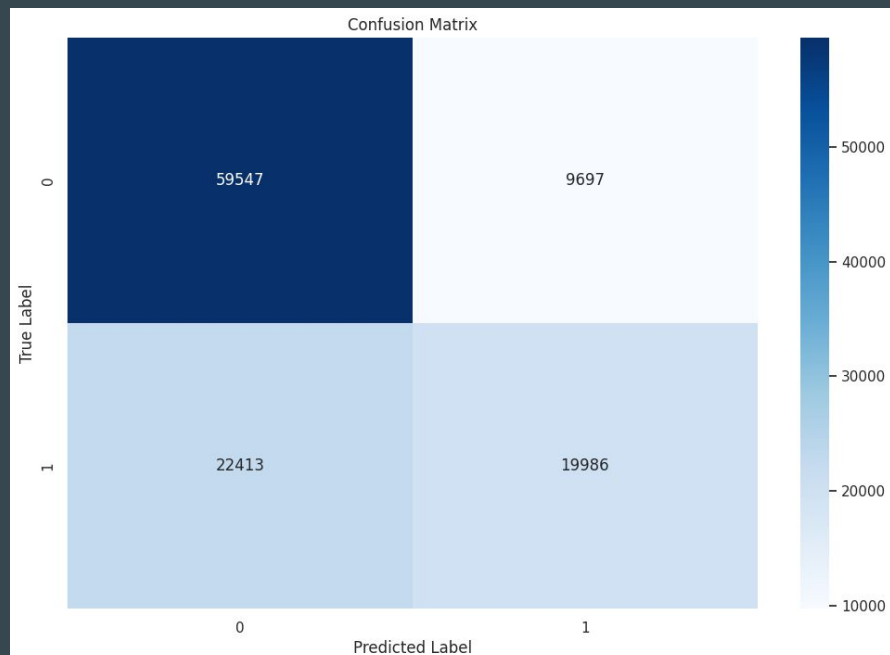
# Logistic Regression

Accuracy: 0.6311457055077345

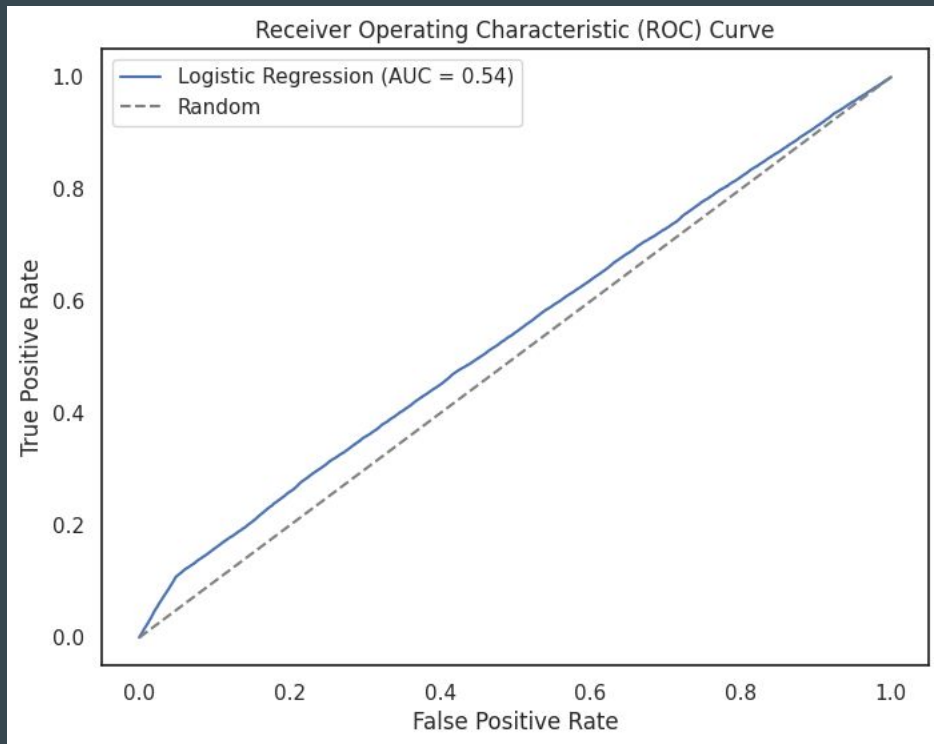


# Random Forest Classifier ✓

Accuracy: 0.7123868043674928



# Evaluation



# Data Preprocessing 2

1. StandardScaler
2. Imblearn.over\_sampling
  - > RandomOverSampler
  - >> SMOTE

Logistic Regression	Random Forest Classifier
Accuracy: 0.5517049882213843	Accuracy: 0.6877457610418924



# Summary

## 1. From EDA:

- Overall fraud rate in inpatient and outpatient data is about 38.12%
- Older beneficiaries have a slight higher risk to be involved in fraud
- Most common codes those also involved in fraud are 4019, 4011

## 2. Data processing & Model selection

	Logistic Regression	Random Forest Classifier
Data Processing 1	Accuracy: 0.6311457055077345	Accuracy: 0.7123868043674928
Data Processing 2	Accuracy: 0.5517049882213843	Accuracy: 0.6877457610418924

# Limitation

1. Not sure the source and date of the dataset
2. Missing definition of some features, ex. Race, State, Chronic diseases
3. Could drop more features provider and Has\_ChronicCondition
4. Could test more models like xgboostclassifier

**Thank You**

# Reference

1. [Health Care Fraud — FBI](#)
2. [Medicare Fraud & Abuse: Prevent, Detect, Report](#)
3. [HEALTHCARE PROVIDER FRAUD DETECTION ANALYSIS | Kaggle](#)
4. [GitHub - Atharvak19/Big-Data-Medicare-Fraud-Detection](#)