# **Big-Data-Medicare-Fraud-Detection**

#### INTRODUCTION

Healthcare is a significant industry in the U.S. with both private and government run programs. The expenses of medicinal services keep on ascending, to some degree because of the expanding populace of the old. U.S. health insurance spending from 2012 to 2014 has expanded by 6.7% to reach \$3 trillion and Medicare spending represents 20% of all human services spending in the U.S. at about \$600 billion. This rising older populace, joined with the expanded expenses of Medicare, need cost-cutting arrangements, where the decrease in misrepresentation is one approach to help recuperate costs and diminish by and large installments. The effect of health care extortion is evaluated to be between 3% to 10% of the country's absolute human services spending and proceeding to antagonistically affect the Medicare program and its recipients (NHCAA 2017). Government has introduced the projects, for example, the Medicare Fraud Strike Force (OIG 2017), authorized to help battle misrepresentation, however proceeded with endeavours are expected to more likely alleviate the impacts of extortion. These are our task's essential presentation and foundation; likewise, it shows that there are immense business open doors for Medicare Fraud Detection frameworks. This project is committed to building big data solutions with substantial applications at the crossing point of the health care and protection industry. This Capstone undertaking will construct a Medicare Fraud Detection model to investigate open information and foresee/identify the false Medicare supplier's dependent on misrepresentation designs, oddity examination, and geo-segment measurements. Our aim is to create a model that predicts the fraud in the health insurance industry using anomaly analysis and geo-demographic metrics. There will be several advantages if we are able to find out the fraud detection accurately. Here are the following advantages: • Protection and Cost deduction Recognizing extortion in advance methods no all the more pursuing the lawbreakers after installments have been made. alongside the insightful and lawful expenses. It's only an a lot more intelligent approach to work together. Furthermore, regardless of whether that cost reserve funds sway your own business primary concern or that of an open element with whom you are gotten, the outcomes must be acceptable ones. • Fraud Prevention becomes easy As our model learns we can help in providing additional information for accessing any fraud. It would provide a helping hand in finding possible patterns by comparing one another. • Payers Can Maintain Compliance with Government Prompt Payment Regs The measure of manual "pursuing" down of cases, regardless of whether pre-or post-installment, is essentially diminished. Examination can occur progressively, as cases are introduced, and goals accomplished a lot quicker. Regardless of whether you are a steward of open dollars or your very own steward business' monetary wellbeing, a prescient examination programming framework just bodes well.

#### PROBLEM STATEMENT

Building an innovative data science model that help in predicting fraud in medical insurance industry by using real time analysis and classification algorithm. This tool can be used by the government so that it benefits the patients, pharmacy, doctors which eventually help in gaining credibility to industry, tackle the increasing costs of healthcare and handle the impact of fraud. Medicinal services misrepresentation is a principle issue that causes generous fiscal misfortune in Medicare/Medicaid and protection industry. The Centres for Medicare and Medicaid Services (CMS) have arrangement Medicare Part D programs since 2006. CMS depends on it to identify and forestall extortion, waste and maltreatment in Part D program. Be that as it may, utilizing the customary techniques, the misrepresentation location is directed on arbitrary examples by human specialists. The results are the examples may be deluding or manual identification is exorbitant. As per Office of Inspector General report: Since 2006, the Medicare Fraud has quickly expanded. The extortion designs incorporate the accompanying four sorts: • Fraud by Service Providers (Doctors, hospitals, pharmacies) • Fraud by Insurance subscribers (patient or patient's employers) • Fraud by insurance carriers • Conspiracy Frauds (involved with all parties)

The main objectives of this project are: • Build a basic Data Model to show the connections among the distinctive datasets and distinguish the key capabilities for extortion recognitions • Build a thorough AI model to recognize misrepresentation design dependent on the various highlights: Service Providers (Doctors, Pharmacies), Insurance supporters (patients), Geo-segment and usually misuse drugs medicines • Setup a benchmark measurement to quantify and assess the test result • Market-prepared item

DATASETS The following public datasets were used: • Part D Prescriber Dataset • Excluded (LEIE) dataset • Payment Received dataset Dataset Link: • CMS Part D datasets: <a href="https://www.cms.gov/Research-Statistics-Data-and-Systems/Statistics-Trends-and-Reports/Medicare-Provider-Charge-Data/Part-D-Prescriber.html">https://www.cms.gov/Research-Statistics-Data-and-Systems/Statistics-Trends-and-Reports/Medicare-Provider-Charge-Data/Part-D-Prescriber.html</a> • LEIE Datasets: <a href="https://oig.hhs.gov/exclusions/exclusions\_list.asp">https://oig.hhs.gov/exclusions/exclusions\_list.asp</a> Office of Inspector General Reports <a href="https://oig.hhs.gov/reports-and-publications/portfolio/index.asp">https://oig.hhs.gov/reports-and-publications/portfolio/index.asp</a> • FDA datasets <a href="https://www.fda.gov/Drugs/InformationOnDrugs/ucm079750.htm#collapseOne">https://www.fda.gov/Drugs/InformationOnDrugs/ucm079750.htm#collapseOne</a>

# 1. Medicare Provider Utilization and Payment Data (Part D Prescriber)

This dataset contains 21 sections and incorporates the accompanying key data • All medicines accumulated at the doctor and medication level • All data on usage, installments, and submitted charges by National Provider Identifier (NPI), Healthcare Common Procedure Code, and Place of Service • All data on the doctor (NPI, Name, City, Practice, and so forth.) The general highlights will be utilized to explore the general misrepresentation discoveries • The quantity of various medications endorsed • The aggregate of the quantity of remedies, • The whole of the quantity of days recommended, • The total of the all out expense • The change of these three entireties amounts • The limit of these three wholes amounts The Prescribing design highlights to make sense of the endorse related extortion • The quantity of various medications recommended • The aggregate of the quantity of solutions • The entirety of the quantity of days recommended • The total of the complete expense • The change of these three entireties amounts • The limit of these three whole amounts The medication design highlights to group the

regularly misuse sedate related extortion • Average Number of Prescriptions per Beneficiary • Number of Pharmacies related with each prescriber • Average number of sorts of medications per recipient • Percentage of Schedule II/III medications (contain the potential maltreatment tranquilize fixings) • Percentage of Brand-Name drugs (costly medications) • Percentage of recipients with an over the top flexibly of a medication Different highlights, for example, the supplier area or claim to fame will be joined with the above highlights to identify the misrepresentation.

# 2. List of Excluded Individuals and Entities (LEIE) database

This database contains a rundown of people and substances that are prohibited from taking an interest in governmentally financed social insurance programs (for example Medicare) because of past medicinal services extortion. We could treat the LEIE dataset as the semi-named information, on the grounds that LEIE is the fraudster-based objective however not a misrepresentation one. We will investigate the LEIE information through the NPI highlight to get together with Part D or Payment dataset, the other way is that we should outline Part D dataset with LEIE rejection rules. In the LEIE dataset, the Column EXCLTYPE is prohibition rules. On the off chance that EXCLTYPE esteem rises to 1128(a)\*, at that point rejection years are five, 1128(c)(g)(i) is ten years and 1128(c)(g)(ii) is changeless prohibition. We may check on the off chance that any cases from avoidance NPI showed up in Part D, at that point we can mark it as misrepresentation. That is the single direction to move LEIE information into the named information. Capabilities interface Part D and LEIE • NPI, Unique supplier recognizable proof number • Medical supplier's strength (or practice) • Number of methodology/benefits the supplier performed • Number of particular Medicare recipients accepting the administration • Number of unmistakable Medicare recipient/every day administrations performed • Average of the charges that the supplier submitted for the administration • Average installment made to a supplier for each guarantee for the administration performed • EXCLUSION-name: Mapped misrepresentation marks from the LEIE database

# 3. Payments Received by Physician from Pharmaceuticals

Doctors in the US are required to proclaim all installments got from pharmaceutical organizations. Notwithstanding "installments got" essential data (ex: sum, date, type) this dataset contains numerous other helpful information components, for example, doctor possession in organization, counselling charges, good cause marker, debate status and so on. The entire datasets are incorporating three sections: General Payment, Research Payment and Physician Ownership Details. Key Features • The whole of general installment. • Name of medication related the installments.

# **METHODOLGY**

The methodology of this project will follow the below procedures: • Data exploration, cleansing and preparation • Build a simple data model to join all datasets • Feature engineering to choose

the effective feature sets for the different fraud patterns • Build a machine learning model to detect the different fraud patterns

#### CONCLUSION

 With the increasing number of populations of over 65 in USA, Medicare Fraud Detection is essential
 All types of Fraud Patterns have been Covered.
 Most Fraud Cases committed are in bay area
 Out of 5 Models Performed, best resulting model is Random Forest with AUC 72 %

#### REFERENCES:

CMS Part D datasets:

https://www.cms.gov/Research-Statistics-Data-and-Systems/Statistics-Trends-and-Reports/Medicare-Provider-Charge-Data/Part-D-Prescriber.html • LEIE Datasets:
https://oig.hhs.gov/exclusions/exclusions\_list.asp • FDA datasets
https://www.fda.gov/Drugs/InformationOnDrugs/ucm079750.htm#collapseOne • Dataset
Downloads. (n.d.). Retrieved June 23, 2020, from
https://www.cms.gov/OpenPayments/Explore-the-Data/Dataset-Downloads