

**Planned Time taken to speak : 25 mins**

**Actual time taken including Questions and answer session : 31 mins**

**Per head speaking time : 5 mins**

**Extra time : 5 mins extra for Surya to emphasize on overall summary**

### Things changed from the time that we showed :



1. Created progressive bar to indicate the status of the CRISP -DM framework and placed it at top left corner

2. Selected graphs were displayed to align with the story
3. News articles have been added. It starts from the Medicare Fraud Suit and ends with wheelchair story
4. Selected features got picked up and explained instead of all the features.

## Slide 1 :



### FRAUD DETECTION IN HEALTHCARE INSURANCE

#### Team Members:

Krithikaa Madhumitha  
Lakshmi Sruthi  
Monica Chandrasekar  
Samyuktha Mahesh  
Surya Narayanan



#### Guided by:

Mr. Animesh Tiwari

Batch : Chennai August 2020  
Course : Post Graduate Program in Data Science and Engineering  
Group : Group 1  
Domain : Finance & Risk Analytics

1

**Speaker : Monica**

**Transcript :**

A very good afternoon to one and all present here. After a 7 month long journey here we are today presenting the capstone project. Health is wealth and we have realized this more in the wake of the ongoing pandemic. Many insurance companies in the early days of the pandemic, halted their support for the people who were getting affected by the COVID-19. This thought made us search for a dataset on insurance companies related to healthcare. And therefore we have picked up a healthcare fraud analytics dataset from kaggle.

## Slide 2 :

Business Understanding

### Industry Review

greatlearning  
Learning for Life

#### HEALTH CARE QUALITY CRITERIA:

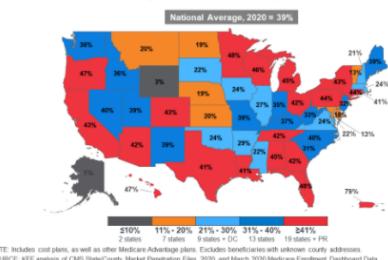
- Expenditure
- Quality
- Availability
- Population's Health
- Upfront Cost



Fig Source: Business Process Incubator



Medicare Advantage Penetration, by State, 2020



People who are 65 or older,  
Certain younger people with disabilities

Fig Source : KFF

2

**Speaker : Monica**

**Transcript:**

We have started the project with the first step in CRISP-DM which is Business Understanding. Just like how a country's wealth is measured with GDP in the same way a countries health is also measured by certain criteria:

- Expenditure
- Quality
- Availability
- Population's Health
- Upfront Cost

Healthcare expenditure is a major concern for the US. Around 530,000 families go bankrupt a year. Therefore, US government has three programs to provide support for the people. Medicare is one such program in which the importance is given to people who are above 65 years old or Certain younger people with disabilities.

## Slide 3:

Business Understanding

### Industry Review

ANNALS OF HEALTH CARE FEBRUARY 4, 2019 ISSUE

## THE PERSONAL TOLL OF WHISTLE-BLOWING

Why one physician took the risk of becoming an F.B.I. informant to expose alleged Medicare fraud.

By Sheelah Kolhatkar January 26, 2019

Nolan Auerbach & White

## End-Stage Renal Disease Fraud

READ MORE

washingtonpost.com/healthcare/fraud/end-stage-renal-disease-fraud

### A Medicare scam that just kept rolling

The government has paid billions to buy power wheelchairs. It has no idea how many of the claims are bogus.

Sections The Washington Post

#### BREAKING POINTS | 2014 UPDATE

Medicare continued to crack down on the lucrative wheelchair scam. But new statistics showed how badly the government may have been duped, during the scam's long-running heyday. Read the full update. ↗

3

**Speaker : Monica & Krithikaa**

**Article 1 Speaker : Monica**  
**Transcript:**

Here we handpicked 3 of the most important articles in medicare fraud.



### Article 1:

This first article is about the person, Darren Sewell, who created history by becoming a whistleblower. He became an undercover agent who filed a suit against his own company and one other company. In Spite of being Vice-President for "Freedom Health" , an health-insurance company , he teamed with an FBI agent Ed Ortega , and filed the biggest crime of America in the year 2009. It is then the Government started taking actions against all the crimes happening in this area.

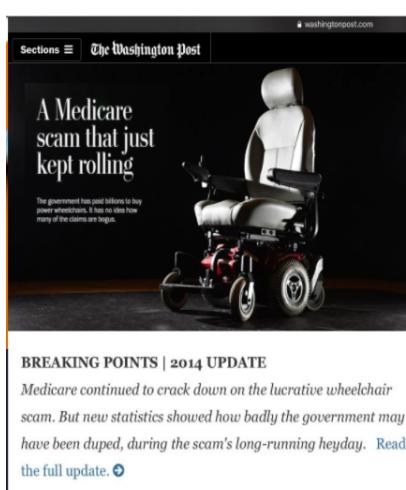
On the morning of September 17, 2009, Darren Sewell left his office at Freedom Health, the Tampa health-insurance company where he was a vice-president, and climbed into his Chevy Tahoe. He drove to Sonny's BBQ, a restaurant nearby, where he picked up barbecue sauce for the sandwich he had brought for lunch. Then he continued on until he reached a drab stretch of road lined with gas stations and scraggly palm trees, where he pulled into a parking lot and waited, as discreetly as possible, to meet his F.B.I. handler.

## Article 2 & 3 Speaker : Krithikaa



### Article 2:

One of the crimes that was discovered with medicare fraud related cases was the crime related to dialysis. According to the National Institute of Diabetes and Digestive and Kidney Diseases, More than 661,000 Americans have kidney failure, resulting in the 9th leading cause of death in the U.S. The treatment that costs about \$90,000 per patient, per year with medicare as their insurance provider. Lot of Elderly Outpatients who had the need for dialysis belonging to the age group of 70 to 85 were booked under false claims indicating renal disease. The patient's information was obtained through various resources. The fraudulent claims involve filing for medications, dialysis treatments etc.,



### Article 3:

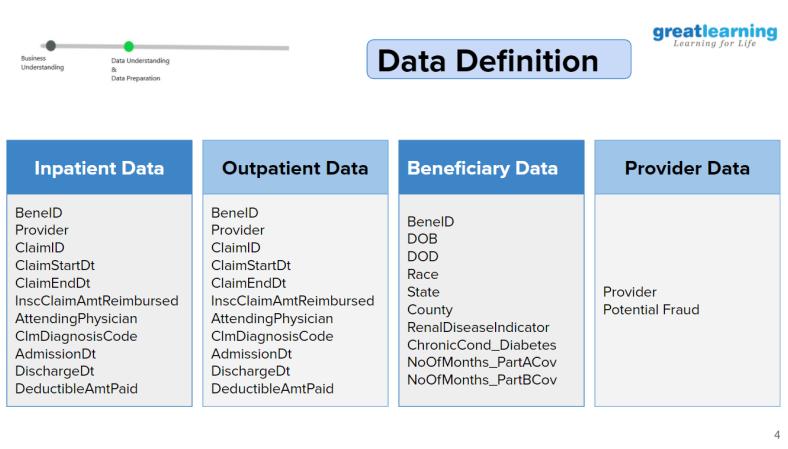
The power wheelchair scam provided a painful and expensive example of why Medicare fraud works so often. The scam started in Miami in the late 1990's and spread to Houston in the early 2000's. The government realised the wheelchair scam, when one of the patients said he had to climb the stairs to procure his wheelchair. Since 1999, Medicare has spent \$8.2 Billion to procure the power wheelchairs. It has taken 15 years for the government to realise the scam. Till now, the government is not able to account for the amount of money that it has lost in this scam.

## Conclusion for Articles :

### Speaker :Krithikaa

Based on the facts by The National Health Care Anti-Fraud Association estimates conservatively that health care fraud ranges from 3 to 10 percent of the nation's total health care expenditures. According to the above percentages, fraud would range from \$102 billion to \$340 billion. Experts project that health care expenditures will soar as high as \$5.5 trillion by 2025 .

## Slide 4 :



**Speaker : Krithikaa**

**Transcript:**

We now move to the process of Data Understanding and Data Processing from the CRISP-DM.

The dataset in its original form consists of 8 different CSV files out of which 4 are Labelled and 4 are Unlabelled. We will be using the labelled dataset having a little over half a million claims and 5410 providers

A person who gets treated by getting admitted in the hospital becomes an inpatient whose insurance is covered under Medicare PartA, whereas a person who gets only medical services becomes an outpatient whose insurance is covered under Medicare PartB.

Based on the Diagnosis, each patient will be assigned to a set of attending physicians, operating physician and other physician.

The Diagnosis Codes here are categorized under ICD Code which is maintained by WHO. The version ICD Code 9 is used in this dataset. Diagnosis Code 4019 indicates Unspecified essential hypertension, 1970 indicates Secondary malignant neoplasm of lung.

Procedures undergone by a patient are marked by the procedural codes defined under ICD-9-CM. Procedure Code 7092 indicates Scar conditions and fibrosis of skin, 331 indicates Incision of lung.

Every person who gets covered under the insurance policy becomes a Beneficiary who can raise a request for claim from their insurance policy through their Provider. The ClaimID will be unique even though the BenelID and ProviderID are repetitive.

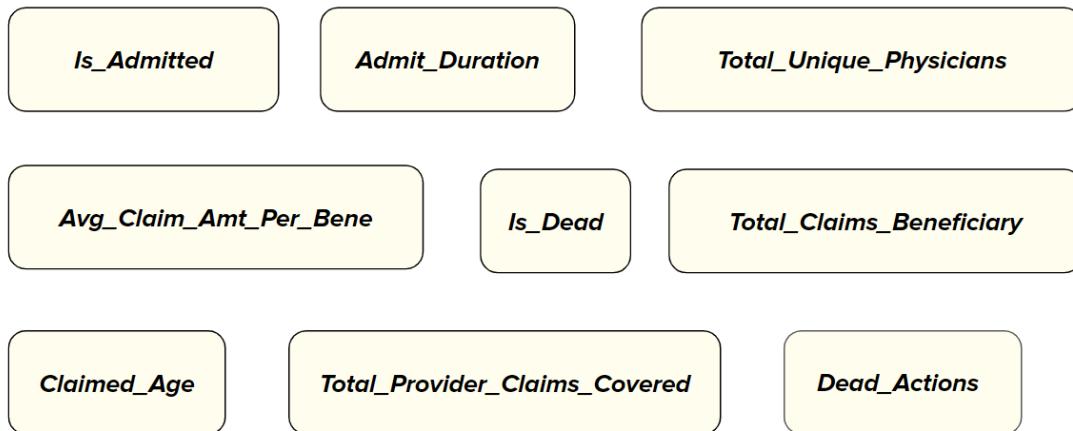
Health insurance companies are often faced with reimbursement as a patient needs to undergo an urgent medical procedure for which he pays from his pocket, this amount will later be reimbursed by the companies. Every Beneficiary would also have paid a part of a claim beforehand and the rest of the amount is paid by the insurance companies when the claim is raised.

Having identified the target variable as Potential Fraud, the goal of the project is to save enormous amounts of money from fraud which would aid those in real need.

# Slide 5 :



## Feature Engineering



5

**Speaker : Krithikaa**

**Transcript:**

In our dataset out of the 56 features, we noticed that there were 25 features with missing values, so we proceeded with Feature Engineering to draw more insights from the Dataset.

The first feature that we added was, IsAdmitted which helps identify if the patient is from the Impatient File or from the Outpatient File.

With the Date of Death column we created a new column called IsDead which works on the Self Comparison property of Nan. We also checked if any actions were taken after the Death of the Person with the DeadActions Column.

Columns such as ClaimDuration, AdmitDuration and ClaimedAge were inferred from their respective date columns. Using the columns InscClaimReimbursementAmt and DeductibleAmtPaid we calculated the Claim Amount for Each Claim for each Beneficiary and also the Average Claim Amount per Beneficiary.

We also calculated the Total Physicians and Total Unique Physicians for each Beneficiary. Using the ClaimID we calculated the Total Claims made by Beneficiary and Total Claims Covered by each Provider.

Post performing the Feature Engineering, we observed 72 columns.

# Slide 6 :



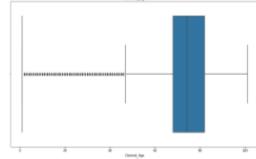
## Missing Value Treatment

```
1 missing_value_info[missing_value_info != 0.0]
AttendingPhysician          0.270149
OperatingPhysician           79.497538
OtherPhysician               64.218548
ClmAdmitDiagnosisCode      73.863109
DeductibleAmtPaid           0.161050
DiagnosisGroupCode          92.749337
C1mDiagnosisCode_1          1.872589
C1mDiagnosisCode_2          35.041588
C1mDiagnosisCode_3          56.458221
C1mDiagnosisCode_4          70.524407
C1mDiagnosisCode_5          79.949517
C1mDiagnosisCode_6          84.881702
C1mDiagnosisCode_7          88.144805
C1mDiagnosisCode_8          90.425843
C1mDiagnosisCode_9          92.509105
C1mDiagnosisCode_10         99.102490
C1mProcedureCode_1           95.824160
C1mProcedureCode_2           99.016501
C1mProcedureCode_3           99.826410
C1mProcedureCode_4           99.978861
C1mProcedureCode_5           99.998388
C1mProcedureCode_6           100.000000
Claim_Amount_For_Bene_Per_Claim  0.161050
Avg_Claim_Amt_Per_Bene     0.001433
dtype: float64
```

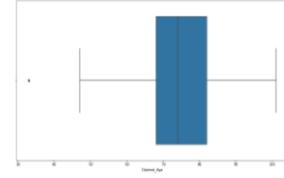
Before Missing Value Treatment

## Outlier Treatment

Before Outlier Treatment



After Outlier Treatment



After Missing Value Treatment

```
missing_value_info3[missing_value_info3 != 0.0]
Series([], dtype: float64)
```

Winsorization Technique

6

Speaker : Surya

Transcript:

We will now be moving to see how the missing values in our dataset were handled. We handled the missing values in three ways: We started by dropping the columns with missing values with more than 40%. We then handled the columns with missing values with imputation.

Three columns namely DeductibleAmtPaid , Claim\_Amount\_For\_Bene\_Per\_Claim and Avg\_Claim\_Amt\_Per\_Bene were filled by taking the mean of their values for each Beneficiary.

It would be unjust to fill these columns Attending Physician, Claim Diagnosis Code1 and Code2, by considering the entire dataset as one single group, therefore we segregated the dataset into two groups. ‘yes\_fraud’ group that will have all the data pertaining to those providers who are fraud, and ‘no\_fraud’ group that will have data pertaining to those providers who are not fraud. The columns were filled in such a fashion that by taking the mode of each group separately.

In our dataset we have 18 true numerical columns out of which 13 columns had outliers which we treated using the Winsorization Technique. As our dataset is quite large, we divided the 13 columns into 5 sets and treated them for outliers separately and then concatenated back together.

## Slide 7 :



### EDA

Potential Fraud  
No  
Yes

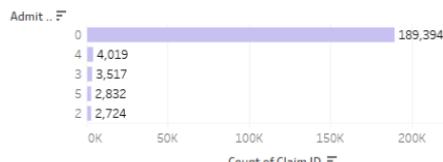
Admit Duration vs Claim Amount (Yes)



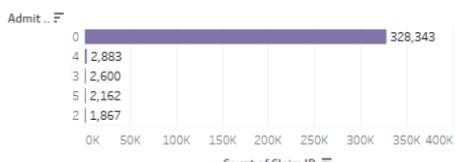
Admit Duration vs Claim Amount (No)



Admit Duration vs Claim Id (Yes)



Admit Duration vs Claim Id (No)



**Speaker : Surya**

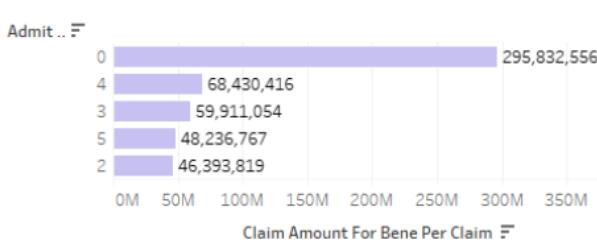
**Transcript:**

We now move to the exploratory data analysis for our project.

All the graphs in this page are based on the Admit Duration.

**Graph : Admit Duration vs Claim Amount (yes)**

Admit Duration vs Claim Amount (Yes)



The first part of the Graphs is between the Admit Duration for Fraudulent and Non-Fraudulent Claim Amount. We can infer that without spending even a single day in the hospital majority of the outpatients are involved in the Fraudulent Claims resulting in the Claim Amount a little over 296M.

**Graph : Admit Duration vs Claim Amount (No) and Admit Duration vs Claim Amount(Yes)**

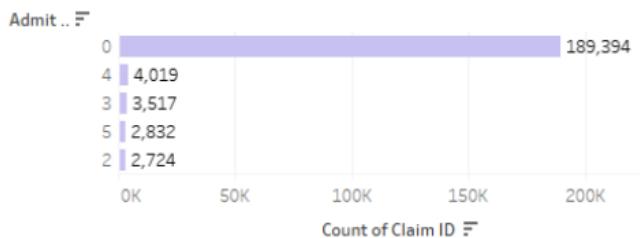
Admit Duration vs Claim Amount (No)



On comparison between the graphs, we can see that the claim amount for the days that are not 0 is higher for Fraudulent Claim Amount while Admit Duration for the non-fraudulent is less comparatively.

## Similarly on comparison of the Admit Duration and Claim Id (Yes)

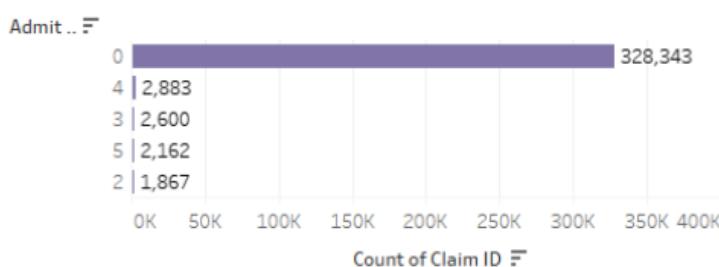
Admit Duration vs Claim Id (Yes)



It can also be seen that there are more ClaimID's registered for the Duration 0 is high compared to other days resulting in more fraudulent activities.

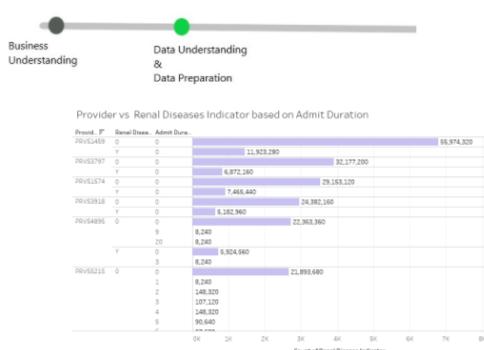
Graph : Admit Duration vs Claim Id (No)

Admit Duration vs Claim Id (No)



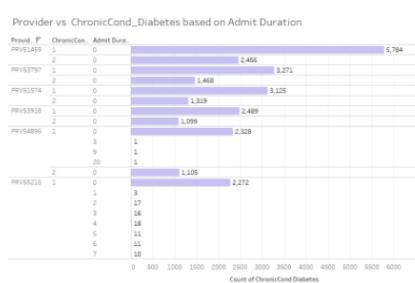
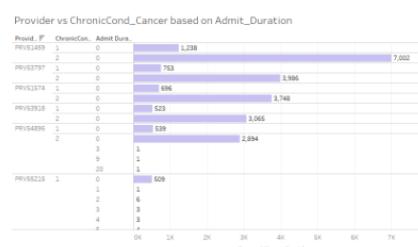
Similarly to the above graph the Claim id for the other days apart from Day 0 is higher than the Claim Id for the Day 0 when compared with Admit Duration vs Claim id (Yes)

## Slide 8 :



## EDA

greatlearning  
Learning for Life



## Speaker : Surya

### Transcript:

Before delving into this graph there are two main terms that have to be known.

**1. Patient Recruiters :** These are the people who get information about a patient, their Date of Birth, Medical Card Id etc., and inform fraudsters about them. These people ensure that there are no previous claims made by the patient for which the fraudsters are planning to make a claim.

### 2. Risk Category:

Medicare has divided a patient into two categories based on their disease.

**High risk category** - those who are having chronic conditions or those whose medical expenses are high

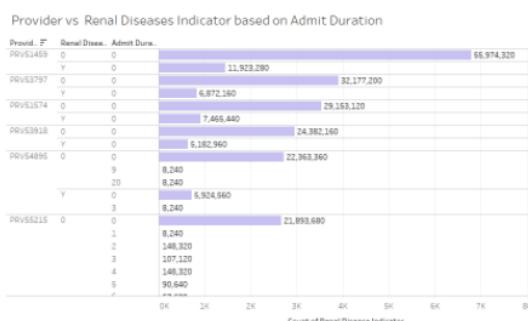
**Low risk category** - those who are not having low expenses for their medication and therapies etc.,

Fraudsters claim a patient as a high risk by entering false details into the medicare system and then obtain the claim money without the knowledge of the patient.

There are times where a physician has also joined this league. One such incident is where a skin doctor falsely reported a patient to have skin cancer. Made the patient undergo false treatments but obtained the money for the skin cancer from medicare.

The graph here talks about such incidents where a patient has been put into a high risk category just to obtain a large sum of money.

### Graph : Provider vs Renal Disease Indicator based on Admit Duration

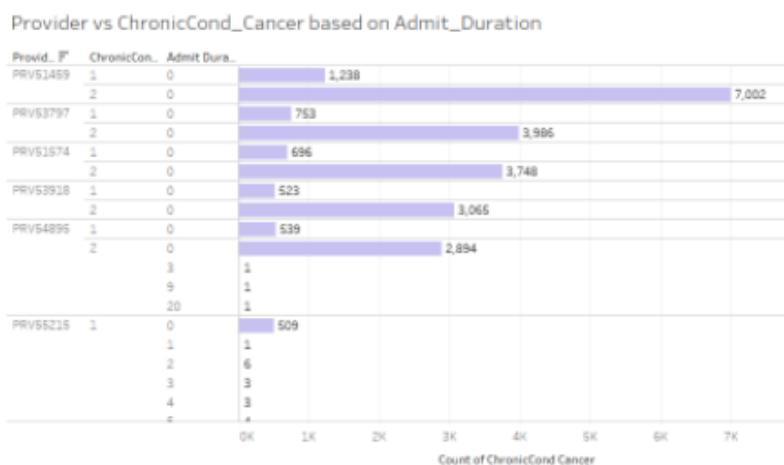


There are two levels in Renal Disease Indicator 0 and Y. Where 'Y' indicates there is a renal issue and 0 indicates there is no renal disease issue.

In this graph we can see that , during admit duration 0 , he has used the false renal disease indicator as a factor to make the claim against the patient. There are more number of providers whose claims are more in during the Admit Duration 0 compared to others

using this indicator.

### Graph : Provider vs ChronicCond\_Cancer based on Admit Duration

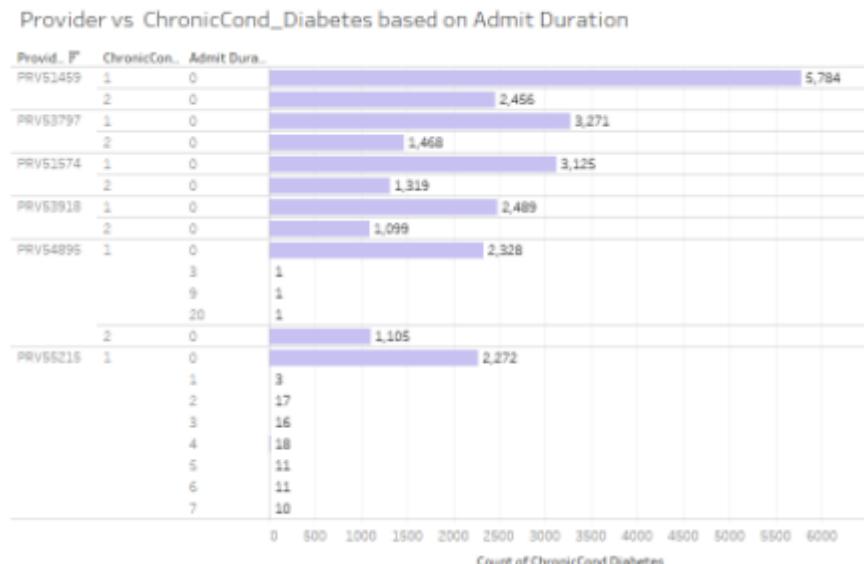


We do know that there are two levels in ChronicCond\_Cancer 1 and 2. But we do not have intelligence as to what 1 and 2 stands for.

In this graph we can see that , during admit duration 0 , he has

used ChronicCond\_Cancer as a factor to make the claim against the patient. There are more number of providers whose claims are more in during the Admit Duration 0 compared to others using this indicator.

## Graph : Provider vs ChronicCond\_Diabetes based on Admit Duration



We do know that there are two levels in ChronicCond\_Diabetes 1 and 2. But we do not have intelligence as to what 1 and 2 stands for.

In this graph we can see that , during admit duration 0 , he has used ChronicCond\_Diabetes as a factor to make the claim against the patient. There are more number of providers whose claims are more in

during the Admit Duration 0 compared to others using this indicator.

\*\*\*\* This point was not spoken but was kept in mind just in case\*\*\*\*

**Another perspective can also be seen by comparing these three graphs together.**

For example : Let's take the first provider.

Let's compare the 0 value of renal disease indicator, 1 value of chroniccond\_diabetes and 1 value of chroniccond\_Cancer and let's assume that 1 stands for no.

If we look at the graph, it can be seen that it is high for both renal indicators and diabetes and less for cancer. Now a question might arise as to why it is high for No indication . The answer for that lies just below . Now let's assume the value 2 as yes in cancer and diabetes,

Renal indicator is less for Y for this provider while it is very high for Cancer and moderately high for Cancer. Now we can see that this provider has used Cancer as a main reason to indicate patients as high risk category and then proceed with their fraudulent activity.

Renal Disease , ChronicCond\_Diabetes,ChornicCond\_Cancer and other ChronicConds comes under high risk category. Patient Recruiters use these indicators and make a claim for those beneficiaries of age group 70 to 85.

## Graph: Admit Duration vs Deduction Amount



From the graph based on the Admit Duration and Deductible Amount, it is noted that there is a small difference between Deductible Amount paid between the Fraudulent and Non-fraudulent claims. As per the federal government's False Claims Act for every false claim \$5,500–11,000. Even Though the difference is quite small, this small difference can make a huge impact on the Amount that can be saved for the Insurance Company.

## Slide 9 :



9

Speaker :Monica

Transcript:

Graph : Is\_Admitted and Potential Fraud



Looking at the graph between Is\_Admitted and Potential Fraud here 0 represents Outpatients and 1 represents Inpatient

We can see that the number of providers in the outpatient is more compared to

inpatient. In general patients are categorized into two categories based on risk score. A patient receives a high amount of money if the risk score is high and low amount of money if the score is low. Certain providers and doctors use this risk score and mark the patient as high risk category. These Outpatients would be falsely marked as high risk category either by stating that they have renal diseases which require dialysis and high cost medications or by stating that they have cancer and getting money for that or by stating other chronic conditions like diabetes, depression etc.

## Graph : Age vs Ins Claim Reimbursement Amount



Now let's look at : Age vs Ins Claim Reimbursement Amount graph

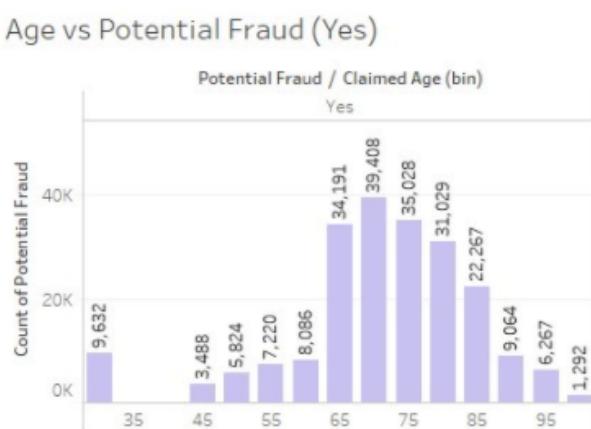
From this graph we can see that the age group of 35 years old has a high insurance claim amount. Most of the Americans love to explore the world and they go on trips to the countries like Vietnam,

Thailand etc. They also love scuba diving, trekking and many exploring activities.

This group of people tends to get affected by different kinds of viruses, bacterial infections, fungal infections and worms. As a result, many have lost their limbs Or they develop some kinds of ambulatory disabilities others may have disability from birth or by accidents.

Now these people require a lot of medicine and a wheelchair. As already mentioned above in the article, wheelchair fraud is another crime that is committed by the providers. And therefore, this age group has a higher amount.

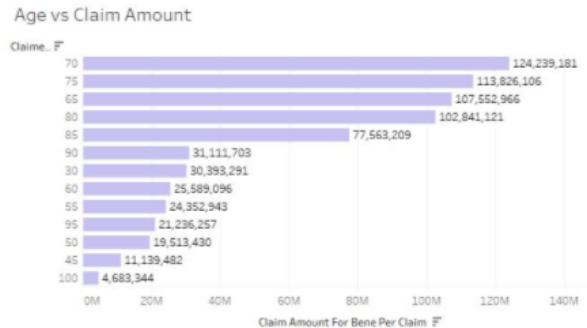
## Graph:Age vs Potential Fraud



This graph Age vs Potential Fraud talks how age plays a major role for fraud activities.

In this we can see that the age group from 65 to 85 have a higher number of fraudulent activities . Medicare being the program that provides support for the age group 65 and above, became the hotspot for the fraudulent activities. More providers concentrated on using these elders as a means to earn more money.

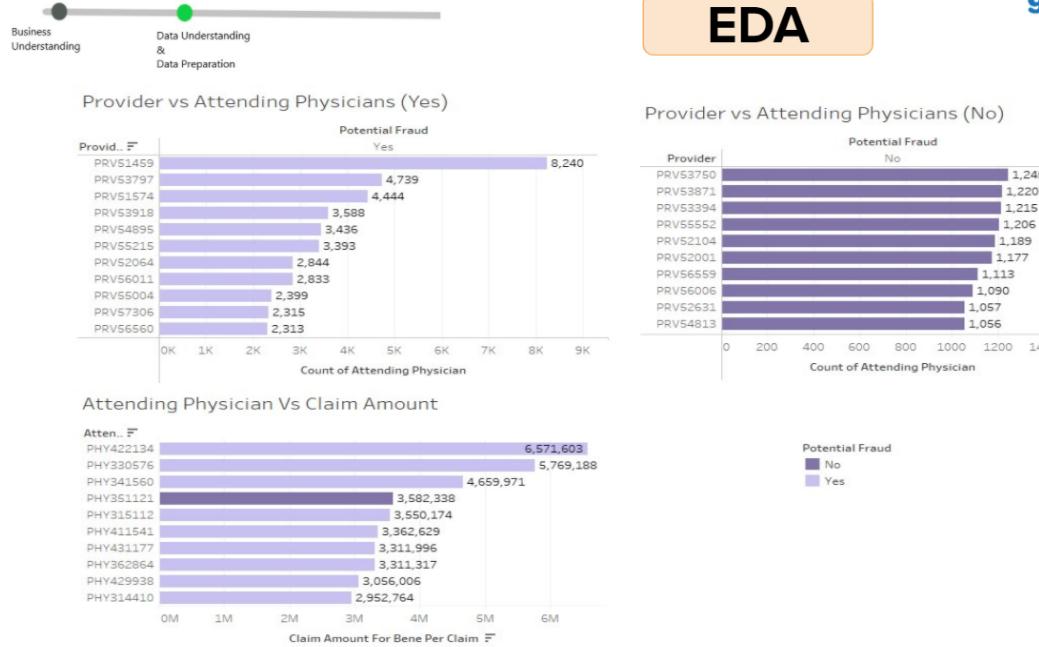
## Graph : Age and Claim Amount



The last graph in this slide is between the Age and Claim Amount which tells us about the amount claimed by a particular age group.

In conjunction with the previous graph , the people in the age group 65 to 85 have the higher number of providers resulting in higher claim amounts as well. Providers book the patients under high risk category without proper documentation and earn all the money.

## Slide 10 :



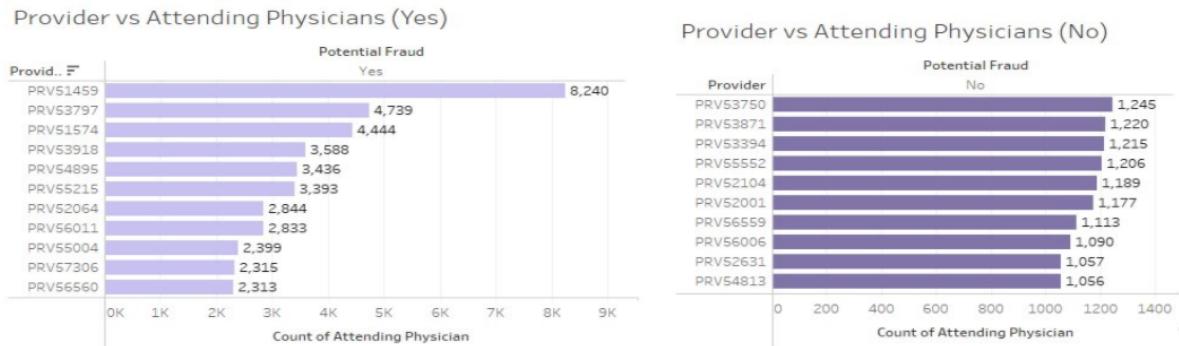
10

Speaker :Monica

Transcript:

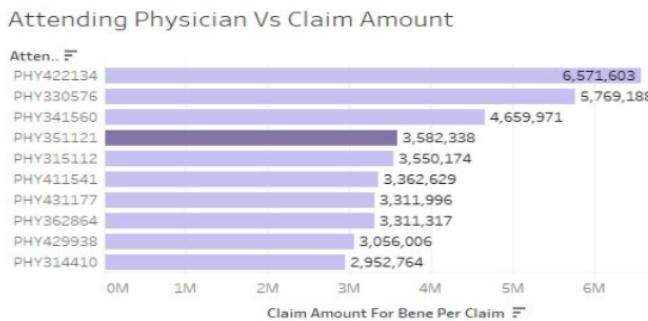
All the graphs in this slide are based on the Attending Physician.

## Graph : Potential Fraud vs Attending Physician



The first two graphs are separated on the basis of potentialFraud. These graphs talk about the relationship between Providers and Attending Physicians. We can notice that the highest number of attending physicians is associated with fraudulent provider PR51459 is 8240. While on the other hand we can see that those providers who are not involved in fraudulent activities have attending physicians associated with them in the range of 1000's only.

## GRAPH : Attending Physician vs Claim Amount



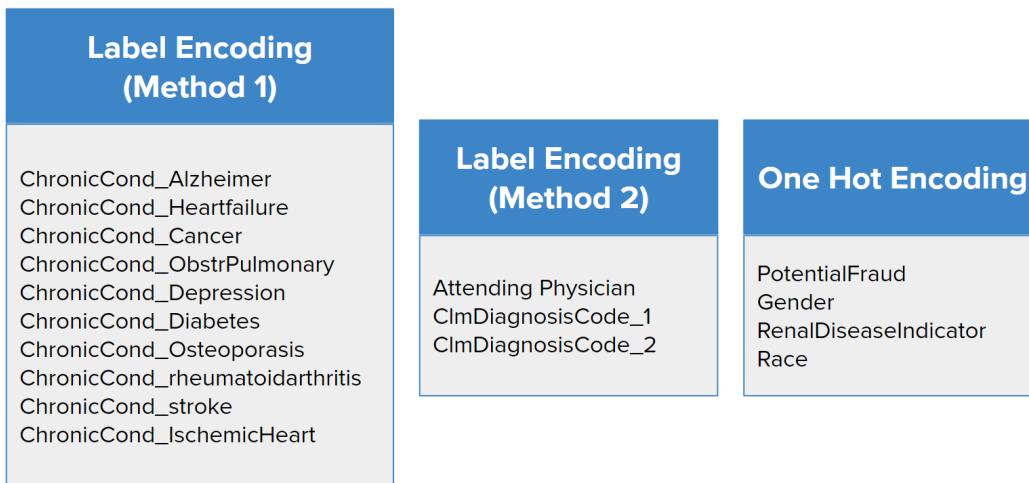
In corresponding to that we can notice that the claim amount is also based on the above facts. Here we can see that only 1 person can be seen who is not involved in fraudulent activity and the rest all have higher claim amounts against them.

This proves the fact that the attending physicians together are involved in the fraudulent activities.

# Slide 11:



## Encoding Categorical Variables



11

**Speaker : Lakshmi**

**Transcript:**

Our dataset has a combination of Numerical and Categorical Variables. The Categorical Variables need to be handled such that the model handles both the data types.

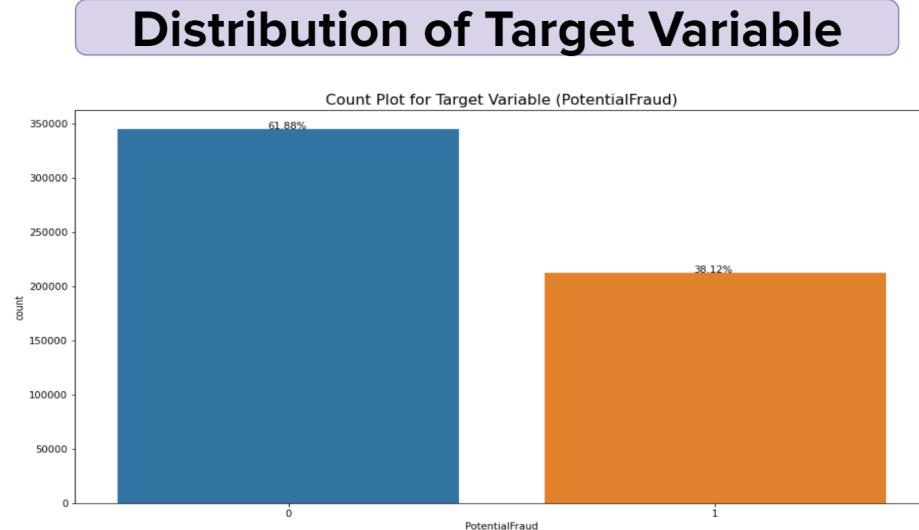
We have handled the Categorical Columns in three ways:

After attempting several methods of encoding for the columns related to the Medical Conditions, we have replaced the 2 with 0 so as to preserve the integrity of the data.

We have performed One Hot Encoding which is the most popular method for handling categorical variables for the Columns Race, Gender, Potential Fraud.

We handled columns AttendingPhysician, ClaimDiagnosisCode1 and ClaimDiagnosisCode2 by encoding it in such a way that if they belong to the Top 25 , they are encoded as 1 and the rest of them are encoded as 0. We have chosen Top 25 for encoding as it has a good distribution.

## Slide 12:



12

Speaker : Lakshmi

Transcript: From the distribution, it is noticed that the dataset is quite balanced.

## Slide 13:



### Statistical Significance of Variables

Default Significant Level | 0.05 |

H0	pval > 0.05	Insignificant i.e. there is no relationship that exists between the features
Ha	pval < 0.05	Significant i.e. there is a relationship that exists between the features

Hypothesis Testing for	Normality Test (Shapiro Test)	Statistical Test Performed	Stats Value	Pval	Significant/Insignificant
Potential Fraud vs State	Not Normal	Chi2_contingency	44912.23342	0	Significant
Potential Fraud vs County	Not Normal	Chi2_contingency	44191.37227	0	Significant
Potential Fraud vs Race	Not Normal	Chi2_contingency	703.9414756	2.93E-152	Significant
Total_Provider_Claims_Covered based on PotentialFraud	Not Normal	Mannwhitneyu	10848645553	0	Significant
Total_Uncode_Physicians based on PotentialFraud	Not Normal	Mannwhitneyu	35513575994	4.00E-100	Significant
Claimed_Age based on PotentialFraud	Not Normal	Mannwhitneyu	36379139593	9.43E-11	Significant
InscClaimAmountReimbursed based on PotentialFraud	Not Normal	Mannwhitneyu	34845178478	2.55E-235	Significant
AdmitDuration baed on PotentialFraud	Not Normal	Mannwhitneyu	34526672195	0	Significant
PotentialFraud vs RenalDiseaseIndicator	Not Normal	Chi2_contingency	31.85690392	1.66E-08	Significant
PotentialFraud vs ChronicCond_Diabetes	Not Normal	Chi2_contingency	15.88789609	6.72E-05	Significant
AdmitDuration vs ChronicCond_Diabetes	Not Normal	Chi2_contingency	2006.879449	0	Significant

13

## Speaker : Lakshmi

### Transcript:

Now we move to the statistical significance of the features.

A statistical test interprets a great deal of information as it gives numerical evidence to draw valid conclusions from test results. Using statistical analysis we can determine the likelihood that a hypothesis should either be accepted or rejected.

By performing the Shapiro test we have noticed that our data is not normal, so we have performed MannwhitneyU test. For statistical tests between two categorical variables we have performed Chi Square.

We have taken the default significance level of 0.05 which indicates a 5% risk of concluding that aims to quantify evidence against a particular hypothesis being true.

In our case the target column is Potential Fraud which is categorical in nature.

If Null Hypothesis is selected, the feature that we used for the test is not significant to build the model, because there is no use to predict the potential fraud. If it satisfies Alternate Hypothesis it is useful in predicting the potential fraud.

We have dropped the predictors “State,’County’,’BeneID’,’ClaimId’,’Provider’ as the level of categories are in thousands and there is no ordinality in the categories among these features therefore label encoding cannot be implemented. We also don't have the infrastructure to do one hot encoding as the data is huge. We now have 40 columns.

## Slide 14:



greatlearning  
Learning for Life

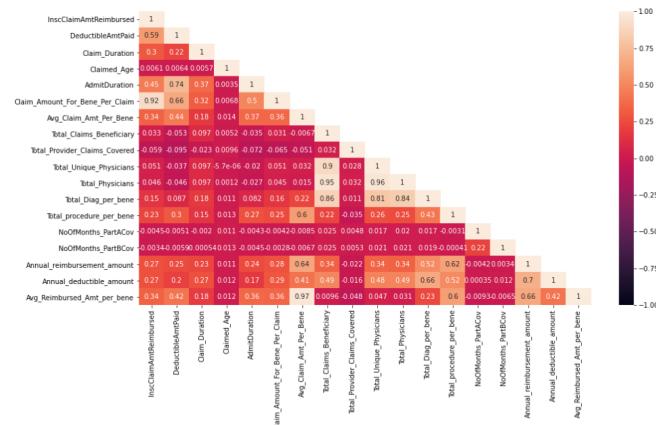
Standardization

### Assumption Checking

Multicollinearity Free

VIF	Features
10 4.131397	Total_Diag_per_bene
9 3.080378	Total_Uncique_Physicians
1 2.986379	DeductibleAmtPaid
12 2.985222	Annual_reimbursement_amount
13 2.658787	Annual_deductible_amount
6 2.504050	AdmitDuration
7 2.331512	Avg_Claim_Amt_Per_Bene
11 1.964363	Total_procedure_per_bene
0 1.661742	InscClaimAmtReimbursed
4 1.288904	Claim_Duration
3 1.049306	NoOfMonths_PartBCov
2 1.049275	NoOfMonths_PartACov
8 1.010136	Total_Provider_Claims_Covered
5 1.000768	Claimed_Age

With Multicollinearity



**Speaker : Samyuktha**

**Transcript:**

We move the next phase of CRISP-DM which is model building and evaluation

We have used Standard Scaler to standardize the true numerical columns.

Post which we check for the existence of multicollinearity. Removing multicollinearity helps to prevent wrong estimation of coefficient and precise calculation for coefficients during model building. Thus it is important for and the first basic step for any regression model.

We check for multicollinearity using Variance Inflation Factor. We drop the columns if the VIF value is high. The features that gets dropped in this process can also be seen having high correlation value i.e. above 0.7

The features that gets dropped are :

1. Total\_Physicians
2. Avg\_Reimbursed\_Amt\_per\_bene
3. Total\_Claims\_Beneficiary
4. Claim\_Amount\_For\_Bene\_Per\_Claim

We then stop the VIF checking process as all the VIF above 5 has been removed. Now the data is multicollinearity free.

We then concatenate both the numerical and categorical features and create a new dataframe.

## Slide 15:



**greatlearning**  
Learning for Life

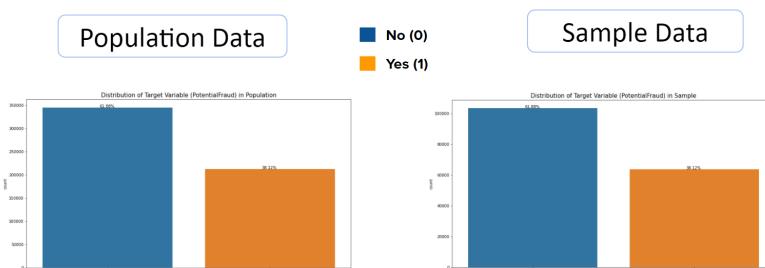
**Speaker : Samyuktha**

**Transcript:**

As the data is too huge for model building , we decided to pick up sample data for model building. After several trial and error we found that the sample size of 30% represented the population data. We did statistical tests to check for the representation as well.

15

### Selecting Sample for Model Building



Distribution of Target Variable

## Slide 16:



### Selecting the Best Model

Model	Train Dataset				Test Dataset			
	Precision	Recall	F1-Score	Accuracy	Precision	Recall	F1-Score	Accuracy
Logistic Regression	0.82	0.67	0.74	0.82	0.82	0.66	0.73	0.82
Decision Tree Classifier	0.95	0.73	0.82	0.88	0.94	0.72	0.82	0.88
Random Forest Classifier	0.95	0.62	0.75	0.84	0.94	0.62	0.75	0.84
<b>XGBoost Classifier</b>	<b>0.99</b>	<b>0.94</b>	<b>0.96</b>	<b>0.93</b>	<b>0.95</b>	<b>0.86</b>	<b>0.91</b>	<b>0.93</b>
MLP Classifier	0.82	0.67	0.74	0.82	0.82	0.66	0.73	0.82

```
----- Classification Report -----
precision    recall   f1-score   support
XGBClassifier
0           0.92    0.97    0.95    31088
1           0.95    0.87    0.91    19152
accuracy          0.94    0.92    0.93    50240
macro avg       0.94    0.92    0.93    50240
weighted avg    0.93    0.93    0.93    50240
```

16

**Speaker : Samyuktha**

**Transcript:**

We then proceed with model building. the logistic model failed. Threshold obtained for it is 1. This is due to the fact that, the number of values for the target representing 0 compared with the target representing 1 is high. For a logistic model, it is quite complicated to form the required equation for the computations when there are more numbers of 0's in the target column . Hence we proceed by building tree based models. It can be seen that xgboost has the highest recall, precision, accuracy and f1 - score. Hence we are choosing xgboost as the best model.

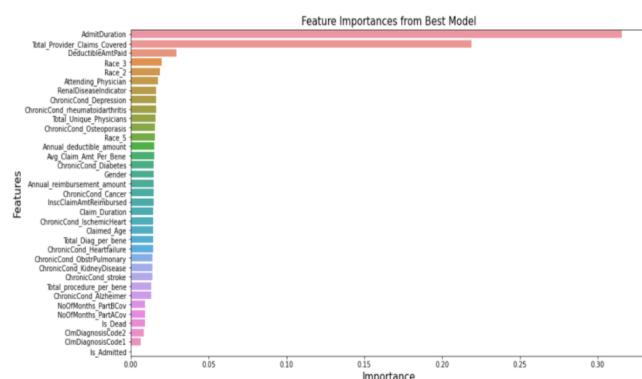
## Slide 17:

### Business Case

XGBClassifier	
----- Confusion Matrix -----	
Predicted:0	Predicted:1
Actual:0	30195 893
Actual:1	2484 16668

No of Fraudulent Claim	19,152
Total Data points	50,240
19152 are fraud out of 50240	0.3812101911
Total Loss calculates is	<b>\$751,255,150</b>
Each fraud costs	<b>\$14,953.32703</b>
Model precision is 95%	18,194
Savings	<b>\$272,060,832</b>
Risk of losing the non fraud customers	<b>\$4,465,000</b>



**Speaker : Samyuktha**

**Transcript:**

Now let's take a look at how xgboost model can help save money in this dataset.

The total number of fraudulent claims is 19 thousand one hundred and fifty two

Total data points is 50 thousand two and forty

Percentage of fraud in this dataset is (given by 19152/50240 )= 0.38%

Total loss = providerfradulent = yes and sum of insurance claim reimbursement amount

= 7 hundred and 51 million, two hundred and 55 thousand and 1 hundred and 50 dollars

Our model can predict 95 percentage of the fraudulent activities which is 18 thousand one hundred and ninety four

Therefore the total savings that can be done is given by 18,194 \* 14593 ie.e. 2 hundred and seventy two million,60 thousand , 8 hundred and thirty two \$

The risk is this model is that there are 893 beneficiaries who might be detected wrongly. If the insurance amount per year is 5000 dollars then this loss will come up to 4million 4 hundred and 65 thousand.

To the right we can see the important features given by the XGBoost classifier. All the analysis that we did in the EDA and statistics correlates with this. Admit Duration gets the high value and where the features from RenalDiseaseIndicator till ChronicCond\_Alzemeris have equal contribution for detection of fraud .

## Slide 18:

**greatlearning**  
Learning for Life

### KEY POINTS SUMMARY

#### From this dataset:

- age group of 35 and 70
- Collaboration between physicians and providers
- More fraudulent providers for outpatients
- Xgboost Model predicted saving \$272,060,832

#### Based on Business Understanding:

- Assigning Wrong chronic disease indicator
- False claims for wheelchair

#### Recommendations for the Industry:

- Prevention of false registration
- Prevent of wrong upcoding of patients' diagnosis code, procedure code and other disease indicators.
- Immediate checking of the patients health history
- Protection of Patient's information

**Speaker : Surya**

**Transcript:**

**From this dataset:**

1. We have inferred that Fraudulent activities occur more in the age group of 35 and 70.
2. From this dataset we have noticed the collaboration with physicians and providers.
3. Outpatient had a lot of fraudulent providers in this dataset.
4. The Xgboost model was able to save \$272060832.

**Based on Business Understanding:**

5. Assignment of Wrong diagnosis code leads to registering the patient in a high risk category resulting in losing more money for Medicare.
6. Lot of fraudulent activities take place when a provider claims money for a wheelchair but the beneficiary does not receive that wheelchair.

**Recommendations for the Industry:**

7. Proper checkpoints have to be established in order to prevent false registration of patients' names to receive the claim amount and to prevent wrong upcoding of patients' diagnosis code, procedure code and other disease indicators.
8. Immediate checking of the patients health history has to be made to speed up the fraud prevention act.
9. Protection of Patient's information has to be given more protection as it will prevent "patient recruiters" from obtaining patient's information for selling into the wrongful agencies.

## Slide 19:

**greatlearning**  
Learning for Life

### REFERENCES

- <https://www.newyorker.com/magazine/2019/02/04/the-personal-toll-of-whistle-blowing>
- <https://spendmenot.com/blog/medical-bankruptcy-statistics/>
- <https://www.whistleblowerfirm.com/healthcare-fraud/end-stage-renal-disease-fraud/>
- <https://www.washingtonpost.com/sf/national/2014/08/16/a-medicare-scam-that-just-kept-rolling/>

## Slide 20:

**THANK YOU**