

ANALYTICS FOR HOSPITALS HEALTH-CARE DATA

REPORT

TEAM ID PNT2022TMID35084

SUBMITTED BY,

JENI ALPHONSA A(963319106036)

JENILA J(963319106039)

JESNEY R(963319106041)

LISHA S(963319106053)

TABLE OF CONTENTS

Project Report Format

1. INTRODUCTION

1.1 Project Overview

1.2 Purpose

2. LITERATURE SURVEY

2.1 Existing problem

2.2 References

2.3 Problem Statement Definition

3. IDEATION & PROPOSED SOLUTION

3.1 Empathy Map Canvas

3.2 Ideation & Brainstorming

3.3 Proposed Solution

3.4 Problem Solution fit

4. REQUIREMENT ANALYSIS

4.1 Functional requirement

4.2 Non-Functional requirements

5. PROJECT DESIGN

5.1 Data Flow Diagrams

5.2 Solution & Technical Architecture

5.3 User Stories

6. PROJECT PLANNING & SCHEDULING

6.1 Sprint Planning & Estimation

6.2 Sprint Delivery Schedule

6.3 Reports from JIRA

7. CODING & SOLUTIONING (Explain the features added in the project along with code)

7.1 Feature 1

7.2 Feature 2

7.3 Database Schema (if Applicable)

8. TESTING

8.1 Test Cases

8.2 User Acceptance Testing

9. RESULTS

9.1 Performance Metrics

10. ADVANTAGES & DISADVANTAGES

11. CONCLUSION

12. FUTURE SCOPE

13. APPENDIX Source Code

GitHub & Project Demo Link

ABSTRACT

The current study performs a systematic literature review (SLR) to synthesise prior research on the applicability of big data analytics (BDA) in healthcare. The SLR examines the outcomes of 41 studies, and presents them in a comprehensive framework. The findings from this study suggest that applications of BDA in healthcare can be observed from five perspectives, namely, health awareness among the general public, interactions among stakeholders in the healthcare ecosystem, hospital management practices, treatment of specific medical conditions, and technology in healthcare service delivery. This SLR recommends actionable future research agendas for scholars and valuable implications for theory and practice.

1. Introduction

1.1 Project overview

Healthcare organizations are under increasing pressure to improve patient care outcomes and achieve better care. While this situation represents a challenge, it also offers organizations an opportunity to dramatically improve the quality of care by leveraging more value and insights from their data. Health care analytics refers to the analysis of data using quantitative and qualitative techniques to explore trends and patterns in the acquired data. While healthcare management uses various metrics for performance, a patient's length of stay is an important one.

Being able to predict the length of stay (LOS) allows hospitals to optimize their treatment plans to reduce LOS, to reduce infection rates among patients, staff, and visitors.

1.2. Purpose

The goal of this project is to accurately predict the Length of Stay for each patient so that the hospitals can optimize resources and function better.

2. Literature survey

2.1 Existing problem

Recent Covid-19 Pandemic has raised alarms over one of the most overlooked

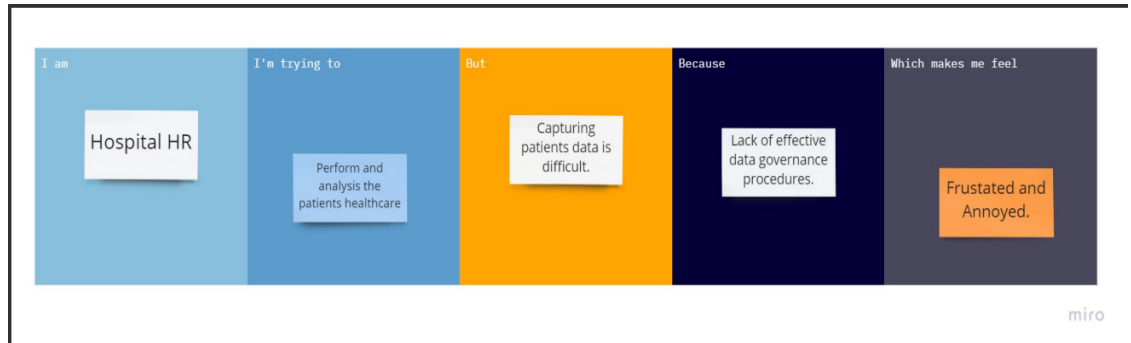
areas to focus: Healthcare Management. While healthcare management has various use cases for using data science, patient length of stay is one critical parameter to observe and predict if one wants to improve the efficiency of the healthcare management in a hospital.

2.2 References

- i. Janatahack: Healthcare AnalyticsII - *Analytics Vidhya* - [Link](#)
- ii. What Is Naive Bayes Algorithm in Machine Learning? - *Rohit Dwivedi* - [Link](#)
- iii. Naïve Bayes for Machine Learning– From Zero to Hero - *Anand Venkataraman* - [Link](#)
- iv. XGBoost Parameters - *XGBoost Documentation* - [Link](#)
- v. Predicting Heart Failure Using Machine Learning, Part 2- *Andrew A Borkowski* - [Link](#)
- vi. How to Tune the Number and Size of DecisionTrees with XGBoostin Python - *JasonBrownlee* - [Link](#)
- vii. Big Data Analytics inHealthcare That Can Save People - *Sandra Durcevic* - [Link](#)
- viii. Learning Process of a Neural Network – *Jordi Torres* - [Link](#)

2.3 Problem statement

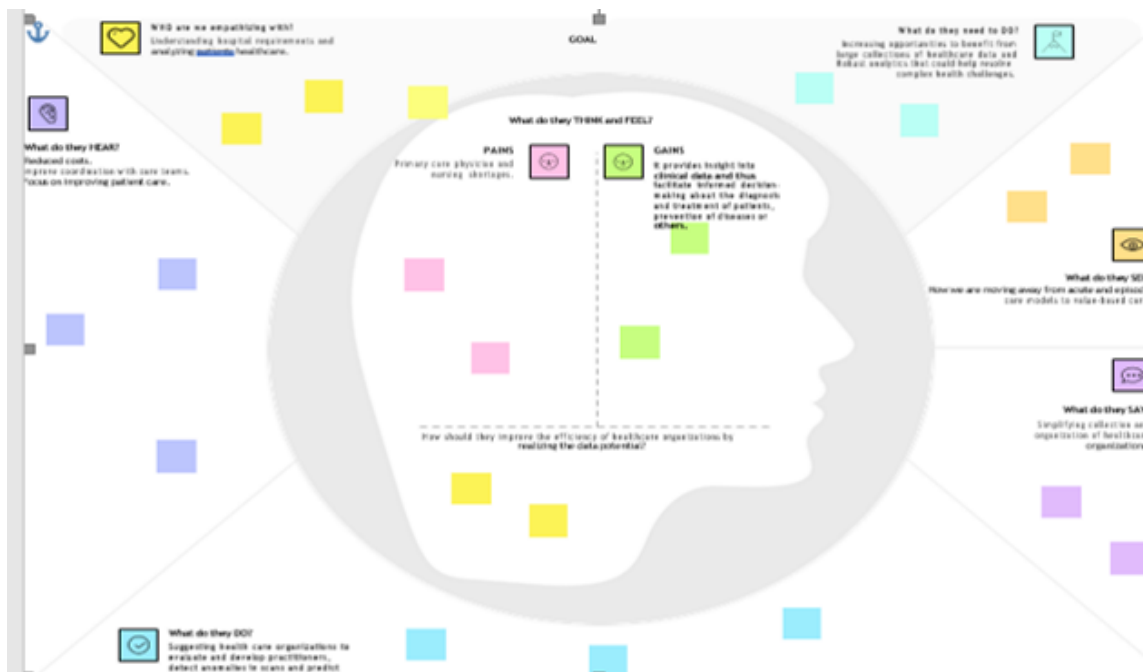
The task is to accurately predict the Length of Stay for each patient on case-by-case basis so that the Hospitals can use this information for optimal resource allocation and better functioning. The length of stay is divided into 11 differentclasses ranging from 0-10 days to more than 100 days.



Problem Statement (PS)	I am (Customer)	I'm trying to	But	Because	Which makes me feel
PS-1	Hospital HR	Perform and analysis the patients healthcare.	Capturing patient data is difficult.	Lack of effective data governance procedures.	Frustrated and Annoyed.

3. Ideation & proposed solution

3.1 Empathy map canvas



3.2 Ideation and Brainstorming

3.3 Proposed solution

S.No.	Parameter	Description
-------	-----------	-------------

1.	Problem Statement (Problem to be solved)	The task is to accurately predict the Length of Stay for each patient on case-by-case basis so that the Hospitals can use this information for optimal resource allocation and better functioning. The length of stay is divided into 11 different classes ranging from 0-10 days to more than 100 days.
2.	Idea / Solution description	Naïve Bayes is a classification technique that works on the principle of Bayes theorem with an assumption of independence among the variables. Here the goal is to predict Length of Stay i.e., “Stay” column (Target Variable) and it is classified into 11 levels. We must find the probability of each patient’s length of stay using feature variables, which contain the patient’s condition and hospital-level information. These feature variables are ordinal and naïve Bayes is a perfect multilevel classifier.
3.	Novelty / Uniqueness	Accurate understanding of the factors associating with the LOS and progressive improvements in processing and monitoring may allow more efficient management of the LOS of inpatients
4.	Social Impact / Customer Satisfaction	A shorter LOS reduces the risk of acquiring staph infections and other healthcare-related conditions, frees up vital bed spaces, and cuts overall medical expenses

5.	Business Model (Revenue Model)	The length of stay (LOS) is an important indicator of the efficiency of hospital management. Reduction in the number of inpatient days results in decreased risk of infection and medication side effects, improvement in the quality of treatment, and increased hospital profit with more efficient bed management
6.	Scalability of the Solution	Remote patient monitoring systems enabling effective distance treatment. Patient portals that allow people to better manage their health themselves;

4. Requirements analysis

4.1 Functional requirements

FR No.	Functional Requirement(Epic)	Sub Requirement (Story/ Sub-Task)
F R-1	User Registration	Registration through Form Registration through Gmail Registration through LinkedIN
F R-2	User Confirmation	Confirmation via Email Confirmation via OTP
F R-3	Operability	Share patient data and make it interoperable among the management
F R-4	Accuracy	The dashboard will be able to predict length of stay based on multiple combinations based on input sources with an accuracy of up to 85%
F R-5	Compliance	The product is to be used within the hospital so any form of data need not be hidden

F R- 6	Productivity	The dashboard is believed to improve the predictions of Length of Stay and thereby creating a scenario of providing better solution
--------------	--------------	---

4.2 Nonfunctional requirements

FR No.	Non-Functional Requirement	Description
NF R-1	Usability	This Dashboards are designed to offer a comprehensive overview of patient's LOS, and do so through the use of data visualization tools like charts and graphs.
NF R-2	Security	General industry level security shall be provided
NF R-3	Reliability	This dashboard will be consistent and reliable to the users and helps the user to use in effective, efficient and reliable manner.
NF R-4	Performance	The dashboard reduces the time needed for analysing data and has an automated system for that which improves the performance

NF R-5	Availability	The dashboard can be available to meet user's demand in timely manner and it also helps to provide necessary information to the user's dataset
NF R-6	Scalability	It is a multi-tenant system which is capable of running on lower-level systems as well.

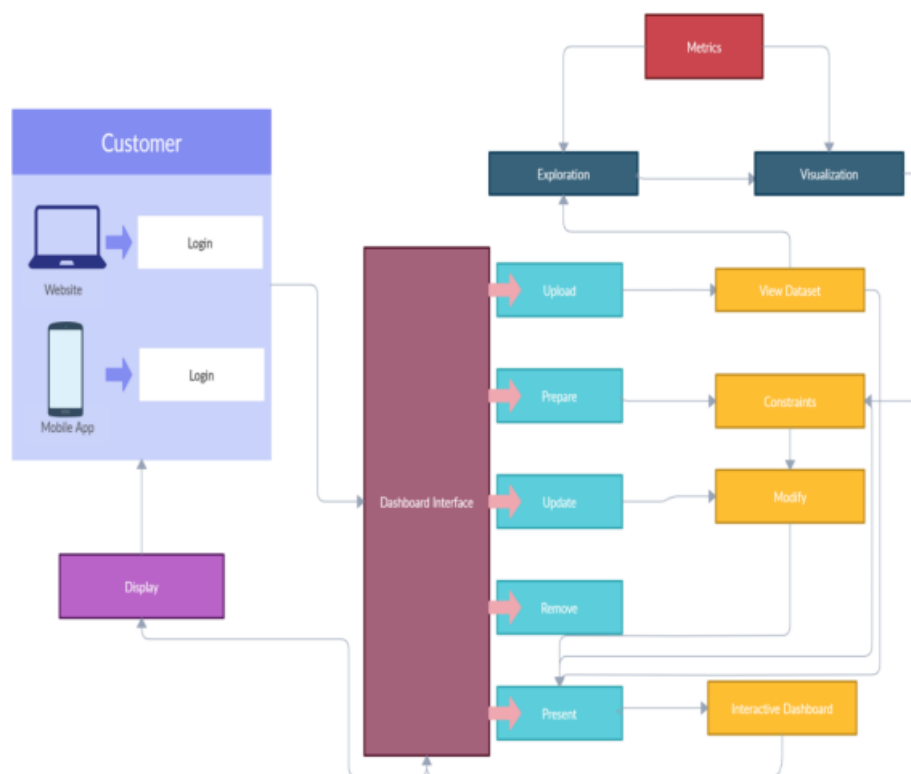
5. PROJECT DESIGN

5.1 Data Flow Diagrams

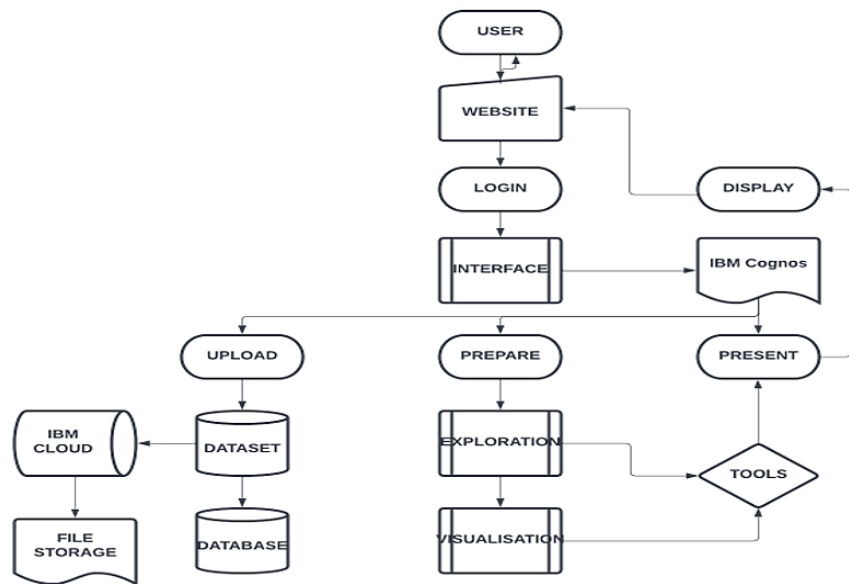
A Data Flow Diagram (DFD) is a traditional visual representation of the information flows within a system. A neat and clear DFD can depict the right amount of the system requirement graphically. It shows how data enters and leaves the system, what changes the information, and where data is stored.

5.2 Solution & Technical Architecture

SOLUTION ARCHITECTURE



TECHNOLOGY ARCHITECTURE



5.3 User Stories

User Type	Functional Requirement (Epic)	User Story Number	User Story / Task	Acceptance criteria	Priority	Release
-----------	-------------------------------	-------------------	-------------------	---------------------	----------	---------

Customer	Dashboard	USN 1	As a user,I can upload the datasets to the dashboard	I can access various operations	Medium	Sprint-4
	View	USN 2	As a user,I can view the patient details	I can view the visual data and the result after the prediction	Medium	Sprint-3
Admin	Analyse	USN 3	As an admin, I will analyse the given dataset	I can analyse the dataset	High	Sprint-2

6. Project planning & scheduling

6.1 Sprint Planning & Estimation

SPRINT 1

- Collection of data
- Data preprocessing
- Upload the dataset

jupyter IBM PROJECT Last Checkpoint: an hour ago (autosaved) Logout

File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3 (ipykernel)

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from scipy import stats

In [2]: import os
os.chdir("C:/Users/Admin/Desktop/DATASETS")

In [3]: df=pd.read_csv('train_data.csv')

In [4]: df
```

Out[4]:

case_id	Hospital_code	Hospital_type_code	City_Code_Hospital	Hospital_region_code	Available Extra Rooms in Hospital	Department	Ward_Type	Ward_Facility_Code	Bed Grade
0	1	8	c	3	Z	3 radiotherapy	R	F	2.0
1	2	2	c	5	Z	2 radiotherapy	S	F	2.0

jupyter IBM PROJECT Last Checkpoint: an hour ago (autosaved) Logout

File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3 (ipykernel)

318438 rows x 18 columns

```
In [5]: #Summary of the dataframe
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 318438 entries, 0 to 318437
Data columns (total 18 columns):
 #   Column                                Non-Null Count  Dtype
---  -
 0   case_id                              318438 non-null  int64
 1   Hospital_code                        318438 non-null  int64
 2   Hospital_type_code                  318438 non-null  object
 3   City_Code_Hospital                  318438 non-null  int64
 4   Hospital_region_code                318438 non-null  object
 5   Available Extra Rooms in Hospital    318438 non-null  int64
 6   Department                          318438 non-null  object
 7   Ward_Type                          318438 non-null  object
 8   Ward_Facility_Code                  318438 non-null  object
 9   Bed Grade                          318325 non-null  float64
10  patientid                           318438 non-null  int64
11  City_Code_Patient                   313906 non-null  float64
12  Type of Admission                   318438 non-null  object
13  Severity of illness                 318438 non-null  object
14  Visitors with Patient               318438 non-null  int64
15  ...
```

jupyter IBM PROJECT Last Checkpoint: an hour ago (autosaved) Logout

File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3 (ipykernel)

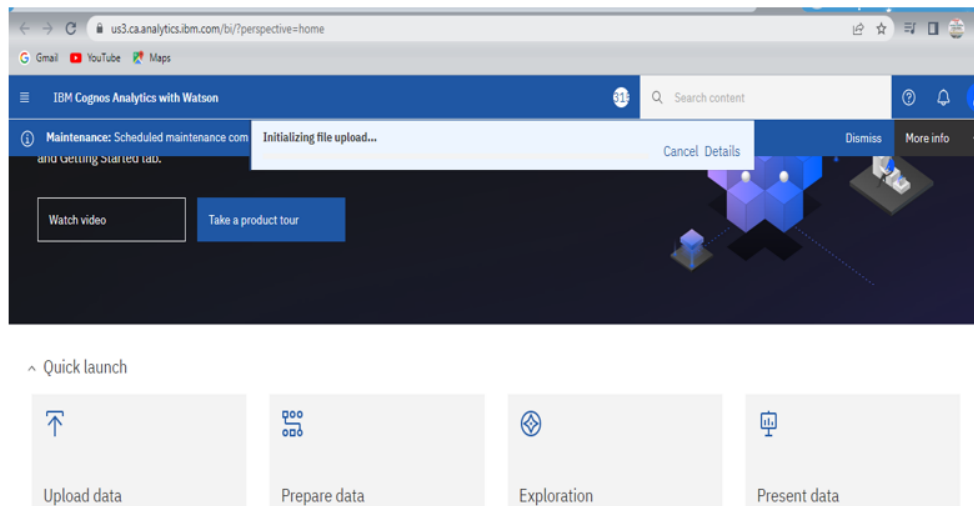
```
In [22]: df.Ward_Type
```

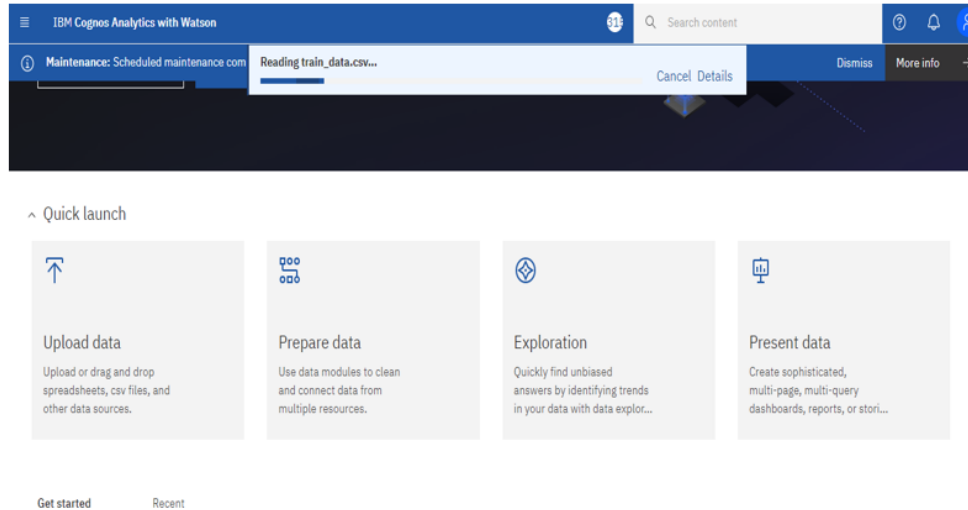
```
Out[22]: 0      R
         1      S
         2      S
         3      R
         4      S
         ..
        318433  Q
        318434  Q
        318435  R
        318436  Q
        318437  Q
        Name: Ward_Type, Length: 318438, dtype: object
```

```
In [23]: df.Ward_Facility_Code
```

```
Out[23]: 0      F
         1      F
         2      E
         3      D
         4      D
         ..
        318433  F
        318434  E
```

Upload the dataset:

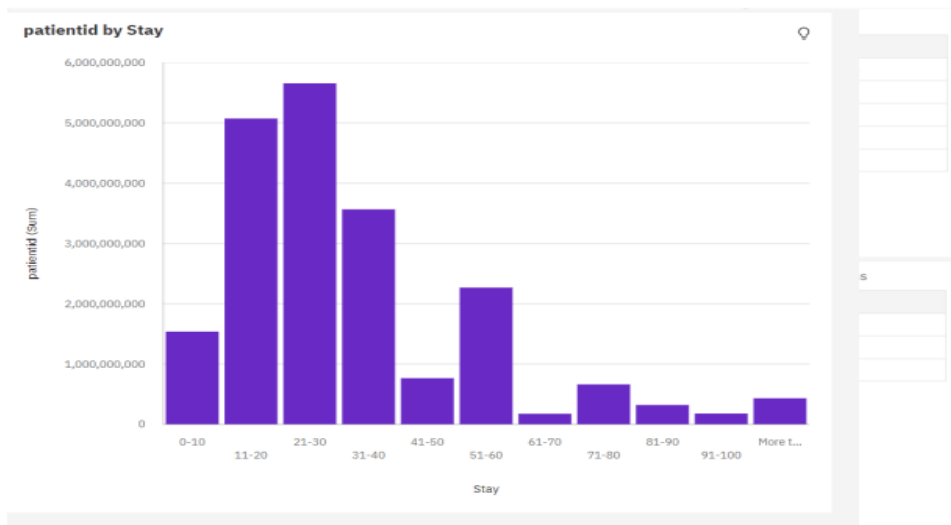




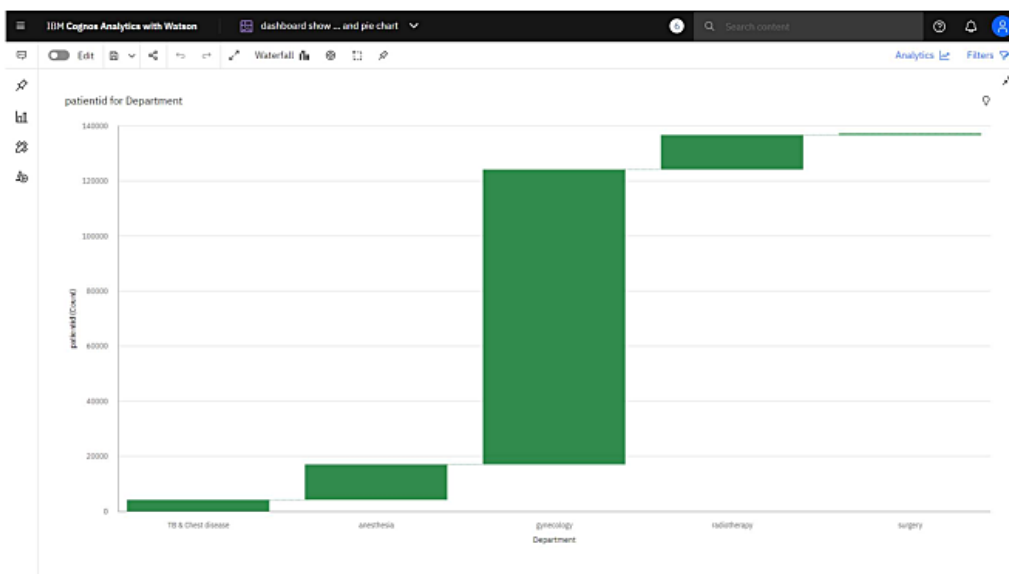
SPRINT 2

DATA EXPLORATION:

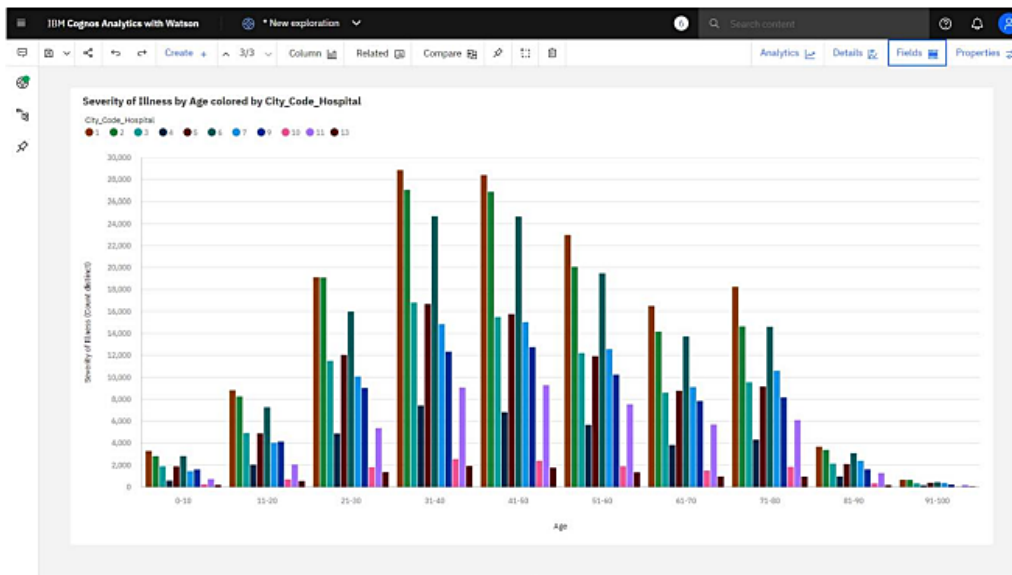
- Patient id by stay
- Patient id for department
- Severity of Illness by Age colored by City Code Hospital
- Case id by Ward Type
- Case id by Department
- Bed Grade by Department
- Case id by Severity of Illness
- Patient by Ward Type
- Available Extra Rooms in Hospital by Ward type
- Stay by Department
- Admission count for Department



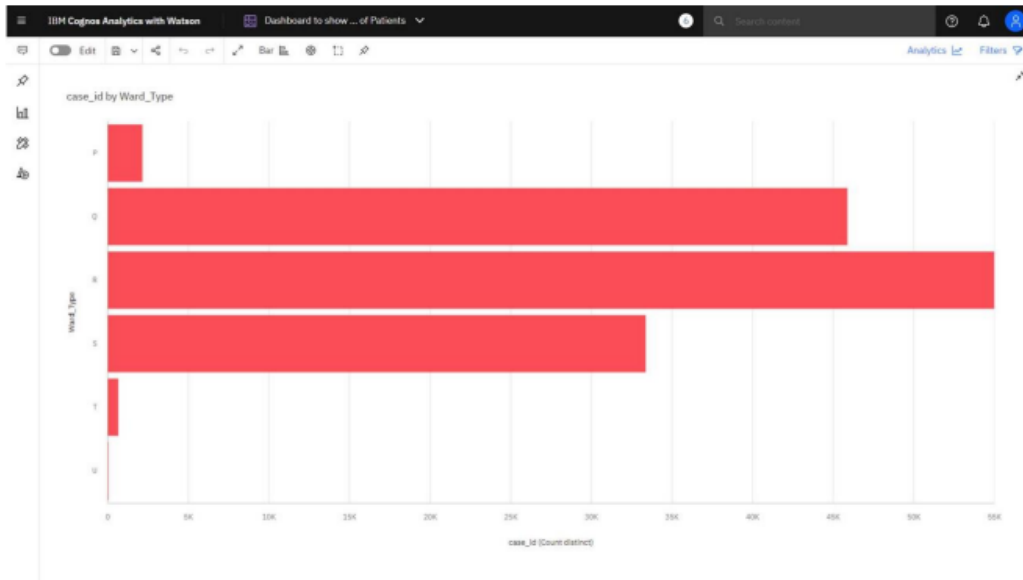
Patient id for department:



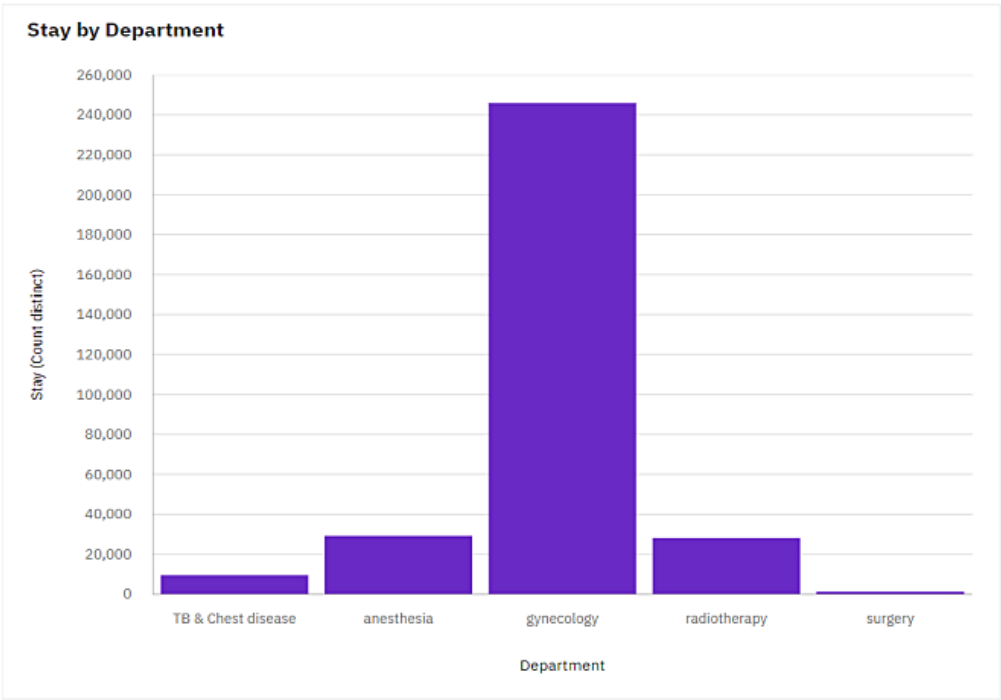
Severity of illness by Age colored by city code Hospital:



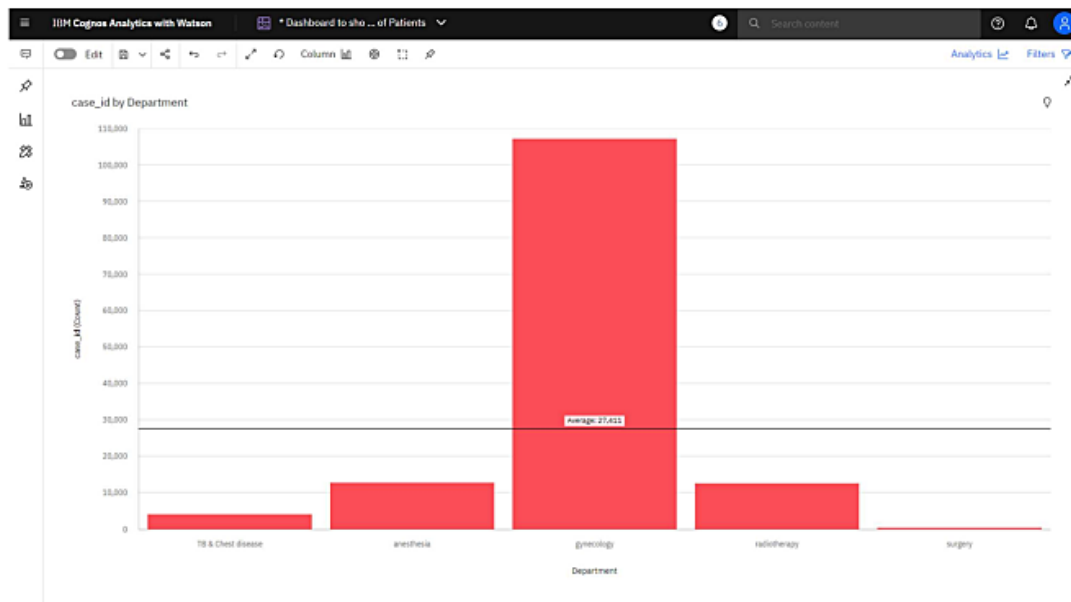
Case id by ward type:



Stay by Department:

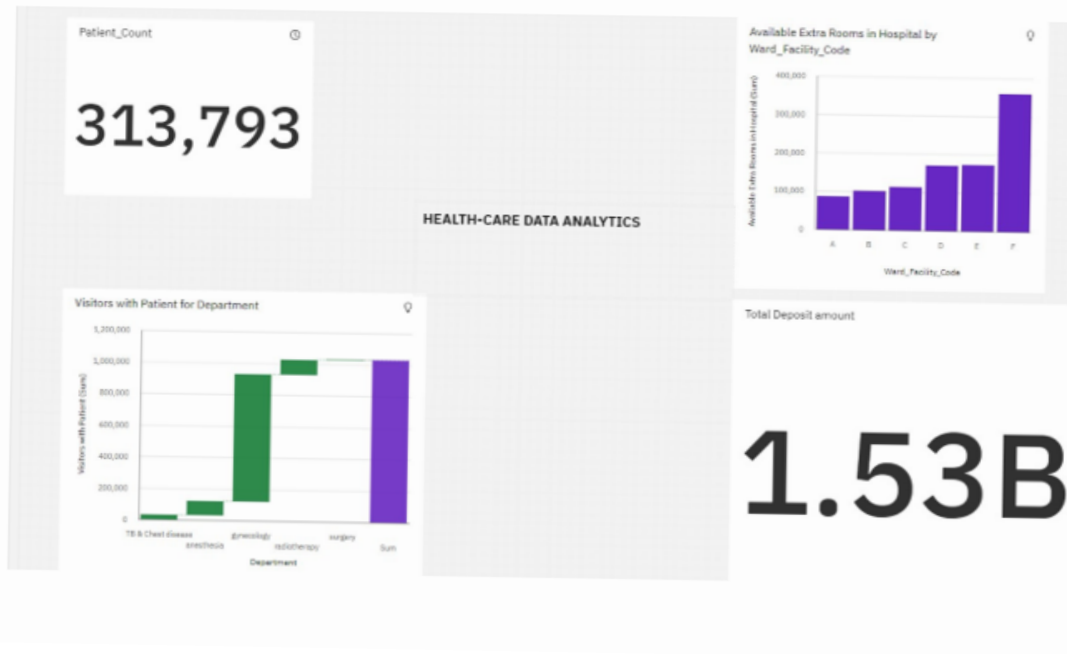


Case id by department:



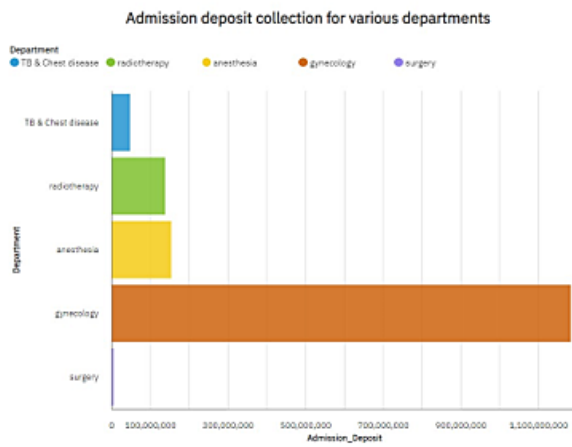
SPRINT 3 CREATION OF DASHBOARD

DASHBOARD

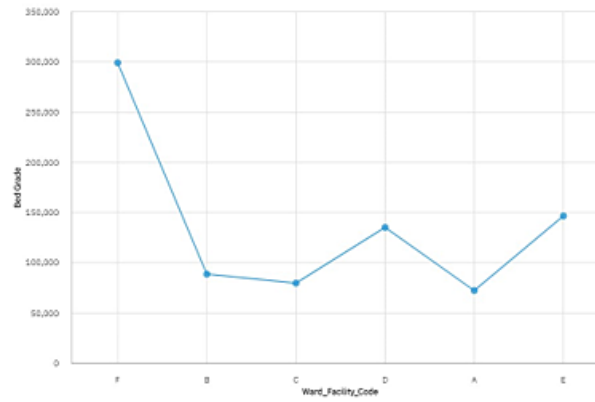


SPRINT 4

Admission deposit and Bed grade

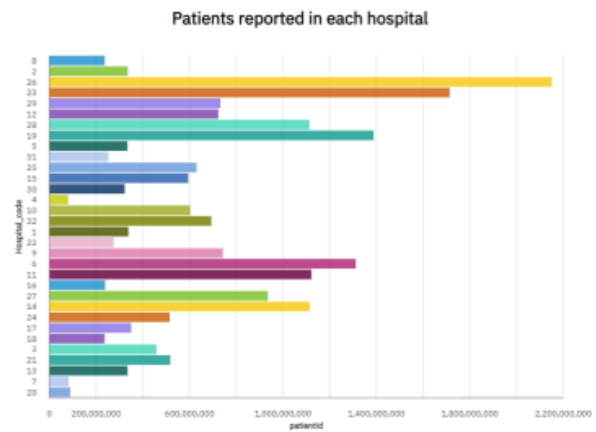


Bed grade compared to ward facility



Available extra beds and Patients reported in each hospital:

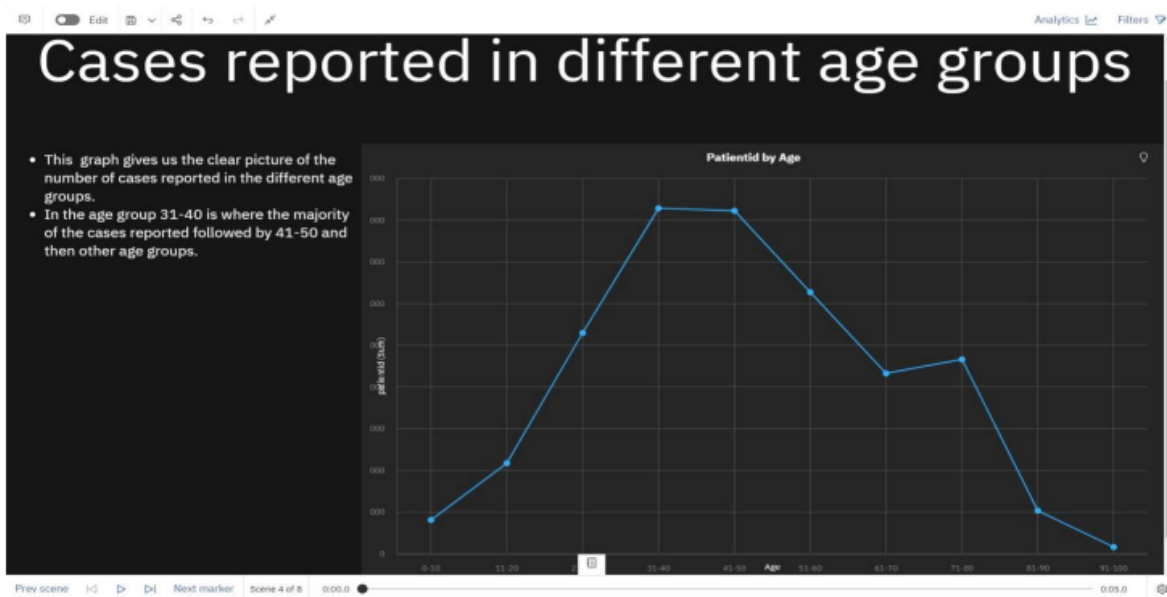
Available extra beds and Patients reported in each hospital



Availability of extra rooms in hospital region code and type code:



Cases reported in different age groups:



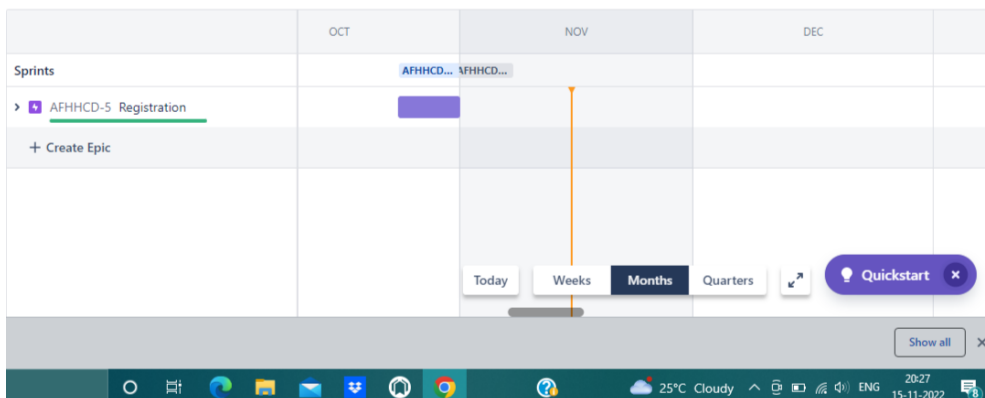
6.2 Sprint DeliverySchedule

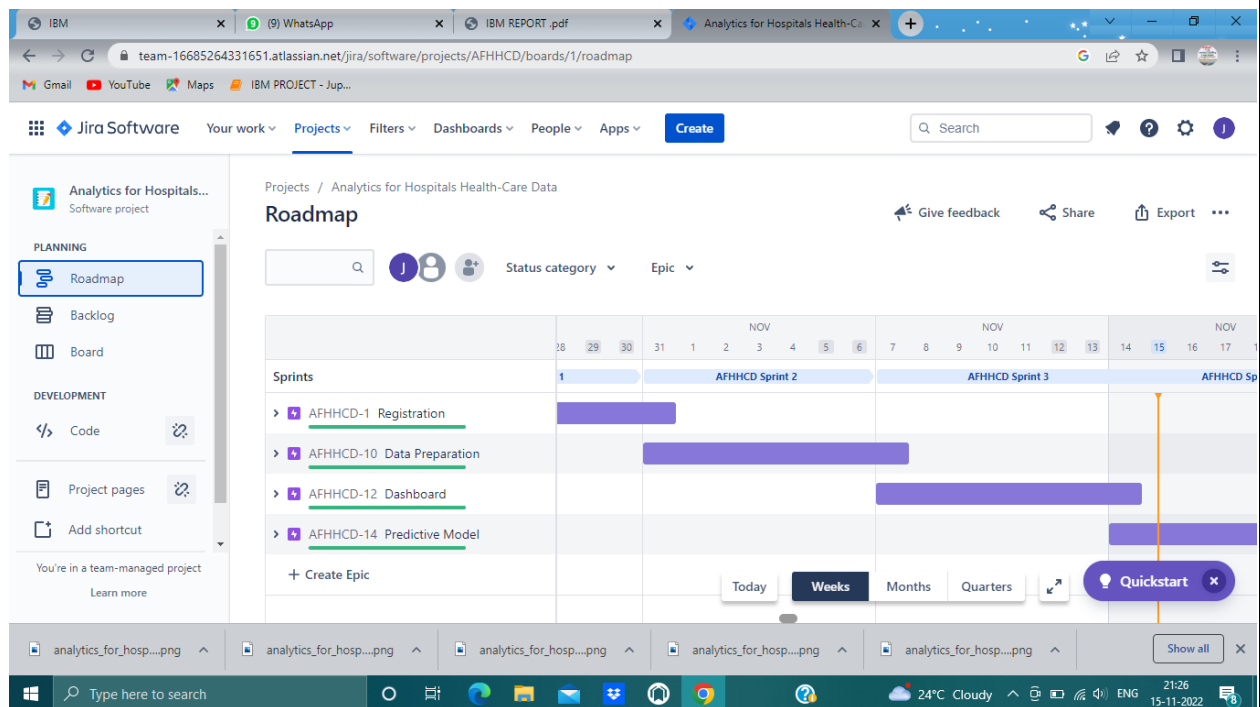
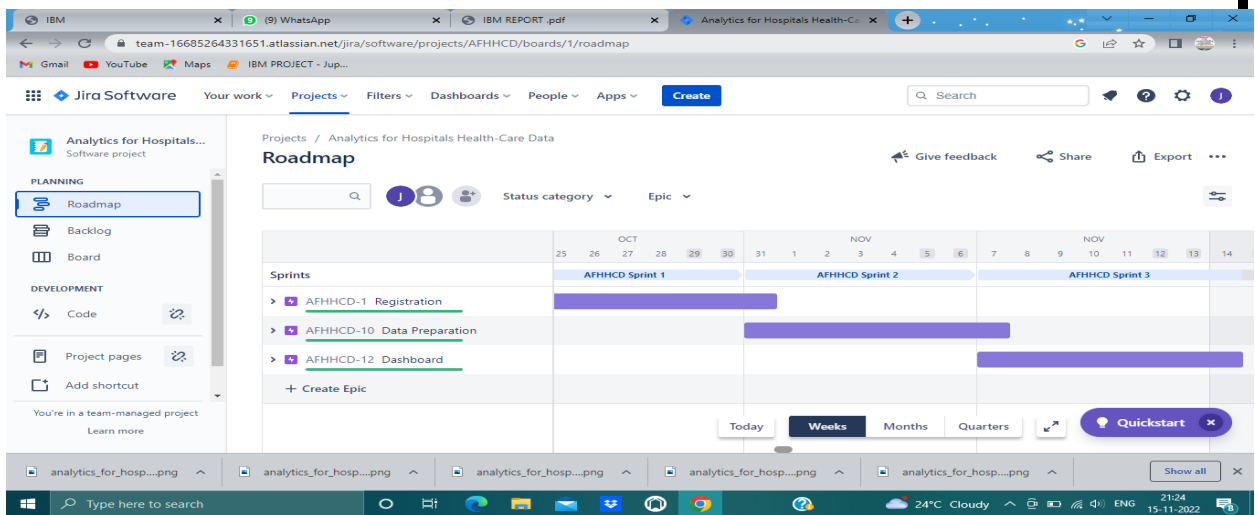
Sprint	Total StoryPoints	Duration	Sprint Start Date	Sprint End Date (Planned)	Story Points Completed (as on Planned)	Sprint Release Date (Actual)

Sprint -1	20	6 Days	24 Oct 2022	29 Oct 2022	20	29 Oct 2022
Sprint -2	20	6 Days	31 Oct 2022	05 Nov 2022	20	05 Nov 2022
Sprint -3	20	6 Days	07 Nov 2022	12 Nov 2022	20	12 Nov 2022
Sprint -4	20	6 Days	14 Nov 2022	19 Nov 2022	20	19 Nov 2022

6.3 Reports from JIRA

Jira Sprints





7. Coding and solutioning Neural Network Model

Neural Networks are built of simple elements called neurons, which take in a real

value, multiply it by weight, and run it through a non-linear activation function. The process records one at a time and learns by comparing their classification of the record with the known actual classification of the record. The errors from the initial classification of the first record are fed back into the network and used to modify the network's algorithm for further iterations. In this neural network model, there are six dense layers, the final layer is an output layer with an activation function "SoftMax". SoftMax is used here because each patient must be classified in one of the 11 levels in the Stay variable. In this model, increasing the number of neurons from each layer to the other layer, will increase the hypothetical space of the model and try to learn more patterns from the data. There are a total of 442,571 trainable parameters. Every layer is activated using "relu" activation function because it overcomes the vanishing gradient problem, allowing models to learn faster and perform better. Finally, evaluating the model with a test set yields an accuracy score of 41.79%. Neural Networks supposedly performs better than any other models. But because of the smaller dataset, it was not able to learn more accurately than the XGBoost model. It nearly took 20 minutes to train the model. In the Naive Bayes model, patients are more likely to be misclassified. This model is biased towards the duration of 21-30 days, it has classified 72,206 patients for this level. Whereas the other two models XGBoost and Neural Networks are predicting mostly similar Length of Stay for the patient. Examining these predictions, many of the patients are staying in the hospital for 21-30 days and very few people are staying for 61-70 days. As far as the distribution of Length of Stay is concerned, 13% of the patients are discharged from the hospital within 20 days and 1% of the overall patients are staying in the hospital for more than 60 days.

XGBoost Model

Boosting is a sequential technique that works on the principle of an ensemble. At any instant T, the model outcomes are weighed based on the outcomes of the previous instant (T -1). It combines the set of weak learners and improves prediction accuracy. Tree ensemble is a set of classification and regression trees. Trees are grown one after another, and they try to reduce the misclassification rate. The final prediction score of the model is calculated by summing up each and individual score. Before feeding train data to the XGB Classifier model, booster parameters must be tuned. Tuning the model can prevent overfitting and can yield higher accuracy. In this XGBoost model, we have used the following parameters for tuning, • learning_rate = 0.1 - step size shrinkage used to prevent overfitting. After each boosting step, we can directly get the weights of new features, and eta shrinks the feature weights to make the boosting process more conservative. • max_depth = 4 – Maximum depth of the tree. This value describes the complexity of the model. Increasing its value results in overfitting. • n_estimators = 800 – Number of gradient boosting trees or rounds. Each new tree

attempts to model and correct for the errors made by the sequence of previous trees. Increasing the number of trees can yield higher accuracy but the model reaches a point of diminishing returns quickly. • objective = 'multi:softmax' – this parameter sets XGBoost to do multiclass classification using the softmax objective because the target variable has 11 Levels.

8.TESTING

Purpose of Document

The product is the responsive dashboard which shows Patient Length of stay in the hospitalbased on the various aspect considering:

1.Based on the severity

of illnesss

2.Based on Departments

3.Type of hospitals

4.Region of Admission

Defect Analysis

This report shows the number of resolved or closed bugs at each severity level, and howthey were resolved

Resolution	Severity 1	Severity 2	Severity 3	Severity 4	Subtotal
By Design	6	3	1	0	10
Duplicate	1	0	0	1	2
External	1	4	1	2	8
Fixed	5	0	6	6	17
Not Reproduced	1	1	0	1	3

Skipped	1	1	0	0	2
Won't Fix	0	1	2	1	4
Totals	15	10	10	11	46

Test Case Analysis

This report shows the number of test cases that have passed, failed, and untested

8.1 Test cases

1. Verify the user is able to get the responsiveness of all the graphs
2. Verify the user get the entire visualization of the dashboard, report.
3. Verify the user get the complete interaction with the website
4. Check if the entire dashboard, Report is visible.
5. User can view pages in the report.
 1. Verify the user is able to access the no of bed based on the region
 2. Verify the user is able to access the bed grade with respect to the severity of illness
 3. Verify the user is able to access the parameters based on the length of stay
 4. Verify the user is able to compare the department based on the Severity of illness.

8.2 USER ACCEPTANCE TESTING

Test case ID	Feature Type	Component	Test Scenario	Actual Result	Status	TC for Automation(Y/N)	BUG ID
--------------	--------------	-----------	---------------	---------------	--------	------------------------	--------

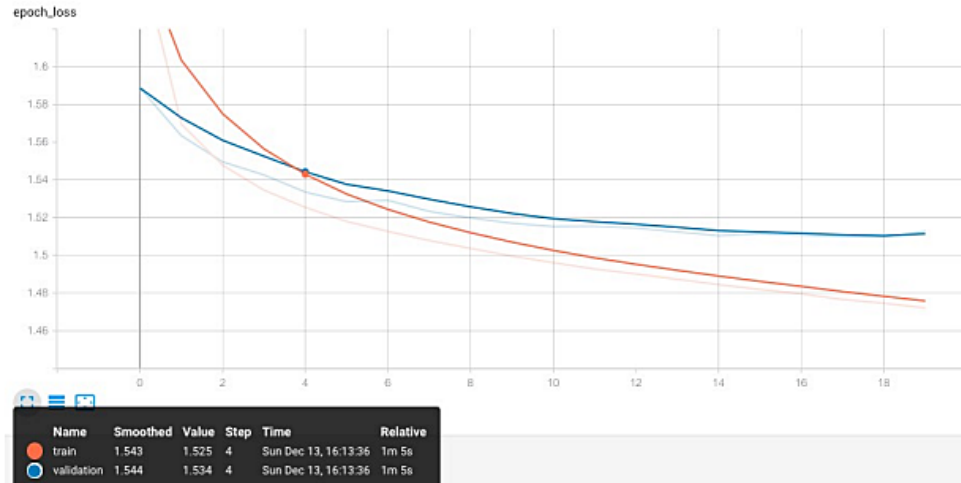
Uploading the data set in the IBM CLOUD	Functional	IBM CLOUD	Loading of all data	UploadedSuccessfully	Pass	Y	-
Responsiveness of dashboard	Functional	Dashboard	Compare the department based on the bed grade	Working as expected	Pass	Y	-
Design	UI	Dashboard	Compatible to the website	Working as expected	Pass	Y	-
Design	Functional	Dashboard	Verify the working of filter	Working as expected	Pass	Y	-
Responsiveness of dashboard	Functional	Dashboard	Verify the user is able to access the bed grade based on the severity of illness	Working as expected	pass	Y	-

Design	Functional	Story	User can view the Report	Working as expected	Pass	Y	-
--------	------------	-------	--------------------------	---------------------	------	---	---

Design	Functional	Report	User can view pages in the report	Working as expected	Pass	Y	-
--------	------------	--------	-----------------------------------	---------------------	------	---	---

9. RESULTS

9.1performance metrices



10. Advantages

1. By predicting a patient's length of stay at the time of admission helps hospitals to allocate resources more efficiently and manage their patients more effectively
2. It helps hospitals in managing resources and in the development of new treatment plans
3. Effective use of hospital resources and reducing the length of stay can reduce overall national medical expenses.

11. Conclusion

In this project, different variables were analyzed that correlate with Length of Stay by using patient-level and hospital-level data. By predicting a patient's length of stay at the time of admission helps hospitals to allocate resources more efficiently and manage their patients more effectively. Identifying factors that associate with LOS to predict and manage the number of days patients stay, could help hospitals in managing resources and in the development of new treatment plans. Effective use of hospital resources and reducing the length of stay can reduce overall national medical expenses.

12. Future Scope

The enormous potential of predictive analysis includes helping identify patients at risk for chronic condition, developing evidence based best practices, proactively spotting potential obstacles to plan adherence.

13. Appendix: Code

Feature engineering:

```
def get_countid_enocde(train, test, cols, name):
    temp = train.groupby(cols)['case_id'].count().reset_index().rename(columns =
{'case_id': name})
    temp2 = test.groupby(cols)['case_id'].count().reset_index().rename(columns =
{'case_id': name}) train = pd.merge(train, temp, how='left', on= cols)
    test = pd.merge(test, temp2, how='left', on= cols) train[name] =
train[name].astype('float')
    test[name] = test[name].astype('float') train[name].fillna(np.median(temp[name]),
inplace = True)
    test[name].fillna(np.median(temp2[name]), inplace = True) return train, test train,
test = get_countid_enocde(train, test, ['patientid'], name = 'count_id_patient') train,
test = get_countid_enocde(train, test, ['patientid', 'Hospital_region_code'], name =
'count_id_patient_hospitalCode') train,
test = get_countid_enocde(train, test, ['patientid', 'Ward_Facility_Code'], name =
'count_id_patient_wardfacilityCode') # Dropping duplicate columns
test1 = test.drop(['Stay', 'patientid', 'Hospital_region_code', 'Ward_Facility_Code'],
axis =1)
train1 = train.drop(['case_id', 'patientid', 'Hospital_region_code',
'Ward_Facility_Code'], axis =1) # Splitting train data for Naive Bayes and XGBoost X1 =
train1.drop('Stay', axis =1)
y1 = train1['Stay'] from sklearn.model_selection import train_test_split X_train,
X_test, y_train,
y_test = train_test_split(X1, y1, test_size =0.20, random_state =100) Models Naïve
bayes Model 24
naive_bayes
import GaussianNB target = y_train.values
features = X_train.values
classifier_nb = GaussianNB()
model_nb = classifier_nb.fit(features, target)
```

```

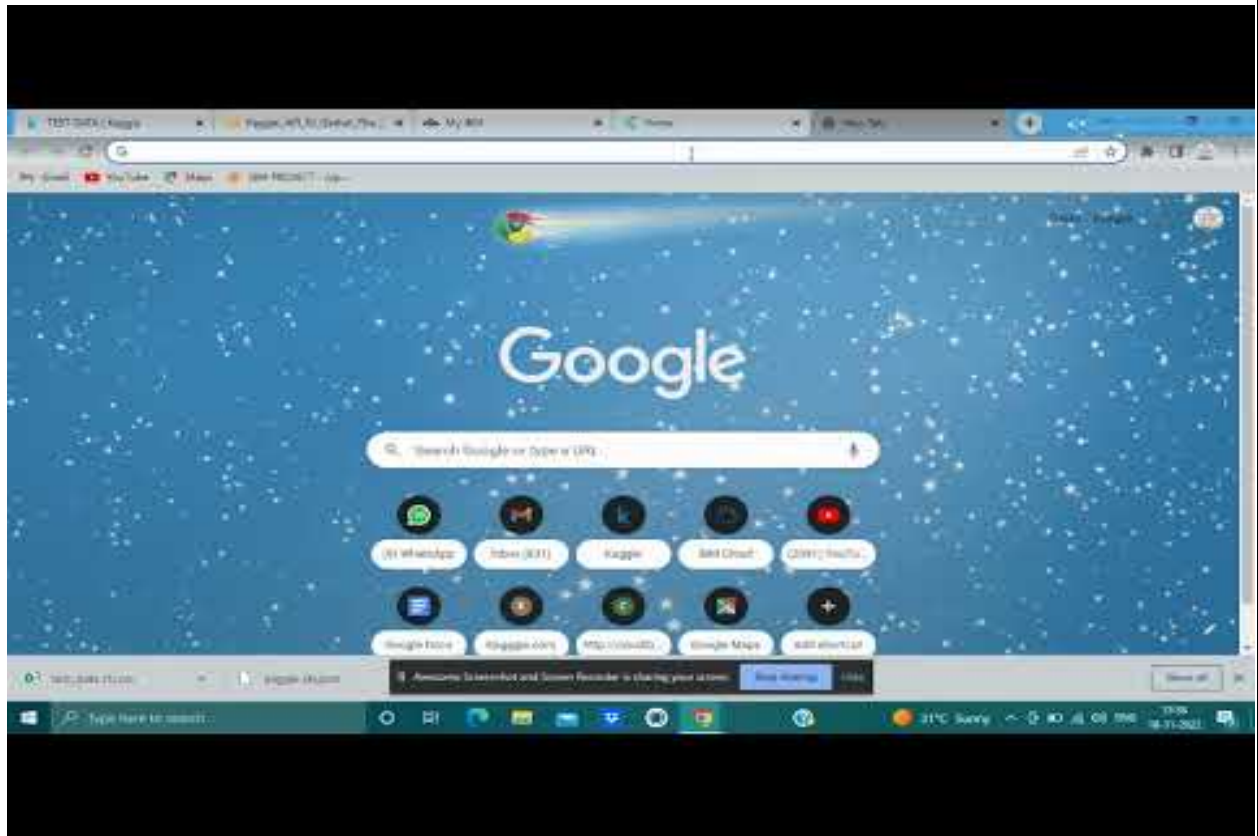
prediction_nb = model_nb.predict(X_test)
from sklearn.metrics
import accuracy_score
acc_score_nb = accuracy_score(prediction_nb,y_test)
print("Accuracy:", acc_score_nb*100)
XGBoost model
import xgboost classifier_xgb = xgboost.XGB
Classifier(max_depth=4, learning_rate=0.1, n_estimators=800,
objective='multi:softmax', reg_alpha=0.5, reg_lambda=1.5, booster='gbtree', n_jobs=4,
min_child_weight=2, base_score= 0.75)
model_xgb = classifier_xgb.fit(X_train, y_train) prediction_xgb =
model_xgb.predict(X_test)
acc_score_xgb = accuracy_score(prediction_xgb,y_test)
print("Accuracy:", acc_score_xgb*100) Neural Network X = train.drop('Stay', axis
=1)
y = train['Stay'] print(X.columns)
z = test.drop('Stay', axis = 1)
print(z.columns)
# Data Scaling from sklearn import preprocessing X_scale =
preprocessing.scale(X)
X_scale.shape X_train, X_test, y_train, y_test = train_test_split(X_scale, y, test_size
=0.20, random_state =100)

```

GITHUB LINK

<https://github.com/IBM-EPBL/IBM-Project-44142-1660722572>

PROJECT DEMO LINK



THANK YOU