# Sprint 3

## TEAM ID - PNT2022TMID21554

## Prediction of Length of Stay

## 1.Uploading Necessary files

## Installing pyspark Libraries



## Ensuring spark is setup and it is running

# Using Machine Learning Algorithm and Principal Component Analysis(PCA)



## Machine Learning

```
[18] input_variable = ['hospital', 'hospital_type', 'hospital_city','hospital_region','available_extra_rooms_in_hospital',
                        'bed_grade','city_code_patient','patient_visitors','admission_deposit',
                        'department_index', 'ward_facility_index', 'ward_type_index', 'illness_severity_index',
                        'type_of_admission_index']

     label = ['stay_days_index']
```

## Principal Component Analysis(PCA)

```
[21] from pyspark.ml.feature import PCA

     pca =PCA(k=10, inputCol="features", outputCol="pcaFeatures")
```

## Standardization

Automatic saving failed. This file was updated remotely or in another tab.   Show diff

## Standardization

```
[22] from pyspark.ml.feature import StandardScaler

     scaler = StandardScaler(inputCol="pcaFeatures", outputCol="scaledFeatures",
                             withStd=True, withMean=False)
```

```
[27] pipeline = Pipeline(stages=[])
```

## Virtualization

## Correlation Matrix

## Decision Tree

```
[37] df.corr().style.background_gradient(cmap='coolwarm').set_precision(2)
```
```
     /usr/local/lib/python3.7/dist-packages/ipykernel launcher.py:1: FutureWarning: this method is deprecated in favour of `Styler.format(precision=..)`
```

# Decision Tree Algorithm for prediction

```
df.corr().style.background_gradient(cmap='coolwarm').set_precision(2)
```
```
/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:1: FutureWarning: this method is deprecated in favour of `Styler.format(precision=..)`
  """Entry point for launching an IPython kernel.
```

| | case_id | Hospital_code | City_Code_Hospital | Available Extra Rooms in Hospital | Bed Grade | patientid | City_Code_Patient | Visitors with Patient | Admission_Deposit |
|---|---|---|---|---|---|---|---|---|---|
| case_id | 1.00 | -0.04 | -0.01 | 0.04 | 0.01 | -0.00 | 0.07 | 0.00 | -0.05 |
| Hospital_code | -0.04 | 1.00 | 0.13 | -0.06 | -0.01 | 0.00 | -0.02 | -0.03 | 0.05 |
| City_Code_Hospital | -0.01 | 0.13 | 1.00 | -0.05 | -0.05 | 0.00 | -0.02 | 0.02 | -0.03 |
| Available Extra Rooms in Hospital | 0.04 | -0.06 | -0.05 | 1.00 | -0.12 | 0.00 | -0.01 | 0.10 | -0.14 |
| Bed Grade | 0.01 | -0.01 | -0.05 | -0.12 | 1.00 | 0.00 | -0.01 | 0.09 | 0.07 |
| patientid | -0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 1.00 | 0.00 | 0.01 | -0.00 |
| City_Code_Patient | 0.07 | -0.02 | -0.02 | -0.01 | -0.01 | 0.00 | 1.00 | -0.01 | 0.03 |
| Visitors with Patient | 0.00 | -0.03 | 0.02 | 0.10 | 0.09 | 0.01 | -0.01 | 1.00 | -0.15 |
| Admission_Deposit | -0.05 | 0.05 | -0.03 | -0.14 | 0.07 | -0.00 | 0.03 | -0.15 | 1.00 |

## Random Forest



```
df.corr().style.background_gradient(cmap='coolwarm').set_precision(2)
```

/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:1: FutureWarning: this method is deprecated in favour of `Styler.format(precision=..)`
  """Entry point for launching an IPython kernel.

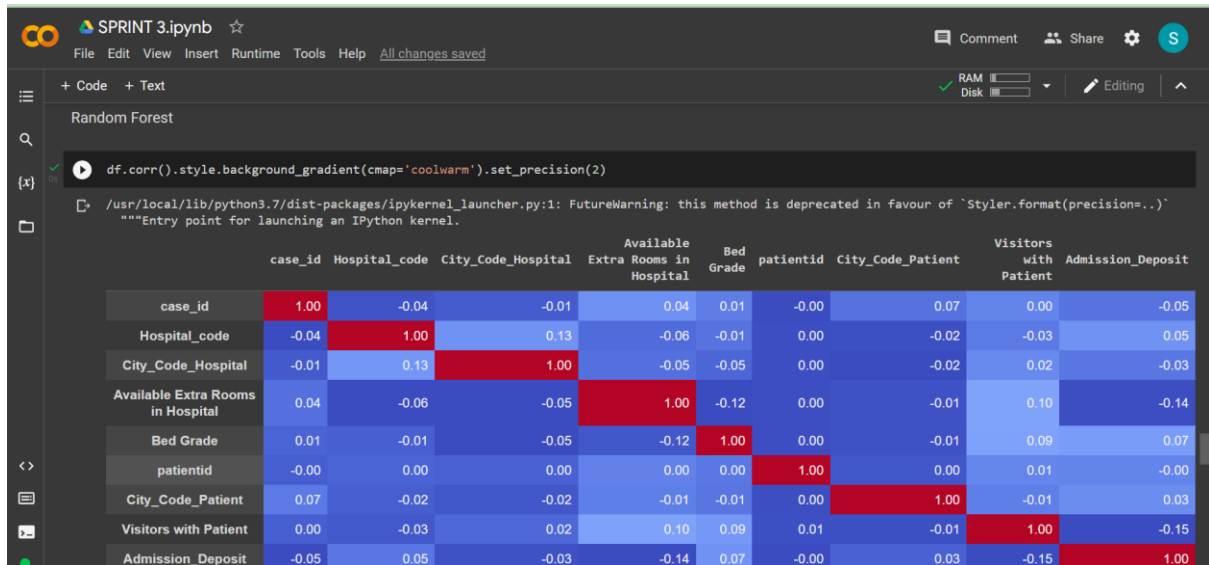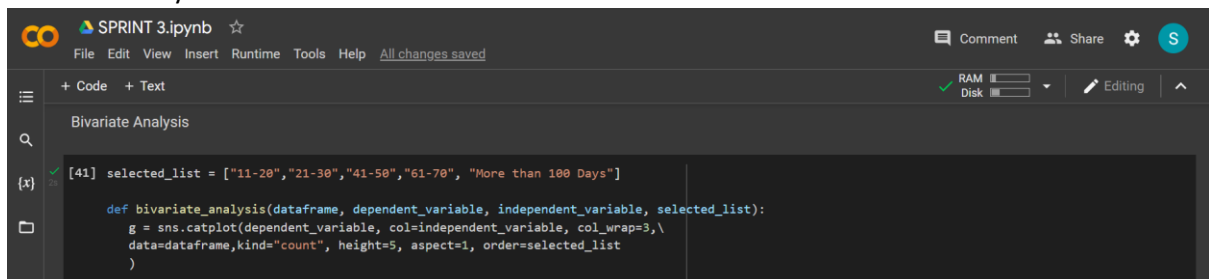| | case_id | Hospital_code | City_Code_Hospital | Available Extra Rooms in Hospital | Bed Grade | patientid | City_Code_Patient | Visitors with Patient | Admission_Deposit |
|---|---|---|---|---|---|---|---|---|---|
| case_id | 1.00 | -0.04 | -0.01 | 0.04 | 0.01 | -0.00 | 0.07 | 0.00 | -0.05 |
| Hospital_code | -0.04 | 1.00 | 0.13 | -0.06 | -0.01 | 0.00 | -0.02 | -0.03 | 0.05 |
| City_Code_Hospital | -0.01 | 0.13 | 1.00 | -0.05 | -0.05 | 0.00 | -0.02 | 0.02 | -0.03 |
| Available Extra Rooms in Hospital | 0.04 | -0.06 | -0.05 | 1.00 | -0.12 | 0.00 | -0.01 | 0.10 | -0.14 |
| Bed Grade | 0.01 | -0.01 | -0.05 | -0.12 | 1.00 | 0.00 | -0.01 | 0.09 | 0.07 |
| patientid | -0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 1.00 | 0.00 | 0.01 | -0.00 |
| City_Code_Patient | 0.07 | -0.02 | -0.02 | -0.01 | -0.01 | 0.00 | 1.00 | -0.01 | 0.03 |
| Visitors with Patient | 0.00 | -0.03 | 0.02 | 0.10 | 0.09 | 0.01 | -0.01 | 1.00 | -0.15 |
| Admission_Deposit | -0.05 | 0.05 | -0.03 | -0.14 | 0.07 | -0.00 | 0.03 | -0.15 | 1.00 |

## Bivariate Analysis



```
[41] selected_list = ["11-20","21-30","41-50","61-70", "More than 100 Days"]

     def bivariate_analysis(dataframe, dependent_variable, independent_variable, selected_list):
       g = sns.catplot(dependent_variable, col=independent_variable, col_wrap=3,\
       data=dataframe,kind="count", height=5, aspect=1, order=selected_list)
       )
```

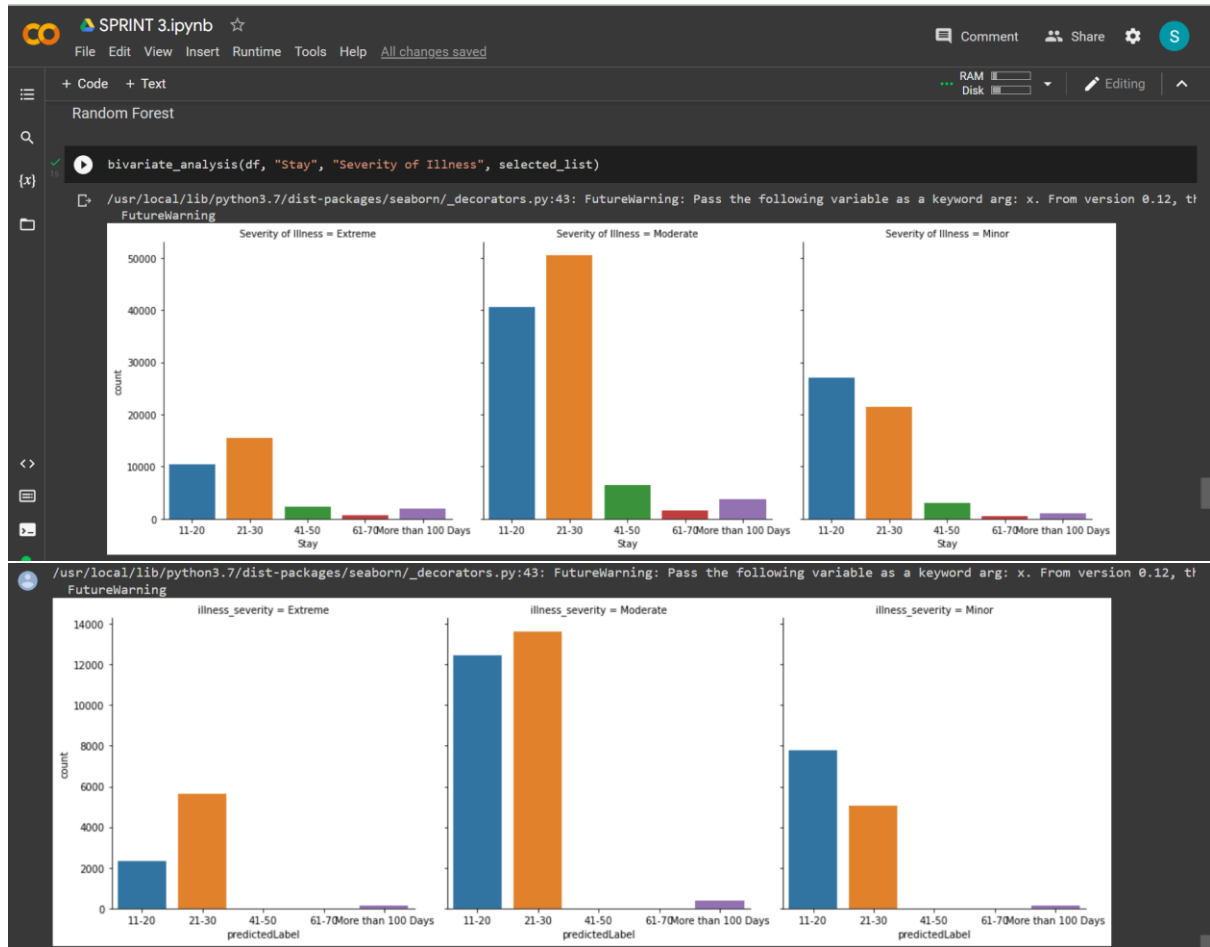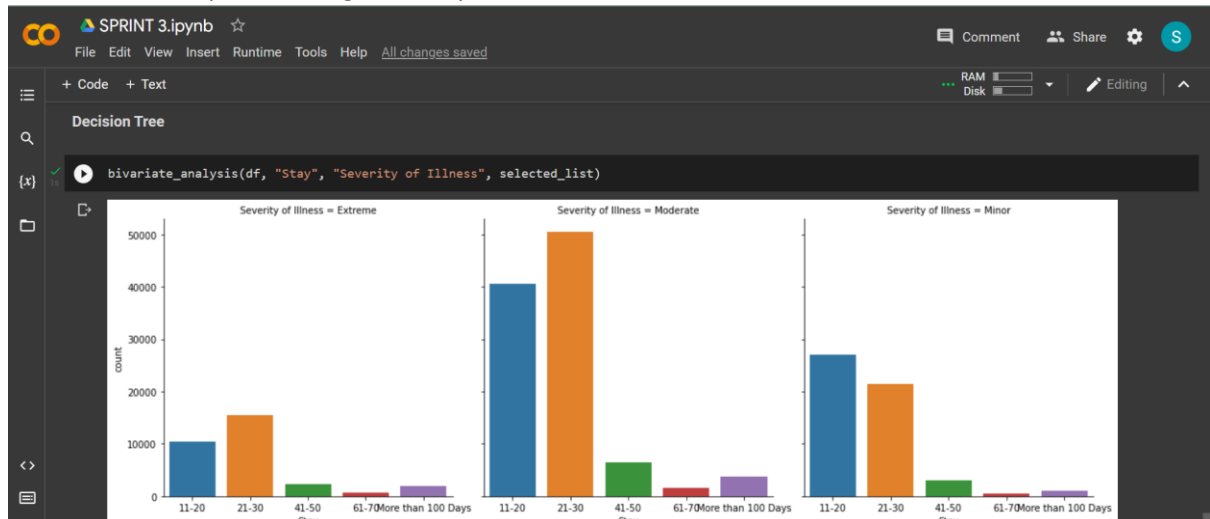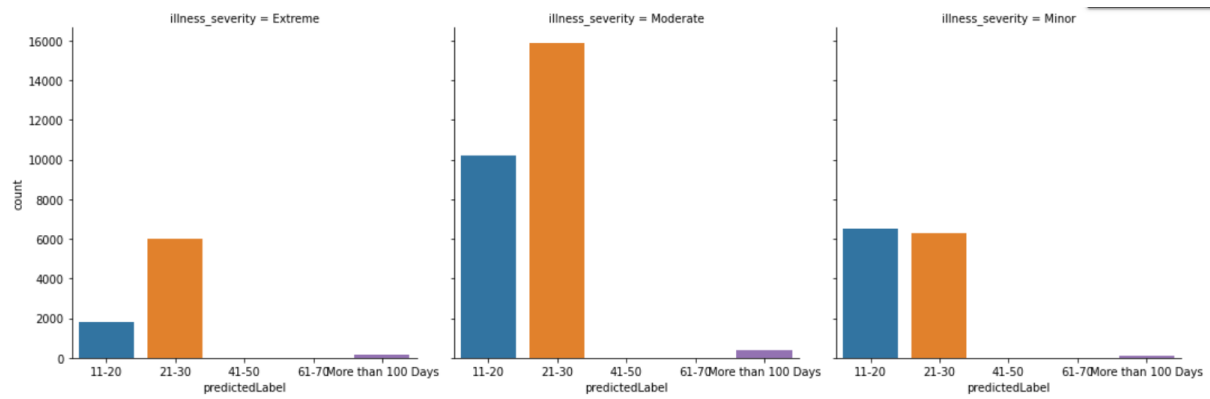## Random Forest to predict Length of Stay



## Decision Tree to predict Length of Stay

Link

https://colab.research.google.com/drive/1Z6ZB9STV8FzW3ry-uNE0jDvupfpJxPLq#scrollTo=UPlncaqRxU8d