

Dashboard Layout

Device Preview

Size

Automatic

Sheets

- Hist_BP
- Hist_BMI
- Hist_Glucose
- Hist_Insulin
- Correlation of ...
- Bubble Chart

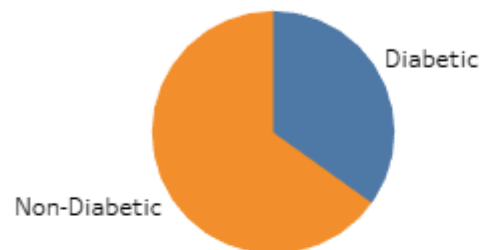
Objects

- Horizontal
- Blank
- Vertical
- Navigation
- Text
- Download
- Image
- Extension
- Web Page
- Ask Data
- Data Story
- Workflow

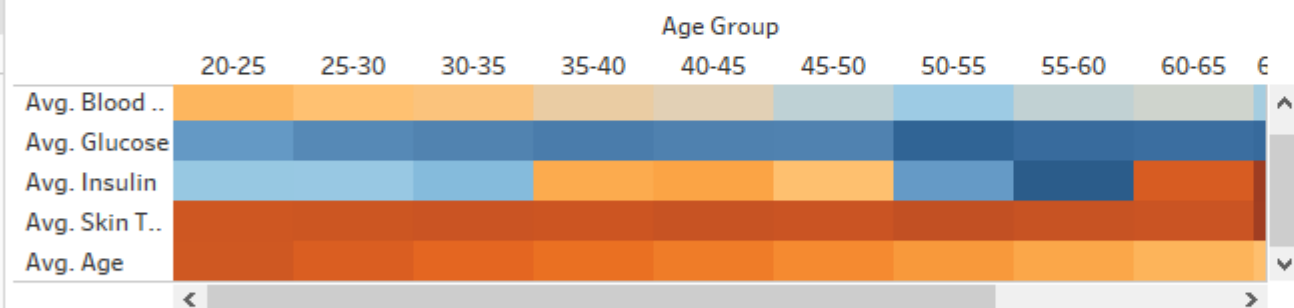
Tiled Floating

☐ Show dashboard title

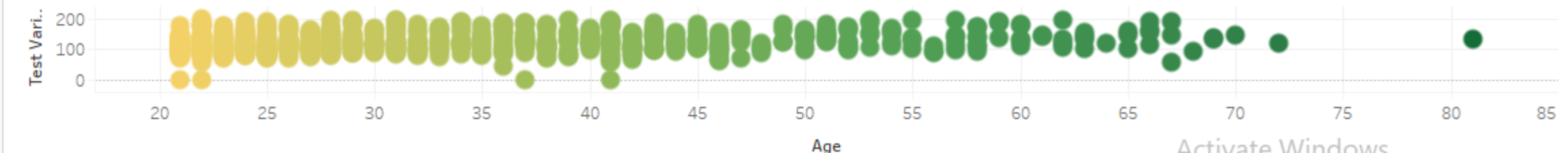
Analysis of Diabetes Population



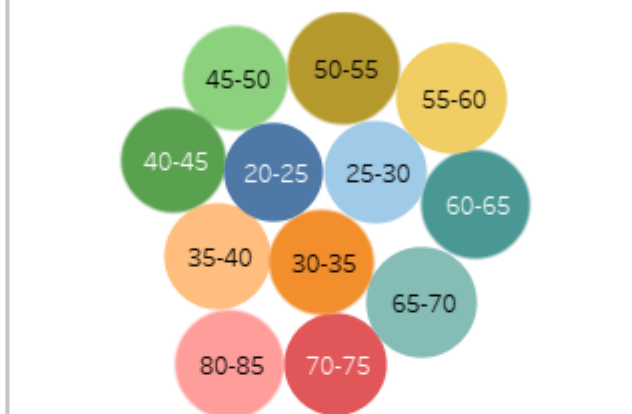
Correlation of Heatmap



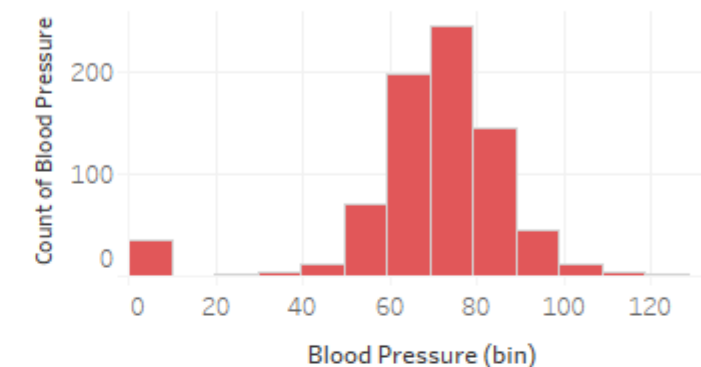
Scatter Chart - Analysis of Variable Relationship



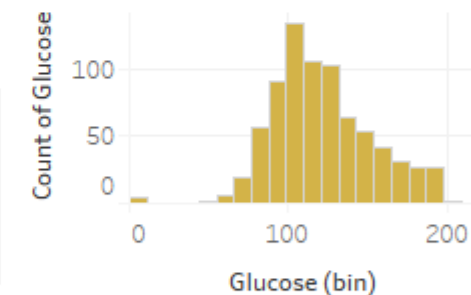
Bubble Chart



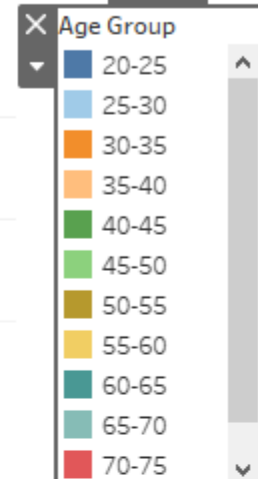
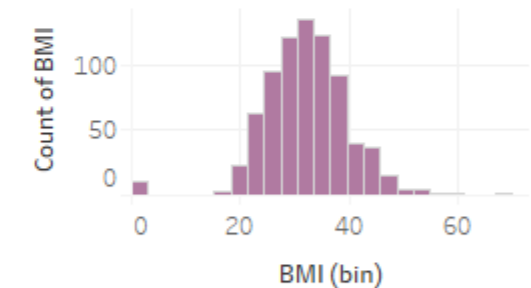
Histogram of Blood Pressure



Hist_Glucose



Hist_BMI



Data Analytics

health care diabetes

Search   

Tables

Age Group
Blood Pressure (bin)
BMI (bin)
Diabetes Pedigree Fun...
Diabetic or Not
Glucose (bin)
Insulin (bin)
Pregnancies (bin)
Skin Thickness (bin)
Measure Names
Age
Blood Pressure
BMI
Bubble
Diabetes Pedigree Fun...
Glucose
Insulin
Outcome
Pregnancies
Skin Thickness

Parameters

Parameter 1

Pages

Filters

Marks

Automatic

Color

Size

Label

Detail

Tooltip

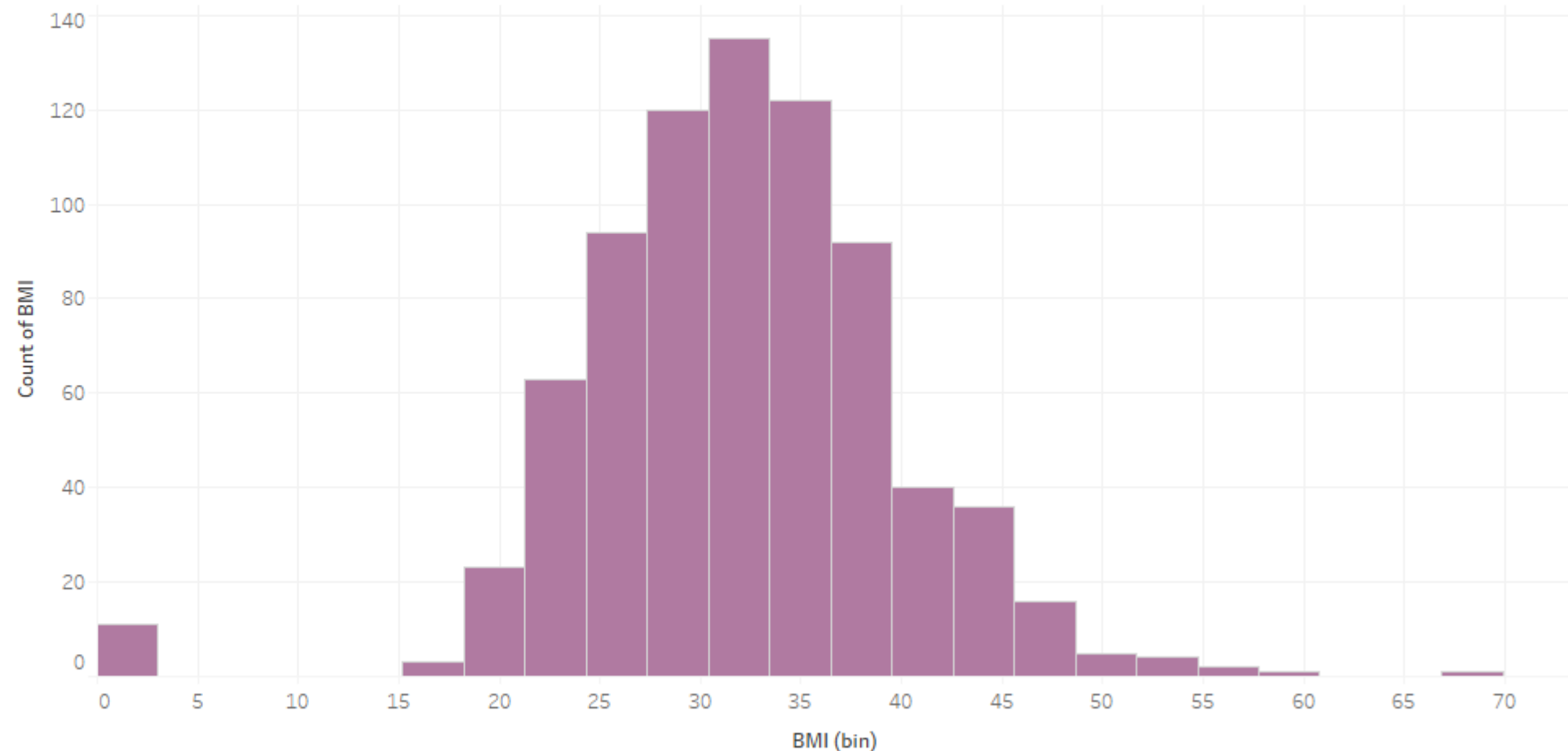
Columns

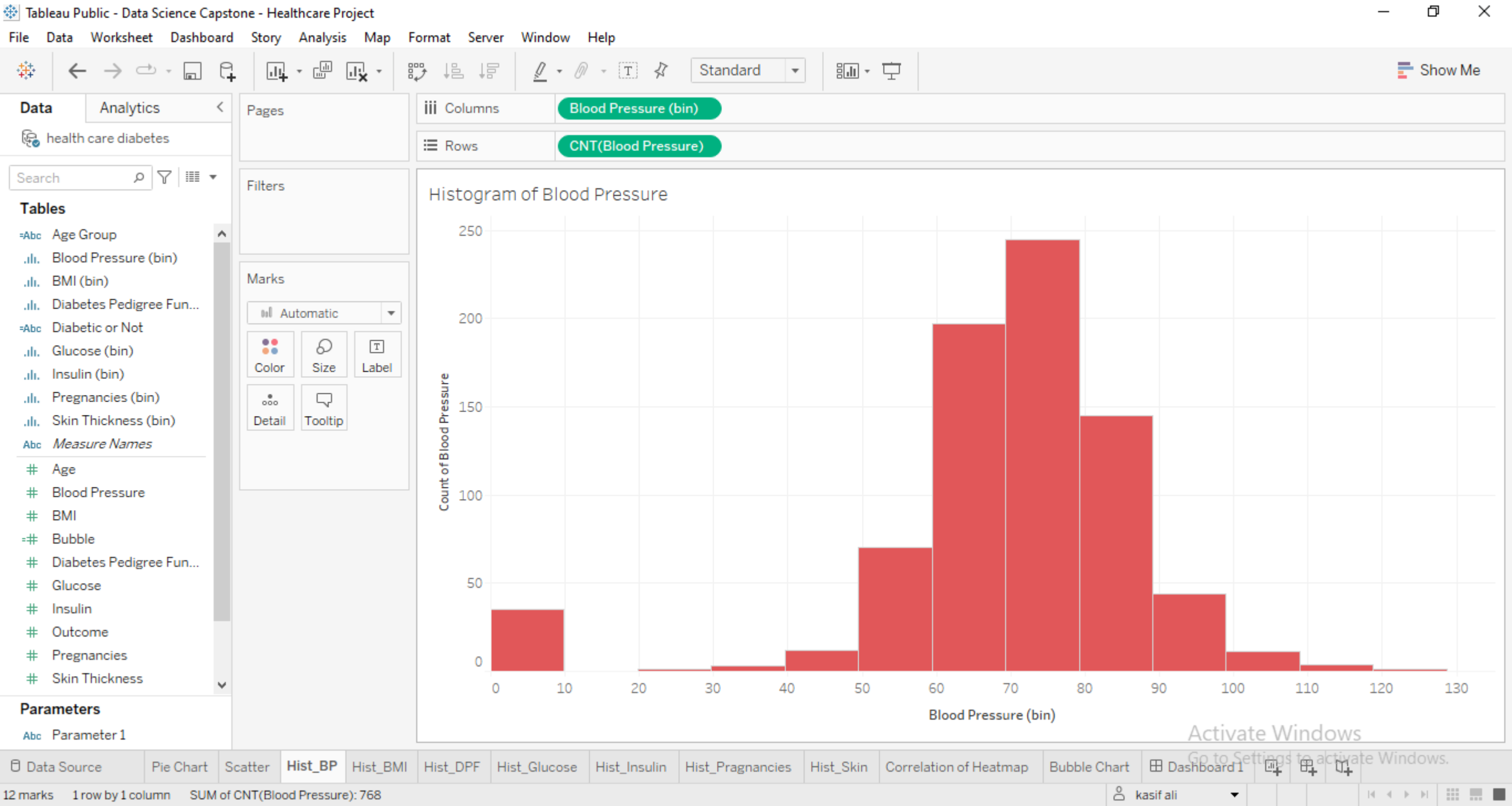
BMI (bin)

Rows

CNT(BMI)

Hist_BMI





Data Analytics

health care diabetes

Search   

Tables

Age Group
Blood Pressure (bin)
BMI (bin)
Diabetes Pedigree Fun...
Diabetic or Not
Glucose (bin)
Insulin (bin)
Pregnancies (bin)
Skin Thickness (bin)
Measure Names

Age
Blood Pressure
BMI
Bubble
Diabetes Pedigree Fun...
Glucose
Insulin
Outcome
Pregnancies
Skin Thickness

Parameters

Parameter 1

Pages

Filters

Marks

Automatic

Color Size Label

Detail Tooltip

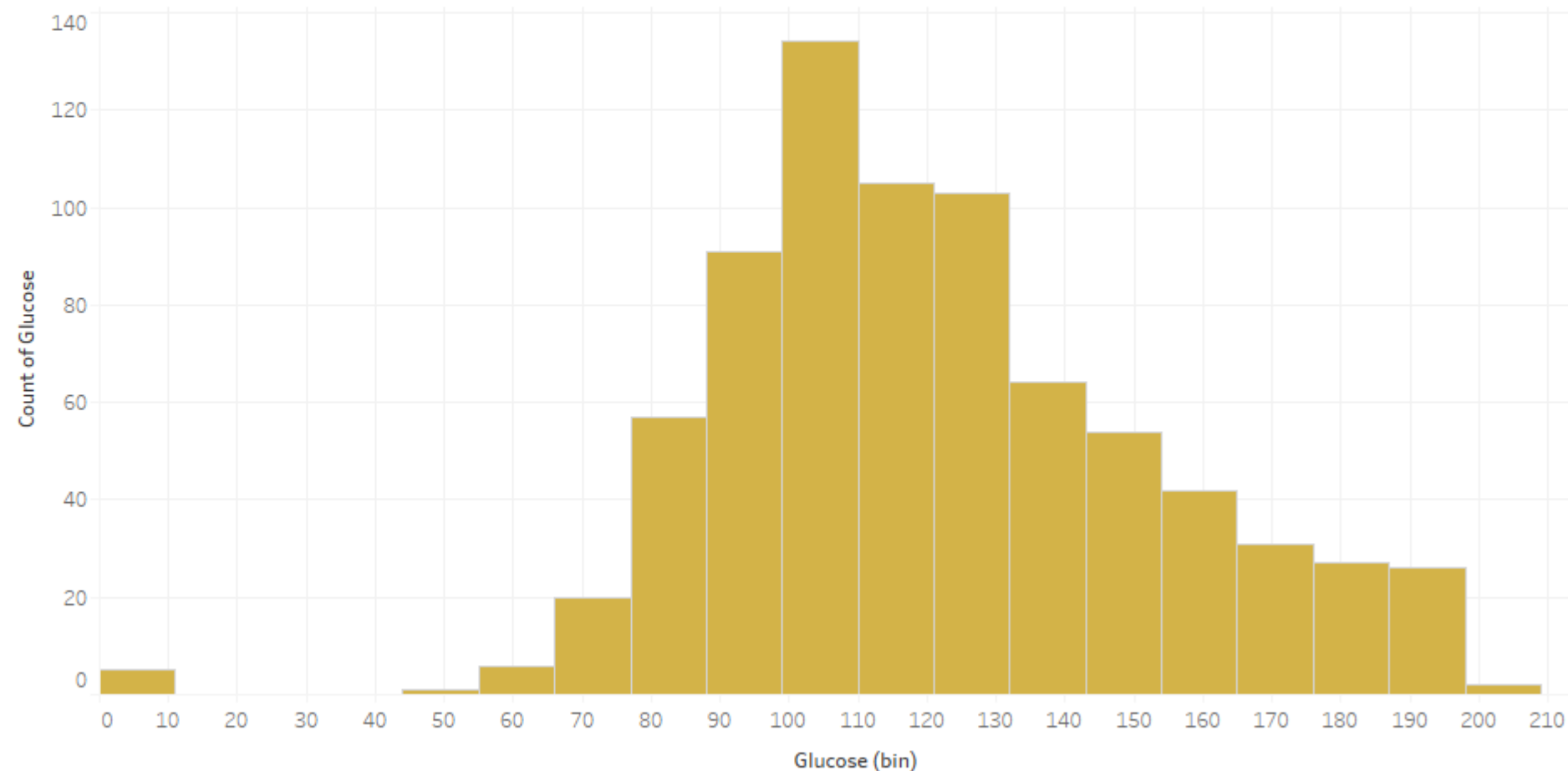
Columns

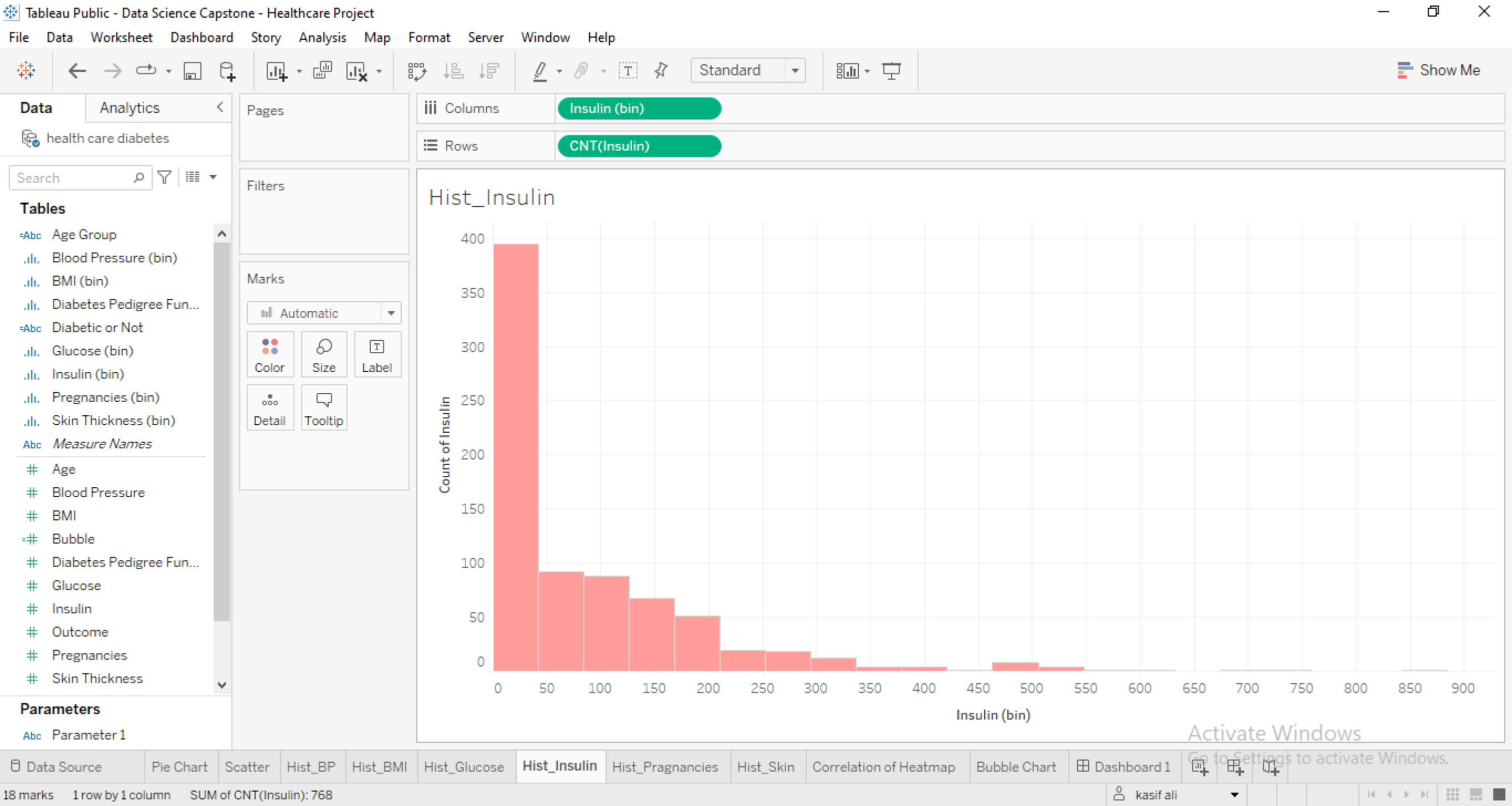
Glucose (bin)

Rows

CNT(Glucose)

Hist_Glucose





Data Analytics

health care diabetes

Search   

Tables

- Age Group
- Blood Pressure (bin)
- BMI (bin)
- Diabetes Pedigree Fun...
- Diabetic or Not
- Glucose (bin)
- Insulin (bin)
- Pregnancies (bin)
- Skin Thickness (bin)
- Measure Names
- Age
- Blood Pressure
- BMI
- Bubble
- Diabetes Pedigree Fun...
- Glucose
- Insulin
- Outcome
- Pregnancies
- Skin Thickness

Parameters

Parameter 1

Pages

Columns

Rows

Filters

Marks

Pie

Color

Size

Label

Detail

Tooltip

Angle

Diabetic or Not

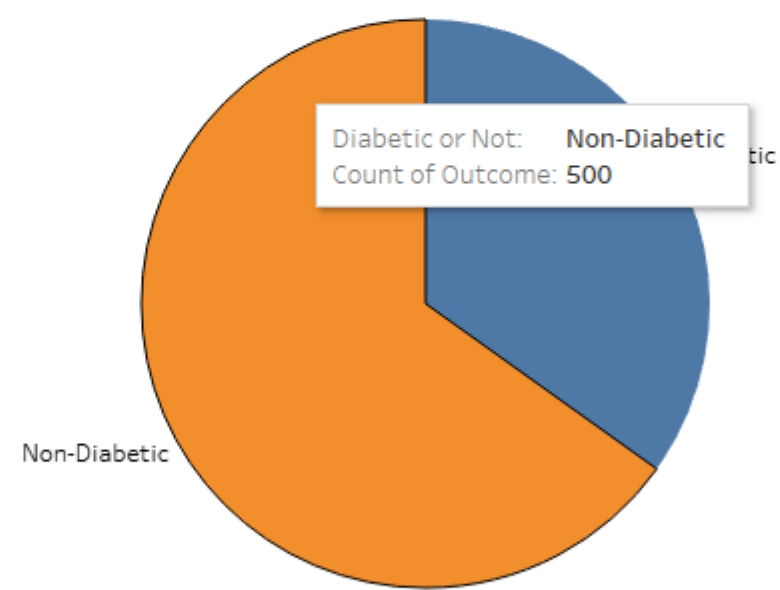
CNT(Outcome)

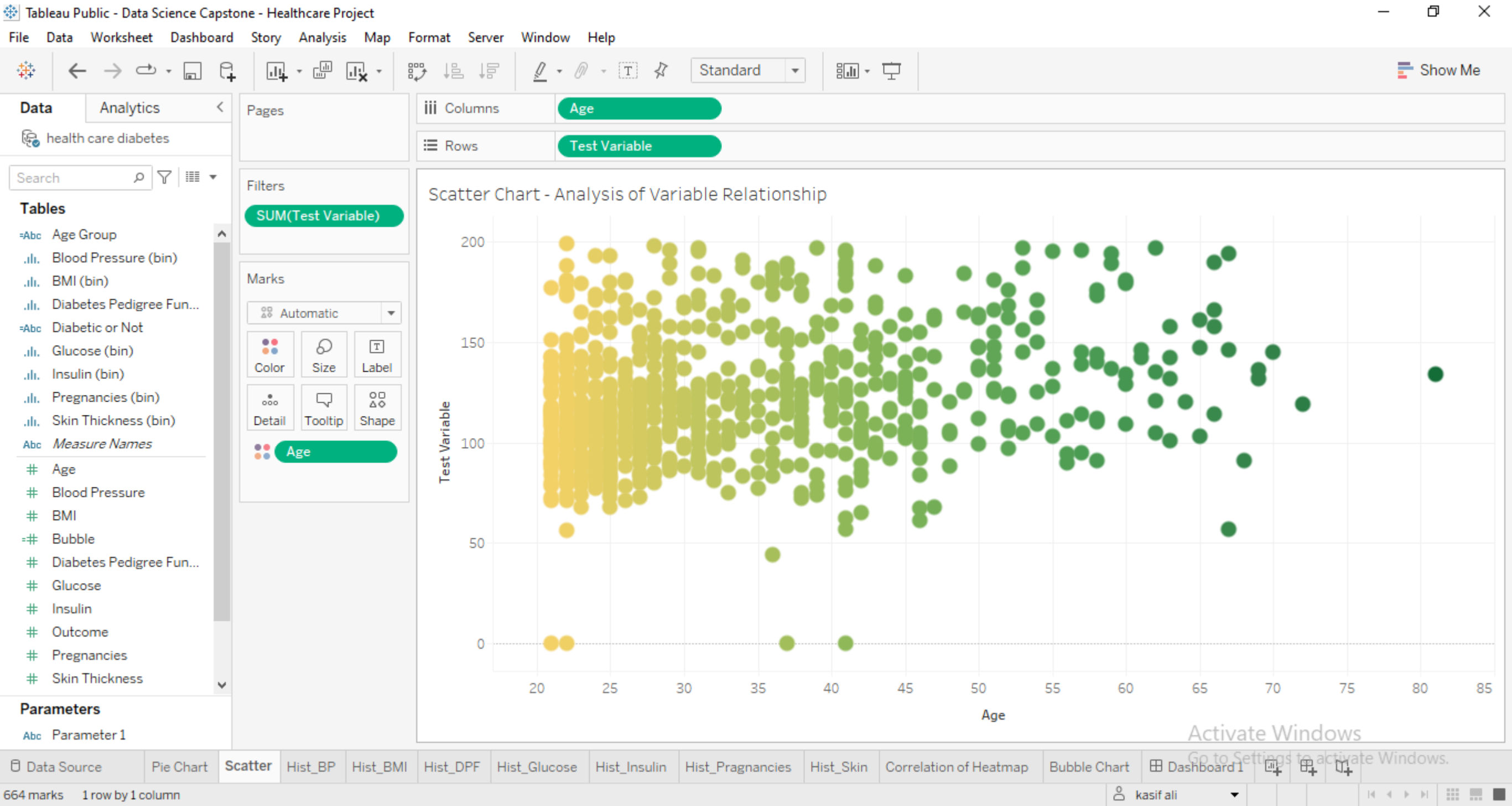
Diabetic or Not

Columns

Rows

Analysis of Diabetes Population





WEEK 1

Reading Data and Performing descriptive analysis

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 from sklearn.model_selection import train_test_split
5 from sklearn.neighbors import KNeighborsClassifier
```

```
1 import os
2 print(os.getcwd())
3 os.listdir(os.getcwd())
```

```
C:\Users\User\Desktop\Capstone Project 2 Healthcare\Project 2

Out[4]: ['_ipynb_checkpoints',
        'Case Study - healthcare Industry.docx',
        'Data Science Capstone Project Healthcare.ipynb',
        'health care diabetes.csv',
        'Healthcare - Diabetes',
        'healthcare appointment data.csv',
        'train.csv']
```

```
1 df=pd.read_csv('health care diabetes.csv')
2 df.head()
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
0	6	148	72	35	0	33.6	0.627	50	1
1	1	85	66	29	0	26.6	0.351	31	0
2	8	183	64	0	0	23.3	0.672	32	1
3	1	89	66	23	94	28.1	0.167	21	0
4	0	137	40	35	168	43.1	2.288	33	1

```
1 # It shows how many null values are present and their data types
2 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 768 entries, 0 to 767
Data columns (total 9 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   Pregnancies            768 non-null    int64
 1   Glucose                768 non-null    int64
 2   BloodPressure          768 non-null    int64
 3   SkinThickness          768 non-null    int64
 4   Insulin                768 non-null    int64
 5   BMI                   768 non-null    float64
 6   DiabetesPedigreeFunction 768 non-null    float64
 7   Age                   768 non-null    int64
 8   Outcome                768 non-null    int64
dtypes: float64(2), int64(7)
memory usage: 54.1 KB
```

```
1 # It displays unique value in each column
2 df.nunique()
```

```
Pregnancies    17
Glucose         136
BloodPressure   47
SkinThickness   51
Insulin        186
BMI             248
DiabetesPedigreeFunction 517
Age             52
Outcome         2
dtype: int64
```

```
1 # Here Min column shows 0 which means missing values
2 df.describe()
```

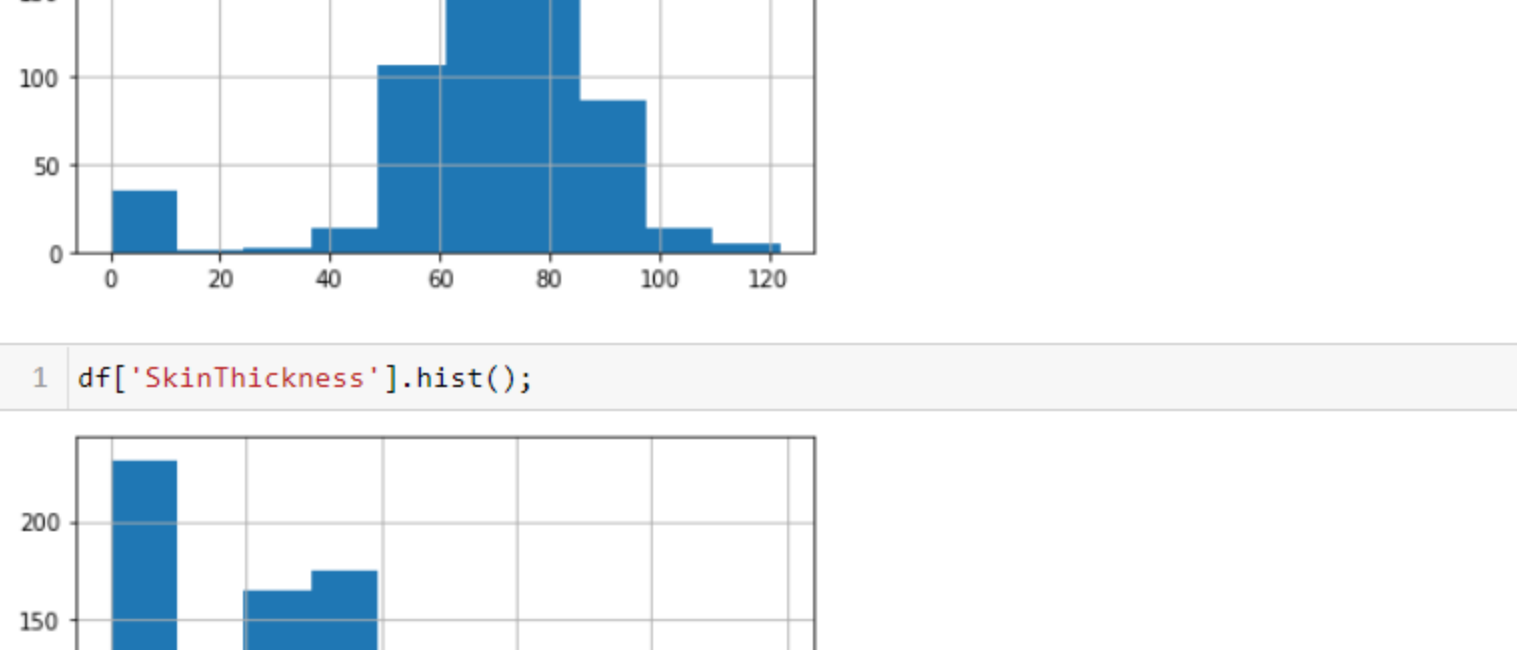
	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
count	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000
mean	3.845052	120.894631	69.155409	29.536450	79.799479	31.902578	0.471876	33.240885	0.348958
std	3.369578	31.972618	19.355807	15.952218	115.244002	7.884160	0.331329	11.760232	0.476951
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.078800	21.000000	0.000000
25%	1.000000	99.000000	62.000000	0.000000	0.000000	27.000000	0.243750	24.000000	0.000000
50%	3.000000	117.000000	72.000000	23.000000	30.500000	32.000000	0.372500	29.000000	0.000000
75%	6.000000	140.250000	80.000000	32.000000	127.250000	36.600000	0.626250	41.000000	1.000000
max	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	2.420000	81.000000	1.000000

```
1 df.corr()
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
Pregnancies	1.000000	0.129459	0.141282	-0.081672	-0.073535	0.017683	-0.033523	0.544341	0.221898
Glucose	0.129459	1.000000	0.152590	0.057320	0.331357	0.221071	0.137327	0.263514	0.466581
BloodPressure	0.141282	0.152590	1.000000	0.207371	0.088933	0.281805	0.041265	-0.239528	0.065688
SkinThickness	-0.081672	0.057328	0.207371	1.000000	0.436783	0.392573	0.183928	-0.139970	0.074752
Insulin	-0.073535	0.331357	0.088933	0.436783	1.000000	0.197859	0.185071	-0.042163	0.130548
BMI	0.017683	0.221071	0.281805	0.392573	0.197859	1.000000	0.140647	0.033561	0.036242
DiabetesPedigreeFunction	-0.033523	0.137337	0.041265	0.183928	0.185071	0.140647	1.000000	0.033561	0.173844
Age	0.544341	0.263514	0.239528	-0.139970	-0.042163	0.036242	0.033561	1.000000	0.238356
Outcome	0.221898	0.466581	0.065688	0.074752	0.130548	0.292695	0.173844	0.238356	1.000000

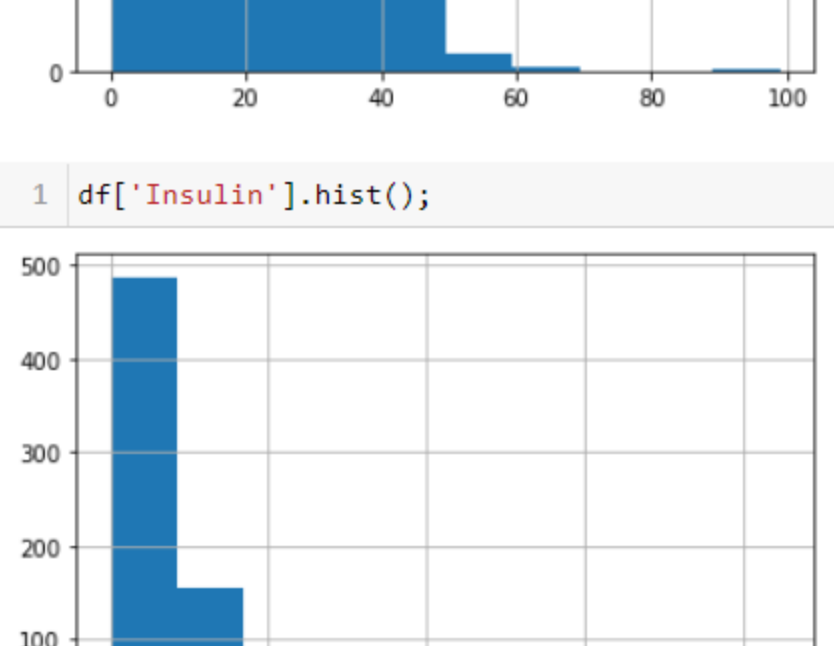
```
1 plt.figure(figsize=(12,8))
2 df.boxplot()
```

<AxesSubplot>

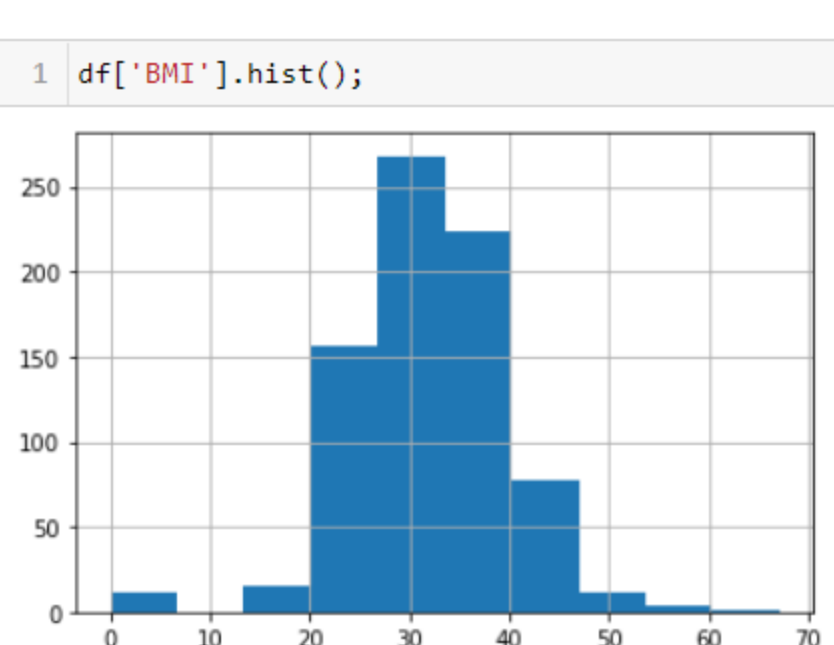


```
1 # Visually exploring these variables using histograms. Treating the missing values accordingly
```

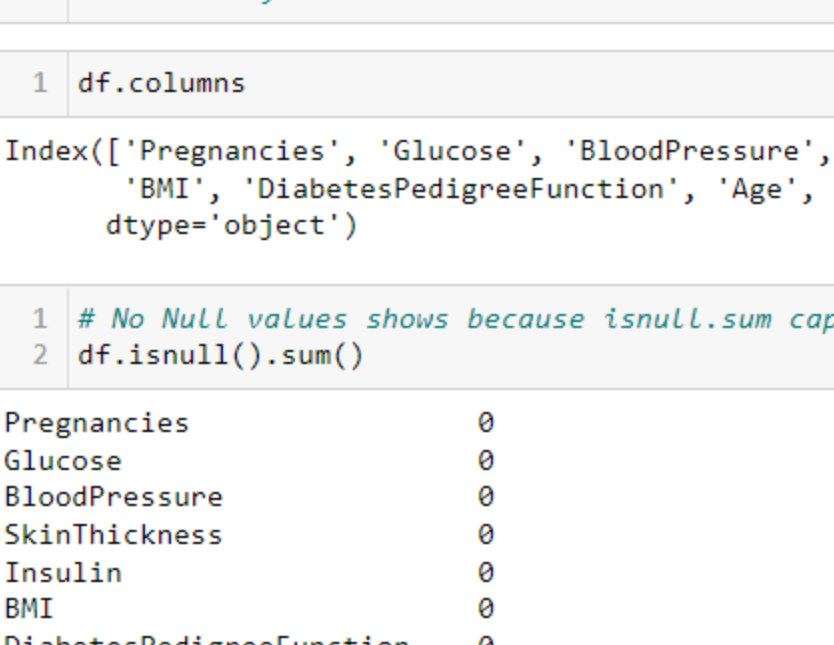
```
1 df['Glucose'].hist()
```



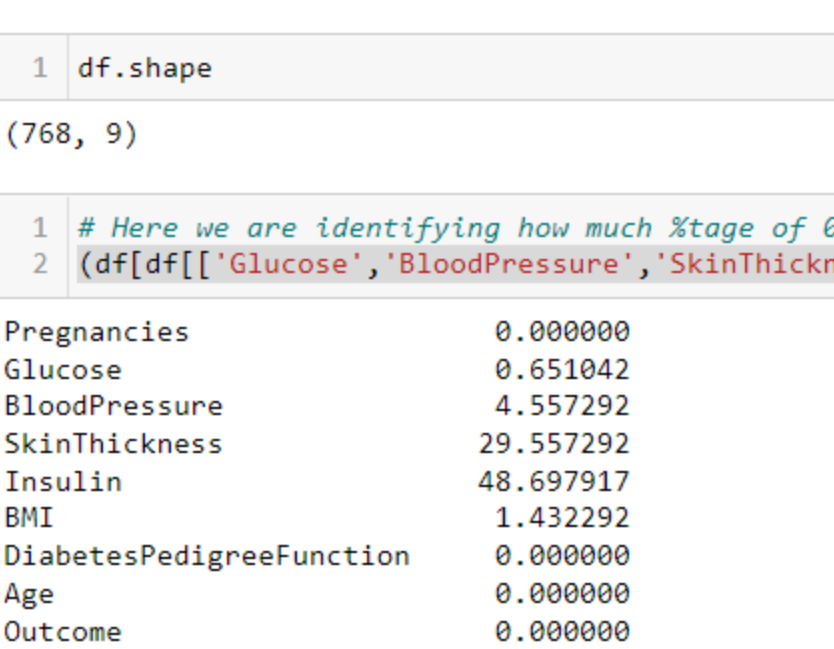
```
1 df['BloodPressure'].hist()
```



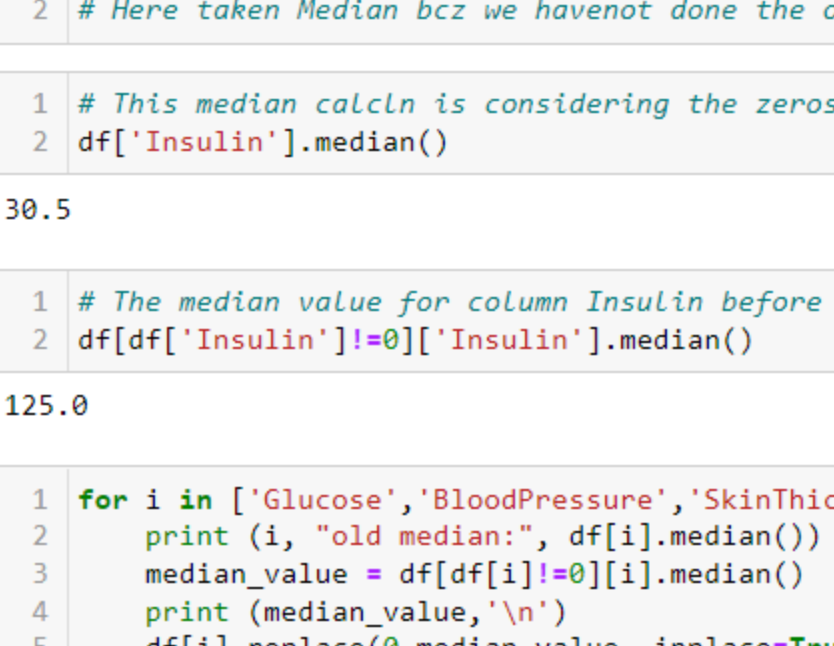
```
1 df['SkinThickness'].hist()
```



```
1 df['Insulin'].hist()
```



```
1 df['BMI'].hist()
```



```
1 # From above histograms, it is clear that Insulin has highly skewed data distribution and remaining 4 variables have
2 # relatively balanced data distribution therefore we will treat missing values in these 5 variables
```

```
1 df.columns
```

```
Out[16]: Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
        'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome'],
        dtype='object')
```

```
1 # No Null values shows because isnull.sum() captures null for NaN values.
2 df.isnull().sum()
```

```
Pregnancies    0
Glucose         0
BloodPressure   0
SkinThickness   0
Insulin         0
BMI             0
DiabetesPedigreeFunction 0
Age             0
Outcome         0
dtype: int64
```

```
1 df.shape
```

```
Out[17]: (768, 9)
```

```
1 # Here we are identifying how much %age of 0 is there in each column.
2 df[df[['Glucose','BloodPressure','SkinThickness','Insulin','BMI']==0]].count()/len(df)*100
```

```
Pregnancies    0.000000
Glucose         0.651842
BloodPressure   4.557292
SkinThickness   29.557292
Insulin        48.697917
BMI             1.432292
DiabetesPedigreeFunction 0.000000
Age             0.000000
Outcome        0.000000
dtype: float64
```

```
1 # Treating the missing value
2 # Here taken Median bcz we haven't done the outlier treatment if we would have done outlier then we can take Mean also.
3 df['Insulin'].median()
```

```
Out[20]: 30.5
```

```
1 # The median value for column Insulin before removing 0 is 30.50 and after removing 0 is 125.00
2 df[df['Insulin']!=0]['Insulin'].median()
```

```
Out[21]: 125.0
```

```
1 for i in ['Glucose','BloodPressure','SkinThickness','Insulin','BMI']:
2     print(i, "old median:", df[i].median())
3     median_value = df[df[i]!=0][i].median()
4     print(median_value, "\n")
5     df[i].replace(0,median_value, inplace=True)
```

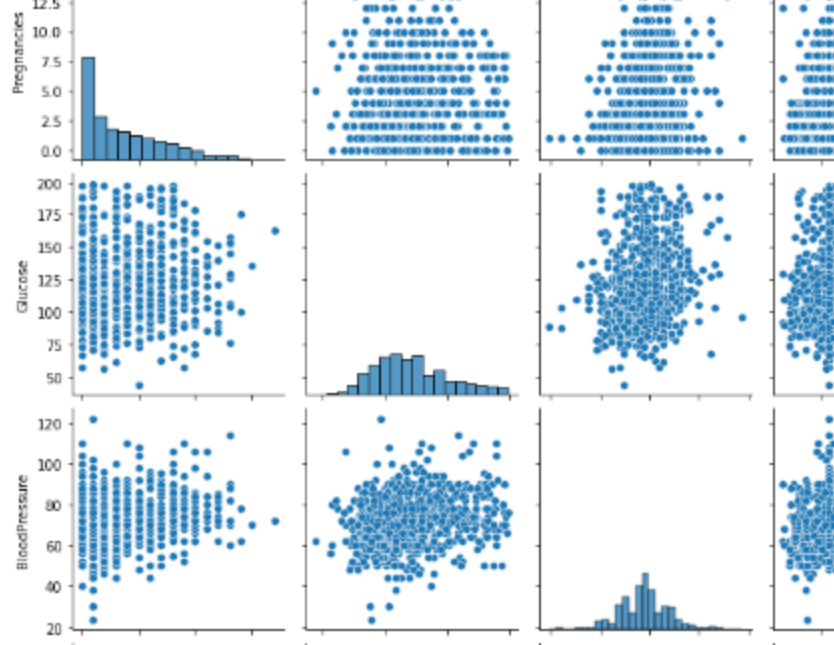
```
Glucose old median: 117.0
117.0
BloodPressure old median: 72.0
72.0
SkinThickness old median: 23.0
29.0
Insulin old median: 30.5
125.0
BMI old median: 32.0
32.3
```

```
1 # here it shows that in all columns wherever the value is 0 the count is shown for all columns
2 df[df[['Glucose','BloodPressure','SkinThickness','Insulin','BMI']==0]].count()
```

```
Pregnancies    0
Glucose         0
BloodPressure   0
SkinThickness   0
Insulin         0
BMI             0
DiabetesPedigreeFunction 0
Age             0
Outcome         0
dtype: int64
```

```
1 # Creating a count (frequency) plot describing the data types and the count of variables:
2 df.describe().value_counts().plot(kind='bar')
```

<AxesSubplot>



WEEK 2

```
1 # DATA EXPLORATION:
2 # Checking the balance of the data by plotting the count of outcomes by their value. Describe your findings and plan
3 # future course of action:
```

```
1 df.columns
```

```
Out[28]: Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
        'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome'],
        dtype='object')
```

```
1 df.shape
```

```
Out[29]: (768, 9)
```

```
1 # It will tell me the number of counts in percentage distribution
2 df['Outcome'].value_counts(normalize=True)*100
```

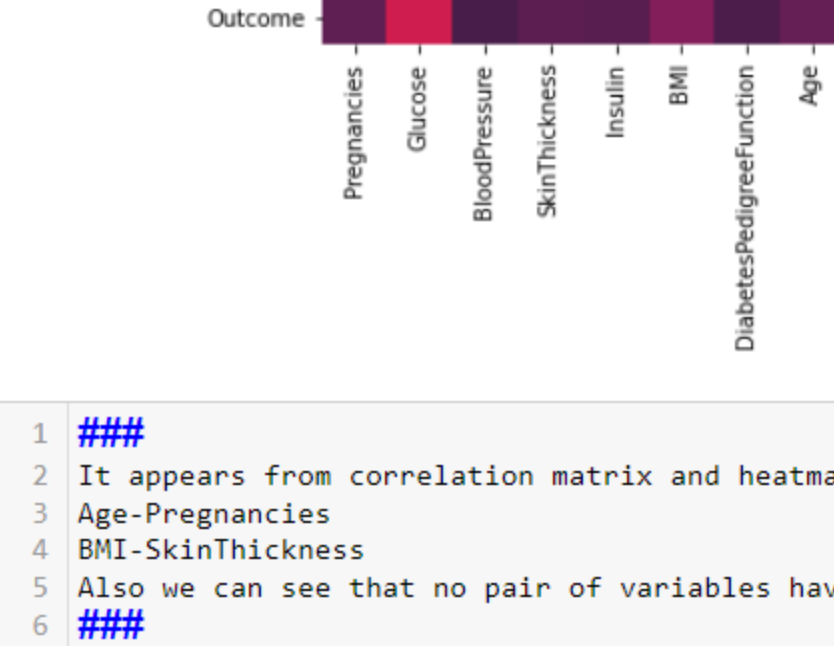
```
Out[30]: 0    65.104167
        1    34.895833
        Name: Outcome, dtype: float64
```

```
1 # It will tell me the number of counts without percentage directly showing no of rows having value 0 & 1.
2 df['Outcome'].value_counts()
```

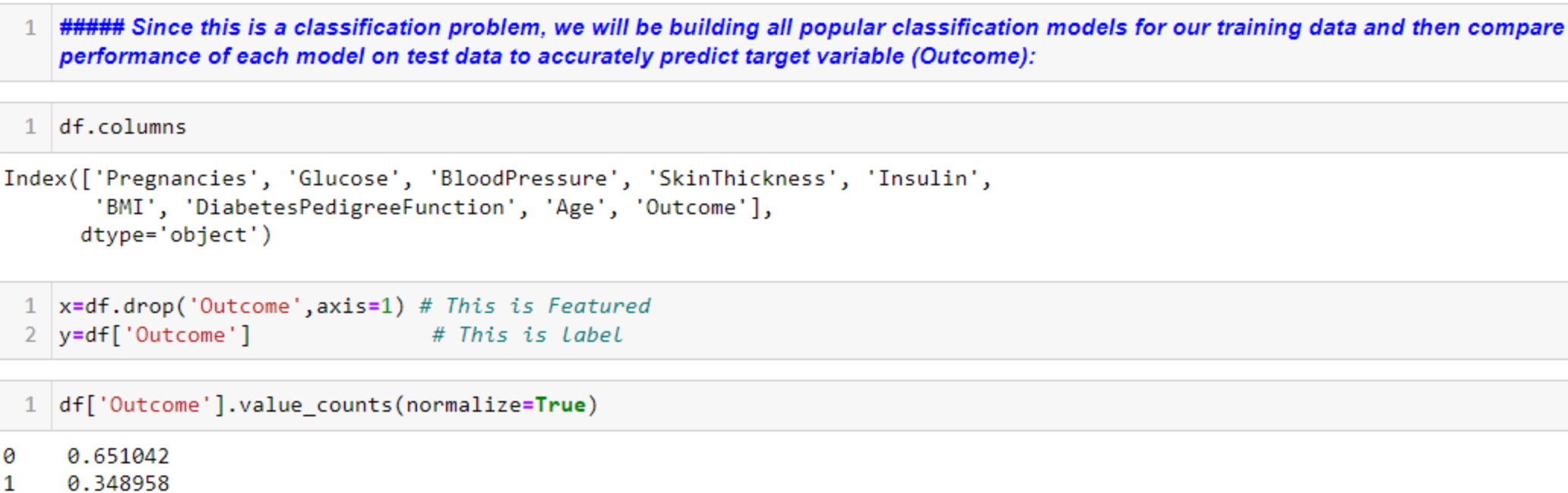
```
Out[31]: 0    500
        1    268
        Name: Outcome, dtype: int64
```

```
1 # This is just showing the above cell in bar plot
2 import warnings
3 warnings.simplefilter('ignore')
4 sns.countplot(df.Outcome)
```

```
Out[33]: <AxesSubplot: xlabel='Outcome', ylabel='count'>
```



```
1 # Creating scatter charts between the pair of variables to understand the relationships.
2 sns.pairplot(df)
```



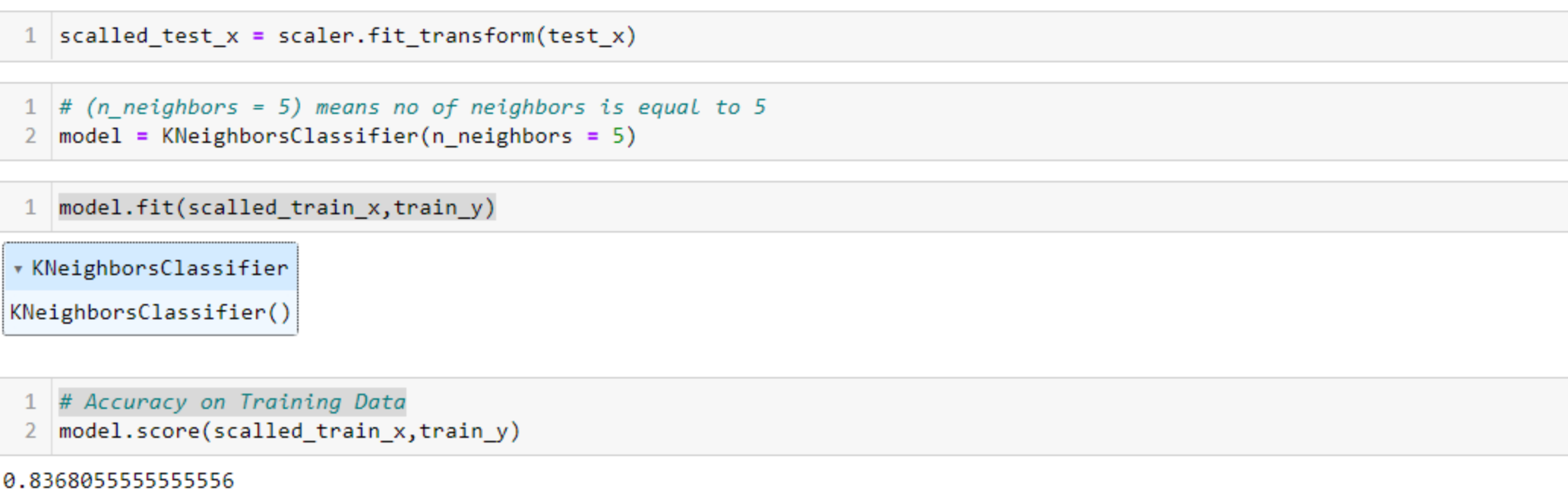
```
1 ##
2 We have some observations from above scatter plot of pairs of features:
3 (1) Glucose alone is impressively good to distinguish between the Outcome classes.
4 (2) Age alone is also able to distinguish between classes to some extent.
5 (3) It seems none of pairs in the dataset is able to clearly distinguish between the Outcome classes.
6 (4) With increase in skin thickness the BMI value also increases.
7 ##
```

```
1 # Performing correlation analysis. Visually exploring it using a heat map:
2 df.corr()
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
Pregnancies	1.000000	0.128213	0.208915	0.081770	0.025047	0.021569	-0.033523	0.544341	0.221898
Glucose	0.128213	1.000000	0.218937	0.192615	0.419451	0.231049	0.137327	0.266909	0.492782
BloodPressure	0.208915	0.218937	1.000000	0.191892	0.045363	0.281267	-0.002378	0.324915	0.165723
SkinThickness	0.081770	0.192615	0.191892	1.000000	0.155610	0.543205	0.102188	0.126107	0.214873
Insulin	0.025047	0.419451	0.045363	0.155610	1.000000	0.180241	0.126503	0.097101	0.203790
BMI	0.021569	0.231049	0.281267	0.543205	0.180241	1.000000	0.153438	0.025597	0.312038
DiabetesPedigreeFunction	-0.033523	0.137327	-0.002378	0.102188	0.126503	0.153438	1.000000	0.033561	0.173844
Age	0.544341	0.266909	0.324915	0.126107	0.097101	0.025597	0.033561	1.000000	0.238356
Outcome	0.221898	0.492782	0.165723	0.214873	0.203790	0.312038	0.173844	0.238356	1.000000

```
1 # In heat map Dark color shows good corr and Light color shows bad corr.
2 sns.heatmap(df.corr())
```

<AxesSubplot>



```
1 ##
2 It appears from correlation matrix and heatmap that there exists significant correlation between some pairs such as -
3 Age-Pregnancies
4 BMI-SkinThickness
5 Also we can see that no pair of variables have negative correlation
6 ##
```

Week 3:

```
1 ##### Since this is a classification problem, we will be building all popular classification models for our training data and then compare
        performance of each model on test data to accurately predict target variable (Outcome):
```

```
1 df.columns
```

```
Out[37]: Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
        'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome'],
        dtype='object')
```

```
1 x=df.drop('Outcome',axis=1) # This is Feature
2 y=df['Outcome'] # This is Label
```

```
Out[41]: 0    0.651042
        1    0.348958
        Name: Outcome, dtype: float64
```

```
1 # Here Stratify do the division b/w train & test so that the balance is maintained of outcome of above line.
2 # In stratify y means the target variable or the label.
3 train_x, test_x, train_y, test_y=train_test_split(x,y,test_size=.25,random_state=42,stratify=y)
```

```
Out[43]: 1 train_x.shape
```

```
Out[44]: (576, 8)
```

```
1 # Here the split or ratio is now same compare to the overall ratio
2 train_y.value_counts(normalize=True)
```

```
Out[45]: 0    0.651042
        1    0.348958
        Name: Outcome, dtype: float64
```

```
1 test_y.value_counts(normalize=True)
```

```
Out[46]: 0    0.651042
        1    0.348958
        Name: Outcome, dtype: float64
```

KNN ALGORITHM

```
1 from sklearn.preprocessing import MinMaxScaler
2 from sklearn.preprocessing import StandardScaler
```

```
1 scaler=StandardScaler()
```

```
1 scaled_train_x = scaler.fit_transform(train_x)
```

```
1 scaled_test_x = scaler.fit_transform(test_x)
```

```
1 # (n_neighbors = 5) means no of neighbors is equal to 5
2 model = KNeighborsClassifier(n_neighbors = 5)
```

```
1 model.fit(scaled_train_x,train_y)
```

```
Out[55]: KNeighborsClassifier()
```

```
1 # Accuracy on Training Data
2 model.score(scaled_train_x,train_y)
```

```
Out[56]: 0.8368055555555556
```

```
1 # Accuracy on Testing Data
2 model.score(scaled_test_x,test_y)
```

```
Out[57]: 0.75
```

```
1 from sklearn.metrics import classification_report,confusion_matrix
```

```
1 y_pred=model.predict(scaled_test_x)
```

```
1 # Checking the counts in prediction
2 pd.DataFrame(y_pred).value_counts()
```

```
Out[61]: 0    137
        1    55
        dtype: int64
```

```
1 # In Healthcare domain Recall is the widely used metric
2 # Sensitivity- Recall of class 1 = out of actual 1 how many we predict correctly
3 # Specificity- Recall of class 0
4 print(classification_report(test_y,y_pred))
```

	precision	recall	f1-score	support
0	0	0.78	0.86	0.82
1	0.67	0.55	0.61	67
accuracy	0.73	0.70	0.75	192
macro avg	0.71	0.65	0.66	192
weighted avg	0.74	0.75	0.74	192

```
1 print(confusion_matrix(test_y,y_pred))
```

```
Out[66]: [[167 16]
        [30 37]]
```

```
1 acc=[]
2 ran=range(2,15)
3 for k in ran:
4     model=KNeighborsClassifier(n_neighbors=k)
5     model.fit(scaled_train_x,train_y)
6     test_acc=model.score(scaled_test_x,test_y)
7     acc.append(test_acc*100)
```

```
1 sns.lineplot(ran,acc)
```

<AxesSubplot>


```
1 # From above line graph we come to know that we are getting the best accuracy at k=8
```

RANDOM FOREST ALGORITHM

```
1 from sklearn.ensemble import RandomForestClassifier
2 RF_model = RandomForestClassifier(n_estimators=500,max_depth=3)
```

```
1 RF_model.fit(scaled_train_x,train_y)
```

```
Out[70]: RandomForestClassifier(
    RandomForestClassifier(max_depth=3, n_estimators=500)
```

```
1 RF_model.score(scaled_train_x,train_y)
```

```
Out[71]: 0.8020833333333334
```

```
1 RF_model.score(scaled_test_x,test_y)
```

```
Out[72]: 0.7239583333333334
```

```
1 rf_y_pred=RF_model.predict(scaled_test_x)
```

```
Out[74]: 0    152
        1    40
        dtype: int64
```

```
1 # In Healthcare domain Recall is the widely used metric
2 # Sensitivity- Recall of class 1 = out of actual 1 how many we predict correctly
3 # Specificity- Recall of class 0
4 print(classification_report(test_y,rf_y_pred))
```

	precision	recall	f1-score	support
0	0	0.74	0.90	0.81
1	0.68	0.40	0.50	67
accuracy	0.73	0.70	0.72	192
macro avg	0.71	0.65	0.66	192
weighted avg	0.72	0.72	0.74	192

```
1 ##### Among all models, RandomForest has given best accuracy and f1_score. Therefore we will build final model using RandomForest.
```