For columns:

```
Glucose,
BloodPressure,
BMI,
DiabetesPedigreeFunction
```

If the column value is 0, then they should be considered as missing data.

1. Firstly, replace all Missing values with relevant figures.

```
import numpy as np
import pandas as pd
df = pd.read csv('Dataset Day7.csv')
print(df.info())
missing value percent = df.isna().sum() / len(df)
* 100
print(missing value percent)
skewness = df.skew()
print(skewness)
df["Glucose"].fillna(df["Glucose"].median(),
inplace=True)
df["BloodPressure"].fillna(df["BloodPressure"].me
an(), inplace=True)
df["BMI"].fillna(df["BMI"].median(),
inplace=True)
df["Outcome"].fillna(df["Outcome"].mean(),
inplace=True)
print(df.info())
```

```
1 C:\Users\tejas\PycharmProjects\pythonProject\venv\Scripts\python.exe C:\Users\
   tejas\PycharmProjects\pythonProject\START\Day8Q1.py
 2 <class 'pandas.core.frame.DataFrame'>
 3 RangeIndex: 768 entries, 0 to 767
 4 Data columns (total 7 columns):
       Column
                                Non-Null Count Dtype
 6 ---
       -----
                                -----
 7 0 Pregnancies
                                768 non-null int64
       Glucose
                               768 non-null int64
       BloodPressure
                               768 non-null int64
10 3
                               768 non-null float64
11 4
       DiabetesPedigreeFunction 768 non-null float64
12 5 Age
                               768 non-null int64
13 6
       Outcome
                                768 non-null int64
14 dtypes: float64(2), int64(5)
15 memory usage: 42.1 KB
16 None
17 Pregnancies
                             0.0
                             0.0
18 Glucose
                             0.0
19 BloodPressure
                            0.0
21 DiabetesPedigreeFunction 0.0
22 Age
                            0.0
23 Outcome
                             0.0
24 dtype: float64
25 Pregnancies
                            0.901674
26 Glucose
                            0.173754
27 BloodPressure
                           -1.843608
28 BMI
                           -0.428982
29 DiabetesPedigreeFunction 1.919911
30 Age
                            1.129597
31 Outcome
                            0.635017
32 dtype: float64
33 <class 'pandas.core.frame.DataFrame'>
34 RangeIndex: 768 entries, 0 to 767
35 Data columns (total 7 columns):
36 # Column
                               Non-Null Count Dtype
37 ---
38 0 Pregnancies
                               768 non-null
                                              int64
                               768 non-null int64
39 1 Glucose
40 2 BloodPressure
                               768 non-null int64
41 3
                                768 non-null float64
42 4
       DiabetesPedigreeFunction 768 non-null float64
                                768 non-null int64
43 5
       Age
                                768 non-null int64
44 6
       Outcome
45 dtypes: float64(2), int64(5)
46 memory usage: 42.1 KB
47 None
48
49 Process finished with exit code 0
```

2. Then remove all existing outliers and get the final data for classification.

```
import numpy as np
import pandas as pd
df = pd.read_csv('Dataset_Day7.csv')
print(df.info())
missing value percent = df.isna().sum() / len(df)
```

```
* 100
print(missing value percent)
skewness = df.skew()
print(skewness)
df["Glucose"].fillna(df["Glucose"].median(),
inplace=True)
df["BloodPressure"].fillna(df["BloodPressure"].me
an(), inplace=True)
df["BMI"].fillna(df["BMI"].median(),
inplace=True)
df["Outcome"].fillna(df["Outcome"].mean(),
inplace=True)
print(df.info())
OutlierData = pd.DataFrame()
temp = df[["Pregnancies", "Glucose",
"BloodPressure", "BMI", "DiabetesPedigreeFunction"
, "Age", "Outcome" ]]
for col in ["Pregnancies", "Glucose",
, "Age", "Outcome"]:
    Q1 = temp[col].quantile(0.25) # Gives 25th
Percentile or Q1
    Q3 = temp[col].quantile(0.75) # Gives 75th
Percentile or Q3
    IQR = Q3 - Q1
    UpperBound = Q3 + 1.5 * IQR
    LowerBound = Q1 - 1.5 * IQR
    OutlierData[col] = temp[col][(temp[col] <</pre>
LowerBound) | (temp[col] > UpperBound)]
    print(len(OutlierData))
    df OutlierFree = df.drop(OutlierData.index,
axis=0)
    df OutlierFree.info()
```

```
tejas\PycharmProjects\pythonProject\START\Day8Q2.py
 2 <class 'pandas.core.frame.DataFrame'>
 3 RangeIndex: 768 entries, 0 to 767
 4 Data columns (total 7 columns):
 5 #
       Column
                                Non-Null Count Dtype
 6 ---
       -----
                                 -----
 7
  Θ
      Pregnancies
                                768 non-null
                                                int64
 8
                                              int64
  1
       Glucose
                                768 non-null
 9
   2
       BloodPressure
                                768 non-null
                                               int64
10
                                 768 non-null
                                                float64
11 4
       DiabetesPedigreeFunction
                                768 non-null
                                                float64
12 5
                                 768 non-null
       Age
                                                int64
13 6
       Outcome
                                 768 non-null
                                                int64
14 dtypes: float64(2), int64(5)
15 memory usage: 42.1 KB
16 None
17 Pregnancies
                             0.0
18 Glucose
                             0.0
19 BloodPressure
                             0.0
20 BMI
                             0.0
21 DiabetesPedigreeFunction
                             0.0
22 Age
                             0.0
23 Outcome
                             0.0
24 dtype: float64
25 Pregnancies
                             0.901674
26 Glucose
                             0.173754
27 BloodPressure
                            -1.843608
28 BMI
                            -0.428982
29 DiabetesPedigreeFunction
                             1.919911
30 Age
                             1.129597
31 Outcome
                             0.635017
32 dtype: float64
33 <class 'pandas.core.frame.DataFrame'>
34 RangeIndex: 768 entries, 0 to 767
35 Data columns (total 7 columns):
36 # Column
                                Non-Null Count Dtype
37 ---
38 0 Pregnancies
                                768 non-null int64
39 1
       Glucose
                                768 non-null int64
40 2
       BloodPressure
                                768 non-null int64
41 3
                                768 non-null
                                                float64
42 4
       DiabetesPedigreeFunction
                                768 non-null
                                                float64
43 5
       Age
                                 768 non-null
                                                int64
44 6 Outcome
                                 768 non-null
                                                int64
45 dtypes: float64(2), int64(5)
46 memory usage: 42.1 KB
47 None
48 4
49 <class 'pandas.core.frame.DataFrame'>
50 Index: 764 entries, 0 to 767
51 Data columns (total 7 columns):
52 #
       Column
                                 Non-Null Count Dtype
53 ---
       -----
                                 -----
54 0
                                764 non-null
       Pregnancies
                                                int64
55 1
                                764 non-null
                                               int64
       Glucose
56 2
                                764 non-null
       BloodPressure
                                               int64
57 3
                                764 non-null
                                              float64
58 4
       DiabetesPedigreeFunction 764 non-null
                                               float64
```

File - Day8Q2

```
59 5 Age
                                  764 non-null
                                                  int64
                                  764 non-null
60 6
       Outcome
                                                 int64
61 dtypes: float64(2), int64(5)
62 memory usage: 47.8 KB
63 4
64 <class 'pandas.core.frame.DataFrame'>
65 Index: 764 entries, 8 to 767
66 Data columns (total 7 columns):
67 # Column
                                  Non-Null Count Dtype
68 ---
69 0
       Pregnancies
                                  764 non-null
                                                 int64
70 1
        Glucose
                                  764 non-null
                                                 int64
        BloodPressure
                                  764 non-null
                                                 int64
72 3
        BMI
                                  764 non-null
                                                 float64
        DiabetesPedigreeFunction
                                 764 non-null
                                                 float64
73 4
74 5
                                  764 non-null
        Age
                                                 int64
75 6
        Outcome
                                  764 non-null
                                                 int64
76 dtypes: float64(2), int64(5)
 77 memory usage: 47.8 KB
78 4
79 <class 'pandas.core.frame.DataFrame'>
80 Index: 764 entries, θ to 767
81 Data columns (total 7 columns):
82 # Column
                                  Non-Null Count Dtype
83 ---
84 0
        Pregnancies
                                  764 non-null
                                                 int64
85 1
        Glucose
                                  764 non-null
                                                 int64
86 2
        BloodPressure
                                  764 non-null
                                                 int64
87
        BMT
                                  764 non-null
                                                  float64
       DiabetesPedigreeFunction
                                  764 non-null
                                                 float64
89
   5
                                  764 non-null
                                                 int64
        Age
90 6
       Outcome
                                  764 non-null
                                                 int64
91 dtypes: float64(2), int64(5)
92 memory usage: 47.8 KB
94 <class 'pandas.core.frame.DataFrame'>
95 Index: 764 entries, 8 to 767
96 Data columns (total 7 columns):
97 # Column
                                  Non-Null Count Dtype
98 ---
99 0
        Pregnancies
                                  764 non-null
                                                  int64
100 1
        Glucose
                                  764 non-null
                                                 int64
101 2
        BloodPressure
                                  764 non-null
                                                 int64
                                                 float64
102 3
        BMT
                                  764 non-null
103 4
        DiabetesPedigreeFunction
                                 764 non-null
                                                 float64
104 5
                                  764 non-null
                                                 int64
        Outcome
                                  764 non-null
                                                 int64
106 dtypes: float64(2), int64(5)
107 memory usage: 47.8 KB
108 4
109 <class 'pandas.core.frame.DataFrame'>
110 Index: 764 entries, 0 to 767
111 Data columns (total 7 columns):
                                  Non-Null Count Dtype
112 # Column
113 ---
114 0
                                  764 non-null
        Pregnancies
                                                 int64
115 1
        Glucose
                                  764 non-null
                                                 int64
116 2
        BloodPressure
                                  764 non-null
                                                  int64
117 3
                                  764 non-null
                                                 float64
```

```
118 4
        DiabetesPedigreeFunction /64 non-null
                                                TL08T64
119 5
                                 764 non-null
                                                int64
        Age
120 6
        Outcome
                                 764 non-null
                                                int64
121 dtypes: float64(2), int64(5)
122 memory usage: 47.8 KB
123 4
124 <class 'pandas.core.frame.DataFrame'>
125 Index: 764 entries, 0 to 767
126 Data columns (total 7 columns):
127 #
       Column
                                 Non-Null Count Dtype
128 ---
       -----
                                 -----
129 0 Pregnancies
                                764 non-null int64
130 1 Glucose
                                764 non-null int64
131 2
       BloodPressure
                                764 non-null
                                               int64
132 3
                                764 non-null float64
133 4
       DiabetesPedigreeFunction 764 non-null float64
134 5
       Age
                                764 non-null
                                               int64
135 6 Outcome
                                 764 non-null
                                                int64
136 dtypes: float64(2), int64(5)
137 memory usage: 47.8 KB
138 4
139 <class 'pandas.core.frame.DataFrame'>
140 Index: 764 entries, 0 to 767
141 Data columns (total 7 columns):
142 # Column
                                Non-Null Count Dtype
143 ---
144 0
       Pregnancies
                                764 non-null
                                                int64
145 1
       Glucose
                                764 non-null
                                                int64
146 2
       BloodPressure
                                764 non-null
                                                int64
147 3
        BMI
                                 764 non-null
                                                float64
148 4
        DiabetesPedigreeFunction 764 non-null
                                                float64
149 5
        Age
                                 764 non-null
                                                int64
150 6
        Outcome
                                 764 non-null
                                                int64
151 dtypes: float64(2), int64(5)
152 memory usage: 47.8 KB
153
154 Process finished with exit code 0
155
```

- 3,Split the data into 70% training and 30% testing data. Then, use a k-Nearest Neighbor algorithm with target variable as 'Outcome'.
 - a. Print the default model performance metrics: Accuracy, Precision, Recall, F1Score
 - b. Plot a Precision & Recall vs k(no. of neighbours) curve (both Prec and Rec on the same graph). Find the k for which F1-score is the highest. **Use any one Distance Metric for this problem.**
 - c. Find the best distance metric, no. of neighbors combination for the kNN algorithm

```
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
```

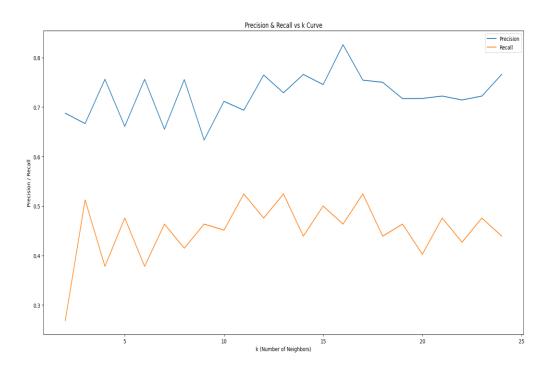
```
from sklearn.model selection import
train test split
from sklearn.neighbors import
KNeighborsClassifier
from sklearn.metrics import accuracy score,
precision score, recall score, f1 score
import math
df = pd.read csv('Dataset Day7.csv')
print(df.info())
missing value percent = df.isna().sum() / len(df)
print(missing value percent)
skewness = df.skew()
print(skewness)
df["Glucose"].fillna(df["Glucose"].median(),
inplace=True)
df["BloodPressure"].fillna(df["BloodPressure"].me
an(), inplace=True)
df["BMI"].fillna(df["BMI"].median(),
inplace=True)
df["Outcome"].fillna(df["Outcome"].mean(),
inplace=True)
print(df.info())
OutlierData = pd.DataFrame()
temp = df[["Pregnancies", "Glucose",
"DiabetesPedigreeFunction", "Age", "Outcome"]]
for col in ["Pregnancies", "Glucose",
"BloodPressure", "BMI",
"DiabetesPedigreeFunction", "Age", "Outcome"]:
   Q1 = temp[col].quantile(0.25) # Gives 25th
Percentile or Q1
   Q3 = temp[col].quantile(0.75) # Gives 75th
Percentile or Q3
    IQR = Q3 - Q1
    UpperBound = Q3 + 1.5 * IQR
    LowerBound = Q1 - 1.5 * IQR
    OutlierData[col] = temp[col][(temp[col] <</pre>
```

```
LowerBound) | (temp[col] > UpperBound)]
df OutlierFree = df.drop(OutlierData.index,
axis=0)
df OutlierFree.info()
X = df OutlierFree.drop('Outcome', axis=1) # all
columns except 'Outcome'
y = df OutlierFree['Outcome']  # target Variable
X train, X test, y train, y test =
train test split(X, y, test size=0.3,
unique classes = np.unique(y)
print('Number of unique classes:',
len(unique classes))
k start = \overline{i}nt(math.sqrt(len(X train)))
print(k start)
metric start = 'euclidean'
knn = KNeighborsClassifier(n neighbors=k start,
metric=metric start)
# fit the model
knn = knn.fit(X train, y train)
y pred = knn.predict(X test)
print("Model Performance metrics are as below:")
print("Accuracy: ", accuracy score(y test,
y pred))
print("Precision: ", precision score(y test,
y pred, zero division=1))
print("Recall: ", recall score(y test, y pred))
print("F1-Score: ", f1 score(y test, y pred))
k values = np.arange(2,25)
metric values = ['euclidean', 'manhattan',
'hamming']
prec = []
rec = []
acc = []
PerfData = pd.DataFrame(columns=['Nearest
Neighbor', 'Distance Metric', 'Precision',
```

```
'Recall', 'Accuracy'])
for dm in metric values:
    for k in k values:
        knn = KNeighborsClassifier(n neighbors=k,
metric=dm)
        knn = knn.fit(X train, y train)
        y pred = knn.predict(X test)
        row = [[k, dm, precision score(y test,
y pred, zero division=1), recall score(y test,
y pred),
                accuracy score(y test, y pred)]]
        df2 = pd.DataFrame(row, columns=['Nearest
Neighbor', 'Distance Metric', 'Precision',
'Recall', 'Accuracy'])
        PerfData = pd.concat([PerfData, df2],
ignore index=True)
print(PerfData.tail())
precision = PerfData['Precision']
recall = PerfData['Recall']
PerfData["F1 Score"] = (2 * PerfData["Precision"]
* PerfData["Recall"]) / (PerfData["Precision"] +
PerfData["Recall"])
print(PerfData[PerfData['F1 Score'] ==
max(PerfData['F1 Score'])])
plt.plot(k values, PerfData[PerfData["Distance
Metric"] == 'euclidean']['Precision'],
label='Precision')
plt.plot(k values, PerfData[PerfData["Distance
Metric"] == 'euclidean']['Recall'], label = 'Recall')
plt.xlabel('k (Number of Neighbors)')
plt.ylabel('Precision / Recall')
plt.title('Precision & Recall vs k Curve')
plt.legend(loc='upper right')
plt.show()
```

```
2 <class 'pandas.core.frame.DataFrame'>
 3 RangeIndex: 768 entries, 0 to 767
 4 Data columns (total 7 columns):
                                 Non-Null Count Dtype
 5 #
       Column
 6 ---
 7 Θ
       Pregnancies
                                 768 non-null
                                                 int64
                                 768 non-null
       Glucose
                                                 int64
 8
   1
       BloodPressure
                                 768 non-null
                                                 int64
10 3
                                 768 non-null
                                                 float64
11 4
       DiabetesPedigreeFunction 768 non-null
                                               float64
12 5
                                 768 non-null
       Age
                                                int64
13 6
       Outcome
                                 768 non-null
                                                 int64
14 dtypes: float64(2), int64(5)
15 memory usage: 42.1 KB
16 None
                              0.0
17 Pregnancies
18 Glucose
                              0.0
19 BloodPressure
                              0.0
20 BMI
                              0.0
21 DiabetesPedigreeFunction
                              0.0
22 Age
                              0.0
23 Outcome
                              0.0
24 dtype: float64
25 Pregnancies
                              0.901674
26 Glucose
                              0.173754
27 BloodPressure
                             -1.843608
28 BMI
                             -0.428982
29 DiabetesPedigreeFunction
                              1.919911
30 Age
                              1.129597
31 Outcome
                              0.635017
32 dtype: float64
33 <class 'pandas.core.frame.DataFrame'>
34 RangeIndex: 768 entries, 0 to 767
35 Data columns (total 7 columns):
                                 Non-Null Count Dtype
36 # Column
37 ---
38 0
       Pregnancies
                                 768 non-null
                                                 int64
39 1
       Glucose
                                 768 non-null
                                                 int64
       BloodPressure
40 2
                                 768 non-null
                                                int64
41 3
                                 768 non-null
                                                float64
42 4 DiabetesPedigreeFunction 768 non-null
43 5 Age
                                 768 non-null
                                                 int64
44 6 Outcome
                                 768 non-null
                                                 int64
45 dtypes: float64(2), int64(5)
46 memory usage: 42.1 KB
47 None
48 <class 'pandas.core.frame.DataFrame'>
49 Index: 764 entries, 0 to 767
50 Data columns (total 7 columns):
51 # Column
                                 Non-Null Count Dtype
52 ---
53 0
       Pregnancies
                                 764 non-null
                                                int64
54
       Glucose
                                 764 non-null
   1
55 2
       BloodPressure
                                 764 non-null
                                                 int64
                                 764 non-null
                                                 float64
56 3
57 4
       DiabetesPedigreeFunction 764 non-null
                                                 float64
58 5
       Age
                                 764 non-null
                                                 int64
```

```
59 6 Outcome
                                  764 non-null
                                                  int64
60 dtypes: float64(2), int64(5)
61 memory usage: 47.8 KB
62 Number of unique classes: 2
63 23
64 Model Performance metrics are as below:
65 Accuracy: 0.7478260869565218
66 Precision: 0.72222222222222
67 Recall: 0.47560975609756095
68 F1-Score: 0.5735294117647057
69
     Nearest Neighbor Distance Metric Precision Recall Accuracy
70 64
                    20
                              hamming
                                            1.0
                                                      0.0 0.643478
71 65
                                                      0.0 0.643478
                    21
                               hamming
                                              1.0
72 66
                    22
                                                      0.0 0.643478
                               hamming
                                             1.0
73 67
                    23
                               hamming
                                             1.0
                                                      0.0 0.643478
74 68
                    24
                               hamming
                                              1.0
                                                      0.0 0.643478
75 Nearest Neighbor Distance Metric Precision
                                                   Recall Accuracy F1 Score
76 15
                             euclidean 0.754386 0.52439 0.769565 0.618705
                    17
77
78 Process finished with exit code 0
79
```



← → | + Q ∓ | B