

# Assignment 4

KNN algorithm on diabetes dataset

```
In [1]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline
import warnings
warnings.filterwarnings('ignore')
from sklearn.model_selection import train_test_split
from sklearn.svm import SVC
from sklearn import metrics
```

```
In [2]: df=pd.read_csv('diabetes.csv')
```

```
In [3]: df.columns
```

```
Out[3]: Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
              'BMI', 'Pedigree', 'Age', 'Outcome'],
              dtype='object')
```

Check for null values. If present remove null values from the dataset

```
In [4]: df.isnull().sum()
```

```
Out[4]: Pregnancies      0
Glucose      0
BloodPressure  0
SkinThickness  0
Insulin      0
BMI          0
Pedigree     0
Age         0
Outcome     0
dtype: int64
```

```
In [ ]:
```

Outcome is the label/target, other columns are features

```
In [7]: X = df.drop('Outcome',axis = 1)
y = df['Outcome']
```

```
In [8]: from sklearn.preprocessing import scale
X = scale(X)
# split into train and test
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random
```

```
In [9]: from sklearn.neighbors import KNeighborsClassifier
knn = KNeighborsClassifier(n_neighbors=7)

knn.fit(X_train, y_train)
y_pred = knn.predict(X_test)
```

```
In [17]: print("Confusion matrix: ")
cs = metrics.confusion_matrix(y_test,y_pred)
print(cs)
```

```
Confusion matrix:
[[123  28]
 [ 37  43]]
```

```
In [12]: print("Accuracy ",metrics.accuracy_score(y_test,y_pred))
```

```
Accuracy  0.7186147186147186
```

Classification error rate: proportion of instances misclassified over the whole set of instances. Error rate is calculated as the total number of two incorrect predictions (FN + FP) divided by the total number of a dataset (examples in the dataset).

Also  $\text{error\_rate} = 1 - \text{accuracy}$

```
In [29]: total_misclassified = cs[0,1] + cs[1,0]
print(total_misclassified)
total_examples = cs[0,0]+cs[0,1]+cs[1,0]+cs[1,1]
print(total_examples)
print("Error rate",total_misclassified/total_examples)
print("Error rate ",1-metrics.accuracy_score(y_test,y_pred))
```

```
65
231
Error rate 0.2813852813852814
Error rate  0.2813852813852814
```

```
In [13]: print("Precision score",metrics.precision_score(y_test,y_pred))
```

```
Precision score 0.6056338028169014
```

```
In [14]: print("Recall score ",metrics.recall_score(y_test,y_pred))
```

```
Recall score  0.5375
```

```
In [15]: print("Classification report ",metrics.classification_report(y_test,y_pred))
```

| Classification report |   |      | precision | recall | f1-score | support |
|-----------------------|---|------|-----------|--------|----------|---------|
|                       | 0 | 0.77 | 0.81      | 0.79   |          | 151     |
|                       | 1 | 0.61 | 0.54      | 0.57   |          | 80      |
| accuracy              |   |      |           | 0.72   |          | 231     |
| macro avg             |   | 0.69 | 0.68      | 0.68   |          | 231     |
| weighted avg          |   | 0.71 | 0.72      | 0.71   |          | 231     |