CS353 ML Lab MidSem

Name: K V Sumanth Reddy

Roll No: 181CO225

Batch: Section-2

Date: 09/03/2021

Q: Build a K-nearest neighbor algorithm to predict whether a patient is having cancer (Malignant tumor) or not (Benign tumor). Use Kaggle dataset from UCI Machine Learning Repository.

Dataset Used: Breast Cancer Wisconsin (Diagnostic) Data Set (https://www.kaggle.com/uciml/breast-cancer-wisconsin-data)

Importing Libraries and Dataset

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

from sklearn.neighbors import KNeighborsClassifier

from sklearn.metrics import confusion_matrix,explained_variance_sconfrom sklearn.metrics import mean_squared_error,r2_score
from sklearn.metrics import accuracy_score, classification_report

dataset = pd.read_csv('data.csv')

#printing 5 sample tuples
dataset.sample(5)
```

		id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean
	312	89382602	В	12.76	13.37	82.29	504.1
	200	877501	В	12.23	19.56	78.54	461.0
<pre>dataset.info()</pre>							

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 569 entries, 0 to 568
Data columns (total 33 columns):
    Column
                              Non-Null Count
                                              Dtype
     _ _ _ _ _
0
    id
                              569 non-null
                                               int64
1
    diagnosis
                              569 non-null
                                              object
2
                              569 non-null
                                               float64
    radius mean
3
    texture mean
                              569 non-null
                                               float64
                                              float64
4
                              569 non-null
    perimeter mean
5
    area mean
                              569 non-null
                                              float64
6
                              569 non-null
                                               float64
    smoothness mean
7
                              569 non-null
                                               float64
    compactness mean
8
    concavity mean
                              569 non-null
                                               float64
9
                              569 non-null
                                               float64
    concave points mean
10
    symmetry mean
                              569 non-null
                                               float64
11
    fractal dimension mean
                              569 non-null
                                               float64
12
                              569 non-null
                                               float64
    radius se
13
                              569 non-null
                                               float64
    texture se
14
                              569 non-null
                                               float64
    perimeter se
15
    area se
                              569 non-null
                                               float64
16
    smoothness se
                              569 non-null
                                               float64
17
                              569 non-null
                                               float64
    compactness se
18
                              569 non-null
                                               float64
    concavity se
                                               float64
19
    concave points se
                              569 non-null
                              569 non-null
                                               float64
20
    symmetry se
21
    fractal dimension se
                              569 non-null
                                               float64
22
                                               float64
    radius worst
                              569 non-null
23
    texture worst
                              569 non-null
                                               float64
24
                                               float64
    perimeter worst
                              569 non-null
25
    area worst
                              569 non-null
                                               float64
26
                                               float64
    smoothness worst
                              569 non-null
27
    compactness worst
                              569 non-null
                                               float64
28
    concavity worst
                              569 non-null
                                               float64
29
    concave points_worst
                              569 non-null
                                               float64
30
    symmetry_worst
                              569 non-null
                                               float64
    fractal_dimension_worst
                              569 non-null
                                               float64
31
    Unnamed: 32
                              0 non-null
                                               float64
dtypes: float64(31), int64(1), object(1)
memory usage: 146.8+ KB
```

Data Preprocessing

```
# Drop unused columns and encoding M to 0 and B to 1
columns = ['Unnamed: 32', 'id', 'diagnosis']
data = {'M': 0, 'B': 1}
```

```
Midsem - Colaboratory
y = uataset[ uiagnosis ].map(uata)
x = dataset.drop(columns, axis=1)
from sklearn.model selection import train test split
x train, x test, y train, y test = train test split(x, y, test size
print('Training dataset size:\nx train -', len(x train), '\ny train
print('Testing dataset size:\nx_test -', len(x_test), '\ny_test -',
    Training dataset size:
    x train - 455
    y_train - 455
    Testing dataset size:
    x_test - 114
    y_test - 114
KNN = []
for i in range(1, 8):
  KNNModel = KNeighborsClassifier(n neighbors = i, metric = 'euclide
  KNN.append(KNNModel)
for i in range(7):
  KNN[i].fit(x_train, y_train)
```

Finding Accuracies for all models using the test dataset

```
train accuracyKNN = [] #store training accuracies
test_accuracyKNN = [] #store testing accuracies
for i in range(7):
  print("-----")
  print('K = {}'.format(i + 1))
  train accuracyKNN.append(KNN[i].score(x train, y train))
  test accuracyKNN.append(KNN[i].score(x test, y test))
  y pred = KNN[i].predict(x test)
  print("Confusion Matrix:\n",confusion matrix(y test, y pred))
 print(classification_report(y_test,y_pred))
  print("\nAccuracy: %.2f" %(accuracy score(y test, y pred)*100))
  print("Mean Squared Error: %.2f" %(mean squared error(y test, y p
  print("Explained Variance: %.2f" %(explained variance score(y test
  print("R2 Score: %.2f" %(r2 score(y test, y pred)*100))
  print("\n\n")
  print("-----
                                  0.94
                         0 9 3 A
                  0 d3
                                  0 d3
      macro avo
```

weighted avg 0.94 0.94 0.94 114

Accuracy: 93.86

Mean Squared Error: 6.14 Explained Variance: 72.40

R2 Score: 72.37

K = 6

Confusion Matrix:

[[35 3] [4 72]]

	precision	recall	f1-score	support
0 1	0.90 0.96	0.92 0.95	0.91 0.95	38 76
accuracy macro avg weighted avg	0.93 0.94	0.93 0.94	0.94 0.93 0.94	114 114 114

Accuracy: 93.86

Mean Squared Error: 6.14 Explained Variance: 72.40

R2 Score: 72.37

K = 7

Confusion Matrix:

[[35 3] [2 74]]

[- / .]]	precision	recall	f1-score	support
0 1	0.95 0.96	0.92 0.97	0.93 0.97	38 76
accuracy macro avg weighted avg	0.95 0.96	0.95 0.96	0.96 0.95 0.96	114 114 114

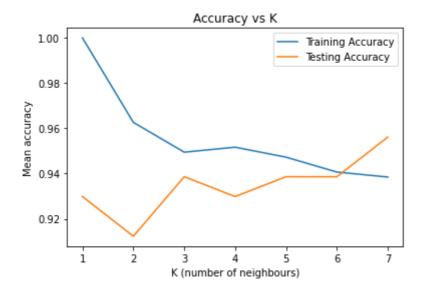
Accuracy: 95.61

Mean Squared Error: 4.39 Explained Variance: 80.30

R2 Score: 80.26

Best accuracy is at neighbours = 7 and the best mean difference between train and test accuracies is at neighbours = 3,5,6

```
fig = plt.figure()
ax = plt.axes()
plt.plot([1,2,3,4,5,6,7], train_accuracyKNN, label = 'Training Accuracy
plt.plot([1,2,3,4,5,6,7], test_accuracyKNN, label = 'Testing Accuracy
plt.xlabel('K (number of neighbours)')
plt.ylabel('Mean accuracy')
plt.title('Accuracy vs K')
plt.legend()
plt.savefig('Graph.png')
plt.show()
```



• ×