**kNN on Iris Dataset** ----Norina Akhtar

```
1 #importing the required libraries
2 import pandas as pd
3 import numpy as np
4 import operator
5 import matplotlib.pyplot as plt
```

```
1 # uplaod file from local drive
2 from google.colab import files
3 uploaded = files.upload()
```

Choose Files   No file chosen          Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.
Saving iris.csv to iris.csv

```
1 import io
2 data = pd.read_csv(io.BytesIO(uploaded['iris.csv']))
3 print(data)
```

```
     sepal.length  sepal.width  petal.length  petal.width    variety
0             5.1          3.5           1.4          0.2     Setosa
1             4.9          3.0           1.4          0.2     Setosa
2             4.7          3.2           1.3          0.2     Setosa
3             4.6          3.1           1.5          0.2     Setosa
4             5.0          3.6           1.4          0.2     Setosa
..            ...          ...           ...          ...        ...
145           6.7          3.0           5.2          2.3  Virginica
146           6.3          2.5           5.0          1.9  Virginica
147           6.5          3.0           5.2          2.0  Virginica
148           6.2          3.4           5.4          2.3  Virginica
149           5.9          3.0           5.1          1.8  Virginica

[150 rows x 5 columns]
```

*part a*

Dividing the dataset as development and test.

```
 1 #randomize the indices
 2 indices = np.random.permutation(data.shape[0])
 3 div = int(0.75 * len(indices))
 4 development_id, test_id = indices[:div], indices[div:]
 5 #dividing the dataset using randomized indices
 6 development_set, test_set = data.loc[development_id,:], data.loc[test_id,:]
 7 print("Development Set:\n", development_set, "\n\nTest Set:\n", test_set)
 8 mean_development_set = development_set.mean(numeric_only=True)
 9 mean_test_set = test_set.mean(numeric_only=True)
10 std_development_set = development_set.std(numeric_only=True)
11 std_test_set = test_set.std(numeric_only=True)
```

```
Development Set:
     sepal.length  sepal.width  petal.length  petal.width    variety
66            5.6          3.0           4.5          1.5  Versicolor
75            6.6          3.0           4.4          1.4  Versicolor
39            5.1          3.4           1.5          0.2     Setosa
34            4.9          3.1           1.5          0.2     Setosa
95            5.7          3.0           4.2          1.2  Versicolor
..            ...          ...           ...          ...        ...
83            6.0          2.7           5.1          1.6  Versicolor
9             4.9          3.1           1.5          0.1     Setosa
40            5.0          3.5           1.3          0.3     Setosa
60            5.0          2.0           3.5          1.0  Versicolor
7             5.0          3.4           1.5          0.2     Setosa

[112 rows x 5 columns]

Test Set:
     sepal.length  sepal.width  petal.length  petal.width    variety
29            4.7          3.2           1.6          0.2     Setosa
3             4.6          3.1           1.5          0.2     Setosa
91            6.1          3.0           4.6          1.4  Versicolor
109           7.2          3.6           6.1          2.5  Virginica
44            5.1          3.8           1.9          0.4     Setosa
137           6.4          3.1           5.5          1.8  Virginica
112           6.8          3.0           5.5          2.1  Virginica
```

```
      38          4.4          3.0          1.3          0.2      Setosa
      10          5.4          3.7          1.5          0.2      Setosa
      24          4.8          3.4          1.9          0.2      Setosa
      35          5.0          3.2          1.2          0.2      Setosa
      94          5.6          2.7          4.2          1.3  Versicolor
     100          6.3          3.3          6.0          2.5   Virginica
      86          6.7          3.1          4.7          1.5  Versicolor
       0          5.1          3.5          1.4          0.2      Setosa
      61          5.9          3.0          4.2          1.5  Versicolor
      88          5.6          3.0          4.1          1.3  Versicolor
     106          4.9          2.5          4.5          1.7   Virginica
     143          6.8          3.2          5.9          2.3   Virginica
      67          5.8          2.7          4.1          1.0  Versicolor
      56          6.3          3.3          4.7          1.6  Versicolor
      36          5.5          3.5          1.3          0.2      Setosa
     133          6.3          2.8          5.1          1.5   Virginica
     121          5.6          2.8          4.9          2.0   Virginica
      76          6.8          2.8          4.8          1.4  Versicolor
      53          5.5          2.3          4.0          1.3  Versicolor
      71          6.1          2.8          4.0          1.3  Versicolor
      89          5.5          2.5          4.0          1.3  Versicolor
      30          4.8          3.1          1.6          0.2      Setosa
     113          5.7          2.5          5.0          2.0   Virginica
      17          5.1          3.5          1.4          0.3      Setosa
      12          4.8          3.0          1.4          0.1      Setosa
      14          5.8          4.0          1.2          0.2      Setosa
     114          5.8          2.8          5.1          2.4   Virginica
      70          5.9          3.2          4.8          1.8  Versicolor
      42          4.4          3.2          1.3          0.2      Setosa
     136          6.3          3.4          5.6          2.4   Virginica
      31          5.4          3.4          1.5          0.4      Setosa
```

**Part b**

Implement kNN using the following hyperparameters:

number of neighbor

```
 * 1,3,5,7
```

distance metric

```
 * euclidean distance
 * normalized euclidean distance
 * cosine similarity
```

Retrieving the 'class' column from the development and test sets and storing it in separate lists. Calculating the mean and standard deviation of the development set and test set for normalizing the data.

```
1 test_class = list(test_set.iloc[:,-1])
2 dev_class = list(development_set.iloc[:,-1])
3 mean_development_set = development_set.mean(numeric_only=True)
4 mean_test_set = test_set.mean(numeric_only=True)
5 std_development_set = development_set.std(numeric_only=True)
6 std_test_set = test_set.std(numeric_only=True)
```

Functions for computing the Euclidean Distance, Normalized Euclidean Distance, Cosine Similarity and k Nearest Neighbor to determine the 'class' for a given input instance.

```
1 def euclideanDistance(data_1, data_2, data_len):
2     dist = 0
3     for i in range(data_len):
4         dist = dist + np.square(data_1[i] - data_2[i])
5     return np.sqrt(dist)
6
7 def normalizedEuclideanDistance(data_1, data_2, data_len, data_mean, data_std):
8     n_dist = 0
9     for i in range(data_len):
10        n_dist = n_dist + (np.square(((data_1[i] - data_mean[i])/data_std[i]) - ((data_2[i] - data_mean[i])/data_std[i])))
11    return np.sqrt(n_dist)
12
```

```
13 def cosineSimilarity(data_1, data_2):
14     dot = np.dot(data_1, data_2[:-1])
15     norm_data_1 = np.linalg.norm(data_1)
16     norm_data_2 = np.linalg.norm(data_2[:-1])
17     cos = dot / (norm_data_1 * norm_data_2)
18     return (1-cos)
19
20 def knn(dataset, testInstance, k, dist_method, dataset_mean, dataset_std):
21     distances = {}
22     length = testInstance.shape[1]
23     if dist_method == 'euclidean':
24         for x in range(len(dataset)):
25             dist_up = euclideanDistance(testInstance, dataset.iloc[x], length)
26             distances[x] = dist_up[0]
27     elif dist_method == 'normalized_euclidean':
28         for x in range(len(dataset)):
29             dist_up = normalizedEuclideanDistance(testInstance, dataset.iloc[x], length, dataset_mean, dataset_std)
30             distances[x] = dist_up[0]
31     elif dist_method == 'cosine':
32         for x in range(len(dataset)):
33             dist_up = cosineSimilarity(testInstance, dataset.iloc[x])
34             distances[x] = dist_up[0]
35 # Sort values based on distance
36     sort_distances = sorted(distances.items(), key=operator.itemgetter(1))
37     neighbors = []
38     # Extracting nearest k neighbors
39     for x in range(k):
40         neighbors.append(sort_distances[x][0])
41     # Initializing counts for 'class' labels counts as 0
42     counts = {"Iris-setosa" : 0, "Iris-versicolor" : 0, "Iris-virginica" : 0}
43     # Computing the most frequent class
44     for x in range(len(neighbors)):
45         response = dataset.iloc[neighbors[x]][-1]
46         if response in counts:
47             counts[response] += 1
48         else:
49             counts[response] = 1
50     # Sorting the class in reverse order to get the most frequest class
51     sort_counts = sorted(counts.items(), key=operator.itemgetter(1), reverse=True)
52     return(sort_counts[0][0])
```

**Part c**

Using the development data set

Iterating all of the development data points and computing the class for each k and each distance metric

```
1   print(development_set)
2
```

```
    sepal.length  sepal.width  petal.length  petal.width     variety
66           5.6          3.0           4.5          1.5  Versicolor
75           6.6          3.0           4.4          1.4  Versicolor
39           5.1          3.4           1.5          0.2      Setosa
34           4.9          3.1           1.5          0.2      Setosa
95           5.7          3.0           4.2          1.2  Versicolor
..           ...          ...           ...          ...         ...
83           6.0          2.7           5.1          1.6  Versicolor
9            4.9          3.1           1.5          0.1      Setosa
40           5.0          3.5           1.3          0.3      Setosa
60           5.0          2.0           3.5          1.0  Versicolor
7            5.0          3.4           1.5          0.2      Setosa

[112 rows x 5 columns]
```

```
1   # Creating a list of list of all columns except 'class' by iterating through the development set
2   row_list = []
3   for index, rows in development_set.iterrows():
4     # print()
5
6     # row_list.append(rows['sepal.length'])
7     my_list =[rows["sepal.length"], rows["sepal.width"], rows["petal.length"], rows["petal.width"]]
8     row_list.append([my_list])
9   # k values for the number of neighbors that need to be considered
10  k_n = [1, 3, 5, 7]
11  # Distance metrics
```

```
12    distance_methods = ['euclidean', 'normalized_euclidean', 'cosine']
13    # Performing kNN on the development set by iterating all of the development set data points and for each k and each distance
14    obs_k = {}
15    for dist_method in distance_methods:
16        development_set_obs_k = {}
17        for k in k_n:
18            development_set_obs = []
19            for i in range(len(row_list)):
20                development_set_obs.append(knn(development_set, pd.DataFrame(row_list[i]), k, dist_method, mean_development_set,
21            development_set_obs_k[k] = development_set_obs
22        # Nested Dictionary containing the observed class for each k and each distance metric (obs_k of the form obs_k[dist_methc
23        obs_k[dist_method] = development_set_obs_k
24        print(dist_method.upper() + " distance method performed on the dataset for all k values!")
```

```
EUCLIDEAN distance method performed on the dataset for all k values!
NORMALIZED_EUCLIDEAN distance method performed on the dataset for all k values!
COSINE distance method performed on the dataset for all k values!
```

Computing the accuracy for the development data set and finding the optimal hyperparametes

```
1  # Calculating the accuracy of the development set by comparing it with the development set 'class' list created earlier
2  accuracy = {}
3  for key in obs_k.keys():
4      accuracy[key] = {}
5      for k_value in obs_k[key].keys():
6          #print('k = ', key)
7          count = 0
8          for i,j in zip(dev_class, obs_k[key][k_value]):
9              if i == j:
10                 count = count + 1
11             else:
12                 pass
13         accuracy[key][k_value] = count/(len(dev_class))
14
15 # Storing the accuracy for each k and each distance metric into a dataframe
16 df_res = pd.DataFrame({'k': k_n})
17 for key in accuracy.keys():
18     value = list(accuracy[key].values())
19     df_res[key] = value
20 print(df_res)
21
22 # Plotting a Bar Chart for accuracy
23 draw = df_res.plot(x='k', y=['euclidean', 'normalized_euclidean', 'cosine'], kind="bar", colormap='YlGnBu')
24 draw.set(ylabel='Accuracy')
25
26 # Ignoring k=1 if the value of accuracy for k=1 is 100%, since this mostly implies overfitting
27 df_res.loc[df_res['k'] == 1.0, ['euclidean', 'normalized_euclidean', 'cosine']] = np.nan
28
29 # Fetching the best k value for using all hyper-parameters
30 # In case the accuracy is the same for different k and different distance metric selecting the first of all the same
31 column_val = [c for c in df_res.columns if not c.startswith('k')]
32 col_max = df_res[column_val].max().idxmax()
33 best_dist_method = col_max
34 row_max = df_res[col_max].argmax()
35 best_k = int(df_res.iloc[row_max]['k'])
36 if df_res.isnull().values.any():
37     print('\n\n\nBest k value is\033[1m', best_k, '\033[0mand best distance metric is\033[1m', best_dist_method, '\033[0m. Ig:
38 else:
39     print('\n\n\nBest k value is\033[1m', best_k, '\033[0mand best distance metric is\033[1m', best_dist_method, '\033[0m.')
```

```
   k  euclidean  normalized_euclidean   cosine
0  1   1.000000              1.000000  1.000000
1  3   0.973214              0.964286  0.973214
2  5   0.964286              0.964286  0.964286
3  7   0.973214              0.973214  0.973214
```

**Using the test dataset**

```
1 print('\n\n\nBest k value is\033[1m', best_k, '\033[0mand best distance metric is\033[1m', best_dist_method, '\033[0m')
```

```
Best k value is 3 and best distance metric is euclidean
```

Using the best k value and best distance metric to determine the class for all rows in the test dataset

```
 1 # Creating a list of list of all columns except 'class' by iterating through the development set
 2 row_list_test = []
 3 for index, rows in test_set.iterrows():
 4     my_list =[rows["sepal.length"], rows["sepal.width"], rows["petal.length"], rows["petal.width"]]
 5     row_list_test.append([my_list])
 6 test_set_obs = []
 7 for i in range(len(row_list_test)):
 8     test_set_obs.append(knn(test_set, pd.DataFrame(row_list_test[i]), best_k, best_dist_method, mean_test_set, std_test_set))
 9 #print(test_set_obs)
10
11 count = 0
12 for i,j in zip(test_class, test_set_obs):
13     if i == j:
14         count = count + 1
15     else:
16         pass
17 accuracy_test = count/(len(test_class))
18 print('Final Accuracy of the Test dataset is ', accuracy_test)
```

```
Final Accuracy of the Test dataset is  0.9473684210526315
```

**References**

https://hc.labnet.sfbu.edu/~henry/sfbu/course//machine_learning_with_r_cookbook/practice_ml_with_r/slide/iris.html

https://ai.plainenglish.io/understanding-confusion-matrix-and-applying-it-on-knn-classifier-on-iris-dataset-b57f85d05cd8

https://stackoverflow.com/questions/16476924/how-to-iterate-over-rows-in-a-dataframe-in-pandas

✓  1s    completed at 3:26 PM                                                                ● ✕