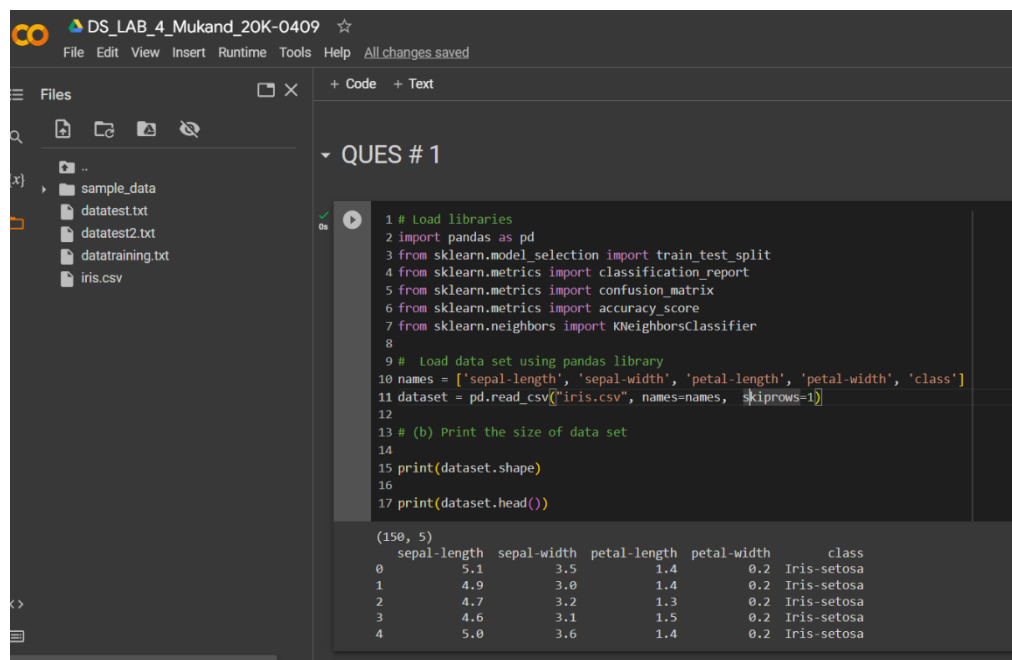


DATA SCIENCE LAB 4

Roll no: 20K-0409

Screen Shots

Ques # 1



```
1 # Load libraries
2 import pandas as pd
3 from sklearn.model_selection import train_test_split
4 from sklearn.metrics import classification_report
5 from sklearn.metrics import confusion_matrix
6 from sklearn.metrics import accuracy_score
7 from sklearn.neighbors import KNeighborsClassifier
8
9 # Load data set using pandas library
10 names = ['sepal-length', 'sepal-width', 'petal-length', 'petal-width', 'class']
11 dataset = pd.read_csv('iris.csv', names=names, skiprows=1)
12
13 # (b) Print the size of data set
14
15 print(dataset.shape)
16
17 print(dataset.head())
```

(150, 5)

	sepal-length	sepal-width	petal-length	petal-width	class
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3.0	1.4	0.2	Iris-setosa
2	4.7	3.2	1.3	0.2	Iris-setosa
3	4.6	3.1	1.5	0.2	Iris-setosa
4	5.0	3.6	1.4	0.2	Iris-setosa

▼ (c) Display the class distribution

```
1 print(dataset.groupby('class').size())
```

```
class
Iris-setosa      50
Iris-versicolor  50
Iris-virginica   50
dtype: int64
```

▼ (d) Now, divide your data using hold out approach (80% for training and 20% for testing) train / test dataset

```
[45] 1 array = dataset.values
     2 X = array[:,0:4]
     3 Y = array[:,4]
     4 t_size = 0.20
     5 seed = 7
     6 X_train_iris, X_test_iris, Y_train_iris, Y_test_iris = train_test_split(X, Y, test_size=t_size, random_state=seed)
     7
     8 print(X_test_iris)
```

```
[[5.9 3.0 5.1 1.8]
 [5.4 3.0 4.5 1.5]
 [5.0 3.5 1.3 0.2]]
```

make predictions

```
[46] 1 from sklearn.model_selection import train_test_split
2 from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
3 from sklearn.neighbors import KNeighborsClassifier
4
5 knn = KNeighborsClassifier()
6 knn.fit(X_train_iris, Y_train_iris)
7
8 predictions = knn.predict(X_test_iris)
9
10 print("Accuracy Score: \n")
11 print(accuracy_score(Y_test_iris, predictions))
12
13 print("\nConfusion Matrix: \n")
14 print(confusion_matrix(Y_test_iris, predictions))
15
16 print("\nClassification : \n")
17 print(classification_report(Y_test_iris, predictions))
```

Accuracy Score:

0.9

Confusion Matrix:

```
[[ 7  0  0]
 [ 0 11  1]
 [ 0  2  9]]
```

Classification :

	precision	recall	f1-score	support
Iris-setosa	1.00	1.00	1.00	7
Iris-versicolor	0.85	0.92	0.88	12
Iris-virginica	0.90	0.82	0.86	11

(f) Repeat (e) by changing the value of k (k=1, 2, 3,..., 10). Print only accuracy

```
[47] 1 print("Accuracy for different k values:")
2 for k in range(1, 11):
3     knn = KNeighborsClassifier(n_neighbors=k)
4     knn.fit(X_train_iris, Y_train_iris)
5     predictions = knn.predict(X_test_iris)
6     accuracy = accuracy_score(Y_test_iris, predictions)
7     print(f"k = {k}: Accuracy = {accuracy:.4f}")
8
```

Accuracy for different k values:

```
k = 1: Accuracy = 0.9000
k = 2: Accuracy = 0.9333
k = 3: Accuracy = 0.9000
k = 4: Accuracy = 0.9333
k = 5: Accuracy = 0.9000
k = 6: Accuracy = 0.8667
k = 7: Accuracy = 0.8667
k = 8: Accuracy = 0.9000
k = 9: Accuracy = 0.9000
k = 10: Accuracy = 0.9000
```

(g) Repeat (e) by changing the value of seed (seed = 1, 2, 3, ..., 10). Print only accuracy

```
[48] 1 print("\nAccuracy for different seed values:")
2 for seed_val in range(1, 11):
3     X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=t_size, random_state=seed_val)
4     knn = KNeighborsClassifier()
5     knn.fit(X_train, Y_train)
6     predictions = knn.predict(X_test)
7     accuracy = accuracy_score(Y_test, predictions)
8     print(f"Accuracy = {accuracy:.4f}")
9
```



Accuracy for different seed values:

```
Accuracy = 1.0000
Accuracy = 1.0000
Accuracy = 0.9667
Accuracy = 0.9667
Accuracy = 0.9333
Accuracy = 0.9667
Accuracy = 0.9000
Accuracy = 0.9000
Accuracy = 1.0000
Accuracy = 0.9667
```

Ques # 2

QUES # 2

```
[59] 1 # Load the training dataset
2 columns = ['date', 'temperature', 'humidity', 'light', 'co2', 'humidity_ratio', 'occupancy']
3 train_data = pd.read_csv("data\training.txt", names=columns, header=0)
4
5 # Dropping the 'date' attribute
6 train_data = train_data.drop('date', axis=1)
7
8 print(train_data.head())
9
```

```
temperature humidity light co2 humidity_ratio occupancy
1      23.18    27.2720  426.0  721.25      0.004793         1
2      23.15    27.2675  429.5  714.00      0.004783         1
3      23.15    27.2450  426.0  713.50      0.004779         1
4      23.15    27.2000  426.0  708.25      0.004772         1
5      23.10    27.2000  426.0  704.50      0.004757         1
```

```
[60] 1 # (b) size of the dataset
2 print(train_data.shape)
3
4 # (c) class distribution
5 print("\nClass Dist:")
6 print(train_data.groupby('occupancy').size())
7
```

(8143, 6)

Class Dist:
occupancy
0 6414
1 1729
dtype: int64

(e) Apply kNN classifier

```
[66] 1 knn = KNeighborsClassifier()
2 knn.fit(X_train0c, Y_train0c)
3
4 predictions = knn.predict(X_test0c)
5
6 print("Accuracy Score: \n")
7 print(accuracy_score(Y_test0c, predictions))
8
9 print("\nConfusion Matrix: \n")
10 print(confusion_matrix(Y_test0c, predictions))
11
12 print("\nClassification : \n")
13 print(classification_report(Y_test0c, predictions))
```

Accuracy Score:

0.9579608601111379

Confusion Matrix:

```
[[9030  366]
 [ 156 2865]]
```

Classification :

	precision	recall	f1-score	support
0	0.98	0.96	0.97	9396
1	0.89	0.95	0.92	3021
accuracy			0.96	12417
macro avg	0.93	0.95	0.94	12417

(f) calculating accuracy

```
1
2 print("Accuracy for different k values:")
3 for k in range(1, 11):
4     knn = KNeighborsClassifier(n_neighbors=k)
5     knn.fit(X_train0c, Y_train0c)
6     predictions = knn.predict(X_test0c)
7     accuracy = accuracy_score(Y_test0c, predictions)
8     print(f"k = {k}: Accuracy = {accuracy:.4f}")
9
```

Accuracy for different k values:

k = 1: Accuracy = 0.9473
k = 2: Accuracy = 0.9462
k = 3: Accuracy = 0.9530
k = 4: Accuracy = 0.9483
k = 5: Accuracy = 0.9580
k = 6: Accuracy = 0.9564
k = 7: Accuracy = 0.9641
k = 8: Accuracy = 0.9635
k = 9: Accuracy = 0.9648
k = 10: Accuracy = 0.9644

Ques#3

QUES # 3

Counter : for counting occurrences of items in a list. In kNN algorithm, Counter is used to identify the most common label among the top k neighbors.

```
1 import numpy as np
2 import pandas as pd
3 from collections import Counter
4
5 # Chi-squared distance function
6 def chi_squared_distance(A, B):
7     return 0.5 * np.sum([(a - b) ** 2] / (a + b)
8                         for (a, b) in zip(A, B)])
9
10 # kNN classifier
11 def knn_classify(X_train, y_train, test_point, k=3):
12     distances = [(y, chi_squared_distance(X, test_point)) for X, y in zip(X_train, y_train)]
13     distances.sort(key=lambda x: x[1])
14     top_k_labels = [item[0] for item in distances[:k]]
15     most_common = Counter(top_k_labels).most_common(1)
16     return most_common[0][0]
```

```
1 # using Iris dataset
2 names = ['sepal-length', 'sepal-width', 'petal-length', 'petal-width', 'class']
3 dataset = pd.read_csv("iris.csv", names=names, skiprows=1)
4
5 array = dataset.values
6 X = array[:,0:4]
7 Y = array[:,4]
8 t_size = 0.20
9 seed = 7
10 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=t_size, random_state=seed)
11
12 # Testing kNN classifier on Iris
13 print("Accuracy for different k values:")
14 for k in range(1, 11):
15     y_pred = [knn_classify(X_train, Y_train, x, k=k) for x in X_test]
16     accuracy = np.mean(y_pred == Y_test)
17     print(f"k = {k}: Accuracy = {accuracy:.4f}")
```

Accuracy for different k values:

k = 1: Accuracy = 0.9000
k = 2: Accuracy = 0.9000
k = 3: Accuracy = 0.9000
k = 4: Accuracy = 0.8667
k = 5: Accuracy = 0.8667
k = 6: Accuracy = 0.8667
k = 7: Accuracy = 0.8667
k = 8: Accuracy = 0.8667
k = 9: Accuracy = 0.9000
k = 10: Accuracy = 0.9000

for occupancy

```
1 # Using a smaller subset for quick verification
2 X_test_sample = X_test0c[:100]
3 Y_test_sample = Y_test0c[:100]
4
5 print("Accuracy for different k values for kNN classifier:")
6
7 for k in range(1, 11):
8     y_pred = [knn_classify(X_train0c, Y_train0c, x, k=k) for x in X_test_sample]
9     accuracy = np.mean(y_pred == Y_test_sample)
10     print(f"k = {k}: Accuracy = {accuracy:.4f}")
11
```

Accuracy for different k values for kNN classifier:

k = 1: Accuracy = 0.9200
k = 2: Accuracy = 0.9200
k = 3: Accuracy = 0.9000
k = 4: Accuracy = 0.9100
k = 5: Accuracy = 0.9400
k = 6: Accuracy = 0.9500
k = 7: Accuracy = 0.9500
k = 8: Accuracy = 0.9600
k = 9: Accuracy = 0.9500
k = 10: Accuracy = 0.9600