



UNIVERSITY OF MISSOURI-KANSAS CITY

Python and Deep Learning

Lab 2 Assignment

Team Members:

Pragathi Thammaneni

Sridevi Mallipudi

Introduction:

The key objective of this assignment is to focus on all the regression/ prediction models by using python machine learning and execute the code successfully with provided data sets.

Objectives:

To code for the 4 questions by implementing the below concepts.

- Implemented scikit-learn
- Linear Discernment Analysis
- Kernel method/Support vector machines
- Analysis of measured accuracy
- Bi-grams|
- Applied lemmatization
- Used Knn algorithm
- Prediction model analysis

Approaches /Methods:

Using Python 3.6, PyCharm (Community edition)

Workflow &Datasets/Parameters and Evaluation:

The below each question will follow different approaches to solve. Coding is done to perform the evaluation of each individual snippet to execute the datasets which are provided as the input parameters.

Question 1:

1. Pick any dataset from the dataset sheet in the class sheet and make one prediction model using your imagination with **Linear Discriminant Analysis***.

- a. In the report provide convincing explanations about the difference between Logistic Regression and Linear Discriminant Analysis.
- b. You can also pick dataset of your own.

*Logistic Regression is a classification algorithm traditionally limited to only two-class classification problems. If you have more than two classes, then the **Linear Discriminant Analysis** algorithm is the preferred linear classification technique.

Solution:

This snippet code provides the above implementation which is linear discriminant analysis and logistic regression. Required packages are downloaded and Iris dataset has been taken in to consideration for predicting the accuracy of each model. The model is trained with the train data and the test data is used to predict the model accuracy. Split ratio of 80%&20% over train and test is taken for building the model.

Below are the differences between the logistic regression and Linear Discriminant Analysis

- LDA and LR are two widely used multivariate statistical methods for data analysis with categorical outcome variables
- LDA makes a few assumptions such as the explanatory or predictor variables must be normally distributed
- LR does not need any underlying assumption made on the distribution of the data but it takes much longer time than LDA for computing.
- From the below code execution, we can say that the Accuracy is more in LDA than compared to LR

Data Set : Iris Data set

Code Snippet:

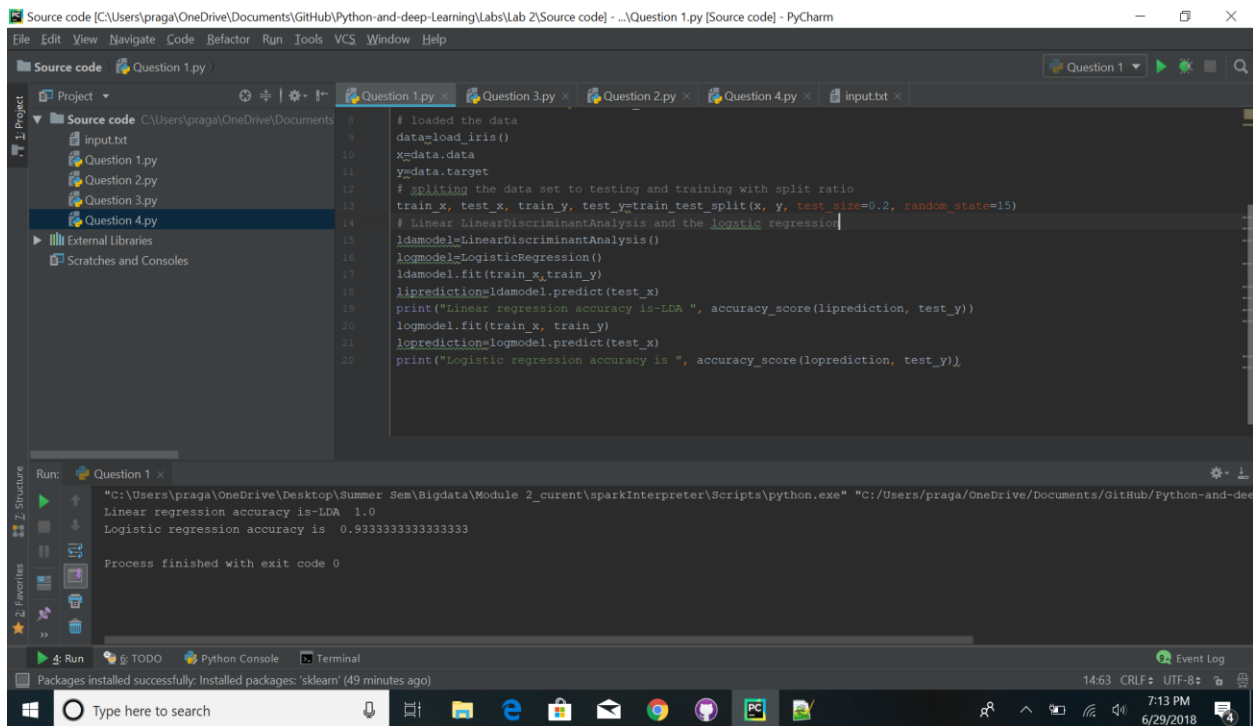
The below is the code for executing the above workflow

```

22 lines (22 sloc) | 1.05 KB
Raw Blame History
1  #e prediction model using with Linear Discriminant Analysis*.
2  # imported the sklearn packages to perform the LinearDiscriminantAnalysis,LogisticRegression
3  from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
4  from sklearn.linear_model import LogisticRegression
5  from sklearn.metrics import accuracy_score
6  from sklearn.model_selection import train_test_split
7  from sklearn.datasets import load_iris
8  # loaded the data
9  data=load_iris()
10 x=data.data
11 y=data.target
12 # splitting the data set to testing and training with split ratio
13 train_x, test_x, train_y, test_y=train_test_split(x, y, test_size=0.2, random_state=15)
14 # Linear LinearDiscriminantAnalysis and the logstic regression
15 ldamodel=LinearDiscriminantAnalysis()
16 logmodel=LogisticRegression()
17 ldamodel.fit(train_x,train_y)
18 lyprediction=ldamodel.predict(test_x)
19 print("Linear regression accuracy is-LDA ", accuracy_score(liprediction, test_y))
20 logmodel.fit(train_x, train_y)
21 loprediction=logmodel.predict(test_x)
22 print("Logistic regression accuracy is ", accuracy_score(loprediction, test_y))

```

Output Screenshot:



The screenshot shows the PyCharm IDE with a Python script in the main editor. The script loads the Iris dataset, splits it into training and testing sets (20% test, 80% train), and applies both Linear Discriminant Analysis (LDA) and Logistic Regression models. The output window shows the results of the training and testing process.

```
# loaded the data
data=load_iris()
x=data.data
y=data.target
# splitting the data set to testing and training with split ratio
train_x, test_x, train_y, test_y=train_test_split(x, y, test_size=0.2, random_state=15)
# Linear LinearDiscriminantAnalysis and the logistic regression
ldamodel=LinearDiscriminantAnalysis()
logmodel=LogisticRegression()
ldamodel.fit(train_x,train_y)
l1prediction=ldamodel.predict(test_x)
print("Linear regression accuracy is-LDA ", accuracy_score(l1prediction, test_y))
logmodel.fit(train_x, train_y)
logprediction=logmodel.predict(test_x)
print("Logistic regression accuracy is ", accuracy_score(logprediction, test_y))
```

Run: Question 1 x

```
"C:\Users\praga\OneDrive\Desktop\Summer Sem\Bigdata\Module 2_curent\sparkinterpreter\Scripts\python.exe" "C:/Users/praga/OneDrive/Documents/GitHub/Python-and-dee
Linear regression accuracy is-LDA 1.0
Logistic regression accuracy is 0.9333333333333333

Process finished with exit code 0
```

14:63 CRLF UTF-8 7:13 PM 6/29/2018

Question 2:

2. Implement Support Vector Machine classification,
 - a. Choose one of the datasets using the datasets features in the scikit-learn
 - b. Load the dataset
 - c. According to your dataset, split the data into 20% testing data, 80% training data (you can also use any other number)
 - d. Apply SVC with Linear kernel
 - e. Apply SVC with RBF kernel
 - f. Report the accuracy of the model on both models separately and report their differences if there are any
 - g. Report your view how can you increase the accuracy and which kernel is the best for your dataset and why

Solution:

This snippet code provides the above implementation of Support vector classification. Imported the required modules from the python packages. Iris dataset has been loaded to perform SVM. The data is splitted in to 20% testing and 80% training data. Then Applied both the Linear kernel and RBF kernel and reported the below screen shots for the predicted accuracy.

- Linear kernel is best compared to RBF kernel but it also depends on the data set and the data size. It varies mainly when the random state is changed. So it is difficult to figure out which is a better kernel, it mainly depends and varies accordingly to the data.
- **Data Set : Iris Data set**

Code Snippet:

The below is the code for executing the above workflow

```
33 lines (33 sloc) | 1.89 KB
Raw Blame History
1 #Implement Support Vector Machine classification,a. Choose one of the datasetsusing the datasets features in the scikit-learnb. Load the d
2 #imported the packages from sklearn module for support vector , accuracy, split train data and datasets
3 from sklearn.svm import SVC
4 from sklearn.metrics import accuracy_score
5 from sklearn.model_selection import train_test_split
6 from sklearn import datasets
7 #load the iris datasets
8 data=datasets.load_iris()
9 #load x and y data
10 x=data.data
11 y=data.target
12 #splitting the training and test data for both x and y for linear kernel
13 train_x, test_x, train_y, test_y=train_test_split(x, y, test_size=0.2, random_state=21)
14 #splitting the training and test data for both x and y for rbf kernel
15 train_x1, test_x1, train_y1, test_y1=train_test_split(x, y, test_size=0.2, random_state=19)
16 #defining the model for linear kernel
17 lmodel=SVC(kernel='linear')
18 #defining the model for rbf kernel
19 rmodel=SVC(kernel='rbf')
20 #fitting the training data into linear kernel
21 lmodel.fit(train_x, train_y)
22 #predicting the test data using linear kernel
23 prediction=lmodel.predict(test_x)
24 #calculate the accuracy score for linear kernel
25 print("Linear kernel Accuracy score is", accuracy_score(prediction, test_y))
26 print(prediction)
27 #fitting the training data into rbc kernel
28 rmodel.fit(train_x1, train_y1)
29 #predicting the test data for rbc kernel
30 pred=lmodel.predict(test_x1)
31 #calculated theb accuracy for rbc kernel
32 print("RBF kernel accuracy score is", accuracy_score(pred, test_y1))
33 print(pred)
```

Output Screenshot:

Variations in the accuracy depending up on the random state of each kernel.

This screenshot shows a PyCharm window with a Python script for SVM classification. The script imports necessary libraries, loads the Iris dataset, and splits it into training and testing sets. It then defines two models: a linear kernel SVM and an RBF kernel SVM. The output shows the accuracy scores for both models.

```
1 #Implement Support Vector Machine classification.a. Choose one of the datasets using the datasets features in the sklearn
2 #Imported the packages from sklearn module for support vector accuracy, split train data and datasets
3 from sklearn.svm import SVC
4 from sklearn.metrics import accuracy_score
5 from sklearn.model_selection import train_test_split
6 from sklearn import datasets
7 #load the iris datasets
8 d=datasets.load_iris()
9 #load x and y data
10 x=data.data
11 y=data.target
12 #splitting the training and test data for both x and y for linear kernel
13 train_x, test_x, train_y, test_y=train_test_split(x, y, test_size=0.2, random_state=19)
14 #splitting the training and test data for both x and y for rbf kernel
15 train_x1, test_x1, train_y1, test_y1=train_test_split(x, y, test_size=0.2, random_state=21)
16 #defining the model for linear kernel
17 lmodel=SVC(kernel='linear')
18 #defining the model for rbf kernel
19 rmodel=SVC(kernel='rbf')
```

Run: Question 2 x

```
"C:\Users\praga\OneDrive\Desktop\Summer Sem\Bigdata\Module 2_curent\sparkInterpreter\Scripts\python.exe" "C:/Users/praga/OneDrive/Documents/GitHub/Python-and-dee
Linear kernel Accuracy score is 1.0
[0 2 1 1 0 0 0 0 1 2 1 0 1 0 2 0 2 0 1 0 1 1 1 1 2 1 2 2 1 2]
RBF kernel accuracy score is 0.9666666666666667
[1 0 0 0 1 1 0 2 0 0 1 1 2 2 0 1 2 1 2 2 1 2 1 0 1 0 0 1 2]
Process finished with exit code 0
```

This screenshot shows a PyCharm window with a Python script for SVM classification. The script imports necessary libraries, loads the Iris dataset, and splits it into training and testing sets. It then defines two models: a linear kernel SVM and an RBF kernel SVM. The output shows the accuracy scores for both models.

```
1 #Implement Support Vector Machine classification.a. Choose one of the datasets using the datasets features in the sklearn
2 #Imported the packages from sklearn module for support vector accuracy, split train data and datasets
3 from sklearn.svm import SVC
4 from sklearn.metrics import accuracy_score
5 from sklearn.model_selection import train_test_split
6 from sklearn import datasets
7 #load the iris datasets
8 d=datasets.load_iris()
9 #load x and y data
10 x=data.data
11 y=data.target
12 #splitting the training and test data for both x and y for linear kernel
13 train_x, test_x, train_y, test_y=train_test_split(x, y, test_size=0.2, random_state=21)
14 #splitting the training and test data for both x and y for rbf kernel
15 train_x1, test_x1, train_y1, test_y1=train_test_split(x, y, test_size=0.2, random_state=19)
16 #defining the model for linear kernel
17 lmodel=SVC(kernel='linear')
18 #defining the model for rbf kernel
19 rmodel=SVC(kernel='rbf')
```

Run: Question 2 x

```
"C:\Users\praga\OneDrive\Desktop\Summer Sem\Bigdata\Module 2_curent\sparkInterpreter\Scripts\python.exe" "C:/Users/praga/OneDrive/Documents/GitHub/Python-and-dee
Linear kernel Accuracy score is 0.9666666666666667
[1 0 0 0 1 1 0 2 0 0 1 1 2 2 0 1 2 1 2 2 1 2 1 0 1 0 0 1 2]
RBF kernel accuracy score is 1.0
[0 2 1 1 0 0 0 0 1 2 1 0 1 0 2 0 2 0 1 0 1 1 1 1 2 1 2 2 1 2]
Process finished with exit code 0
```

Question 3:

3. Write a program

Take an Input file. Use the simple approach below to summarize a text file:

- a. Read the file
- b. Apply lemmatization on the words
- c. Apply the bigram on the text
- d. Calculate the word frequency (bi-gram frequency) of the words (bi-grams)
- f. Choose top five bi-grams that have been repeated most
- g. Go through the original text that you had in the file
- h. Find all the sentences with those most repeated bi-grams
- i. Extract those sentences and concatenate
- j. Enjoy the summarization

Solution:

This snippet code provides the above implementation of lemmatization, bigram, top rated 5-words and the concatenation of the sentences. Input file is loaded and below code is used for performing the above steps.

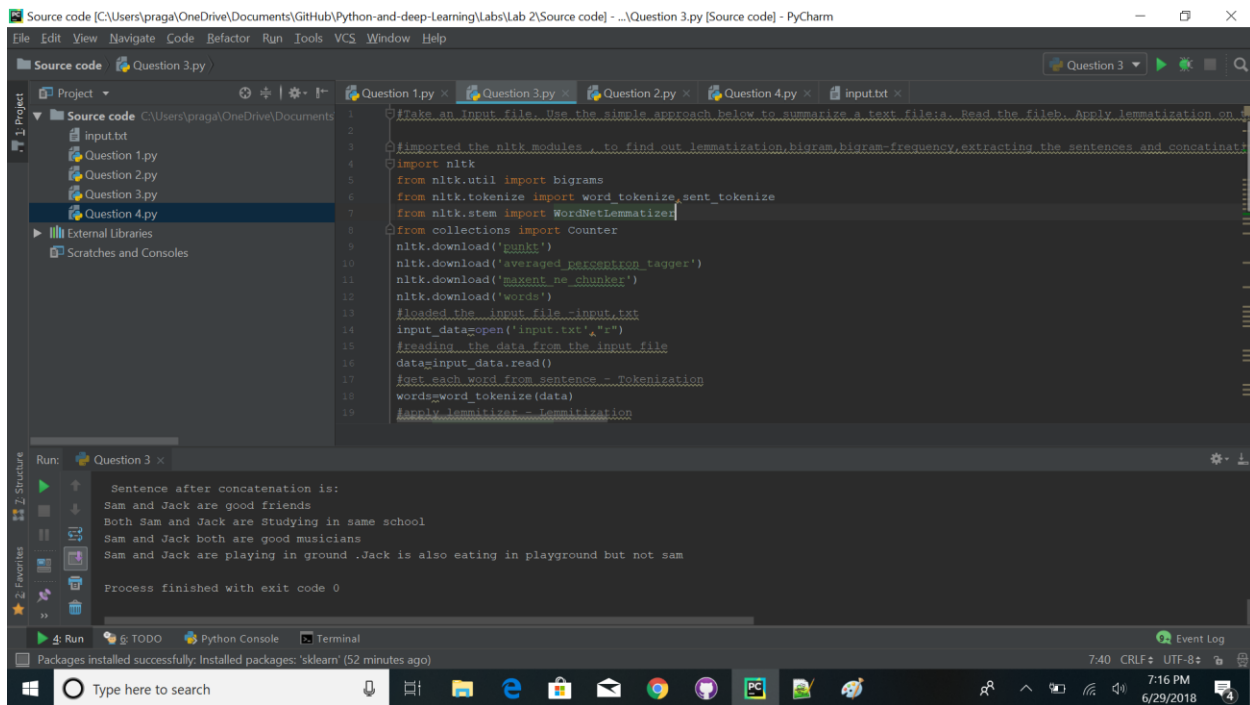
Data Set : A sample text file with the paragraph

Code Snippet:

The below is the code for executing the above workflow

```
1 #Take an Input file. Use the simple approach below to summarize a text file:a. Read the file. Apply lemmatization on the wordsc. Apply the
2
3 #imported the nltk modules , to find out lemmatization,bigram,bigram-frequency,extracting the sentences and concatenation
4 import nltk
5 from nltk.util import bigrams
6 from nltk.tokenize import word_tokenize,sent_tokenize
7 from nltk.stem import WordNetLemmatizer
8 from collections import Counter
9 nltk.download('punkt')
10 nltk.download('averaged_perceptron_tagger')
11 nltk.download('maxent_ne_chunker')
12 nltk.download('words')
13 #loaded the input file -input.txt
14 input_data=open('input.txt',"r")
15 #reading the data from the input file
16 data=input_data.read()
17 #get each word from sentence - Tokenization
18 words=word_tokenize(data)
19 #apply lemmatizer - Lemmatization
20 lemm=WordNetLemmatizer()
21 list=[]
22 list1=[]
23 print("Lemmatizer:"+ "\n")
24 #printing the lemmatizer for each word
25 for k in words:
26     print("Lemmatizer for word "+k+" is "+lemm.lemmatize(k))
27
28 print("\n")
29 print("Bigrams:"+ "\n")
```

Output Screenshot:



The screenshot shows the PyCharm IDE with a Python script in the main editor and its output in the Run console. The script uses NLTK for text processing, including downloading corpora, tokenizing, and lemmatizing text. The output in the Run console shows the concatenated sentences and the process finished with exit code 0.

```
1 #Take an input file. Use the simple approach below to summarize a text file:- Read the file, Apply lemmatization on
2
3 #Imported the nltk modules , to find out lemmatization,bigram,bigram-frequency,extracting the sentences and concatenat
4 import nltk
5 from nltk.util import bigrams
6 from nltk.tokenize import word_tokenize,sent_tokenize
7 from nltk.stem import WordNetLemmatizer
8 from collections import Counter
9 nltk.download('punkt')
10 nltk.download('averaged_perceptron_tagger')
11 nltk.download('maxent_ne_chunker')
12 nltk.download('words')
13 #loaded the input file -input.txt
14 input_data=open('input.txt','r')
15 #reading the data from the input file
16 data=input_data.read()
17 #get each word from sentence - Tokenization
18 words=word_tokenize(data)
19 #Apply lemmatizer = Lemmatization
```

Run: Question 3 x

Sentence after concatenation is:
Sam and Jack are good friends
Both Sam and Jack are studying in same school
Sam and Jack both are good musicians
Sam and Jack are playing in ground .Jack is also eating in playground but not sam

Process finished with exit code 0

Question 4:

4. Report your views on the k nearest neighbor algorithm when we change the K how it will affect the accuracy. Provide a good justification for the changes of the accuracy when we change the amount of K. For example: compare the accuracy when K=1 and K is a big number like 50, why the accuracy will change

Solution:

This snippet code provides the above implementation KNN algorithm with different k values. When the k value increases, the model will better fit into the training data and get the accurate results while compared to the less k value. The K value has more effect on the test data So there will be variation in predicting the accuracy.

Data Set : Iris Data set

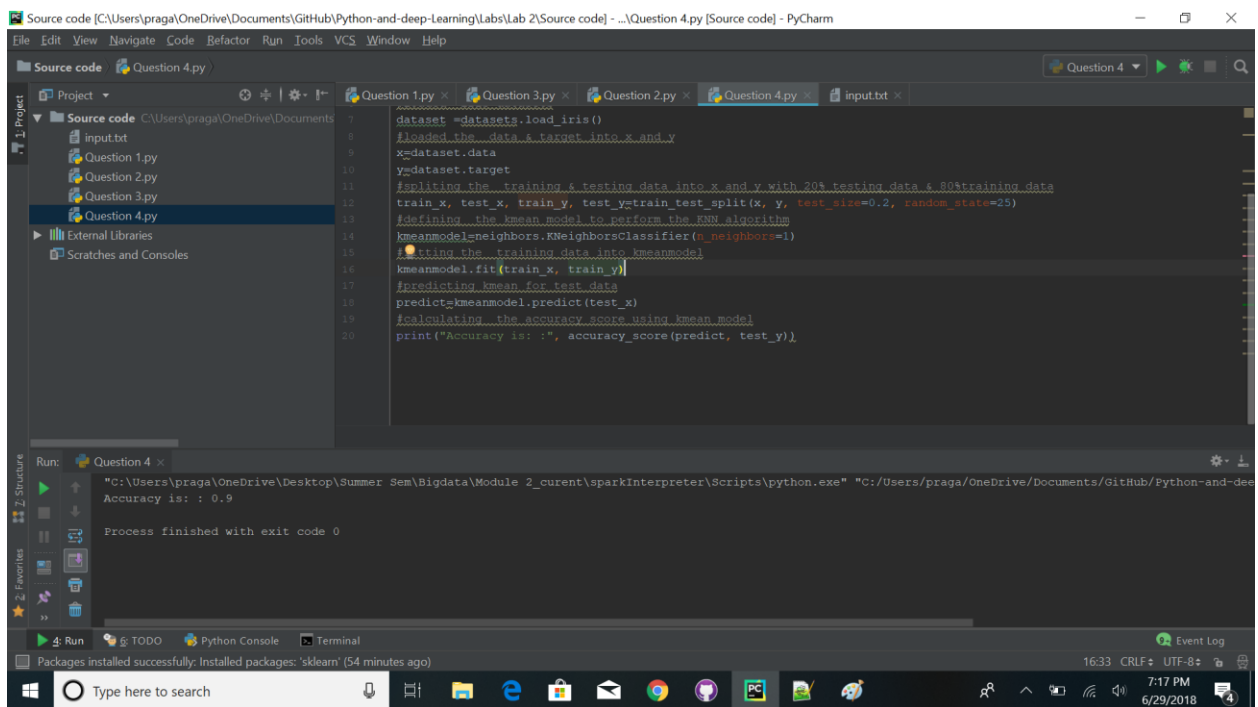
Code Snippet:

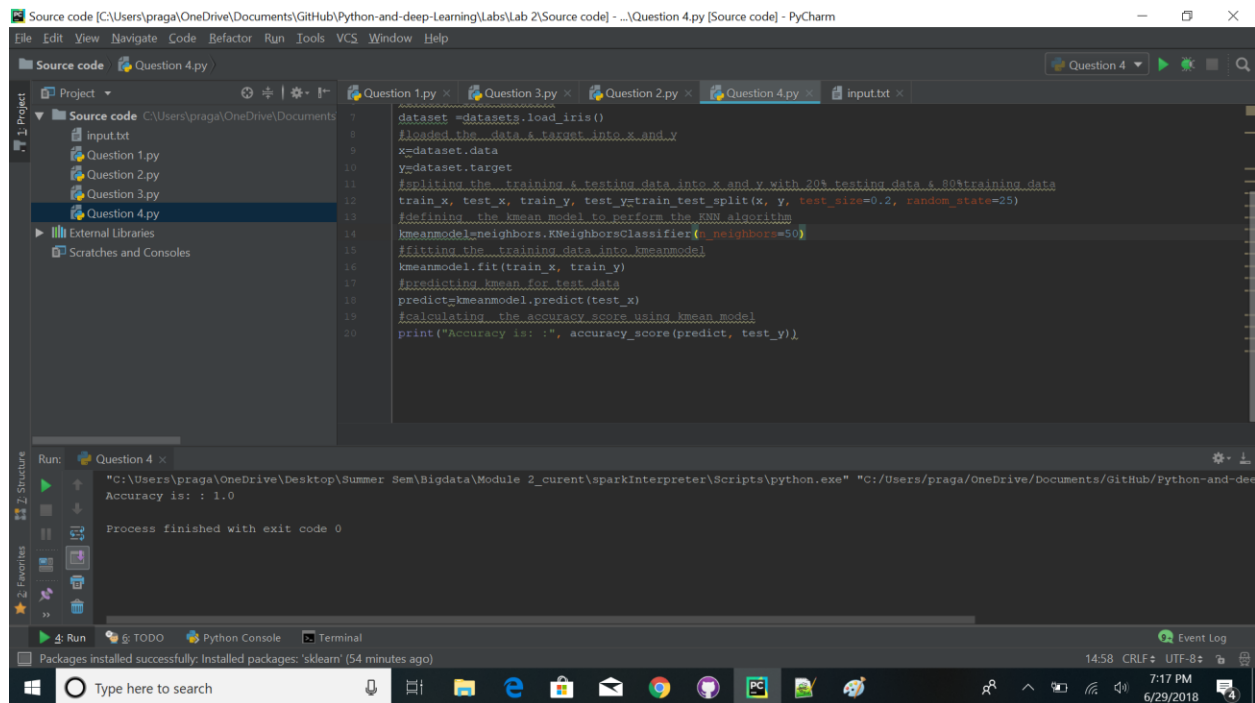
The below is the code for executing the above workflow


```
20 lines (20 sloc) | 941 Bytes
Raw Blame History

1 # k nearest neighbor algorithm
2 #imported the sklearn package for accuracy,metrics,datasets,split train & test data
3 from sklearn.metrics import accuracy_score
4 from sklearn.model_selection import train_test_split
5 from sklearn import neighbors,datasets
6 #loaded iris datasets
7 dataset =datasets.load_iris()
8 #loaded the data & target into x and y
9 x=dataset.data
10 y=dataset.target
11 #splitting the training & testing data into x and y with 20% testing data & 80%training data
12 train_x, test_x, train_y, test_y=train_test_split(x, y, test_size=0.2, random_state=25)
13 #defining the kmean model to perform the KNN algorithm
14 kmeanmodel=neighbors.KNeighborsClassifier(n_neighbors=50)
15 #fitting the training data into kmeanmodel
16 kmeanmodel.fit(train_x, train_y)
17 #predicting kmean for test data
18 predict=kmeanmodel.predict(test_x)
19 #calculating the accuracy score using kmean model
20 print("Accuracy is: :", accuracy_score(predict, test_y))
```

Output Screenshot:





Conclusion:

As stated the above workflow with certain set of parameters is followed in solving the execution by implementing the core and basic concepts of the python programming.

Source code link <https://github.com/PragathiThammaneni/Python-and-deep-Learning/tree/master/Lab%202/Source%20code>

Video Link: <https://youtu.be/OdIXO6z8gCE> (made a short video as it should be less time)

Wiki Link:

<https://github.com/PragathiThammaneni/Python-and-deep-Learning/wiki/Lab-2-Assignment>