# CS5590 APS -Python Programming

# LAB1

Dukkipati, Sri Sai Nithin Chowdary – 4

Inakollu, Sri Naga Bhuvaneswari - 9

Kolluri, Nikhitha – 12

**INTRODUCTION:**

In this lab we will go through some basics of python, use BeautifulSoup package for web scrapping, plotting patterns using some classification algorithms, applying tokenization, lemmatization and trigrams on some file, creating some multiple regression and evaluating the model using RMSE and R2 techniques. Cleaning all the dataset before its usage.
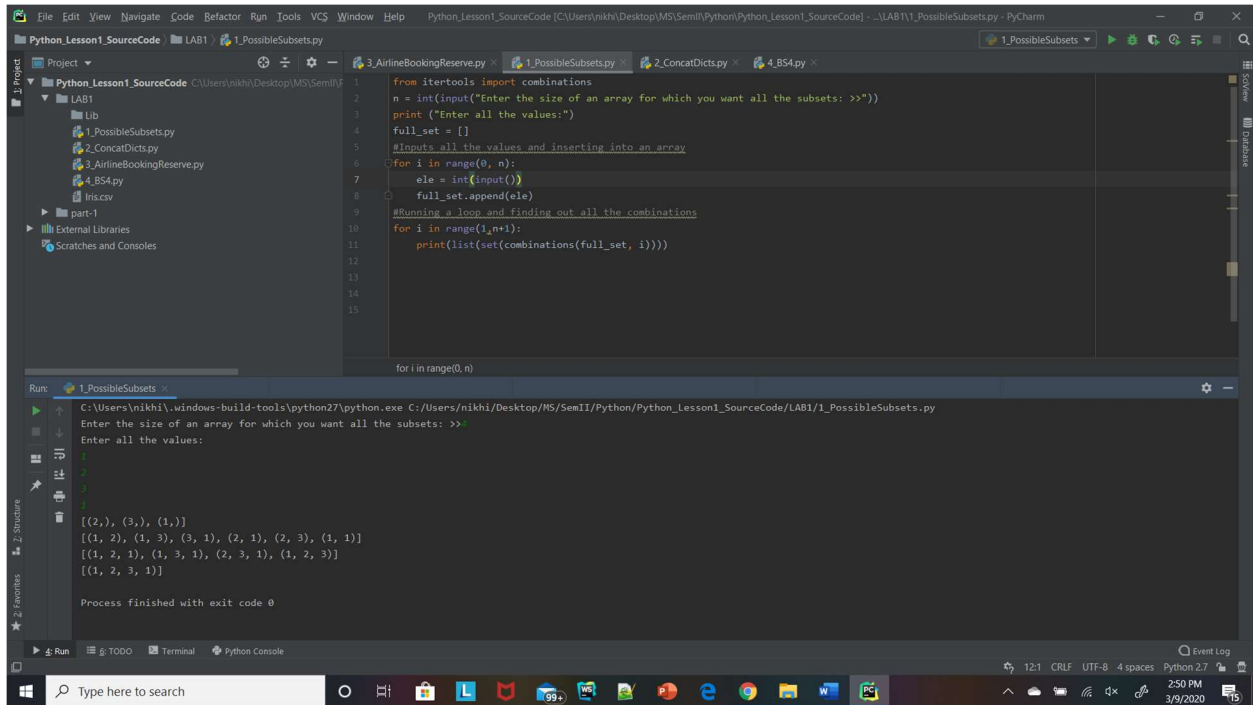
**OBJECTIVES:**

Some objectives of this lab are:

- Basics of python such as performing some operations on subsets, concatenating two dictionaries and some operations on it.
- Airline Booking Reservation System which takes in the details of customer, his preferences and then display the available airlines from which the customer can select source, destination, the type of airlines he wish to fly in, the class (business or economy) and number of bags he wish to check-in . The system should display the fare for the ticket and prints the ticket. The system should display all the details of the employee that is booking the ticket.
- Use the BeautifulSoup package and from the school catalogue site get all the course names and their description.
- Take a dataset which will have both numeric and non-numeric data, perform data analysis on the chosen dataset and plot patterns. Before plotting patterns remove all the null values from the data set, remove all the features which are unrelated to the target class and all the categorical features should be encoded. Report the classification of the data set based on some classification algorithms such as Naïve Bayes, SVM and KNN.
- Choose any dataset you wish and them perform K-means on that and then visualize the data using matplotlib or seaborn. Select the best K by applying elbow method, evaluate the silhouette score. Before evaluating the score, clean the data set with EDA.
- Take an input file which has few sentences and apply tokenization, lemmatization, trigrams, most repeated trigrams, extract few sentences and concatenate them and print the result.
- Choose a data set of your choice, create multiple regression. Use RMSE and R2 to evaluate the model, display the improved results after EDA.

**APPROACHES/METHODS:**

1. Return all possible subsets for a given collection of integers without including the null subset.

a. First we will be accepting the length of the integers from the user.
b. Then we will ask them to input the integers for the set using a for loop and append all these integers to a list.
c. Then using itertools.combinations we will get all the subsets.
d. If you apply set on this combinations we can eliminate the duplicates.
e. We use a loop to get all the subsets of a given set.



2. Concatenate two dictionaries and sort them based on the values.
    a. First declare and define two dictionaries.
    b. Merge the dictionaries using the update method.
    c. Sorting the merged dictionary by using sorted method which takes lambda function and sorts based on the value field.
    d. We use a for loop to print all the key value pairs in sorted order.

3. Airline Booking Reservation System
   a. First we have taken a flight class which will take some inputs from user, like source, destination, airlines, class and will generate the flight number. It also has a method that displays the details.
   b. Employee class will display all the details of the employee. It uses the concept of inheritance and overrides the printing method of the Flight class.
   c. Passenger class will take all the details of the passenger.
   d. Baggage class would take the number of check-in bags from the passenger and calculate the baggage fare.
   e. TicketCost class takes the inputs from all the classes and would calculate the fare of the ticket. A method in TicketCost would display all the ticket details.

```python
#This class inputs the journey details and gives the flight number
class Flight:
    def __init__(self, flightNumber):
        Flight.source = input('Enter Source : ')
        Flight.destination = input('Enter Destination : ')
        print('Available Airlines :')
        print('SouthWest')
        print('American')
        print('Spirit')
        Flight.airlinesName = input('Which Airlines do you prefer : ')
        self.flightNumber = flightNumber

    #This method prints the airlines details
    def print_details(self):
        print('Airlines : ', Flight.airlinesName)
        print('Flight Number : ', self.flightNumber)
        print(Flight.source, '-->', Flight.destination)
        print('*******************************************************************')

#This class gives all the employee details (Inheritance)
class Employee(Flight):
    def __init__(self, emp_id, emp_name, emp_gender):
        self.emp_name = emp_name
        self.emp_id = emp_id
        self.emp_gender = emp_gender

    #Method Overriding
    def print_details(self):
        print("Name of employee: ", self.emp_name)
        print('Employee id: ', self.emp_id)
```

Employee

---

```python
        print('Employee gender: ', self.emp_gender)

#This class inputs all the passenger details
class Passenger:
    def __init__(self):
        Passenger.passenger_fname = input('Enter first name : ')
        Passenger.passenger_lname = input('Enter last name : ')
        Passenger.passenger_passportNo = input('Enter passport number: ')
        Passenger.passenger_gender = input('Enter gender : ')
        Passenger.passenger_class = input('Business or Economy class? : ')

#This class calculates the baggage fare based on the number of bags
class Baggage:
    def __init__(self):
        Baggage.noOfBags = int(input('Number of bags you want to checkin : '))
        Baggage.totBagFare = 0;
        Baggage.noOfBags = Baggage.noOfBags
        if(Baggage.noOfBags > 2):
            for i in range(Baggage.noOfBags-2):
                Baggage.totBagFare += 60
        print('You can take two bags for free !!! Total bag fare is for ', Baggage.noOfBags,'is ',Baggage.totBagFare)

#This class calculates the ticket cost based on the class, flight and bags (Inheritance)
class TicketCost(Baggage, Passenger, Flight):
    def __init__(self):
        TicketCost.baseCost = 100
        TicketCost.baseCost = TicketCost.baseCost + Baggage.totBagFare
        if(Passenger.passenger_class == 'business'):
            TicketCost.baseCost = TicketCost.baseCost + 100
        print('Total ticket cost is : ', TicketCost.baseCost)
```

Employee

4. WebScrapping using BeautifulSoup package
   a. First we would get the html of the given URL using BS4 package
   b. Then convert all the html into plain text.
   c. We find all the div tags with class='courseblock'.
   d. We then run a loop to get all the course names and course descriptions.
   e. We then get the text only without the tags using the text option.
   f. Finally we display the result.

File   Edit   View   Navigate   Code   Refactor   Run   Tools   VCS   Window   Help

Python_Lesson1_SourceCode > LAB1 > 4_BS4.py

```python
import requests
from bs4 import BeautifulSoup

def get_links():
    url = "https://catalog.umkc.edu/course-offerings/graduate/comp-sci/"
    sourceCode = requests.get(url)
    plainText= sourceCode.text
    soup=BeautifulSoup(plainText,"html.parser")
    print('******************************************************************')
    # Finding all the divs with class courseblock
    result = soup.find_all('div', {'class': 'courseblock'})
    # getting and printing all the course names and corse descriptions from the scrapped code
    for c in result:
        result1 = c.find('span', {'class': 'code'}).text
        result2 = c.find('p', {'class': 'courseblockdesc'}).text
        print(result1)
        print(result2)

get_links()
```

get_links() > for c in result

Run: 4_BS4

```
Advanced research by a group of doctoral students based on intensive readings from the current research literature under the direction of one or more doctoral faculty.  Original research results of
COMP-SCI 5699A

Doctoral research in computer science.
COMP-SCI 5899


Process finished with exit code 0
```

4: Run    6: TODO    Terminal    Python Console

15:65    CRLF    UTF-8    4 spaces    Python 2.7

Type here to search

3:15 PM
3/9/2020

---

Run: 4_BS4

```
C:\Users\nikhi\.windows-build-tools\python27\python.exe C:/Users/nikhi/Desktop/MS/SemII/Python/Python_Lesson1_SourceCode/LAB1/4_BS4.py
******************************************************************
COMP-SCI 5101

A review of mathematical logic, sets, relations, functions, mathematical induction, and algebraic structures with emphasis on computing applications. Recurrence relations and their use in the analys
COMP-SCI 5102

This course covers concurrency and control of asynchronous processes, deadlocks, memory management, processor and disk scheduling, parallel processing, and file system organization in operating syst
COMP-SCI 5103

A review of linear and hierarchical data structures, including stacks, queues, lists, trees, priority queues, advanced tree structures, hashing tables, dictionaries and disjoint-sets. Asymptotic ana
COMP-SCI 5514

Fiber optic cable and its characteristics, optical sources and transmitters, optical detectors and receivers, optical components such as couplers and connectors, WDM and OFDM techniques, modulation
COMP-SCI 5525

Cloud computing systems operate in a very large scale, and are impacting the economics and the assumptions behind computing significantly. This special topics course provides a comprehensive overvie
COMP-SCI 5531

Components of an operating system, scheduling/routing mechanisms, process control blocks, design and test various operating system components.
COMP-SCI 5540

This course will introduce the essential characteristics of Big Data and why it demands rethinking how we store, process, and manage massive amounts of structured and unstructured data. It will cove
COMP-SCI 5542

Big Data analytics focus on analyzing large amounts of data to find useful information and to make use of the information for better business decisions. This course introduces students to the practi
COMP-SCI 5543

This course teaches students fundamental theory and practice in the field of big data analytics and real time distributed systems for real time big data applications. In this course, students will h
COMP-SCI 5551
```

4: Run    6: TODO    Terminal    Python Console    Event Log

```
Readings in an area selected by the graduate student in consultation with a faculty member.  Arrangements must be made prior to registration.
COMP-SCI 5598

Graduate research based on intensive readings from the current research literature under the direction of a faculty member.  Arrangements must be made prior to registration.
COMP-SCI 5599

A project investigation leading to a thesis, or written report under the direction of a faculty member. A prospectus must be accepted prior to registration.
COMP-SCI 5690

A lecture course presenting advanced research level topics. This course is intended to allow faculty and visiting scholars to offer special courses in selected research areas.
COMP-SCI 5690ND


COMP-SCI 5697

Readings in an area selected by the doctoral student in consultation with a doctoral faculty member.  Arrangements must be made prior to registration.
COMP-SCI 5698

Advanced research by a group of doctoral students based on intensive readings from the current research literature under the direction of one or more doctoral faculty.  Original research results of
COMP-SCI 5699A

Doctoral research in computer science.
COMP-SCI 5899



Process finished with exit code 0
```

## 5. SVM, KNN, Naïve Bayes classifier

### 5a.

a. Took the data glass_type.cvs file which consists of both numeric and non-numberic data

b. Used label label encoder to change non numeric data to numeric one

c. Used fillna to d=fill all the null values with mean data(adding noise)

d. Used spilt function for splitting the data into test and training data



```python
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report
from sklearn.preprocessing import LabelEncoder


train_df = pd.read_csv('/content/drive/My Drive/Colab Notebooks/glass_type.csv')
Label_Encoder = LabelEncoder()
print("Before transforming the data using labelencoder",train_df)
train_df["Type"] = Label_Encoder.fit_transform(train_df["Type"])
train_df=train_df.apply(lambda x: x.fillna(x.mean()),axis=0)
print("After transforming the data using labelencoder",train_df)

print(train_df.isnull().sum())
X_train = train_df.drop("Type",axis=1)
Y_train = train_df["Type"]


X_train, X_test, Y_train, y_test= train_test_split(X_train, Y_train, test_size=0.4, random_state=0)
```

```
Before transforming the data using labelencoder          RI    Na    Mg    Al    Si    K    Ca   Ba  Fe            Type
0     1.52101  13.64  4.49  1.10  71.78  0.06  8.75  0.00  0.0  Patterned glass
1     1.51761  13.89  3.60  1.36  72.73  0.48  7.83  0.00  0.0  Patterned glass
2     1.51618  13.53  3.55  1.54  72.99  0.39  7.78  0.00  0.0  Patterned glass
3     1.51766  13.21  3.69  1.29  72.61  0.57  8.22  0.00  0.0  Patterned glass
4     1.51742  13.27  3.62  1.24  73.08  0.55  8.07  0.00  0.0  Patterned glass
..        ...    ...   ...   ...    ...   ...   ...   ...  ...       ...
209   1.51623  14.14  0.00  2.88  72.61  0.08  9.18  1.06  0.0   silicate glass
210   1.51685  14.92  0.00  1.99  73.06  0.00  8.40  1.59  0.0   silicate glass
211   1.52065  14.36  0.00  2.02  73.42  0.00  8.44  1.64  0.0   silicate glass
212   1.51651  14.38  0.00  1.94  73.61  0.00  8.48  1.57  0.0   silicate glass
213   1.51711  14.23  0.00  2.08  73.36  0.00  8.62  1.67  0.0   silicate glass

[214 rows x 10 columns]
After transforming the data using labelencoder          RI    Na    Mg    Al    Si    K    Ca   Ba  Fe  Type
0     1.52101  13.64  4.49  1.10  71.78  0.06  8.75  0.00  0.0     2
1     1.51761  13.89  3.60  1.36  72.73  0.48  7.83  0.00  0.0     2
```

```
Before transforming the data using labelencoder          RI    Na   Mg    Al     Si    K    Ca    Ba   Fe         Type
0    1.52101  13.64  4.49  1.10  71.78  0.06  8.75  0.00  0.0  Patterned glass
1    1.51761  13.89  3.60  1.36  72.73  0.48  7.83  0.00  0.0  Patterned glass
2    1.51618  13.53  3.55  1.54  72.99  0.39  7.78  0.00  0.0  Patterned glass
3    1.51766  13.21  3.69  1.29  72.61  0.57  8.22  0.00  0.0  Patterned glass
4    1.51742  13.27  3.62  1.24  73.08  0.55  8.07  0.00  0.0  Patterned glass
..       ...    ...   ...   ...    ...   ...   ...   ...  ...             ...
209  1.51623  14.14  0.00  2.88  72.61  0.08  9.18  1.06  0.0   silicate glass
210  1.51685  14.92  0.00  1.99  73.06  0.00  8.40  1.59  0.0   silicate glass
211  1.52065  14.36  0.00  2.02  73.42  0.00  8.44  1.64  0.0   silicate glass
212  1.51651  14.38  0.00  1.94  73.61  0.00  8.48  1.57  0.0   silicate glass
213  1.51711  14.23  0.00  2.08  73.36  0.00  8.62  1.67  0.0   silicate glass

[214 rows x 10 columns]
After transforming the data using labelencoder          RI    Na   Mg    Al     Si    K    Ca    Ba   Fe  Type
0    1.52101  13.64  4.49  1.10  71.78  0.06  8.75  0.00  0.0     2
1    1.51761  13.89  3.60  1.36  72.73  0.48  7.83  0.00  0.0     2
2    1.51618  13.53  3.55  1.54  72.99  0.39  7.78  0.00  0.0     2
3    1.51766  13.21  3.69  1.29  72.61  0.57  8.22  0.00  0.0     2
4    1.51742  13.27  3.62  1.24  73.08  0.55  8.07  0.00  0.0     2
..       ...    ...   ...   ...    ...   ...   ...   ...  ...   ...
209  1.51623  14.14  0.00  2.88  72.61  0.08  9.18  1.06  0.0     5
210  1.51685  14.92  0.00  1.99  73.06  0.00  8.40  1.59  0.0     5
211  1.52065  14.36  0.00  2.02  73.42  0.00  8.44  1.64  0.0     5
212  1.51651  14.38  0.00  1.94  73.61  0.00  8.48  1.57  0.0     5
213  1.51711  14.23  0.00  2.08  73.36  0.00  8.62  1.67  0.0     5

[214 rows x 10 columns]
RI      0
Na      0
Mg      0
Al      0
Si      0
K       0
Ca      0
Ba      0
Fe      0
Type    0
dtype: int64
```

**5b.**

e.  Called the class svm, gaussianNb, KNeihborsClassifier.

f. Fitted the model with training data

g.  Predicted the accuracy with score function with test data

**We got more accuracy with KNN where k=3**

## 6. K-means

### 6a.

  a. Used cc.csv data set

  b. Visualized data using bar graph with seaborn library

  c. Extracted columns that are required

  d. Cleaned the data, used fillna to fill all null values(fitting noise data)

  e. Used StandardSCaler function for converting the data into numeric form

  f.  Used  datafram method to convert the data again back to dataframe
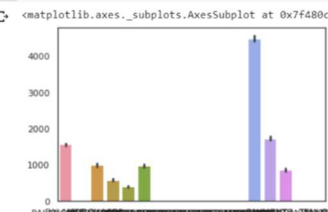
  g. Used KMeans clustering function to get best k value

**Files**

- CC.csv
- Copy of LAB1_6.ipynb
- Copy of LAB1_7.ipynb
- Copy of Lab1_8.ipynb
- Copy of test.ipynb
- Heart.csv
- ICP7_1st.ipynb
- ICP7_2nd.ipynb
- Iris.csv
- LAB1_5.ipynb
- LAB1_6.ipynb
- LAB1_7.ipynb
- Lab1_8.ipynb
- big.txt
- glass.csv
- glass_type.csv
- nlp_input.txt
- spelling_corrector.ipynb
- test.ipynb
- webscrap.txt
- winequality-red.csv
- Iotmini project video
- Video for increment 3
- personal documents
- tweets

Disk 79.35 GB available

```python
import pandas as pd
import seaborn as sns
sns.set(style="white", color_codes=True)
import warnings
warnings.filterwarnings("ignore")

#Reading the Data
data = pd.read_csv('/content/drive/My Drive/Colab Notebooks/CC.csv')
sns.barplot(data=data)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f480c2b72e8>
```



```python
[21] #Removing the unrequired columns
     x = data.iloc[:,[1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17]]
     #Filling all the null values with mean
     x=x.apply(lambda x: x.fillna(x.mean()),axis=0)
     print(x.isnull().sum())
```

```
BALANCE                0
BALANCE_FREQUENCY      0
```

**CC.csv**

1 to 10 of 8

| CUST_ID | BALANCE | BALANCE_FREQUENCY | PUR |
|---------|---------|-------------------|-----|
| C10001 | 40.900749 | 0.818182 | 95.4 |
| C10002 | 3202.467416 | 0.909091 | 0 |
| C10003 | 2495.148862 | 1 | 773. |
| C10004 | 1666.670542 | 0.636364 | 1499 |
| C10005 | 817.714335 | 1 | 16 |
| C10006 | 1809.828751 | 1 | 1333 |
| C10007 | 627.260806 | 1 | 7091 |
| C10008 | 1823.652743 | 1 | 436. |
| C10009 | 1014.926473 | 1 | 861. |
| C10010 | 152.225975 | 0.545455 | 1281 |

Show 10 per page   1   2   10   100

---

**LAB1_6.ipynb**
File Edit View Insert Runtime Tools Help   All changes saved

**Files**

- drive
  - My Drive
    - AICHAT PPT
    - AICHAT VIDEO
    - Colab Notebooks
      - CC.csv
      - Copy of test.ipynb
      - Heart.csv
      - ICP7_1st.ipynb
      - ICP7_2nd.ipynb
      - Iris.csv
      - LAB1_5.ipynb
      - LAB1_6.ipynb
      - big.txt
      - glass.csv
      - glass_type.csv
      - spelling_corrector.ipynb
      - test.ipynb
      - webscrap.txt
    - Iotmini project video
    - Video for increment 3
    - personal documents
    - tweets
    - video
    - video increment 4
    - AI BASED SURVEY CHATBO...
    - AI BASED SURVEY CHATBO...

Disk 79.35 GB available

```python
[4] import pandas as pd
    import seaborn as sns
    sns.set(style="white", color_codes=True)
    import warnings
    warnings.filterwarnings("ignore")

    #Reading the Data
    data = pd.read_csv('/content/drive/My Drive/Colab Notebooks/CC.csv')
```

```python
[5] #Removing the unrequired columns
    x = data.iloc[:,[1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17]]
    #Filling all the null values with mean
    x=x.apply(lambda x: x.fillna(x.mean()),axis=0)
    print(x.isnull().sum())
```

```
BALANCE                           0
BALANCE_FREQUENCY                 0
PURCHASES                         0
ONEOFF_PURCHASES                  0
INSTALLMENTS_PURCHASES            0
CASH_ADVANCE                      0
PURCHASES_FREQUENCY               0
ONEOFF_PURCHASES_FREQUENCY        0
PURCHASES_INSTALLMENTS_FREQUENCY  0
CASH_ADVANCE_FREQUENCY            0
CASH_ADVANCE_TRX                  0
PURCHASES_TRX                     0
CREDIT_LIMIT                      0
PAYMENTS                          0
MINIMUM_PAYMENTS                  0
PRC_FULL_PAYMENT                  0
TENURE                            0
dtype: int64
```
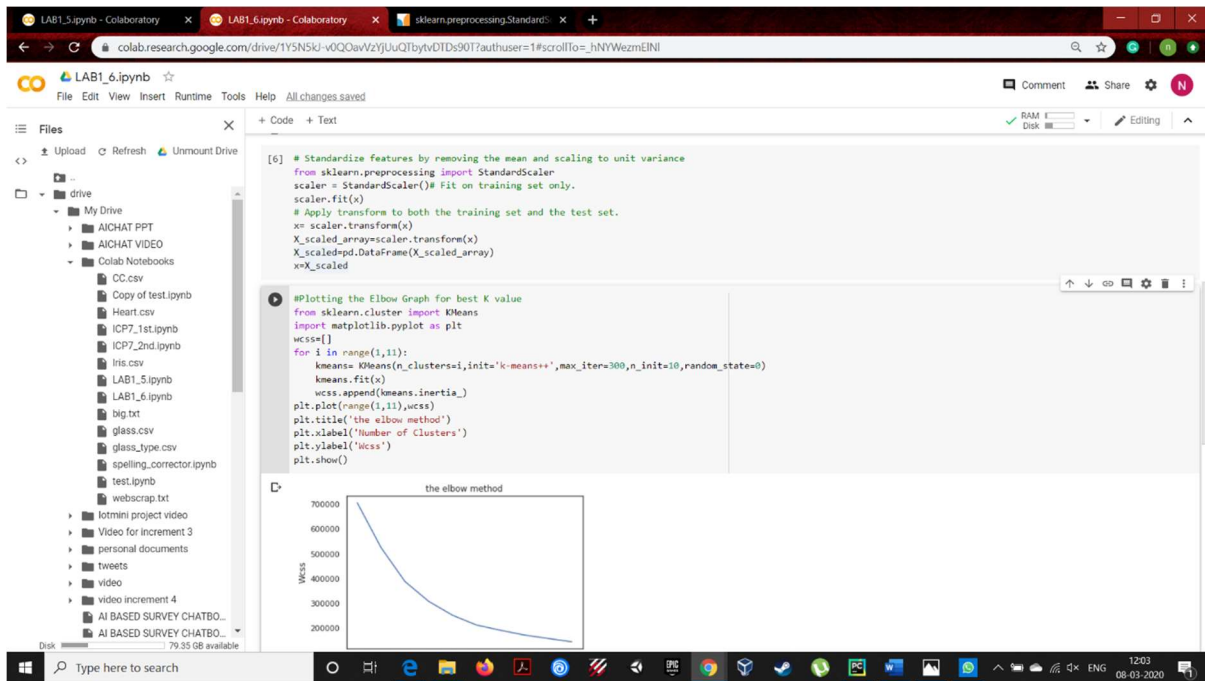
```python
[6] # Standardize features by removing the mean and scaling to unit variance
    from sklearn.preprocessing import StandardScaler
```
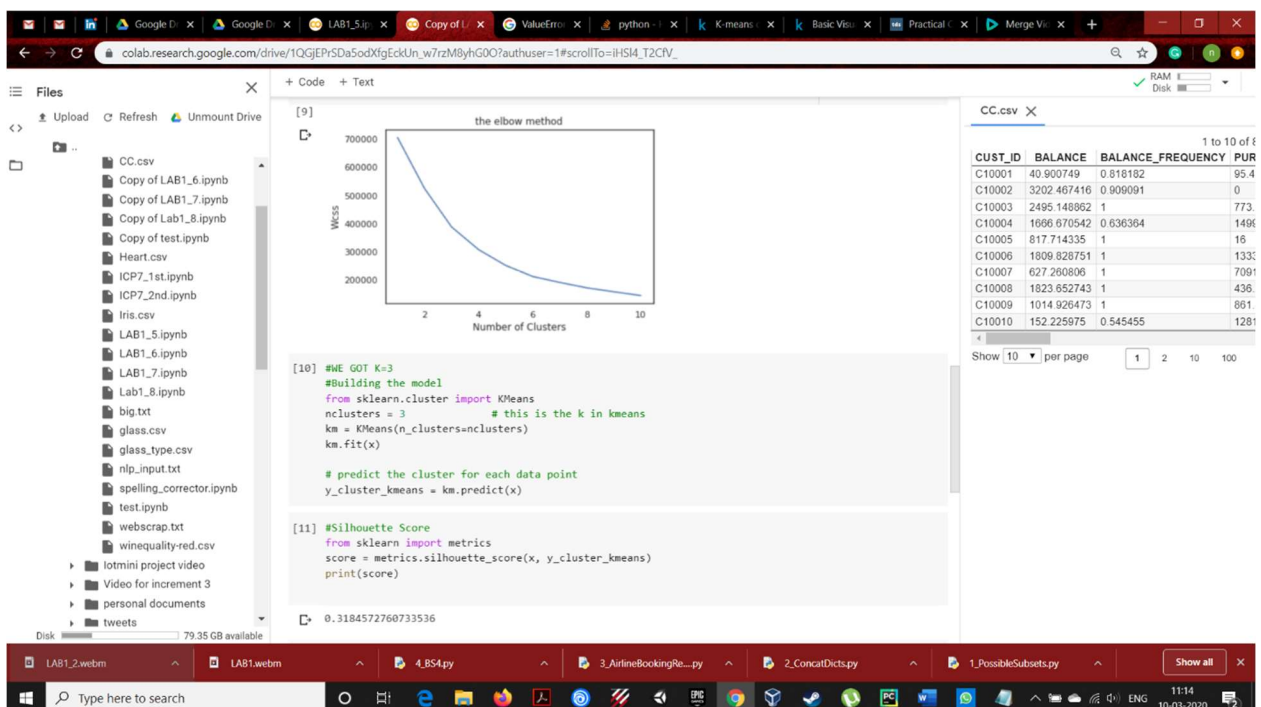
**6b.**

h. plotted an elbow graph for visualizing the elbow where k lies

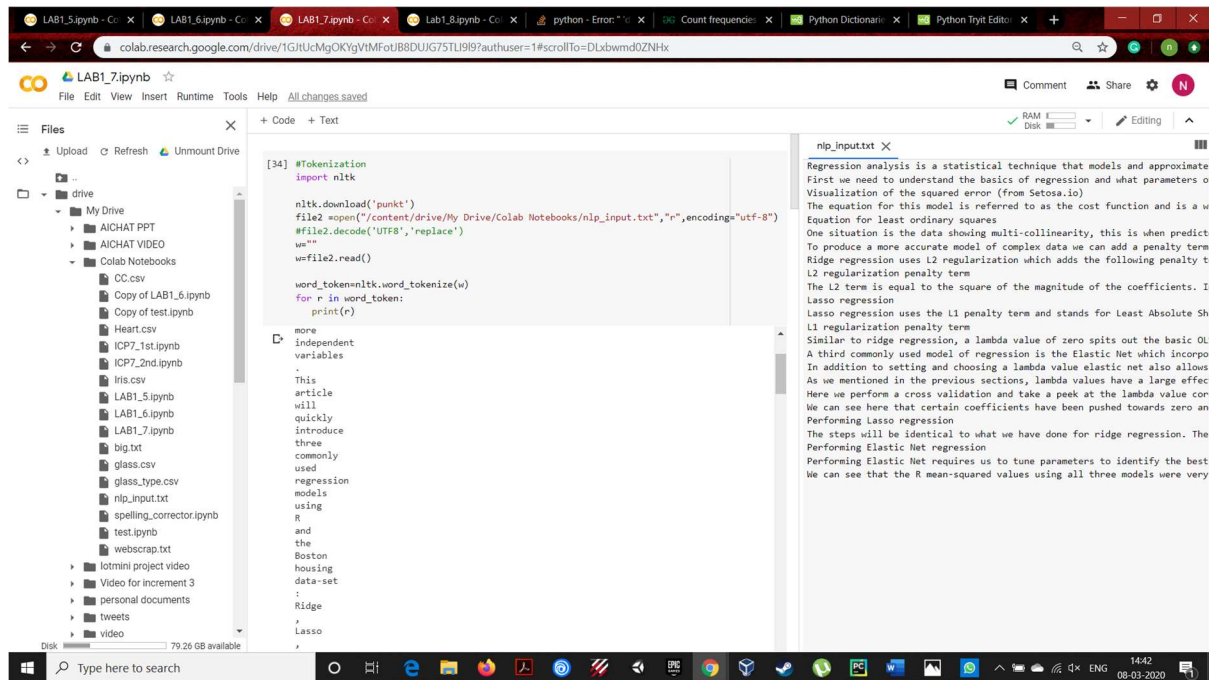i. Fitted a Knn model with the best k value ( which is n=3)

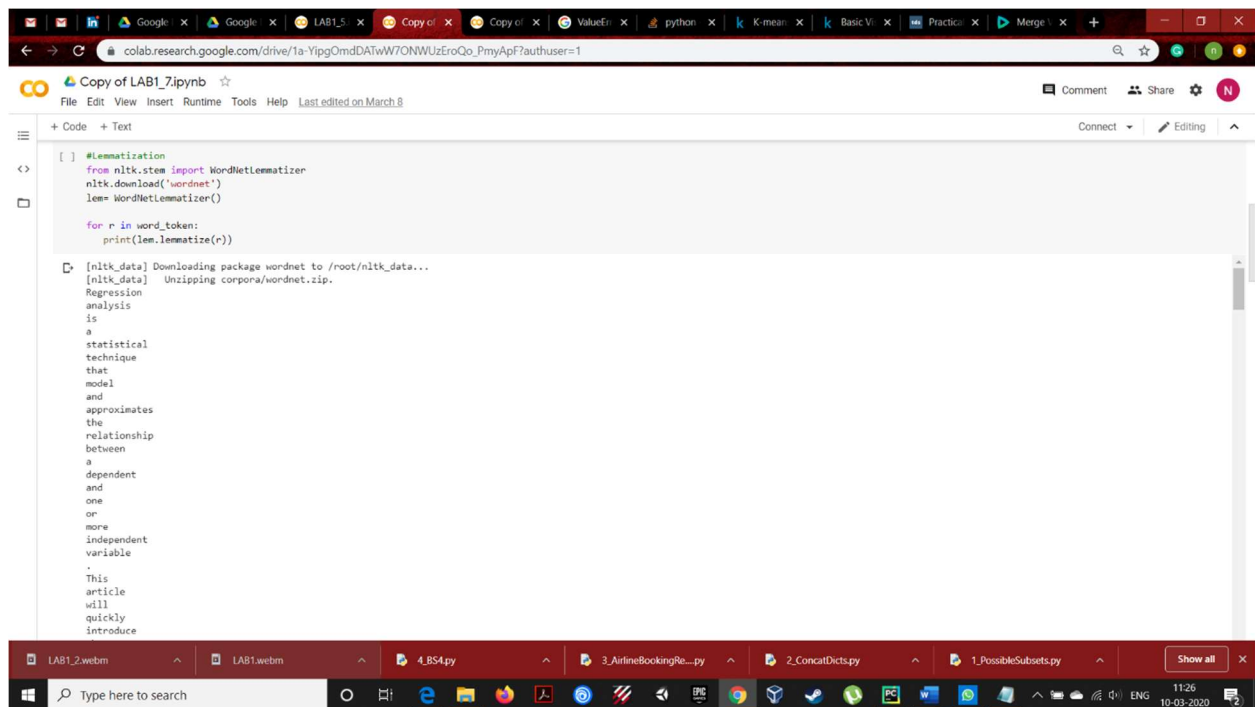j. finally calculated the silhouette score

**7Ans.**

**7a,b**

    a**.** Read the given file nlp_input.txt with encoding UTF-8

    b. Used nltk library

    c. Used word_tokenize function for tokenizing each word

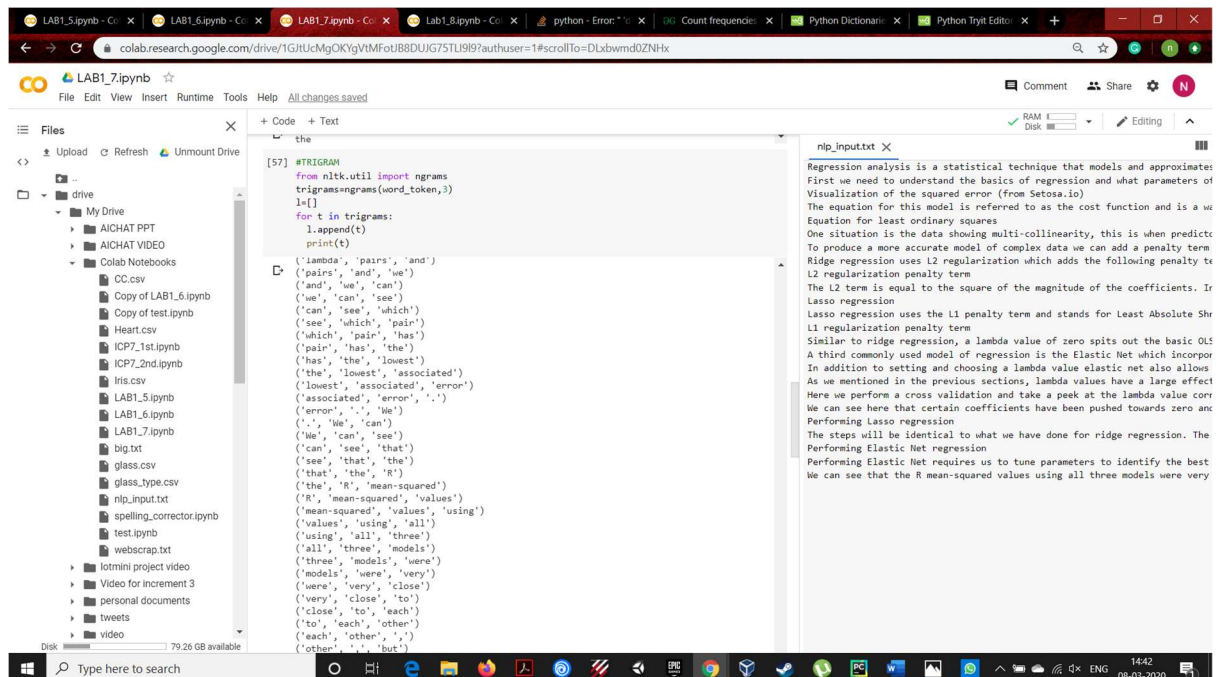    d. Used lemmatize function to apply lemmatization on each word

**7c.**

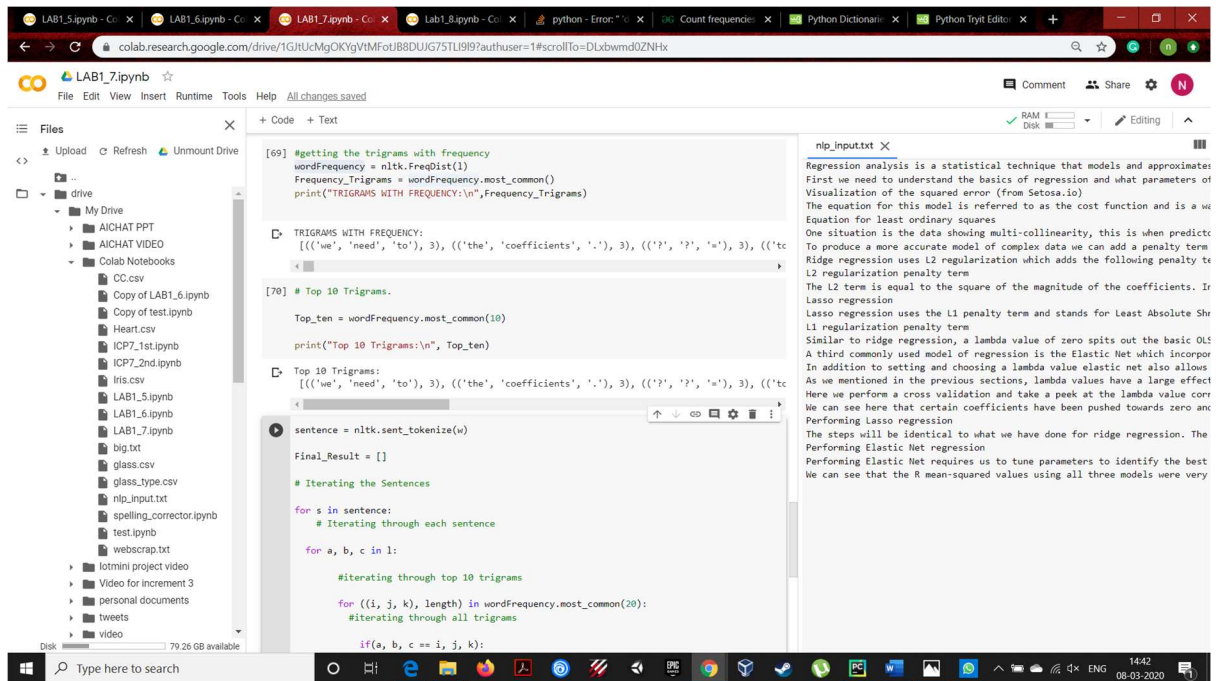> e. Used ngram fuction for n=3 to get all trigrams

> f. Used FreqDist() class to find the all trigrams frequency

**7d.**

g. Used most_common(10) function for extracting the top 10 trigrams based on their frequency



**7e,f,g,h .**

h. Used sent_tokenize to get each sentence

i. Iterated  sentence with trigrams  to get sentence with most trigrams with most frequency

## 8. Multiple Regression (Calculating RMSE, and R2)

a. Used winequality-red.csv file

b. Dropped predicting quality from training data

c. Used LinearRegression class to for fitting the model

d. Calculated the RMSE score and R^2 score (before cleaning the data)

e. Cleaned the data and filled all the null values with value

f. Used fillna to fill all null values with mean value(adding noise)

g. Fitted the model again with cleaned data

h. calculated the RMSE and R^2 score (After cleaning data)

## Note:

Before cleaning the data

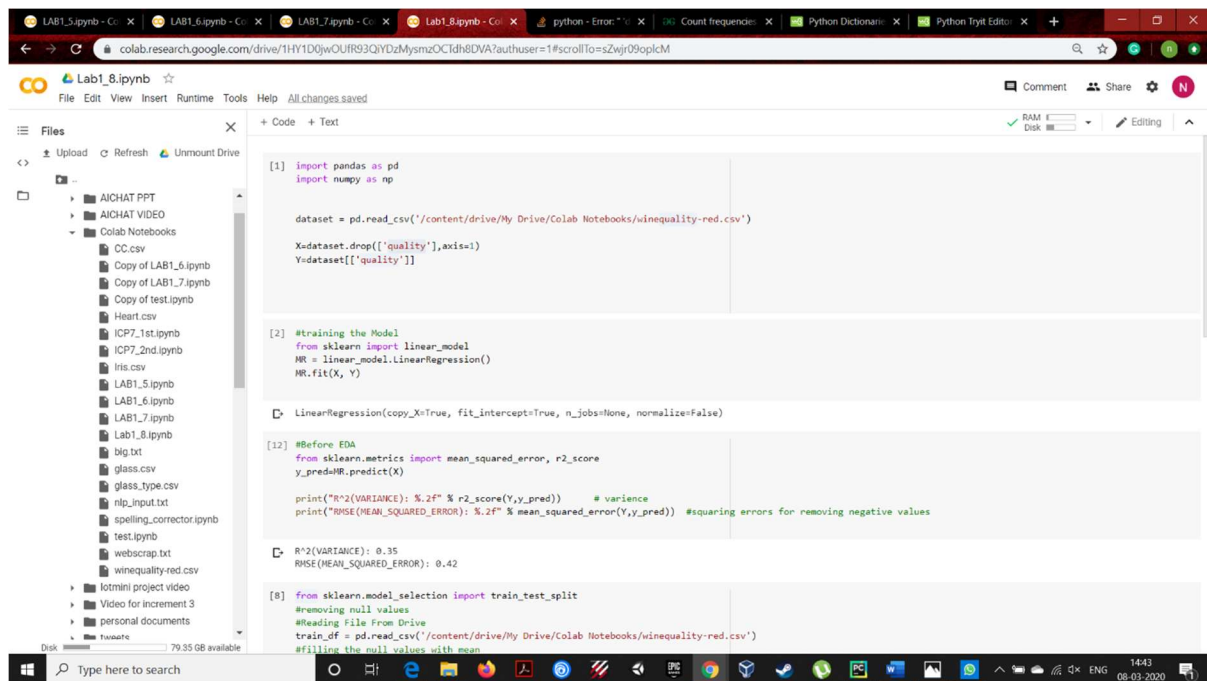**RMSE** : 0.42

 **R^2 score**: 0.35

After Cleaning the data

**RMSE** : 0.42

 **R^2 score**: 0.36

**Lab1_8.ipynb**

File Edit View Insert Runtime Tools Help  All changes saved

Comment   Share

+ Code  + Text

```
[8]  from sklearn.model_selection import train_test_split
     #removing null values
     #Reading File From Drive
     train_df = pd.read_csv('/content/drive/My Drive/Colab Notebooks/winequality-red.csv')
     #filling the null values with mean
     train_df=train_df.apply(lambda x: x.fillna(x.mean()),axis=0)
     print(train_df.isnull().sum())
     #dropping the predicting value from training data
     X_train = train_df.drop("quality",axis=1)
     Y_train = train_df["quality"]

     #using the inbuild function for splitting train data and test data
     X_train, X_test, Y_train, y_test= train_test_split(X_train, Y_train, test_size=0.4, random_state=0)
```

```
fixed acidity         0
volatile acidity      0
citric acid           0
residual sugar        0
chlorides             0
free sulfur dioxide   0
total sulfur dioxide  0
density               0
pH                    0
sulphates             0
alcohol               0
quality               0
dtype: int64
```

```
[9]  #training the Model
     from sklearn import linear_model
     MR = linear_model.LinearRegression()
     MR.fit(X_train, Y_train)
```

```
LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)
```

```
#After EDA
from sklearn.metrics import mean_squared_error, r2_score
y_pred=MR.predict(X_train)
```

Disk ▭▭▭▭    79.35 GB available

**Copy of Lab1_8.ipynb**

File Edit View Insert Runtime Tools Help  Last edited on March 8

Comment   Share

+ Code  + Text

```
Y_train = train_df["quality"]
[ ]
     #using the inbuild function for splitting train data and test data
     X_train, X_test, Y_train, y_test= train_test_split(X_train, Y_train, test_size=0.4, random_state=0)
```

```
fixed acidity         0
volatile acidity      0
citric acid           0
residual sugar        0
chlorides             0
free sulfur dioxide   0
total sulfur dioxide  0
density               0
pH                    0
sulphates             0
alcohol               0
quality               0
dtype: int64
```

```
[ ]  #training the Model
     from sklearn import linear_model
     MR = linear_model.LinearRegression()
     MR.fit(X_train, Y_train)
```

```
LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)
```

```
[ ]  #After EDA
     from sklearn.metrics import mean_squared_error, r2_score
     y_pred=MR.predict(X_train)

     print(" R^2(VARIANCE): %.2f" % r2_score(Y_train,y_pred))       # varience
     print("RMSE(Mean Squared Error): %.2f" % mean_squared_error(Y_train,y_pred))  #squaring errors for removing negative values
```

```
R^2(VARIANCE): 0.36
RMSE(Mean Squared Error): 0.42
```

**DATASETS USED:**

- Glass_type.csv
- Cc.csv

- Winequality-red.csv

**EVALUATION & DISCUSSION:**

1. A set of integers were taken and all the possible subsets without the empty set were printed.
2. Two dictionaries were taken and concatenation and sorting were applied to the merged dictionary.
3. Airline Booking Reservation System was build using the concepts of class, inheritance and method overriding.
4. By using the BeautifulSoup, webscrapping was performed, all the course names and their descriptions were displayed.
5. Take random data set, clean data, split the data, fit and evaluate the accuracy for each of three classifiers (KNN, Naive Bayes, SVM).
6. Take random data set, clean the data, use kmeans function for iterating through all values and get best k value, plot a graph to find elbow, fit model with best k value, and calculate the silhouette score
7. Take the given text file, use the in-build functions to tokenize the data, lemmatize, Find all trigrams using ngrams function where n=3, freqDest() to get frequency, most_common(10) to get top ten trigrams and sentence tokenize to tokenize each sentence and get sentence with words with more frequency of trigrams.
8. Take the random dataset clean the data, remove predicting column from training data set, fit the model, calculate the RMSE, R^2 score, and find the score without cleaning the data.

**CONSLUSION:**

All the programs were successfully executed according to the objectives specified and have met the evaluation criteria.