



Python and Deep Learning Programming

Lab Assignment1

Team Members:

Vishnu vardhan Reddy Manne

Sai Jaswanth Gattidi

Vandith Thotla

Department of Computer Science
University of Missouri, Kansas City

Introduction:

This Lab assignment will help us to get experienced with some of the concepts in python like tuples, dictionaries, OOPS concepts and BeautifulSoup package. This lab work also contains the application of some of the machine learning algorithms like Navie Bayes, SVM, KNN, Multiple regression and K-Means clustering.

Objectives:

The objectives of this Lab are:

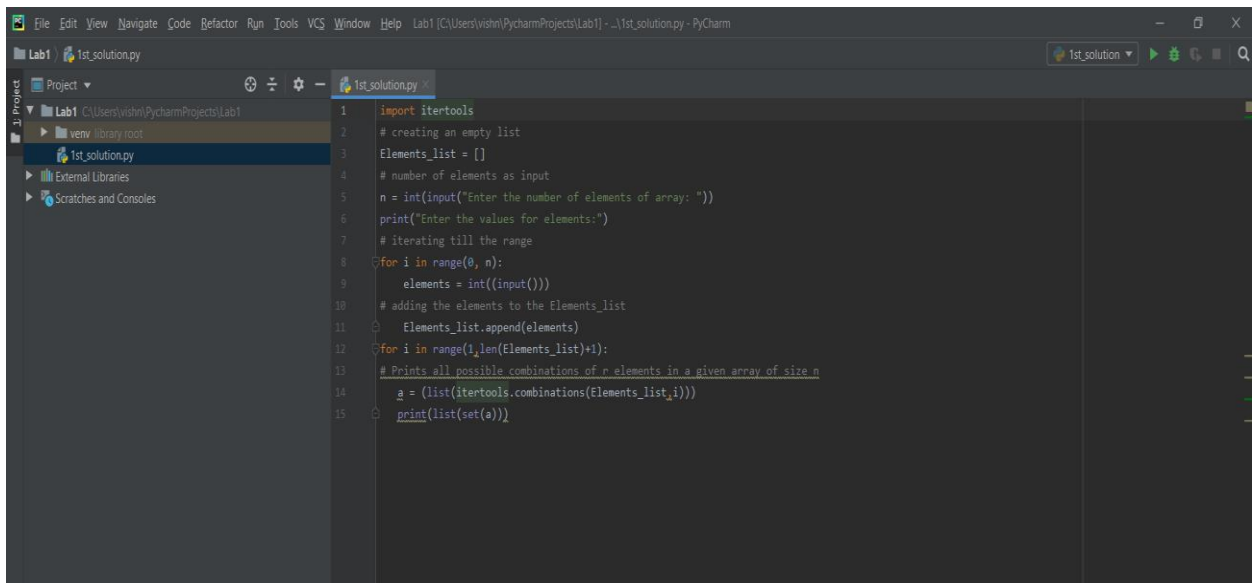
- To get all the possible subsets without any null or duplicate combinations in a given set.
- Creating multiple dictionaries and merging them and then sorting the merged dictionary content based on value.
- Writing a python program to create Airline Booking Reservation system with classes Flight, Person, Employee, Passenger by using Inheritance and Method Overriding concepts
- Using beautifulsoup4 and requests package to get the contents of a webpage and then printing all the course codes and their respective course descriptions from the webpage.
- Picking any dataset that includes both numeric and non-numeric data and performing exploratory data analysis (EDA) on the chosen dataset like handling null values, finding the correlation between the features, replacing the null values with mean, encoding the non-numeric features. Applying the classification Algorithms: Naïve Bayes, KNN, and SVM on the chosen dataset and reporting the classification algorithm with the best result.
- Picking the dataset of our choice and performing exploratory data analysis and then applying the K-means on the dataset and visualizing the clusters using matplotlib or seaborn. Reporting the K value using Elbow method and evaluating the k value using silhouette score.
- To take an input file and read it and apply Tokeniation, Lemmatization and finding the trigrams for the words, top10 most repeated trigrams, extracting sentences with most repeated trigrams and concatenating them and printing the result.

- To pick a dataset and apply Multiple regression and reporting the R squared, and RMSE values before and after performing the Exploratory data analysis on the chosen dataset.

Approaches/Methods:

1.

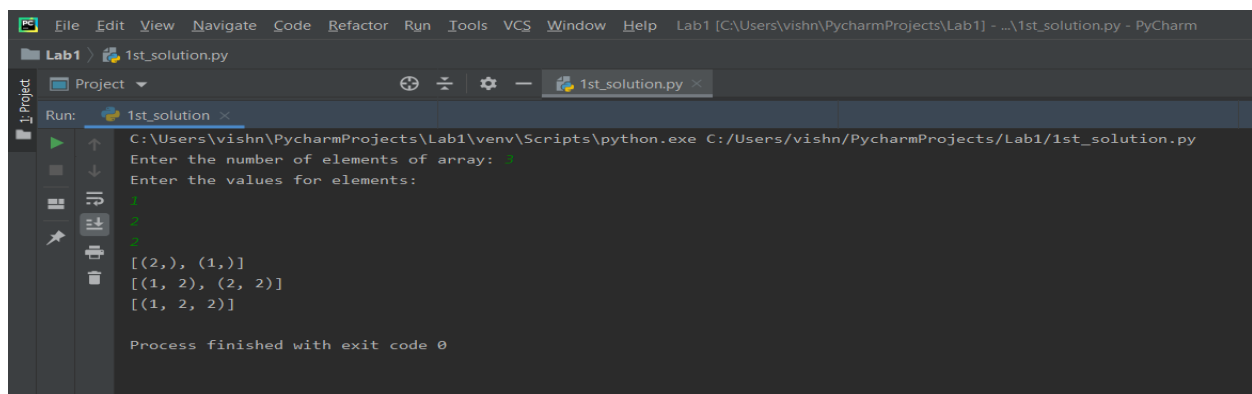
Code:



```

1 import itertools
2 # creating an empty list
3 Elements_list = []
4 # number of elements as input
5 n = int(input("Enter the number of elements of array: "))
6 print("Enter the values for elements:")
7 # iterating till the range
8 for i in range(0, n):
9     elements = int(input())
10    # adding the elements to the Elements_list
11    Elements_list.append(elements)
12    for i in range(1, len(Elements_list)+1):
13        # Prints all possible combinations of r elements in a given array of size n
14        a = (list(itertools.combinations(Elements_list, i)))
15        print(list(set(a)))
  
```

Output:



```

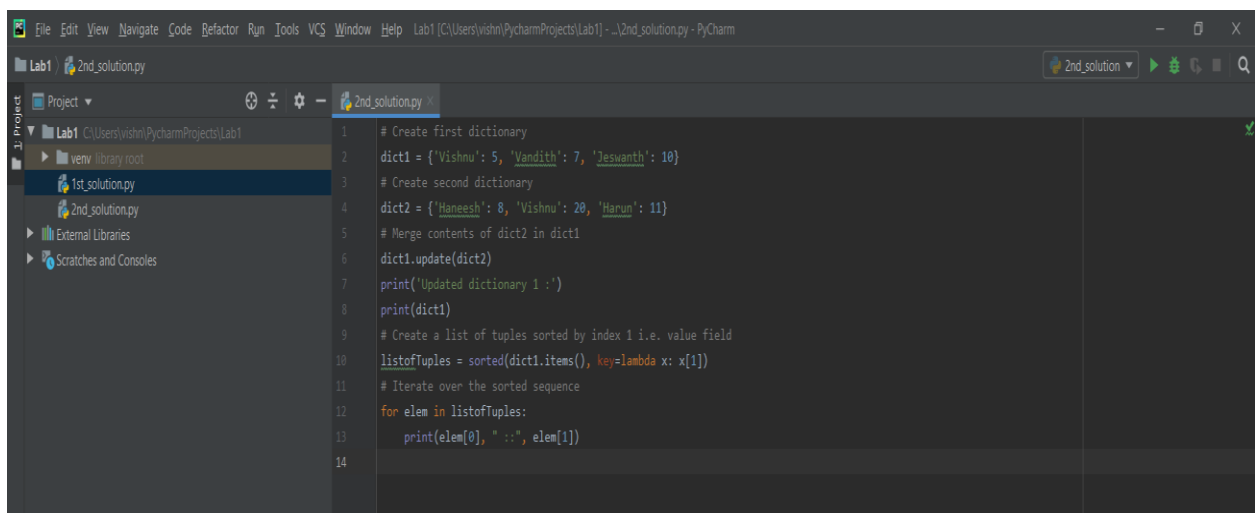
Run: 1st_solution x
C:\Users\vishn\PycharmProjects\Lab1\venv\Scripts\python.exe C:/Users/vishn/PycharmProjects/Lab1/1st_solution.py
Enter the number of elements of array: 2
Enter the values for elements:
1
2
[(2, ), (1, )]
[(1, 2), (2, 2)]
[(1, 2, 2)]
Process finished with exit code 0
  
```

- Here itertools is a module that implements a number of iterator building blocks.

- We create an empty list and then, we take the number of elements as input to the array from the terminal. Then we will append the elements taken as input from the terminal to the empty list created.
- Then we will iterate for i in the range of 0 to n (which is size of the list) and print all the combinations of elements in the given array of size n.
- Here we get some of the duplicate combinations so we use list(set(a)) in order to eliminate them. Then we will get all the possible subsets without any null and duplicate combinations.

2.

Code:

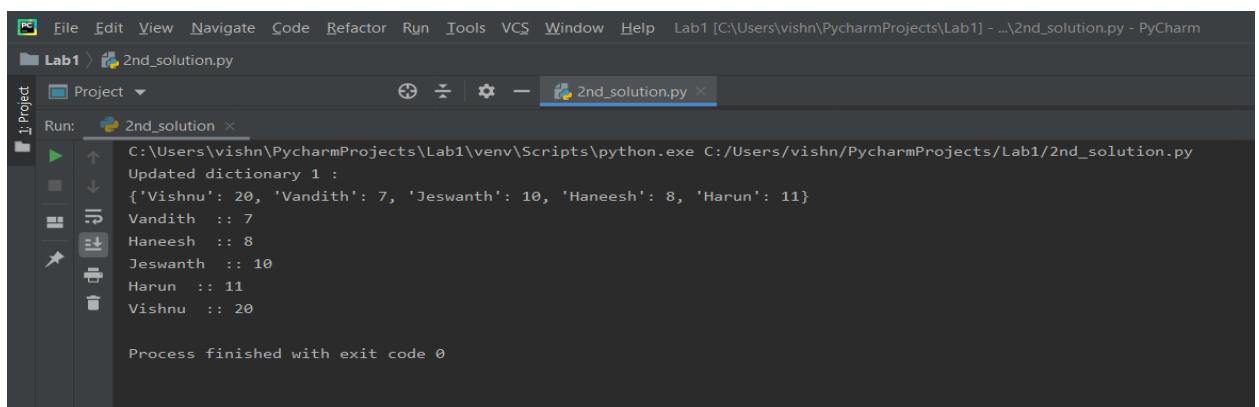


```

1 # Create first dictionary
2 dict1 = {'Vishnu': 5, 'Vandith': 7, 'Jeswanth': 10}
3 # Create second dictionary
4 dict2 = {'Haneesh': 8, 'Vishnu': 20, 'Harun': 11}
5 # Merge contents of dict2 in dict1
6 dict1.update(dict2)
7 print('Updated dictionary 1 :')
8 print(dict1)
9 # Create a list of tuples sorted by index 1 i.e. value field
10 listofTuples = sorted(dict1.items(), key=lambda x: x[1])
11 # Iterate over the sorted sequence
12 for elem in listofTuples:
13     print(elem[0], " :: ", elem[1])
14

```

Output:



```

Run: 2nd_solution
C:\Users\vishn\PycharmProjects\Lab1\venv\Scripts\python.exe C:/Users/vishn/PycharmProjects/Lab1/2nd_solution.py
Updated dictionary 1 :
{'Vishnu': 20, 'Vandith': 7, 'Jeswanth': 10, 'Haneesh': 8, 'Harun': 11}
Vandith :: 7
Haneesh :: 8
Jeswanth :: 10
Harun :: 11
Vishnu :: 20

Process finished with exit code 0

```

- Here we create two dictionaries dict1 & dict2 and we merge the contents of dictionary2(dict2) to dictionary1(dict1) using the update function.

- Then in order to sort the dictionary content based on value we create a sorted list of tuples and then will iterate over this sorted list of values and print the key value pairs in sorted order of values.

3(a).

Code:

```

1 class Airline:
2     def __init__(self, flightNo):
3         Airline.source = input('Enter Src : ')
4         Airline.destination = input('Enter Dest : ')
5         print('Airlines that are available :')
6         print('DELTA')
7         print('SPICEJET')
8         print('INDIGO')
9         print('QATAR')
10        Airline.airlinesName = input('Name the Airlines that you prefer to travel: ')
11        self.flightNo = flightNo
12
13        #This method prints the details of the airlines
14        def print_details(self):
15            print('Airlines : ', Airline.airlinesName)
16            print('Flight Number : ', self.flightNo)
17            print(Airline.source, ' -> ', Airline.destination)
18
19        #Inheritance
20        class Employee(Airline):
21            def __init__(self, employee_id, employee_name, employee_gender):
22                self.employee_name = employee_name
23                self.employee_id = employee_id
24                self.employee_gender = employee_gender
25
26            #Method Overriding
27            def print_details(self):
28                print('Name of employee: ', self.employee_name)
29                print('Employee id: ', self.employee_id)
30                print('Employee gender: ', self.employee_gender)
31
32            #This class inputs all the Traveller details
33            class Traveller:
34                def __init__(self):
35                    Traveller.traveller_fn = input('Enter first name : ')
36                    Traveller.traveller_ln = input('Enter last name : ')
37                    Traveller.traveller_PNo = input('Enter passport number: ')
38                    Traveller.traveller_gender = input('Enter gender : ')
39
40            def __init__(self):
41                Traveller.traveller_PNo = input('Enter passport number: ')
42                Traveller.traveller_gender = input('Enter gender : ')
43                Traveller.traveller_class = input('Business or Economy class? : ')
44
45            #This class is used to calculate the baggage fare
46            class Baggage:
47                def __init__(self):
48                    Baggage.numberOfBags = int(input('Number of bags you want to checkin : '))
49                    Baggage.totalBagfare = 0
50                    Baggage.numberOfBags = Baggage.numberOfBags
51                    if(Baggage.numberOfBags > 3):
52                        for i in range(Baggage.numberOfBags-3):
53                            Baggage.totalBagfare += 80
54                    print('You can take two bags for free !!! Total bag fare is for ', Baggage.numberOfBags, 'is ', Baggage.totalBagfare)
55
56            #This class calculates the ticket cost based on the class, Flight and bags
57            class TicketCost(Baggage, Traveller, Airline):
58                def __init__(self):
59                    TicketCost.baseCost = 250
60                    TicketCost.baseCost = TicketCost.baseCost + Baggage.totalBagfare
61                    if(Traveller.traveller_class == 'business'):
62                        TicketCost.baseCost = TicketCost.baseCost + 250
63                    print('Total ticket cost is : ', TicketCost.baseCost)
64
65            #These details are displayed on the ticket by using this method
66            def ticketDisplay(self):
67                print('Ticket Details')
68                print('*****')
69                print('Traveller Name : ', Traveller.traveller_fn, ' ', Traveller.traveller_ln)
70                print('Traveller Passport Number : ', Traveller.traveller_PNo)
71                print('Gender : ', Traveller.traveller_gender)
72                print('Class : ', Traveller.traveller_class)
73                print('Total number of bags checked in : ', Baggage.numberOfBags)
74                print('Total fare for the trip is : ', TicketCost.baseCost)
75
76            employee = Employee(5555, 'haneesh', 'male')
77            employee.print_details()
78
79            #Traveller = Traveller()
80            Airline.__init__()
  
```

```

lab3.py
48 class TicketCost(Airline):
49     def __init__(self):
50         TicketCost.baseCost = 250
51         TicketCost.baseCost = TicketCost.baseCost + Baggage.totalBagfare
52         if(Traveller.traveller_class == "business"):
53             TicketCost.baseCost = TicketCost.baseCost + 250
54         print("Total ticket cost is : ", TicketCost.baseCost)
55         #These details are displayed on the ticket by using this method
56     def ticketDisplay(self):
57         print("Ticket Details")
58         print("=====")
59         print("Traveller Name : ", Traveller.traveller_fn, " ", Traveller.traveller_ln)
60         print("Traveller Passport Number : ", Traveller.traveller_PNo)
61         print("Gender : ", Traveller.traveller_gender)
62         print("Class : ", Traveller.traveller_class)
63         print("Total number of bags checked in : ", Baggage.numberofBags)
64         print("Total Fare for the trip is : ", TicketCost.baseCost)
65 employee = Employee(5555, "haneesh", "male")
66 employee.print_details()
67 flight = Airline("RX1006")
68 traveller = Traveller()
69 bags = Baggage()
70 ticket = TicketCost()
71 ticket.ticketDisplay()
72 flight.print_details()
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100

```

Output:

```

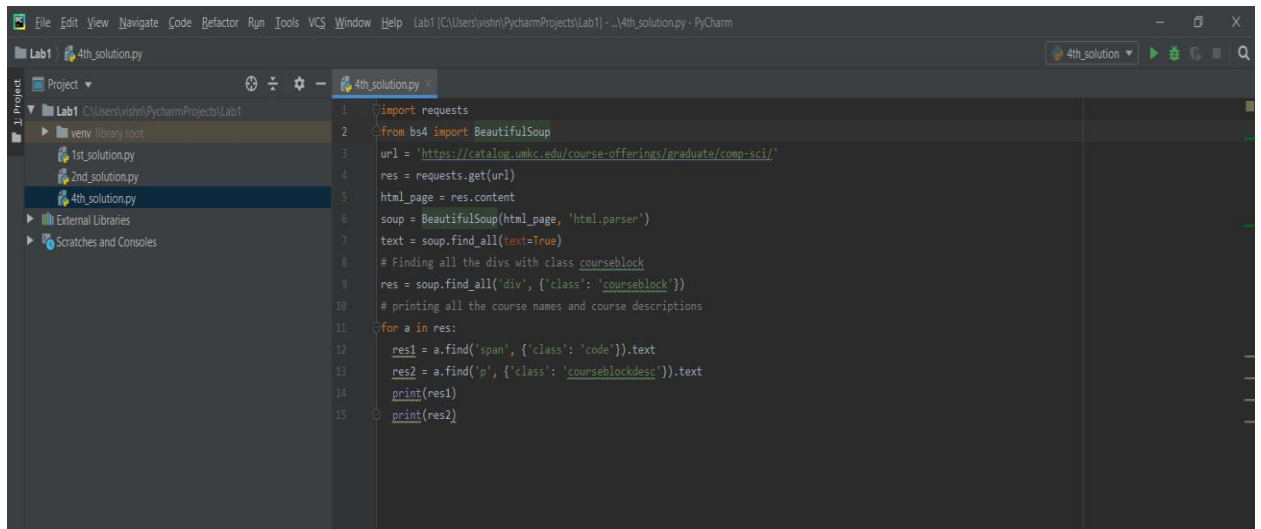
Run
lab3
C:\Users\haneesh\AppData\Local\Temp\3_inheritance.py - C:\Users\haneesh\Downloads\lab3.py - PyCharm
Name of employee: haneesh
Employee id: 5555
Employee gender: male
Enter Src : hyderabad
Enter Dest : vijayavada
Airlines that are available :
DELTA
SPICEJET
INDIGO
QATAR
Name the Airlines that you prefer to travel: DELTA
Enter first name : sai jashwanth
Enter last name : gattidi
Enter passport number: RS12345
Enter gender : male
Business or Economy class?: business
Number of bags you want to checkin : 4
You can take two bags for free !!! Total bag fare is for 4 is 80
Total ticket cost is : 330
Ticket Details
=====
Traveller Name : SAI JASHWANTH GATTIDI
Traveller Passport Number : RS12345
Gender : MALE
Class : Business
Total number of bags checked in : 4
Total Fare for the trip is : 330
Airlines : DELTA
Flight Number : RX1006
vijayavada --> hyderabad
Process finished with exit code 0

```

- We have created an Airline class that will take input from the terminal like source, destination, and preferred airlines then airline class will generate the flight number for the provided details.
- Employee class will display the details of the employee. Employee class will use inheritance concept and overrides the printing method of the Airline class.
- Traveller class will take all the details of the traveller. Baggage class would take the number of check-in bags from the traveller and calculate the baggage fare. TicketCost class takes the inputs from all the classes and calculates the fare of the ticket. A method in TicketCost would display all the ticket details.

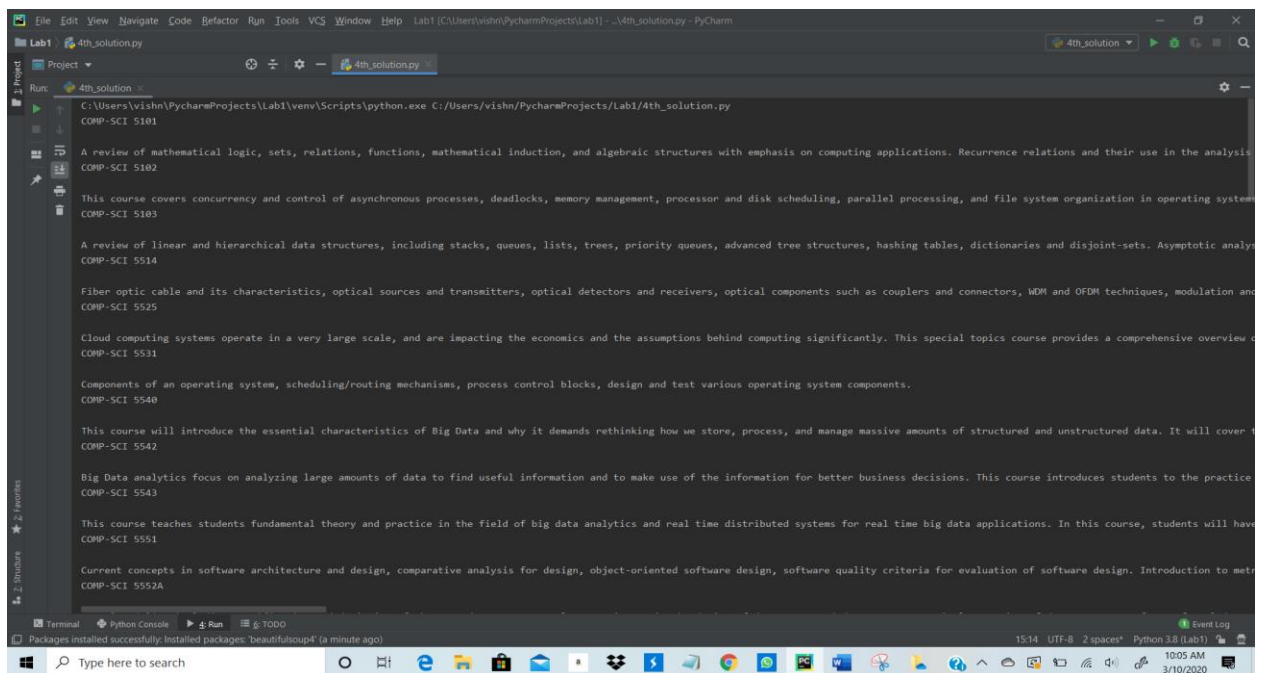
4.

Code:

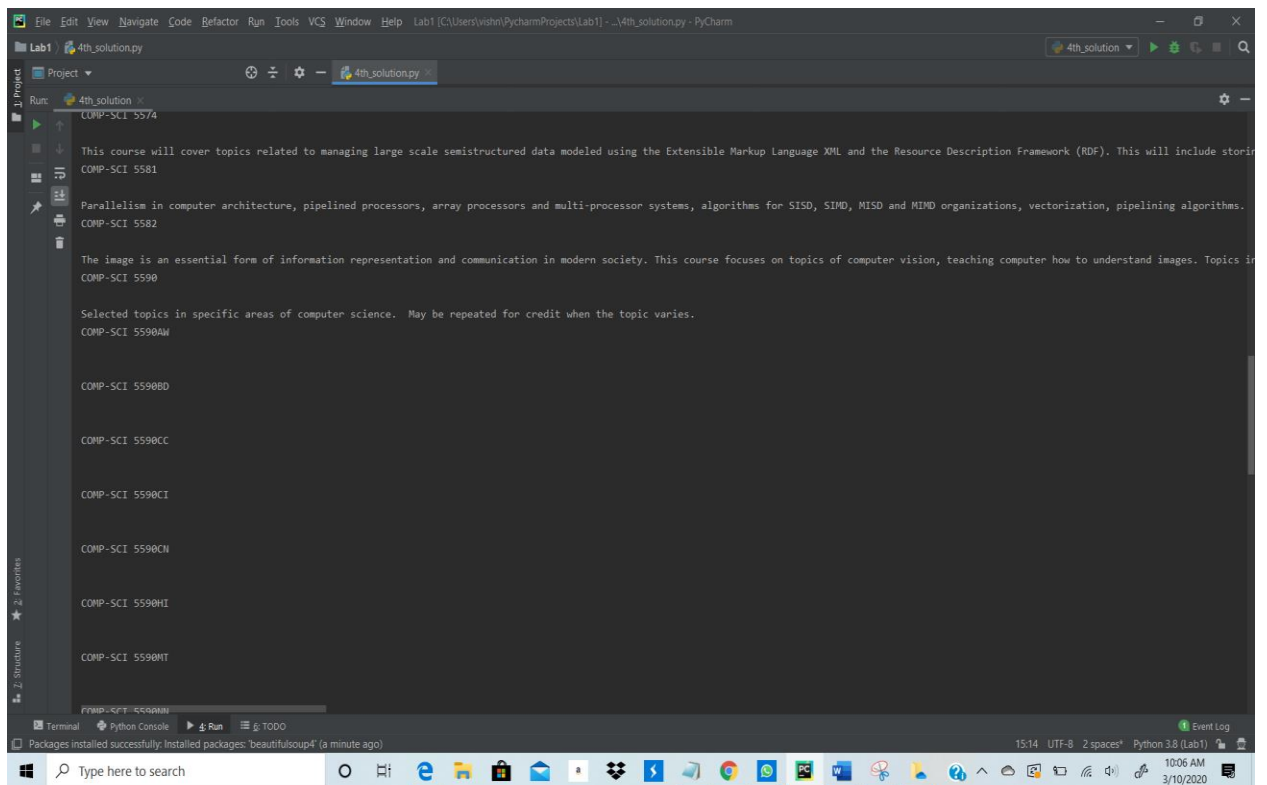
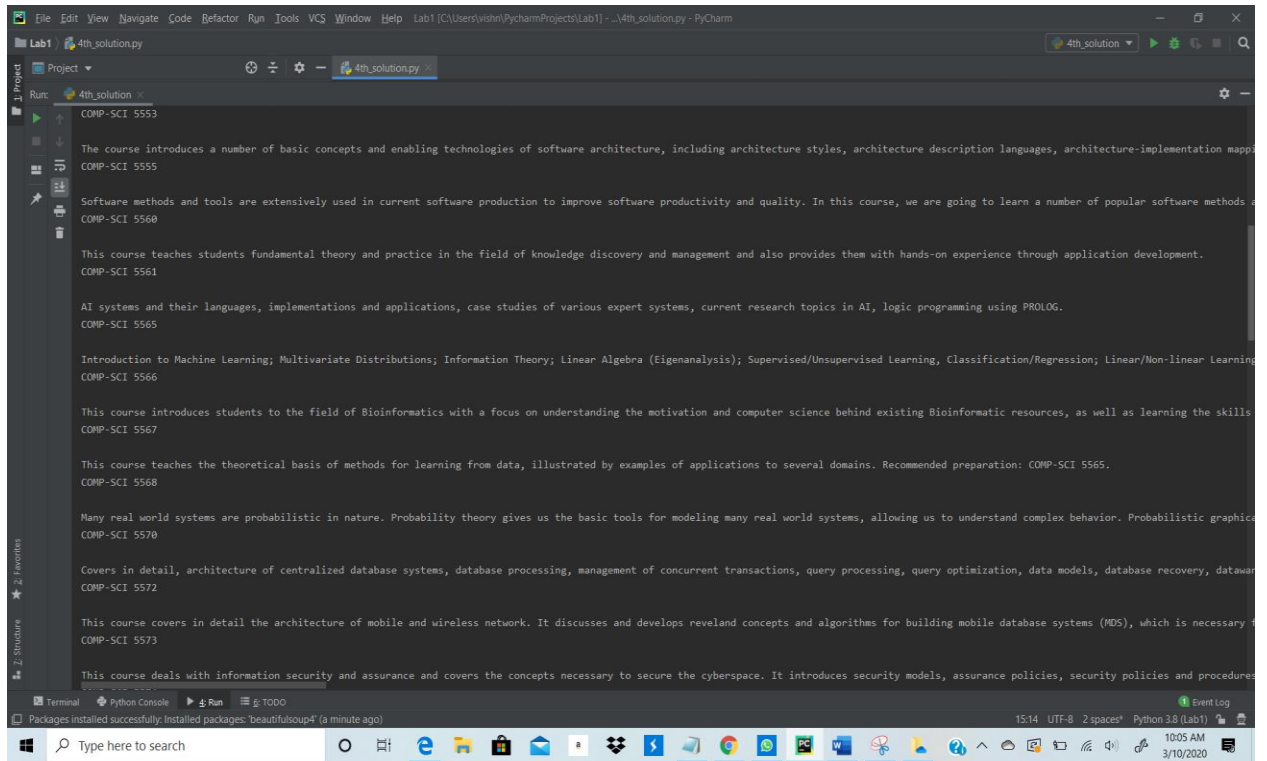


```
1 import requests
2 from bs4 import BeautifulSoup
3 url = 'https://catalog.unm.edu/course-offerings/graduate/comp-sci/'
4 res = requests.get(url)
5 html_page = res.content
6 soup = BeautifulSoup(html_page, 'html.parser')
7 text = soup.find_all(text=True)
8 # Finding all the divs with class courseblock
9 res = soup.find_all('div', {'class': 'courseblock'})
10 # printing all the course names and course descriptions
11 for a in res:
12     res1 = a.find('span', {'class': 'code'}).text
13     res2 = a.find('p', {'class': 'courseblockdesc'}).text
14     print(res1)
15     print(res2)
```

Output:



```
Run: 4th_solution.py
C:\Users\vishn\P\chamProjects\Lab1\venv\Scripts\python.exe C:\Users\vishn\PchamProjects\Lab1\4th_solution.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 analysis
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 systems.
COMP-SCI 5103
A review of linear and hierarchical data structures, including stacks, queues, lists, priority queues, advanced tree structures, hashing tables, dictionaries and disjoint-sets. Asymptotic analysis.
COMP-SCI 5114
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 and
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 overview of
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 cover
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 practice
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 have
COMP-SCI 5551
Current concepts in software architecture and design, comparative analysis for design, object-oriented software design, software quality criteria for evaluation of software design. Introduction to met
COMP-SCI 5552A
```




```

Application of the algorithmic techniques learned in COMP-SCI 5596A to provide suitable security countermeasures to the variety of security threats across the spectrum of computing.
COMP-SCI 5597
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
  
```

- Using BeautifulSoup we will first get the content of the webpage and converted all the html content to plain text.
- Then we have inspected the webpage and we get to know that courseblock represents the subject name, code represents the subject code and courseblockdesc gives the description of the subject.
- Then we find all the div tags with class:courseblock. Then we have printed all the course code and their respective descriptions from the given URL.

5(a).

Code & Output:

```

from google.colab import drive
drive.mount('/content/drive/')

Go to this URL in a browser: https://accounts.google.com/o/oauth2/auth?client_id=947318989803-6b6gk8d9df4d4g3pfee6491hc@rc4i.apps.googleusercontent.com&redirect_uri=urn%3Aietf%3Awg%3A...
Enter your authorization code:
.....
Mounted at /content/drive/

[ ] import warnings
warnings.filterwarnings("ignore")
import pandas as pd
import numpy as np
data=pd.read_csv('drive/My Drive/Lab1/cars.csv', delimiter=',', header=None, skiprows=1, names=['mpg','cylinders','cubicinches','hp','weightlbs','time-to-60','year','brand'])

# Nulls Handling
nulls = pd.DataFrame(data.isnull().sum().sort_values(ascending=False))
nulls.columns = ['Features']
nulls.index.name = 'Nulls_count'
print(nulls)
x = data.select_dtypes(include=[np.number]).interpolate().dropna()
print(sum(x.isnull().sum() != 0))

# Encoding non-numeric features
from sklearn.preprocessing import LabelEncoder
data = data.apply(LabelEncoder().fit_transform)

# Here we are filling the null values with mean value
data=data.apply(lambda x: x.fillna(x.mean()),axis=0)
print(data["brand"])
print(data.isnull().sum())
  
```

colab.research.google.com/drive/1pyRUESb4MQnj_VY8jTelo6yy6X9UfUez

Copy of 5th_Solution.ipynb

File Edit View Insert Runtime Tools Help Last saved at 10:41 AM

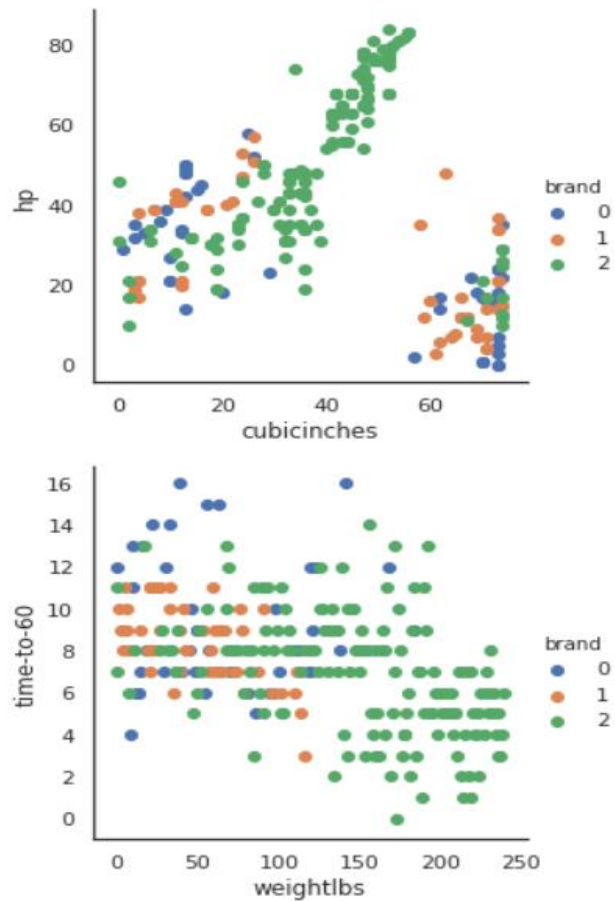
Code + Text

```
[ ] # Visualize data to analyze our features correlations
import seaborn as sns
sns.set(style="white", color_codes=True)
import matplotlib.pyplot as plt
sns.FacetGrid(data, hue='brand', height=4).map(plt.scatter, 'mpg', 'cylinders').add_legend()
plt.show()
sns.FacetGrid(data, hue='brand', height=4).map(plt.scatter, 'cubicinches', 'hp').add_legend()
plt.show()
sns.FacetGrid(data, hue='brand', height=4).map(plt.scatter, 'weightlbs', 'time-to-60').add_legend()
plt.show()

# Split data into train and test
from sklearn.model_selection import train_test_split
x = data.drop(['brand'], axis=1)
y = data['brand']
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=0)
```

Features

Nulls_count	brand
0	0
0	1
0	2
0	3
0	4
0	5
0	6
0	7
0	8
0	9
0	10
0	11
0	12
0	13
0	14
0	15
0	16
0	17
0	18
0	19
0	20
0	21
0	22
0	23
0	24
0	25
0	26
0	27
0	28
0	29
0	30
0	31
0	32
0	33
0	34
0	35
0	36
0	37
0	38
0	39
0	40
0	41
0	42
0	43
0	44
0	45
0	46
0	47
0	48
0	49
0	50
0	51
0	52
0	53
0	54
0	55
0	56
0	57
0	58
0	59
0	60
0	61
0	62
0	63
0	64
0	65
0	66
0	67
0	68
0	69
0	70
0	71
0	72
0	73
0	74
0	75
0	76
0	77
0	78
0	79
0	80
0	81
0	82
0	83
0	84
0	85
0	86
0	87
0	88
0	89
0	90
0	91
0	92
0	93
0	94
0	95
0	96
0	97
0	98
0	99
0	100
0	101
0	102
0	103
0	104
0	105
0	106
0	107
0	108
0	109
0	110
0	111
0	112
0	113
0	114
0	115
0	116
0	117
0	118
0	119
0	120
0	121
0	122
0	123
0	124
0	125
0	126
0	127
0	128
0	129
0	130
0	131
0	132
0	133
0	134
0	135
0	136
0	137
0	138
0	139
0	140
0	141
0	142
0	143
0	144
0	145
0	146
0	147
0	148
0	149
0	150
0	151
0	152
0	153
0	154
0	155
0	156
0	157
0	158
0	159
0	160
0	161
0	162
0	163
0	164
0	165
0	166
0	167
0	168
0	169
0	170
0	171
0	172
0	173
0	174
0	175
0	176
0	177
0	178
0	179
0	180
0	181
0	182
0	183
0	184
0	185
0	186
0	187
0	188
0	189
0	190
0	191
0	192
0	193
0	194
0	195
0	196
0	197
0	198
0	199
0	200
0	201
0	202
0	203
0	204
0	205
0	206
0	207
0	208
0	209
0	210
0	211
0	212
0	213
0	214
0	215
0	216
0	217
0	218
0	219
0	220
0	221
0	222
0	223
0	224
0	225
0	226
0	227
0	228
0	229
0	230
0	231
0	232
0	233
0	234
0	235
0	236
0	237
0	238
0	239
0	240
0	241
0	242
0	243
0	244
0	245
0	246
0	247
0	248
0	249
0	250
0	251
0	252
0	253
0	254
0	255
0	256
0	257
0	258
0	259
0	260
0	261
0	262
0	263
0	264
0	265
0	266
0	267
0	268
0	269
0	270
0	271
0	272
0	273
0	274
0	275
0	276
0	277
0	278
0	279
0	280
0	281
0	282
0	283
0	284
0	285
0	286
0	287
0	288
0	289
0	290
0	291
0	292
0	293
0	294
0	295
0	296
0	297
0	298
0	299
0	300
0	301
0	302
0	303
0	304
0	305
0	306
0	307
0	308
0	309
0	310
0	311
0	312
0	313
0	314
0	315
0	316
0	317
0	318
0	319
0	320
0	321
0	322
0	323
0	324
0	325
0	326
0	327
0	328
0	329
0	330
0	331
0	332
0	333
0	334
0	335
0	336
0	337
0	338
0	339
0	340
0	341
0	342
0	343
0	344
0	345
0	346
0	347
0	348
0	349
0	350
0	351
0	352
0	353
0	354
0	355
0	356
0	357
0	358
0	359
0	360
0	361
0	362
0	363
0	364
0	365
0	366
0	367
0	368
0	369
0	370
0	371
0	372
0	373
0	374
0	375
0	376
0	377
0	378
0	379
0	380
0	381
0	382
0	383
0	384
0	385
0	386
0	387
0	388
0	389
0	390
0	391
0	392
0	393
0	394
0	395
0	396
0	397
0	398
0	399
0	400
0	401
0	402
0	403
0	404
0	405
0	406
0	407
0	408
0	409
0	410
0	411
0	412
0	413
0	414
0	415
0	416
0	417
0	418
0	419
0	420
0	421
0	422
0	423
0	424
0	425
0	426
0	427
0	428
0	429
0	430
0	431
0	432
0	433
0	434
0	435
0	436
0	437
0	438
0	439
0	440
0	441
0	442
0	443
0	444
0	445
0	446
0	447
0	448
0	449
0	450
0	451
0	452
0	453
0	454
0	455
0	456
0	457
0	458
0	459
0	460
0	461
0	462
0	463
0	464
0	465
0	466
0	467
0	468
0	469
0	470
0	471
0	472
0	473
0	474
0	475
0	476
0	477
0	478
0	479
0	480
0	481
0	482
0	483
0	484
0	485
0	486
0	487
0	488
0	489
0	490
0	491
0	492
0	493
0	494
0	495
0	496
0	497
0	498
0	499
0	500
0	501
0	502
0	503
0	504
0	505
0	506
0	507
0	508
0	509
0	510
0	511
0	512
0	513
0	514
0	515
0	516
0	517
0	518
0	519
0	520
0	521
0	522
0	523
0	524
0	525
0	526



- In this program we will first import the required libraries.
 - We have taken cars.csv as the sample dataset which contains both the numeric and non-numeric data. We have applied exploratory data analysis and cleaned the dataset by handling all the null values by filling them with the mean values and encoded the non-numeric features using Label encoder.
- Then we have visualized all the correlation features and plotted them using seaborn. Then by using Train_test_split method we have divided the dataset into train and test datasets with test dataset of size = 0.2 of the cars.csv dataset.

5(b)

Code & Output:

```
colab.research.google.com/drive/1pyRUsb4MQnj_VY8jTelo6yy6X9UFUez#scrollTo=fQV_U1Vjnh4M
```

Copy of 5th_Solution.ipynb

File Edit View Insert Runtime Tools Help Last saved at 10:41 AM

+ Code + Text

```
[ ] # KNN method
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import classification_report
knn = KNeighborsClassifier(n_neighbors=3)
knn.fit(x_train, y_train)
# Evaluate model
score = round(knn.score(x_train, y_train)* 100, 2)
print('K-Neighbors accuracy training score: ', score)
print('Classification report:')
y_pred = knn.predict(x_test)
print(classification_report(y_test, y_pred))
```

K-Neighbors accuracy training score: 85.1

Classification report:

	precision	recall	f1-score	support
0	0.58	0.88	0.70	8
1	0.75	0.50	0.60	12
2	0.88	0.88	0.88	33
accuracy			0.79	53
macro avg	0.74	0.75	0.73	53
weighted avg	0.81	0.79	0.79	53

```
colab.research.google.com/drive/1pyRUsb4MQnj_VY8jTelo6yy6X9UFUez#scrollTo=fQV_U1Vjnh4M
```

Copy of 5th_Solution.ipynb

File Edit View Insert Runtime Tools Help Last saved at 10:41 AM

+ Code + Text

```
[ ] # Naive Bayes method
from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import classification_report
nb = GaussianNB()
nb.fit(x_train, y_train)
# Evaluate model
score = round(nb.score(x_train, y_train)* 100, 2)
print('Naive Bayes accuracy training score: ', score)
print('Classification report:')
y_pred = nb.predict(x_test)
print(classification_report(y_test, y_pred))
```

Naive Bayes accuracy training score: 72.12

Classification report:

	precision	recall	f1-score	support
0	0.40	0.75	0.52	8
1	0.64	0.58	0.61	12
2	0.89	0.73	0.80	33
accuracy			0.70	53
macro avg	0.64	0.69	0.64	53
weighted avg	0.76	0.70	0.71	53

```
colab.research.google.com/drive/1pyRUsb4MQnj_VY8jTelo6yy6X9UFUez#scrollTo=fQV_U1Vjnh4M
```

Copy of 5th_Solution.ipynb

File Edit View Insert Runtime Tools Help Last saved at 10:41 AM

+ Code + Text

```
[ ] # SVM method
from sklearn.svm import SVC, LinearSVC
svc = SVC()
svc.fit(x_train, y_train)
# Evaluate model
score = round(svc.score(x_train, y_train)* 100, 2)
print('Support Vector Machines score: ', score)
print('Classification report:')
y_pred = svc.predict(x_test)
print(classification_report(y_test, y_pred))
```

Support Vector Machines score: 70.19

Classification report:

	precision	recall	f1-score	support
0	0.00	0.00	0.00	8
1	0.45	0.42	0.43	12
2	0.74	0.94	0.83	33
accuracy			0.68	53
macro avg	0.40	0.45	0.42	53
weighted avg	0.56	0.68	0.61	53

- After Splitting the dataset into train and test datasets we have applied all the three classification algorithms: K-nearest neighbors, Naïve bayes, and support vector machines and fitted the models with the training data.
- Calculated the accuracy score for each of the classification algorithm and found that 'k-Nearest neighbors' with k=3 has the highest accuracy 85.1 when compared with other two classification algorithms 'support vector machines' (70.10), and 'Naïve bayes' (72.12)

6(a).

Code & Output:

The first screenshot shows the initial code in the Colab notebook. It includes imports for Google Drive, pandas, numpy, and sklearn. The code mounts the Google Drive, reads a CSV file, handles nulls, and encodes non-numeric features.

```
[1] from google.colab import drive
drive.mount('/content/drive')

Go to this URL in a browser: https://accounts.google.com/o/oauth2/auth?client_id=947318989803-6b60k8d9d4d3pfee6491hc@rc4i.apps.googleusercontent.com&redirect_uri=urn%3Aietf%3Awg%3A...

Enter your authorization code:
.....
Mounted at /content/drive

import warnings
warnings.filterwarnings("ignore")
import pandas as pd
import numpy as np
data = pd.read_csv('/content/drive/My Drive/Lab1/cars.csv', delimiter=',', header=None, skiprows=1, names=['mpg', 'cylinders', 'cubicinches', 'hp', 'weightlbs', 'time-to-60', 'year', 'brand'])
# Nulls Handling
nulls = pd.DataFrame(data.isnull().sum().sort_values(ascending=False))
nulls.columns = ['Features']
nulls.index.name = 'Nulls_count'
print(nulls)
x = data.select_dtypes(include=[np.number]).interpolate().dropna()
print(sum(x.isnull().sum() != 0))

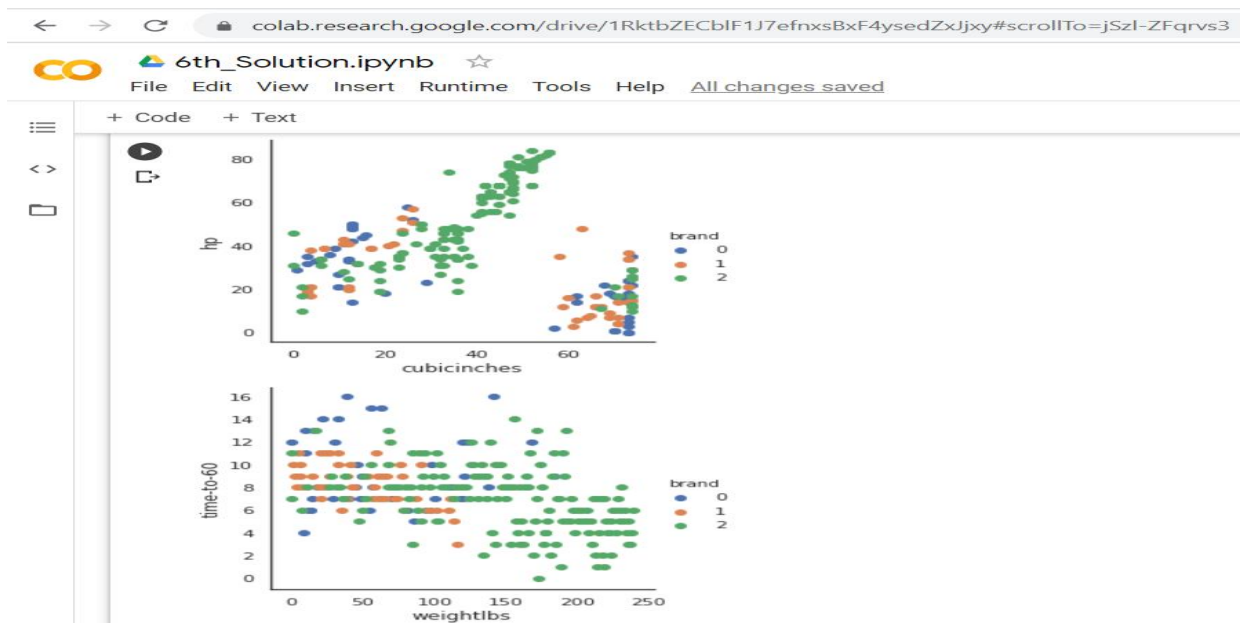
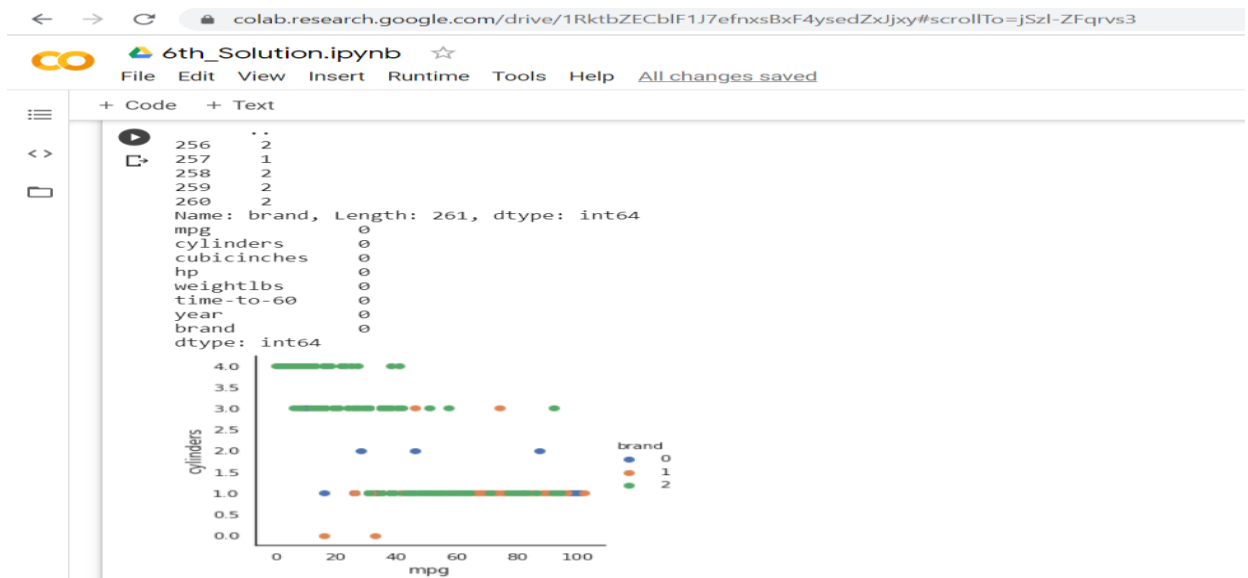
# Encoding non-numeric features
from sklearn.preprocessing import LabelEncoder
data = data.apply(LabelEncoder().fit_transform)
```

The second screenshot shows the continuation of the code. It includes filling null values with the mean, visualizing data with faceted plots, and applying the k-means algorithm.

```
# Here we are filling the null values with mean value
data=data.apply(lambda x: x.fillna(x.mean()),axis=0)
print(data["brand"])
print(data.isnull().sum())

# Visualize data to analyze our features correlations
import seaborn as sns
sns.set(style="white", color_codes=True)
import matplotlib.pyplot as plt
sns.FacetGrid(data, hue='brand', height=4).map(plt.scatter, 'mpg', 'cylinders').add_legend()
plt.show()
sns.FacetGrid(data, hue='brand', height=4).map(plt.scatter, 'cubicinches', 'hp').add_legend()
plt.show()
sns.FacetGrid(data, hue='brand', height=4).map(plt.scatter, 'weightlbs', 'time-to-60').add_legend()
plt.show()

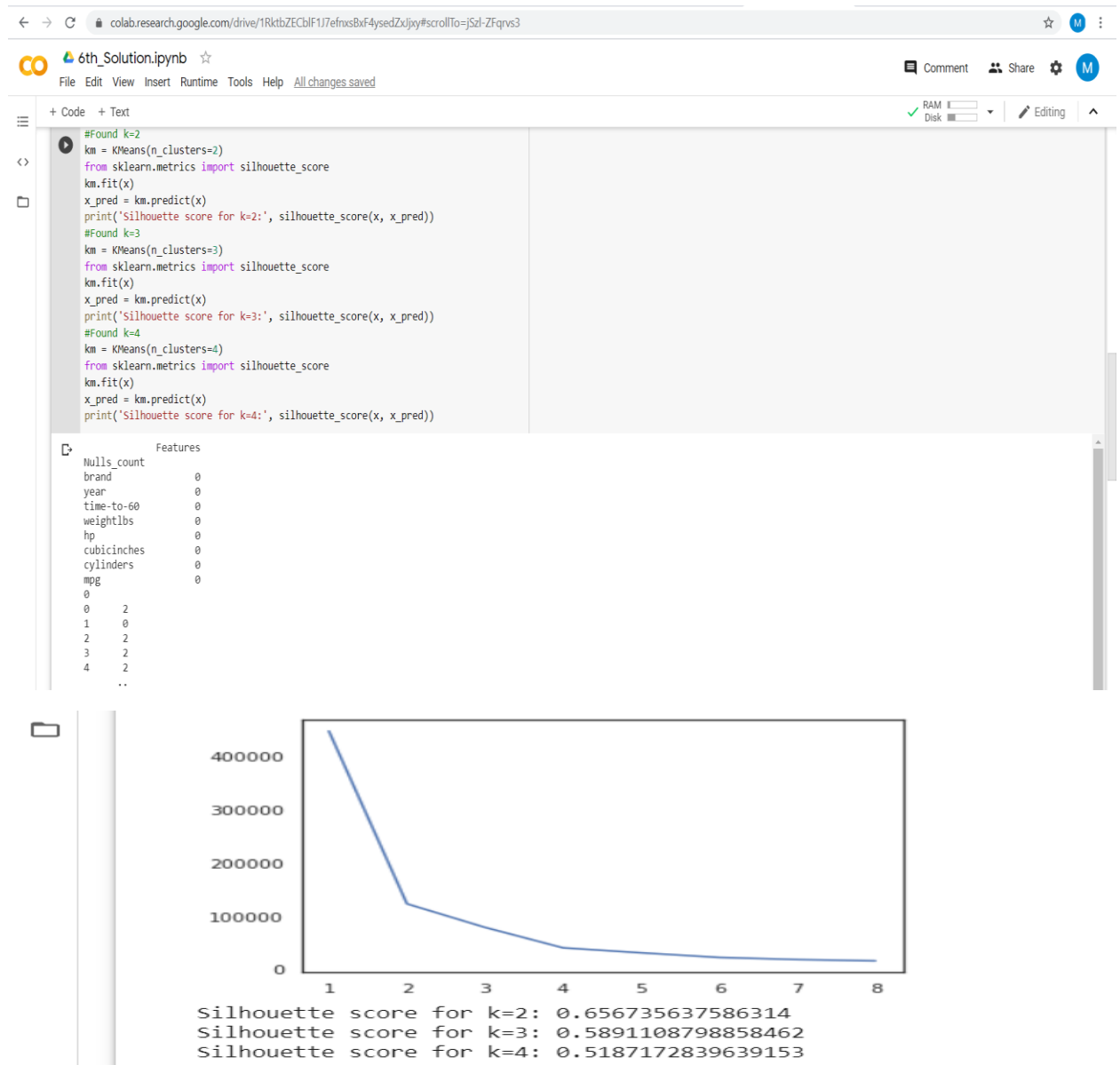
#Apply k-means algorithm
from sklearn.cluster import KMeans
wcss = []
for i in range(1, 9):
    km = KMeans(n_clusters=i, init='k-means++', max_iter=300, n_init=10, random_state=0)
    km.fit(x)
    wcss.append(km.inertia_)
#Visualize elbow method
import matplotlib.pyplot as plt
plt.plot(range(1, 9), wcss)
plt.title = 'The Elbow Method'
plt.xlabel = 'n-clusters'
plt.ylabel = 'wcss'
plt.show()
```



- We have taken cars.csv as the sample dataset which contains both the numeric and non-numeric data. We have applied exploratory data analysis and cleaned the dataset by handling all the null values by filling them with the mean values and encoded the non-numeric features using Label encoder.
- Then we have visualized all the correlation features and plotted them using seaborn.
- Used K-means clusterings to get the best value for k which is 2 by using the elbow method.

6(b).

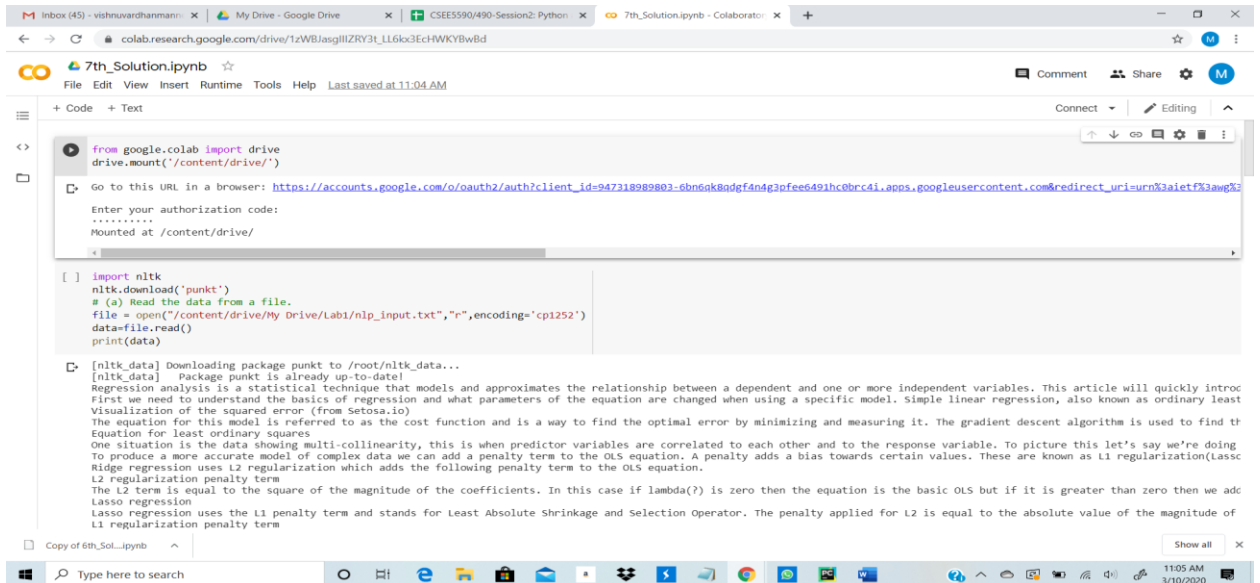
Code & Output:



- Calculated the Silhouette score for k=2,3,4 and when the number of clusters k=2 we got the highest silhouette score which is 0.656

7(a).

Code & Output:



The screenshot shows a Google Colab notebook titled "7th_Solution.ipynb". The code in the first cell is as follows:

```
from google.colab import drive
drive.mount('/content/drive/')

Go to this URL in a browser: https://accounts.google.com/o/oauth2/auth?client\_id=947318989803-6b66nk8d9f4n4g3pfee6491hc0br:ci.apps.googleusercontent.com&redirect\_uri=urn%3Aietf%3Awww%3Ahttp%3A%2F%2Flocalhost%3A8080%2F&scope=https%3A%2F%2Fwww.googleapis.com/auth/drive

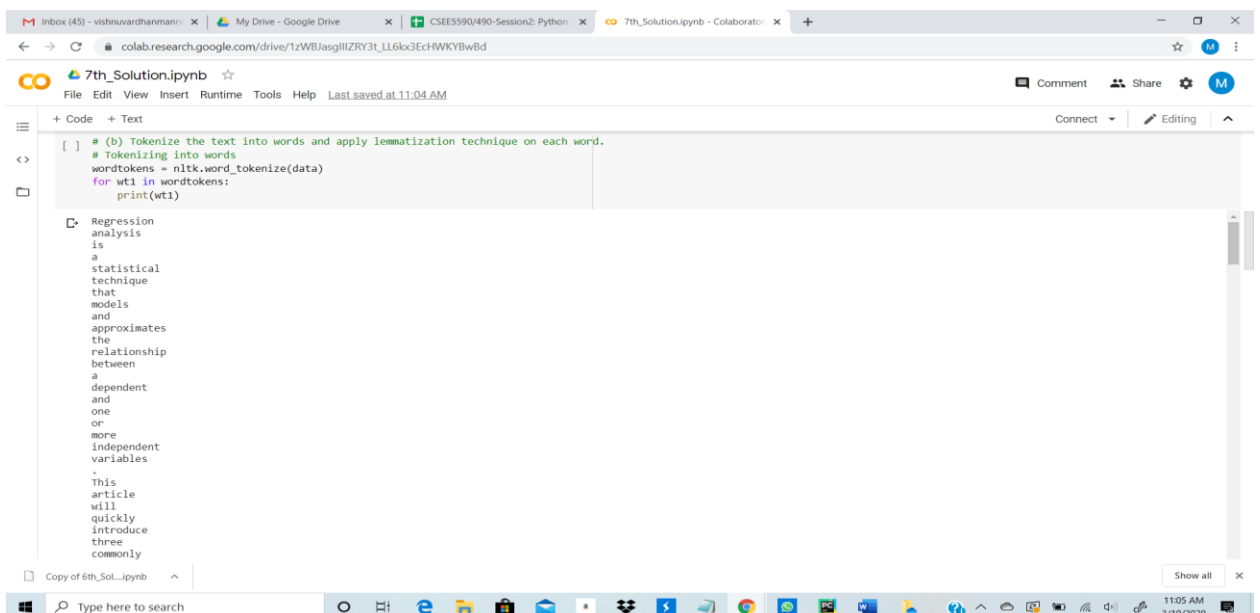
Enter your authorization code:
.....
Mounted at /content/drive/

[ ] import nltk
nltk.download('punkt')
# (a) Read the data from a file.
file = open("/content/drive/My Drive/Lab1/nlp_input.txt", "r", encoding='cp1252')
data=file.read()
print(data)
```

The output shows the package punkt is already up-to-date and the content of the file nlp_input.txt, which is a paragraph about regression analysis.

- Import the nltk library and then read the given nlp_input.txt file with encoding cp1252.
- Used file.read() to get the complete data in the give file and printed the result.

7(b).



The screenshot shows a Google Colab notebook titled "7th_Solution.ipynb". The code in the second cell is as follows:

```
# (b) Tokenize the text into words and apply lemmatization technique on each word.
# Tokenizing into words
wordtokens = nltk.word_tokenize(data)
for wt1 in wordtokens:
    print(wt1)
```

The output shows the tokens of the text from the previous cell, one by one.

- We have used word_tokenize function for tokenizing text into words.


```

[ ] # Lemmatization
from nltk.stem import WordNetLemmatizer
nltk.download('wordnet')
lemmatizer = WordNetLemmatizer()
for wt2 in wordtokens:
    print(lemmatizer.lemmatize(wt2))

[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data] Unzipping corpora/wordnet.zip.
Regression
analysis
is
a
statistical
technique
that
model
and
approximates
the
relationship
between
a
dependent
and
one
or
more
independent
variable
.
This
article
will
quickly

```

- We have imported wordnetlemmatizer and applied lemmatization on each word.

7(c,d).

```

[ ] # (c) Find all the trigrams for the words.
from nltk.util import ngrams
trigramoutput = []
trigrams=ngrams(wordtokens,3)
for t in trigrams:
    trigramoutput.append(t)
print(trigramoutput)

[('Regression', 'analysis', 'is'), ('analysis', 'is', 'a'), ('is', 'a', 'statistical'), ('a', 'statistical', 'technique'), ('statistical', 'technique', 'that'), ('technique', 'that', 'model'), ('model', 'and', 'approximates'), ('and', 'approximates', 'the'), ('approximates', 'the', 'relationship'), ('the', 'relationship', 'between'), ('relationship', 'between', 'a'), ('between', 'a', 'dependent'), ('dependent', 'and', 'one'), ('and', 'one', 'or'), ('one', 'or', 'more'), ('or', 'more', 'independent'), ('more', 'independent', 'variable'), ('independent', 'variable', '.'), ('variable', '.', 'This'), ('.', 'This', 'article'), ('This', 'article', 'will'), ('article', 'will', 'quickly')]

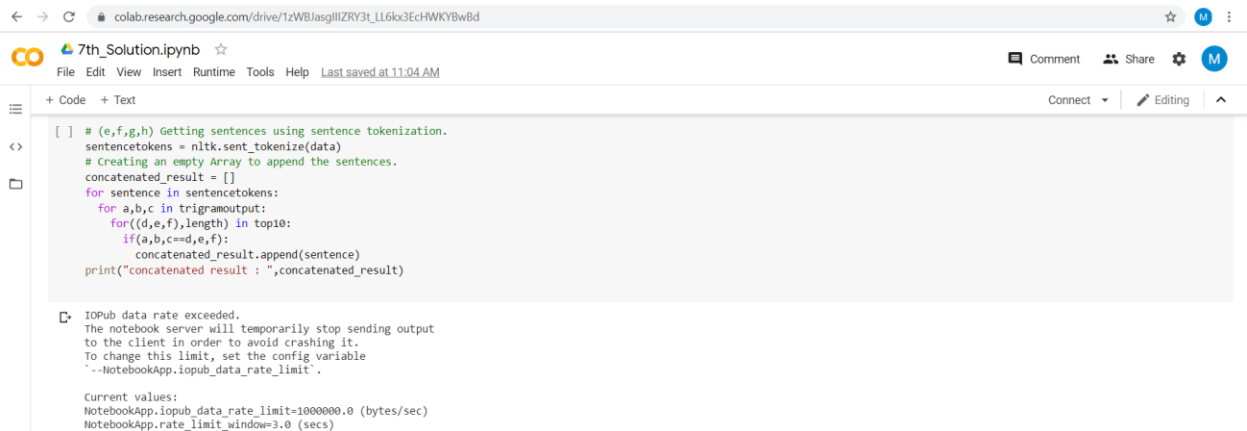
[ ] # (d) Extract the top 10 of the most repeated trigrams based on their count.
# Using trigramoutput we need to fetch the frequency of each word in the file
wordfrequency = nltk.FreqDist(trigramoutput)
# Printing the most common words
commonwords = wordfrequency.most_common()
print("Trigrams Frequency : \n", commonwords)
# Top 10 Trigrams
top10 = wordfrequency.most_common(10)
print("Top 10 Trigrams : \n", top10)

Trigrams Frequency :
[('we', 'need', 'to'), (('the', 'coefficients', '.'), 3), (('?', '?', '='), 3), (('to', 'find', 'the'), 2), (('find', 'the', 'optimal'), 2), (('over', 'a', 'number'), 2), (('a', 'statistical', 'technique'), 2), (('that', 'model', 'and'), 2), (('approximates', 'the', 'relationship'), 2), (('between', 'a', 'dependent'), 2)]
Top 10 Trigrams :
[('we', 'need', 'to'), (('the', 'coefficients', '.'), 3), (('?', '?', '='), 3), (('to', 'find', 'the'), 2), (('find', 'the', 'optimal'), 2), (('over', 'a', 'number'), 2), (('a', 'statistical', 'technique'), 2), (('that', 'model', 'and'), 2), (('approximates', 'the', 'relationship'), 2), (('between', 'a', 'dependent'), 2)]

```

- Imported the ngrams library and the value of n should be 3 in order to obtain all the trigrams.
- To extract the top 10 most repeated trigrams first we used FreqDist function to obtain the word frequency of trigrams. Then have used most_common(10) function to get the top 10 trigrams.

7(e,f,g).



```
[ ] # (e,f,g,h) Getting sentences using sentence tokenization.
sentencetokens = nltk.sent_tokenize(data)
# Creating an empty Array to append the sentences.
concatenated_result = []
for sentence in sentencetokens:
    for a,b,c in trigramoutput:
        for ((d,e,f),length) in top10:
            if(a,b,c==d,e,f):
                concatenated_result.append(sentence)
print("concatenated_result : ",concatenated_result)
```

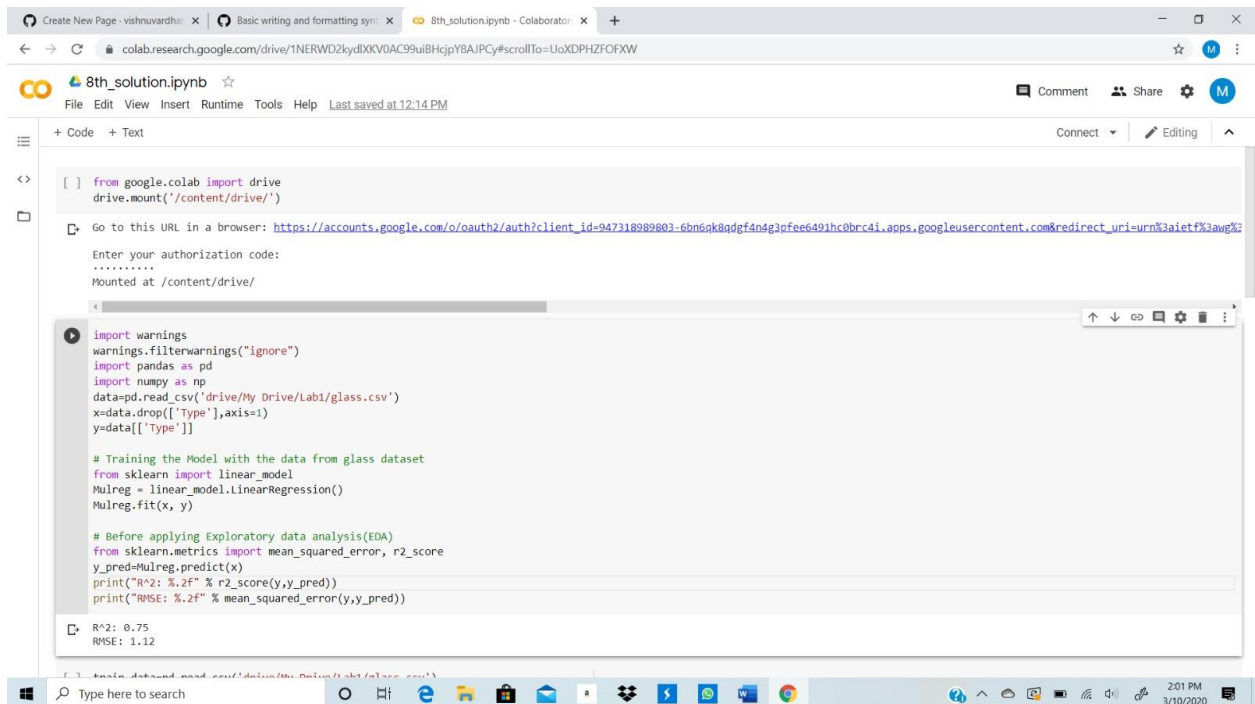
IOPub data rate exceeded.
The notebook server will temporarily stop sending output
to the client in order to avoid crashing it.
To change this limit, set the config variable
`--NotebookApp.iopub_data_rate_limit`.

Current values:
NotebookApp.iopub_data_rate_limit=1000000.0 (bytes/sec)
NotebookApp.rate_limit_window=3.0 (secs)

- We have used sent_tokenize to get every sentence and created an empty array to append each sentence. Then we have iterated sentence with trigrams to get sentence with most trigrams with most recurrence and concatenated the result and printed the output.

8.

Code & Output:



```
[ ] from google.colab import drive
drive.mount('/content/drive/')

Go to this URL in a browser: https://accounts.google.com/o/oauth2/auth?client_id=947318989803-6bn6nk8qdgfmg3pfee6491hc0brca1.apps.googleusercontent.com&redirect_uri=urn%3Aietf%3Awg%3A...
Enter your authorization code:
.....
Mounted at /content/drive/

import warnings
warnings.filterwarnings("ignore")
import pandas as pd
import numpy as np
data=pd.read_csv('drive/My Drive/Lab1/glass.csv')
x=data.drop(['Type'],axis=1)
y=data[['Type']]

# Training the Model with the data from glass dataset
from sklearn import linear_model
Mulreg = linear_model.LinearRegression()
Mulreg.fit(x, y)

# Before applying Exploratory data analysis(EDA)
from sklearn.metrics import mean_squared_error, r2_score
y_pred=Mulreg.predict(x)
print("R^2: %.2f" % r2_score(y,y_pred))
print("RMSE: %.2f" % mean_squared_error(y,y_pred))
```

R^2: 0.75
RMSE: 1.12

The first screenshot shows the initial data loading and cleaning steps in a Google Colab notebook. The code imports pandas and sklearn, reads the 'glass.csv' dataset, handles null values by filling them with the mean of each feature, and splits the data into training and testing sets. It then trains a linear regression model and calculates the RMSE and R-squared scores.

```
train_data=pd.read_csv('drive/My Drive/Lab1/glass.csv')

# Nulls Handling
nulls = pd.DataFrame(train_data.isnull().sum().sort_values(ascending=False))
nulls.columns = ['Features']
nulls.index.name = 'Nulls_count'
print(nulls)
x = train_data.select_dtypes(include=[np.number]).interpolate().dropna()
print(sum(x.isnull().sum() != 0))

# Here we are filling the null values with mean value
train_data=train_data.apply(lambda x: x.fillna(x.mean()),axis=0)
print(train_data['Type'])
print(train_data.isnull().sum())

# Split data into train and test
from sklearn.model_selection import train_test_split
x_train = train_data.drop(['Type'], axis=1)
y_train = train_data['Type']
x_train, x_test, y_train, y_test = train_test_split(x_train, y_train, test_size =0.4, random_state=0)

# Training the Model with the data from glass dataset
from sklearn import linear_model
Mulreg = linear_model.LinearRegression()
Mulreg.fit(x_train, y_train)

# After applying Exploratory data analysis(EDA)
from sklearn.metrics import mean_squared_error, r2_score
y_pred=Mulreg.predict(x_train)
print("R^2: %.2f" % r2_score(y_train,y_pred))
print("RMSE: %.2f" % mean_squared_error(y_train,y_pred))
```

The second screenshot shows the output of the EDA step, displaying the 'Features' column of the training data. The output shows the distribution of the 'Type' variable and the values of the other features.

```
[ ] Nulls_count
Type 0
Fe 0
Ba 0
Ca 0
K 0
Si 0
Al 0
Pb 0
Na 0
RI 0
0 1
1 1
2 1
3 1
4 1
..
209 7
210 7
211 7
212 7
213 7
Name: Type, Length: 214, dtype: int64
RI 0
Na 0
Pb 0
Al 0
Si 0
Ca 0
Ba 0
Fe 0
Type 0
dtype: int64
R^2: 0.78
RMSE: 1.08
```

- Imported the required libraries and we have used glass.csv dataset. Dropped the target variable 'Type' from training data and used linear regression for fitting the model. Before applying the Exploratory data analysis we have calculated the RMSE and R squared scores.
- Now we have cleaned the dataset by applying exploratory data analysis and used linear regression for fitting the model again. Now again we have calculated the RMSE and R squared scores.

- Before applying EDA the R_squared and RMSE values are 0.75 and 1.12.
- After applying EDA the R_squared and RMSE values are 0.78 and 1.08.
- The R_squared value is increased and RMSE value is decreased slightly after cleaning the dataset.

Workflow:

The work flow for all algorithms in this lab assignment is as follows:

- Created Python files using PyCharm IDE and used Google colab.
- Downloaded the datasets from Kaggle.
- Pre-processing of the data.
- Split the data into train and test using train_test_split().
- Trained the model with that data.
- Metrics calculation and plotting the data using matplotlib and seaborn.

Datasets:

- Cars.csv
- Glass.csv

Evaluation & Discussion:

We ensured that all the objectives of the lab assignment are met.

Conclusion:

We have completed the assignment on fundamentals of python and have trained the classifier with models like SVM, KNN, Naive Byes. We have also applied K-means clustering and used Elbow method and silhouette score to find the best k value and also applied Multiple regression technique to dataset and evaluated R2, RMSE scores before and after EDA successfully.