# Python Programming LAB1

## Report

Submitted By:

Dhairya Chandra (6)
Tanvi Jain (22)

Course:

CSEE5590/490

# INTRODUCTION

This lab assignment focused on how to make us familiar with python basic topics. This lab assignment focused on how to make us familiar with the basics of python

The Dataset we used in completing the Lab Assignment was:

1. Heart Disease UCI - Ronit. "Heart Disease UCI." *Kaggle*, 25 June 2018, www.kaggle.com/ronitf/heart-disease-uci.


2. Wine Quality - *UCI Machine Learning Repository: Wine Quality Data Set*, archive.ics.uci.edu/ml/datasets/wine+quality.

# OBJECTIVE

Main objective is to learn the basics of python and using multiple packages like beautifulsoup, nltk, etc. Also learned how to classification on data by splitting data into test and train.

Here is the list of Questions:

1. Given a collection of integers that might contain duplicates, nums, return all possible subsets

2. Concatenate two dictionaries and sort the concatenated dictionary by value.

3. Airline Booking Reservation System

4. fetch the course name and overview of course using beautiful soup

5. Perform 3 classification algorithms on dataset.

6. Apply K-means on the dataset and also find silhouette score, elbow score.

7. Perform Natural Language Processing on an input file

8. Do Multiple Regression on a dataset. Evaluate both RMSE and R2.

# QUESTION 1

Given a collection of integers that might contain duplicates, nums, return all possible subsets. Do not include null subset.

## SOLUTION

```python
def subset(arr, n):
    list = []

    for i in range(2 ** n):
        subset = ""
        # consider each element in the set
        for j in range(n):
            if (i & (1 << j)) != 0:
                subset += str(arr[j]) + " "
        # if subset is encountered for the first time
        if subset not in _list and len(subset) > 0:
            list.append(subset)
            # consider every subset
    for subset in _list:
        # split the subset and print its elements
        arr = subset.split(' ')
        for string in arr:
            print(string, end=" ")
        print()

if __name__ == '__main__':
    arr = [1, 2, 2]
    n = len(arr)
    subset(arr, n)
```
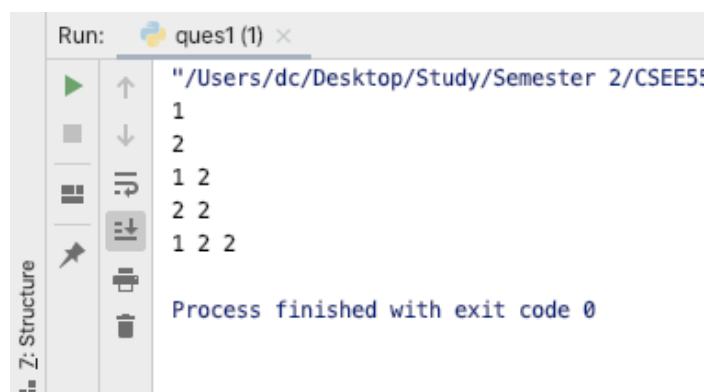
## WORKFLOW/APPROACH

- We are making a function subset which is called in main function
- Now we are considering each element in the set using for loop
- If subset is there then it is appended
- Also we are splitting the subset and printing all elements.
- In main function we are inputting the array and entering the values in subset function

# OUTPUT

```
Run:    ques1 (1)  ×
  ▶   ↑   "/Users/dc/Desktop/Study/Semester 2/CSEE5
  ■   ↓   1
          2
  ≡   ⇥   1 2
      ≡↓  2 2
  ★       1 2 2
      🖶
          Process finished with exit code 0
      🗑
```

# QUESTION 2

Concatenate two dictionaries and sort the concatenated
dictionary by value.

## SOLUTION

```python
# Create first dictionary
dict1 = {'Dhairya': 6, 'Tanvi': 22, 'John': 10}

# Create second dictionary
dict2 = {'Raju': 8, 'Bill': 20, 'Mark': 11}

# Concatenating two dictionaries
dict1.update(dict2)

# Printing the concatenated dictionary
print('Concatenated dictionary :', dict1)

# Sorting dictionaries by values
sortdict =  sorted(dict1.items(), key=lambda x: x[1])

# Printing the sorted dictionary
print('Sorted dictionary :', sortdict)
```
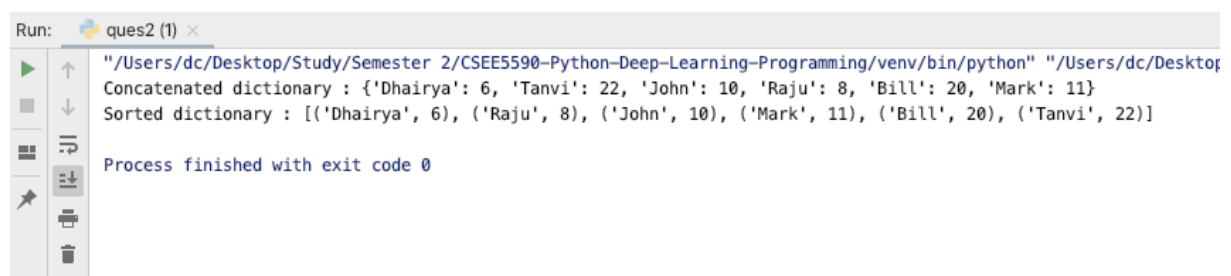
## WORKFLOW/APPROACH

- We are declaring two dictionaries dict1 and dict2
- To concatenate these two dictionaries I am using update function
- After that I am sorting the concatenated dictionaries by values using key lambda
- At end we are printing both dictionaries concatenated and sorted

## OUTPUT

```
Run:    ques2 (1) ×
   ▶  ↑   "/Users/dc/Desktop/Study/Semester 2/CSEE5590-Python-Deep-Learning-Programming/venv/bin/python" "/Users/dc/Desktop
      ↓    Concatenated dictionary : {'Dhairya': 6, 'Tanvi': 22, 'John': 10, 'Raju': 8, 'Bill': 20, 'Mark': 11}
   ■       Sorted dictionary : [('Dhairya', 6), ('Raju', 8), ('John', 10), ('Mark', 11), ('Bill', 20), ('Tanvi', 22)]

           Process finished with exit code 0
```

# QUESTION 3

Airline booking reservation system

## SOLUTION

```python
import random

class Flight:
    # Flight class - Default constructor
    def __init__(self, airline_name, flight_number):
        self.airline_name = airline_name
        self.flight_number = flight_number

    # Displaying flight details
    def flight_details(self):
        print('Airlines : ', self.airline_name)
        print('Flight number : Boeing', self.flight_number)


class employee:
    # Employee class - Default constructor
    def __init__(self, e_id, e_name, e_age, e_gender):
        self.e_name = e_name
        self.e_age = e_age
        self.__e_id = e_id
        self.e_gender = e_gender

    # Displaying employee details
    def e_display(self):
        print("Name of employee: ", self.e_name)
        print('Employee ID: ', self.__e_id)
        print('Employee Age: ', self.e_age)
        print('Employee Gender: ', self.e_gender)


class Passenger:
    # Passenger class - for fetching details of the passenger
    def __init__(self):
        Passenger.__passport_number = input("Enter the passport number of the passenger: ")
        Passenger.name = input('Enter name of the passenger: ')
        Passenger.age = input('Enter age of passenger : ')
        Passenger.gender = input('Enter the gender: ')
        Passenger.class_type = input('Select business or economy class: ')
```

```python
40
41 ●↓ ⊟class Baggage:
42          cabin_bag = 1
43          bag_fare = 0
44
45          # calculating cost for checked in bags more than 2
46 ●↓ ⊟     def __init__(self, checked_bags):
47              self.checked_bags = checked_bags
48              if checked_bags > 2:
49                  for i in checked_bags:
50                      self.bag_fare += 100
51              print("Number of checked bags allowed: ", checked_bags, "bag fare: $", self.bag_fare)
52
53
54 ●↓ ⊟class Fare(Baggage):
55          # Cost is fixed for purchasing at counter
56          counter = 1150
57          # Cost varies with ticket is purchased through online and fair is generated through random function
58          online = random.randint(1110, 2200)
59          total_fare = online
60
61 ●↓ ⊟     def __init__(self):
62              super().__init__(2)  # Super call
63              x = input('Buy ticket through online or counter: ')
64              if x == 'online':
65                  Fare.total_fare = self.online + self.bag_fare
66              elif x == 'counter':
67                  Fare.total_fare = self.counter + self.bag_fare
68              else:
69                  x = input('Enter correct transaction type:')
70              print("Total Fare before class type: $", Fare.total_fare)
71


71
72
73 ⊟class Ticket(Passenger, Fare):  # Multiple inheritence
74 ⊟     def __init__(self):
75              print("Passenger name:", Passenger.name)  # Accessing parent class variable
76              if Passenger.class_type == "business":
77                  Fare.total_fare += 100
78              else:
79                  pass
80              print("Passenger class type:", Passenger.class_type)
81              print("Total fare: $", Fare.total_fare)  # Displaying total fare
82
83
84 print("-------------- Flight Details --------------")
85 f1 = Flight('Air India', 777)
86 f1.flight_details()
87 print("-------------- Employee Details --------------")
88 e0 = employee('4588124', 'Dhairya Chandra', 22, 'Male')
89 e0.e_display()
90
91 print("-------------- Book Ticket --------------")
92 p1 = Passenger()
93
94 fare1 = Fare()
95
96 print("-------------- Passenger Details --------------")
97 t = Ticket()
98
```
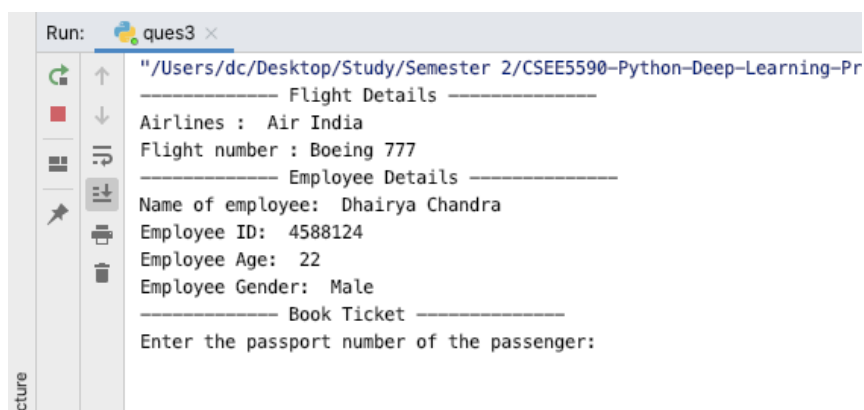
# WORKFLOW/APPROACH

- We have created six classes: flight, employee, passenger, baggage, fare and ticket
- Employers will take input from passenger using passenger class
- Flight cost is fixed if purchased from counter and varies if purchased online which is calculated using random function
- Using multiple inheritance flight details will be shown which is fetched from all classes declared
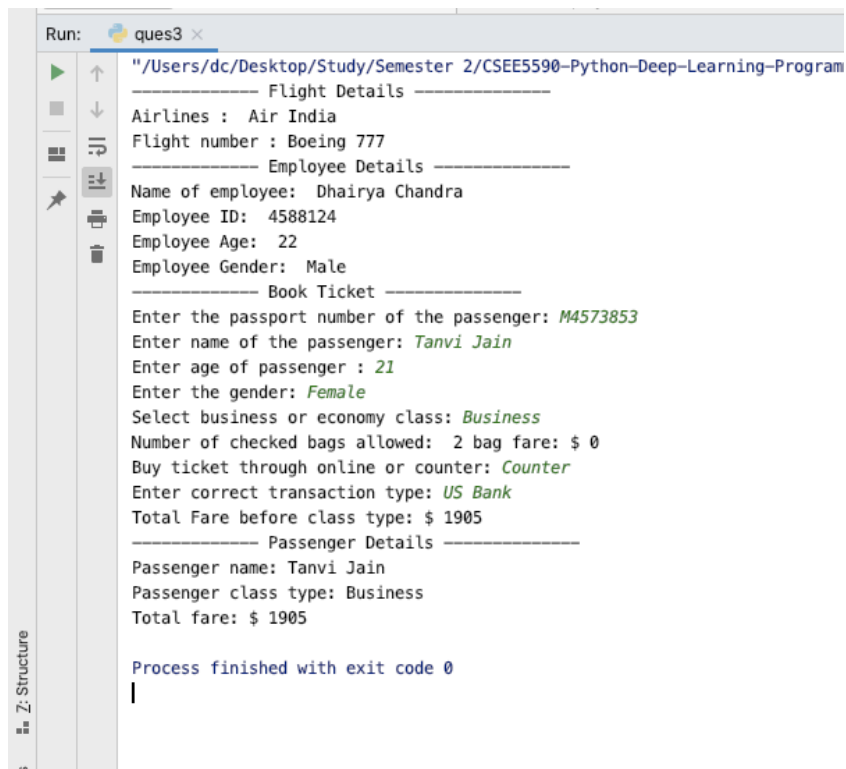- Employee details and flight details will be shown on the console

# OUTPUT

This will be showed by default:



Now employee has to enter information of passenger so that ticket get generated:

# QUESTION 4

fetching the course name and overview of course using Beautiful Soup

## SOLUTION

```python
from bs4 import BeautifulSoup
import requests

# Enter URL from where you have to fetch the data
url  = requests.get("https://catalog.umkc.edu/course-offerings/graduate/comp-sci/")
data = url.text
soup = BeautifulSoup(data, "html.parser")

# Using FindAll to find specific class

for p in soup.findAll('p',{'class':'courseblocktitle'}):
    for title, detail in zip(soup.findAll('span', {'class': 'title'}),
                             soup.findAll('p', {'class': 'courseblockdesc'})):

        # Print each course detail
        print ("Course Name: ", title.string)
        print("Detail: ", detail.string)
        print("--------------------------")
```
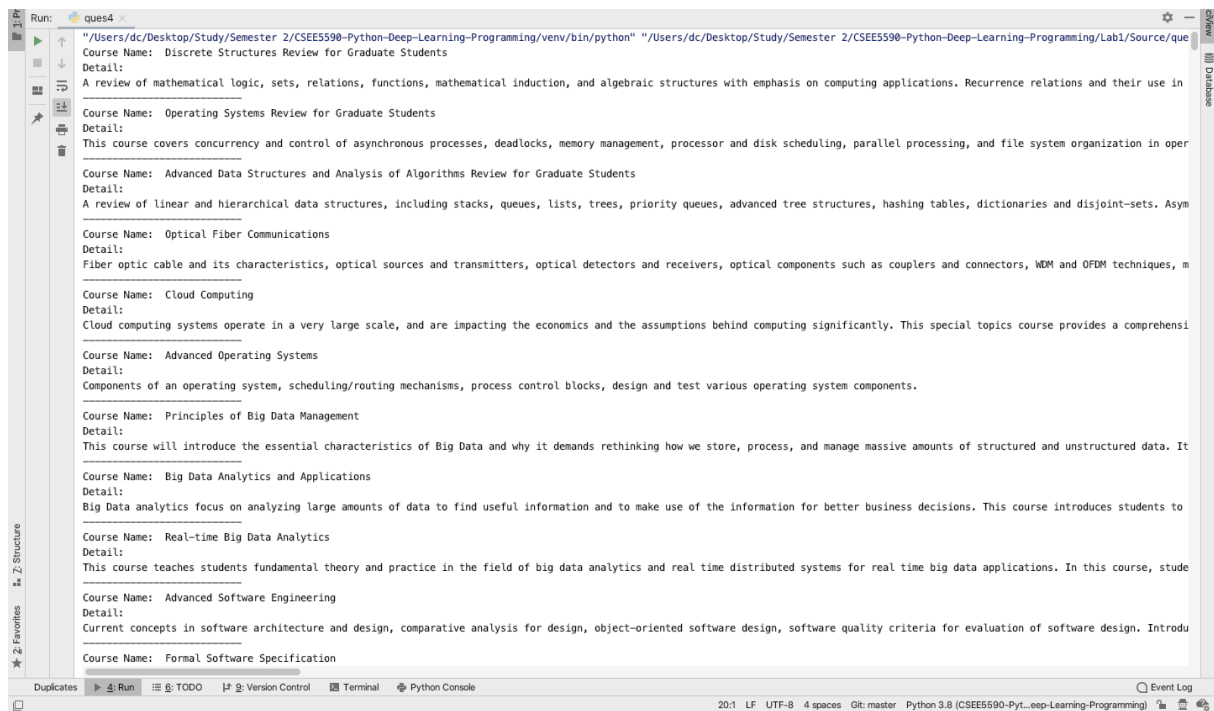
## WORKFLOW/APPROACH

- We have used beautiful soup to fetch contents from a web page
- Webpage is parsed using beautiful soup
- Now using inspect element I have selected the class where the content data is there
- I am using findall function to access each class and then using zip function I am showing course title with its description
- Then I am printing course name and details

# OUTPUT

```
"/Users/dc/Desktop/Study/Semester 2/CSEE5590-Python-Deep-Learning-Programming/venv/bin/python" "/Users/dc/Desktop/Study/Semester 2/CSEE5590-Python-Deep-Learning-Programming/Lab1/Source/ques
Course Name:  Discrete Structures Review for Graduate Students
Detail:
A review of mathematical logic, sets, relations, functions, mathematical induction, and algebraic structures with emphasis on computing applications. Recurrence relations and their use in
------------------------
Course Name:  Operating Systems Review for Graduate Students
Detail:
This course covers concurrency and control of asynchronous processes, deadlocks, memory management, processor and disk scheduling, parallel processing, and file system organization in oper
------------------------
Course Name:  Advanced Data Structures and Analysis of Algorithms Review for Graduate Students
Detail:
A review of linear and hierarchical data structures, including stacks, queues, lists, trees, priority queues, advanced tree structures, hashing tables, dictionaries and disjoint-sets. Asym
------------------------
Course Name:  Optical Fiber Communications
Detail:
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, m
------------------------
Course Name:  Cloud Computing
Detail:
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 comprehensi
------------------------
Course Name:  Advanced Operating Systems
Detail:
Components of an operating system, scheduling/routing mechanisms, process control blocks, design and test various operating system components.
------------------------
Course Name:  Principles of Big Data Management
Detail:
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
------------------------
Course Name:  Big Data Analytics and Applications
Detail:
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
------------------------
Course Name:  Real-time Big Data Analytics
Detail:
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, stude
------------------------
Course Name:  Advanced Software Engineering
Detail:
Current concepts in software architecture and design, comparative analysis for design, object-oriented software design, software quality criteria for evaluation of software design. Introdu
------------------------
Course Name:  Formal Software Specification
```

# QUESTION 5

Perform three classification algorithms Naïve Bayes, SVM and KNN on any dataset.

## SOLUTION

Dataset Used is **Heart Disease UCI**

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.naive_bayes import GaussianNB
from sklearn.neighbors import KNeighborsClassifier
from sklearn import svm
from sklearn.preprocessing import StandardScaler

# Dataset: https://www.kaggle.com/ronitf/heart-disease-uci
train = pd.read_csv('heart.csv')

# Working with Numeric Features
numeric_features = train.select_dtypes(include=[np.number])
corr = numeric_features.corr()
plt.figure(figsize=(20,20));
sns.heatmap(corr, annot=True, cmap="YlGnBu")
plt.show();
print (corr['age'].sort_values(ascending=False)[:4], '\n')

# Null values
nulls = pd.DataFrame(train.isnull().sum().sort_values(ascending=False)[:25])
nulls.columns = ['Null Count']
nulls.index.name = 'Feature'
print(nulls)

# Handling the missing value
data = train.select_dtypes(include=[np.number]).interpolate().dropna()
print(sum(data.isnull().sum() != 0))

# Encoding the categorial feature
data_binary = pd.get_dummies(train)
data_binary.head()
```

```
36
37     # Spliting Test and Train data
38     x_train, x_test, y_train, y_test = train_test_split(data_binary,train['age'])
39     performance = []
40
41
42     # --------------- Using Naive Bayes classification ---------------
43     GNB = GaussianNB()
44
45     # Training Model
46     GNB.fit(x_train,y_train)
47     train_score = GNB.score(x_train,y_train)
48
49     # Predicting Output
50     test_score = GNB.score(x_test,y_test)
51     print(f'Gaussian Naive Bayes : Training score: {train_score}, Test score: {test_score}')
52
53
54     # --------------- Using KNN classification ---------------
55     knn = KNeighborsClassifier(n_neighbors=5)
56
57     # Training Model
58     knn.fit(data_binary,train['age'])
59     knn.score(x_train,y_train)
60
61     # Predicting Output
62     train_score = knn.score(x_train,y_train)
63     test_score = knn.score(x_test,y_test)
64     print(f'K Neighbors : Training score: {train_score}, Test score: {test_score}')
65
66
67     # --------------- creating SVM classification ---------------
68     svc = svm.SVC(kernel='linear')
69
70     # Training Model
71     scaler = StandardScaler()
72     scaler.fit(data_binary,train['age'])
73     x_train_scaled = scaler.transform(x_train)
74     x_test_scaled = scaler.transform(x_test)
75     svc.fit(x_train_scaled,y_train)
76
77     # Predicting Output
78     train_score = svc.score(x_train_scaled,y_train)
79     test_score = svc.score(x_test_scaled, y_test)
80     print(f'SVM : Training score: {train_score}, Test score: {test_score}')
81
```

# WORKFLOW/APPROACH

- We have used heart disease UCI database
- Then we are finding top 3 features
- After that we are print all null values from the dataset
- Then we are encoding the categorical features
- For classification we are splitting the test and train dataset using train_test_split
- Three classification techniques are used on dataset: Naïve bayes, SVM and KNN
- We are getting these results:
  - Naïve Bayes - Train score: 100%, Test score: 86.84%
  - SVM - Train score: 24.22%, Test score: 22.36%
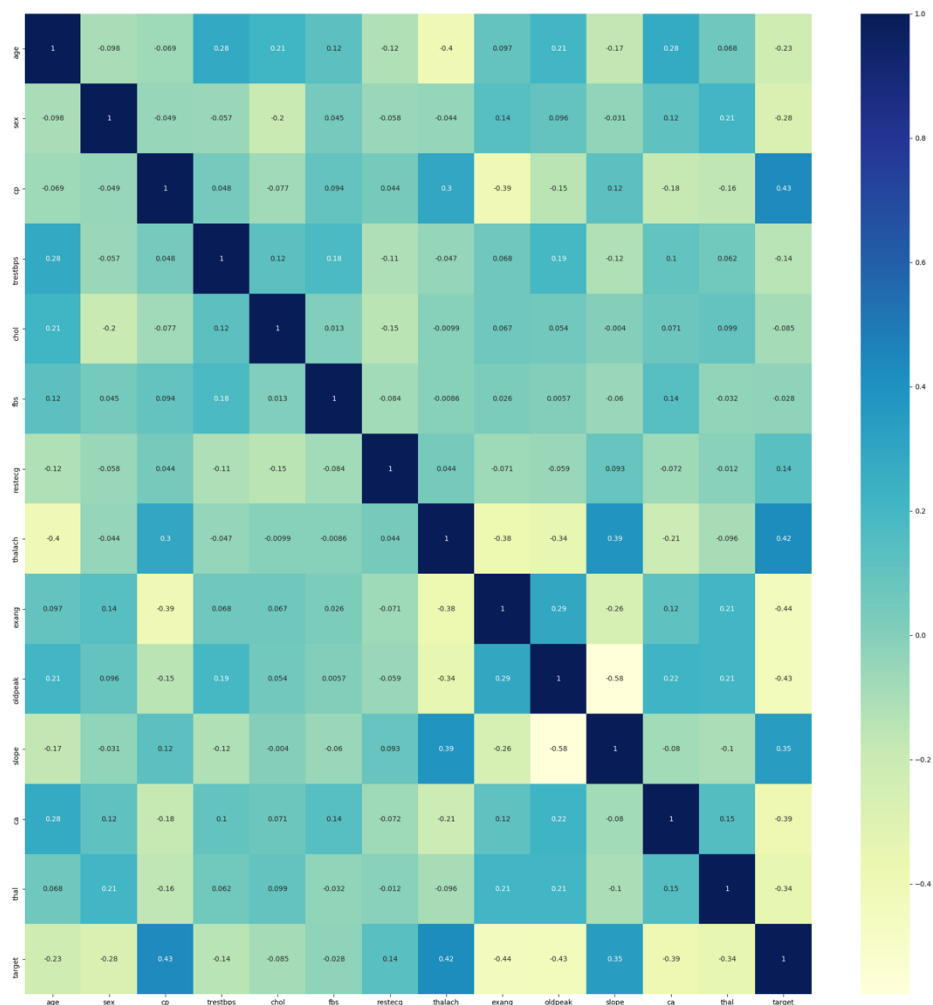  - KNN - Train score: 86.78%, Test score: 10.26%

# OUTPUT

```
Run:    ques5 ×
    "/Users/dc/Desktop/Study/Semester 2/CSEE5590-Python-Deep-Learning-Programming/venv/bin/python" "/Users/dc/
    age           1.000000
    trestbps      0.279351
    ca            0.276326
    chol          0.213678
    Name: age, dtype: float64

                Null Count
    Feature
    target             0
    thal               0
    ca                 0
    slope              0
    oldpeak            0
    exang              0
    thalach            0
    restecg            0
    fbs                0
    chol               0
    trestbps           0
    cp                 0
    sex                0
    age                0
    0
    Gaussian Naive Bayes : Training score: 1.0, Test score: 0.868421052631579
    K Neighbors : Training score: 0.2422907488986784, Test score: 0.2236842105263158
    SVM : Training score: 0.8678414096916299, Test score: 0.10526315789473684

    Process finished with exit code 0
```

# QUESTION 6

Apply K-means on any dataset and visualize the clusters using matplotlib or seaborn

## SOLUTION

Dataset Used is **Wine Quality Data Set**

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.cluster import KMeans
from sklearn import preprocessing
from sklearn.decomposition import PCA

from sklearn.metrics import silhouette_score

dataset = pd.read_csv('winequality-red.csv')

# Null values
nulls = pd.DataFrame(dataset.isnull().sum().sort_values(ascending=False)[:25])
nulls.columns = ['Null Count']
nulls.index.name = 'Feature'
print(nulls)

# handling the missing value
data = dataset.select_dtypes(include=[np.number]).interpolate().dropna()

# find the most correlated features
numeric_features = dataset.select_dtypes(include=[np.number])
corr = numeric_features.corr()
print (corr['quality'].sort_values(ascending=False)[:4], '\n')

# Preprocessing the data
scaler = preprocessing.StandardScaler()
scaler.fit(data)
X_scaled_array = scaler.transform(data)
X_scaled = pd.DataFrame(X_scaled_array, columns = data.columns)


wcss = []
# elbow method to know the number of clusters
for i in range(2,12):
    kmeans = KMeans(n_clusters=i,init='k-means++',max_iter=300,n_init=10,random_state=0)
    kmeans.fit(data)
    wcss.append(kmeans.inertia_)
    score = silhouette_score(data, kmeans.labels_, metric='euclidean')
    print("For n_clusters = {}, silhouette score is {}".format(i, score))

plt.plot(range(1,11),wcss)
plt.title('the elbow method')
plt.xlabel('Number of Clusters')
plt.ylabel('Wcss')
plt.show()
```

# WORKFLOW/APPROACH

- We have used Wine Quality dataset
- We have printed the null values, handled the missing values, and then finding the most correlated features
- Preprocessing is done on the dataset
- Used elbow method on multiple clusters and generating the silhouette score
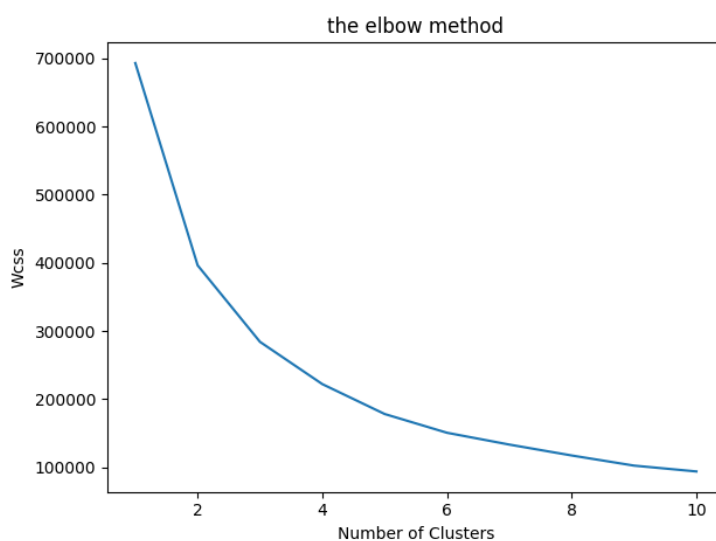
# OUTPUT



```
Run:    ques6 ×
        "/Users/dc/Desktop/Study/Semester 2/CSEE5590-Python-Deep-Learning-Programming/venv/bin/python" "/Users/dc/Desk
                        Null Count
        Feature
        quality                    0
        alcohol                    0
        sulphates                  0
        pH                         0
        density                    0
        total sulfur dioxide       0
        free sulfur dioxide        0
        chlorides                  0
        residual sugar             0
        citric acid                0
        volatile acidity           0
        fixed acidity              0
        quality        1.000000
        alcohol        0.476166
        sulphates      0.251397
        citric acid    0.226373
        Name: quality, dtype: float64

        For n_clusters = 2, silhouette score is 0.6027870469574543
        For n_clusters = 3, silhouette score is 0.5184003155871573
        For n_clusters = 4, silhouette score is 0.48448390096387717
        For n_clusters = 5, silhouette score is 0.4451649073633646
        For n_clusters = 6, silhouette score is 0.4462630538654589
        For n_clusters = 7, silhouette score is 0.3920564307436043
        For n_clusters = 8, silhouette score is 0.3854113495216165
        For n_clusters = 9, silhouette score is 0.3823228452563935
        For n_clusters = 10, silhouette score is 0.3801474576185881
        For n_clusters = 11, silhouette score is 0.3796184204120568

        Process finished with exit code 0
```

# QUESTION 7

Perform Natural Language Processing on an input file

## SOLUTION

```python
import nltk
from nltk.stem import WordNetLemmatizer

# A. Read the data from a file
rfile = open("nlp_input.txt", "r", encoding="utf8", errors='ignore')

# opening a file in write mode to store the tokenized words in it
with open("tokenize.txt", "w") as t:
    # reading each sentence
    for sentence in rfile:
        # Performing word tokenization
        wtokens = nltk.word_tokenize(sentence)

        # Writing it to a file
        for w in wtokens:
            t.write(str("\n"))
            t.write(str(w))

# creating a lemmatization object
lemmatizer = WordNetLemmatizer()

# Opening a new file in write mode
with open("lemmatize.txt", "w") as l:
    # Opening the file consisting of words tokenized
    w = open("tokenize.txt", "r")
    for words in w:
        # B. Performing lematization on each word
        le = lemmatizer.lemmatize(words)
        l.write(str(le))
```

```python
30
31     a = open("nlp_input.txt", "r", encoding="utf8", errors='ignore')
32     with open("trigram.txt", "w") as tri:
33         for sentence in a:
34             # C. performing trigram on each sentence
35             trigram = nltk.trigrams(sentence.split())
36
37             # trigram are written on to a file
38             for ti in trigram:
39                 tri.write(str("\n"))
40                 tri.write(str(ti))
41
42     f1 = open("nlp_input.txt", "r", encoding="utf8", errors='ignore')
43     fileread = f1.read()
44
45     tg = []
46     word_tokens = nltk.word_tokenize(fileread)
47     for t in nltk.ngrams(word_tokens, 3):
48         tg.append(t)
49
50     wordFreq = nltk.FreqDist(tg)
51     mostCommon = wordFreq.most_common()
52
53     # D. Extract the top 10 of the most repeated trigrams based on their count.
54     Top_ten_trigrams = wordFreq.most_common(10)
55     print("Top 10 Trigrams:\n", Top_ten_trigrams, "\n")
56
57     # E. Go through the text in the file
58     # F. Find all the sentences with the most repeated tri-grams
59     # G. Extract those sentences and concatenate
60
```

```python
60
61
62     sent_tokens = nltk.sent_tokenize(fileread)
63     concat_result = []
64     # Iterating the Sentences
65     for s in sent_tokens:
66         # Iterating all the trigrams
67         for a, b, c in tg:
68             # iterating the top 10 trigrams from all the trigrams
69             for ((p, q, r), length) in wordFreq.most_common(20):  # Comparing the each with the top 10 trigrams
70                 if (a, b, c == p, q, r):
71                     concat_result.append(s)
72
73     # H. Print the concatenated result
74     print("Concatenated Array: \n", concat_result)
75     print("Maximum of Concatenated Array: \n ", max(concat_result))
76
```
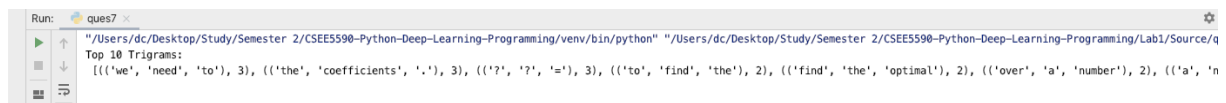
# WORKFLOW/APPROACH

- Reading input from given text file
- Tokenizing the input file and storing the tokenized text in txt file
- Lemmatize the input file to split words into parts of speech and stored in separate txt file
- Performed trigram and split the sentence and stored in a txt file
- Top 10 repeated trigrams are printed
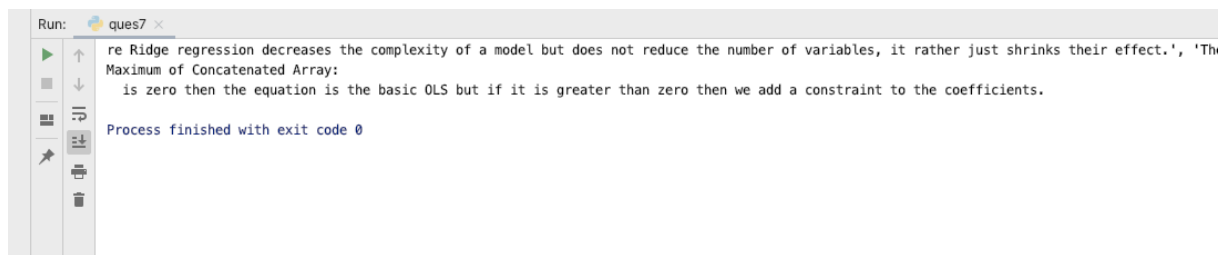- Then we finally printed the most repeated trigram

# OUTPUT

Top 10 Trigrams:



Concatenated Array and Maximum of Concatenated Array:

# QUESTION 8

Do Multiple Regression on a dataset. Evaluate both RMSE and R2

## SOLUTION

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns; sns.set(color_codes=True)

#Dataset: https://www.kaggle.com/ronitf/heart-disease-uci
train = pd.read_csv('heart.csv')

# Null values
nulls = pd.DataFrame(train.isnull().sum().sort_values(ascending=False)[:25])
nulls.columns = ['Null Count']
nulls.index.name = 'Feature'
print(nulls)

# Replacing null values with mean values
data = train.select_dtypes(include=[np.number]).interpolate().dropna()


# Using Pearson Correlation and ploting in the heat map
plt.figure(figsize=(20,20))
cor = data.corr()
sns.heatmap(cor, annot=True, cmap=plt.cm.Reds)
plt.show()

# Printing the correlation with the target feature "quality"
print(cor['target'].sort_values(ascending=False)[:5],'\n')

# Build a multiple linear regression model
y = data['target']
X = data.drop(['target'],axis =1)

print(X.shape)
```

```python
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=42, test_size=.20)
from sklearn import linear_model
lr = linear_model.LinearRegression()
model = lr.fit(X_train, y_train)

# Evaluate the performance and visualize results
print("R2: \n", model.score(X_test, y_test))
predictions = model.predict(X_test)
from sklearn.metrics import mean_squared_error
print('RMSE: \n', mean_squared_error(y_test, predictions))

# visualize
actual_values = y_test
plt.scatter(predictions, actual_values, alpha=.75,
            color='b') # alpha helps to show overlapping data
plt.xlabel('Predicted ')
plt.ylabel('Actual')
plt.title('Linear Regression Model')
plt.show()
```
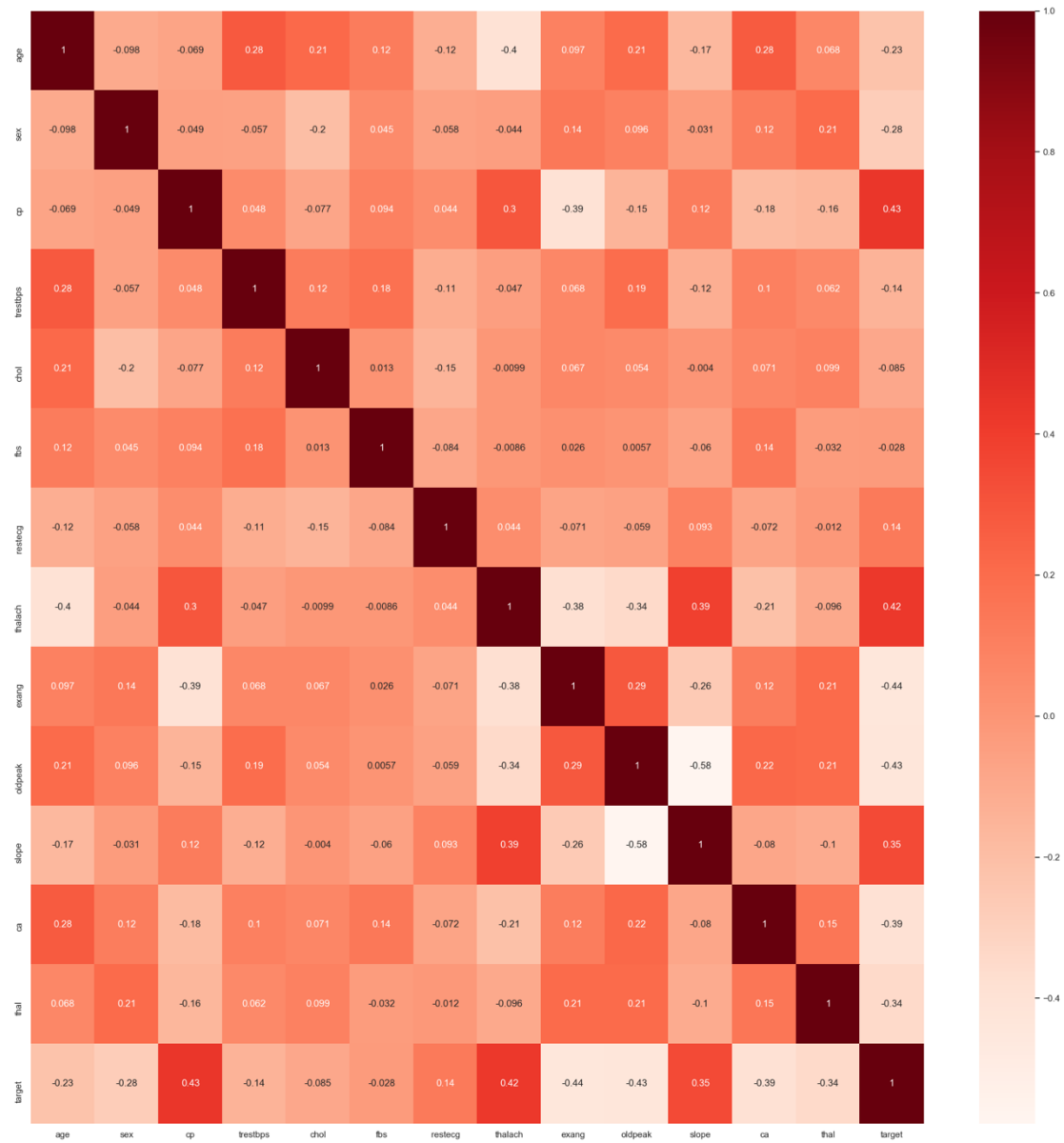
# WORKFLOW/APPROACH

- Using Heart Disease dataset
- Printing the null values
- Plotting the heat map of pearson correlation
- Printing the top 5 correlated features
- Building the linear regression model and splitting the dataset for training and testing
- Evaluated the scores for R2, RMSE
- Finally plotting the linear regression model

# OUTPUT

# CONCLUSION

We have learned how to use python for various tasks like extracting data from website, doing classification such as KNN, SVM and Naïve Bayes.

Also learned all basic usage of python in programming.


YOUTUBE: https://youtu.be/q-cGjCc369Y




GITHUB:  https://github.com/dhairyachandra/CSEE5590-Python-Deep-Learning-Programming/tree/master/Lab1




GITHUB WIKI:
https://github.com/dhairyachandra/CSEE5590-Python-Deep-Learning-Programming/wiki/Lab-1