

# IC 152

## Computing & Data Science

### Lab 6

**Mohit Verma**

**B20215**

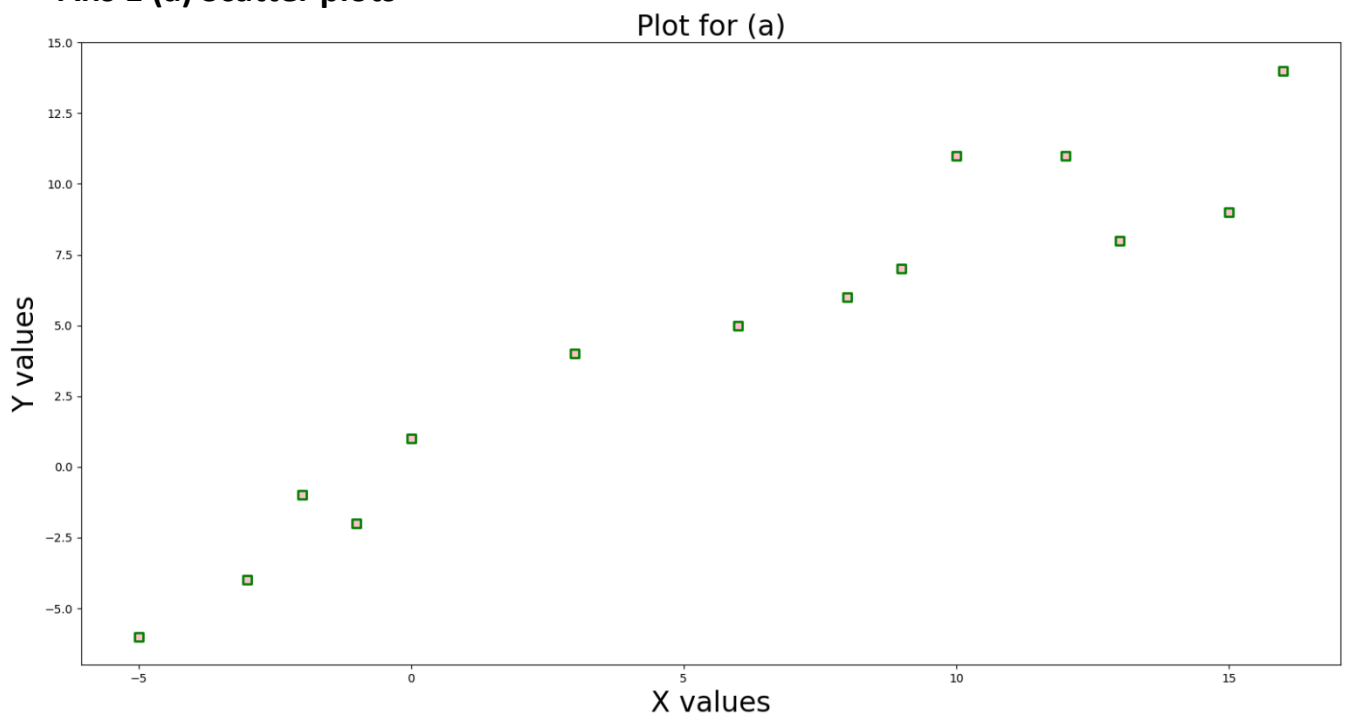
**Instructor- Dr.Padmanabhan Rajan**

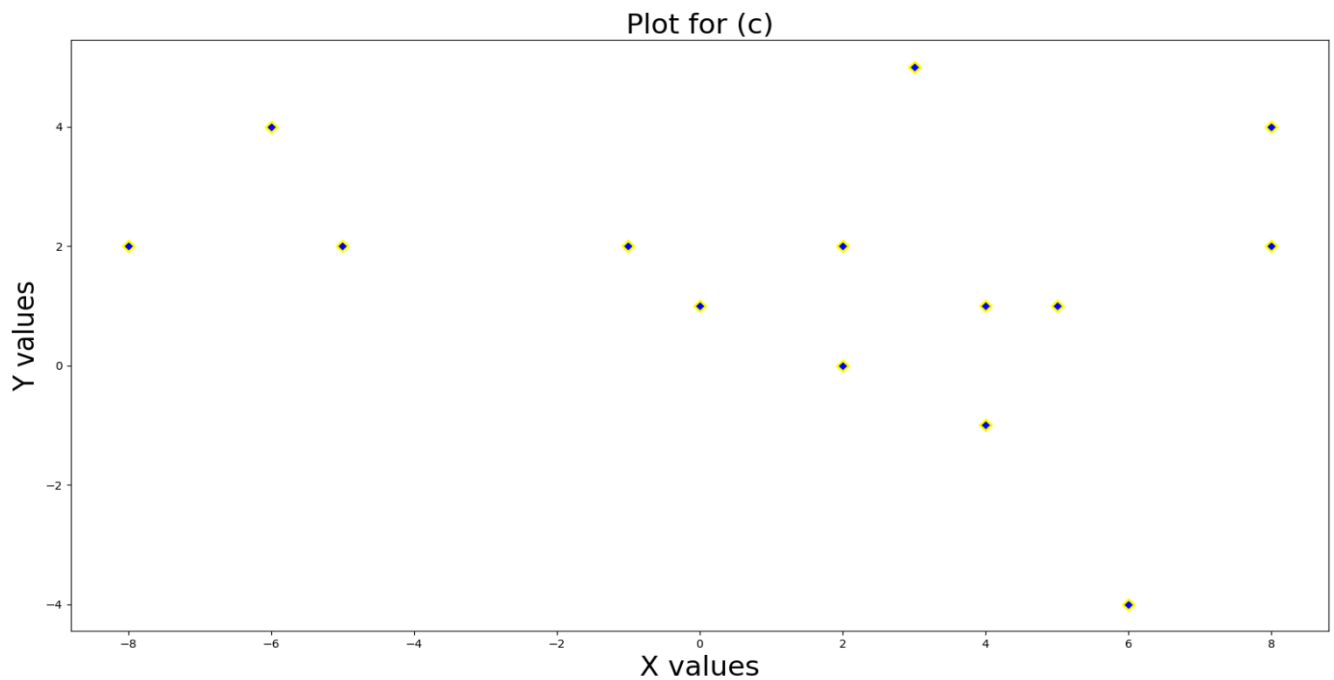
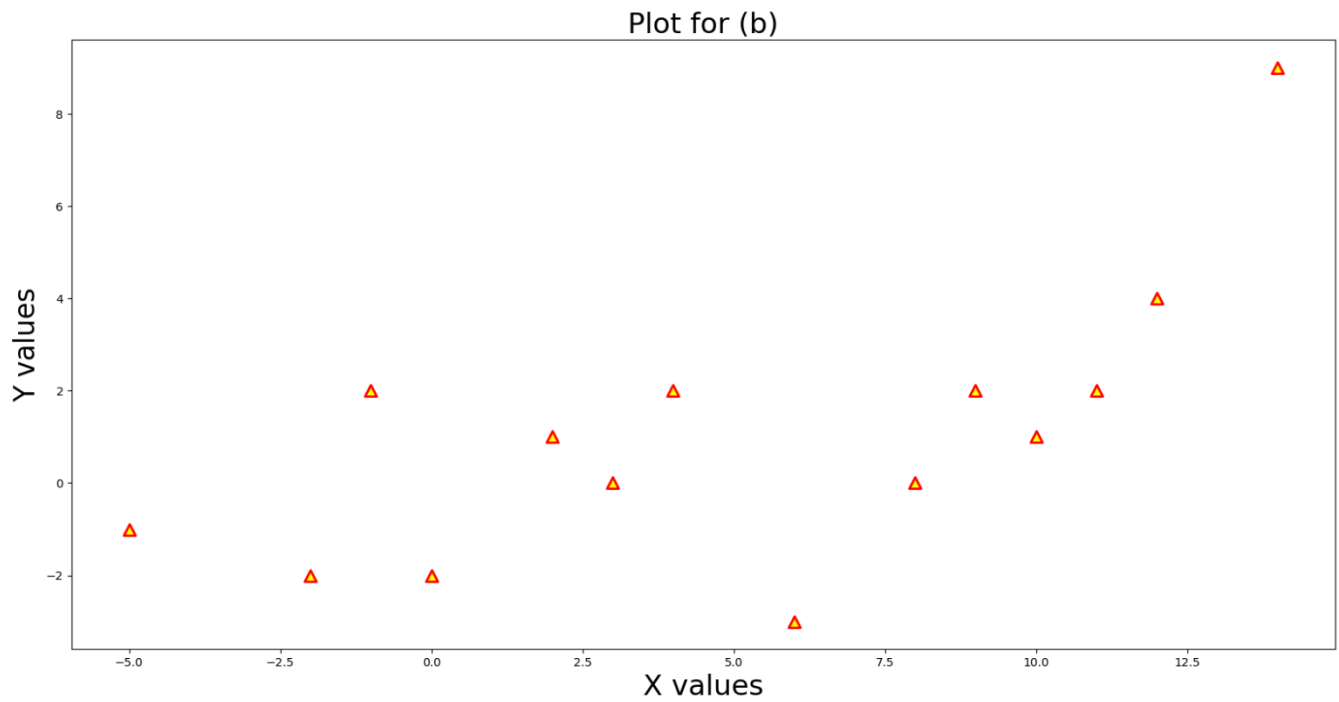
## Questions & Answers

**Q1 (i)** Plot the variables (x and y) as scatter plots given below.

**(ii)** Find out the correlation between the different cases of two variables (x and y) as given below and use the equation (mentioned above). Analyse the numerical value and scatter plots of the variables (i.e., if correlation is positive, y should increase with increase in x).

**Ans 1 (a) Scatter plots**





**Ans 1 (ii)**

```
In [6]: runfile('C:/Users/verma/OneDrive/Desktop/LAB 6/untitled1.py', wdir='C:/Users/verma/OneDrive/Desktop/LAB 6')
Reloaded modules: jupyter_client.session, zmq.eventloop, zmq.eventloop.ioloop, tornado.platform,
tornado.platform.asyncio, tornado.gen, zmq.eventloop.zmqstream, jupyter_client.jsonutil, jupyter_client.adapter,
spyder, spyder.pil_patch, PIL, PIL._version, PIL.Image, PIL.ImageMode, PIL.TiffTags, PIL._binary, PIL._util,
PIL._imaging, cffi, cffi.api, cffi.lock, cffi.error, cffi.model, IPython.lib.guisupport, IPython.external,
IPython.external.qt_for_kernel, IPython.utils.version, IPython.external.qt_loaders
Correlation for (a) 0.9591483545329047
Correlation for (b) 0.6470857165183262
Correlation for (c) -0.2518839194843887
```

```
In [7]:
```

We observe that Correlation for part(a) is 0.9591 which is positive and very high(close to 1) which shows that the value of y increases with x and we have a strong relation between x and y values.

For part(b) Correlation is 0.6470 which is positive which shows that value of y increases as x increases.

For part (c) the Correlation -0.2518 which is negative which indicates that the value of y decreases with increase in x. It has weak relation between x and y.

**Code for 1(i) and 1(ii) is shown –**

```
1  # -*- coding: utf-8 -*-
2  """
3  Created on Thu Apr  8 20:19:11 2021
4
5  @author: verma
6  """
7
8  # Q1
9  import matplotlib.pyplot as plt
10 def correlation(a,b):
11     prod_xy=0
12     sum_x2=0
13     n=len(a)
14     sum_y2=0
15     sum_x=0
16     sum_y=0
17     for x,y in zip(a,b):
18         prod_xy+=(x*y)
19         sum_x2+=(x**2)
20         sum_y2+=(y**2)
21         sum_x+=x
22         sum_y+=y
23     p=((n*prod_xy)-(sum_x*sum_y))/((((n*sum_x2)-(sum_x**2)**0.5)*((n*sum_y2-(sum_y**2)**0.5))
24     return p
25
26
27 #Q1 (A)
28 x1=[3,-1,6,12,8,10,9,13,15,-3,0,-5,-2,16]
29 x2=[3,-1,6,12,8,10,9,-5,-2,4,2,0,14,11]
30 x3=[3,4,2,5,8,-8,-6,-1,0,-5,4,2,6,8]
31 y1=[4,-2,5,11,6,11,7,8,9,-4,1,-6,-1,14]
32 y2=[0,2,-3,4,0,1,2,-1,-2,2,1,-2,9,2]
33 y3=[5,-1,0,1,2,2,4,2,1,2,1,2,-4,4]
34 plt.scatter(x1,y1,c="pink",linewidths=2,marker="s",edgecolor="green",s=50)
35 plt.title("Plot for (a)",fontsize=24)
36 plt.xlabel("X values",fontsize=24)
37 plt.ylabel("Y values",fontsize=24)
38 plt.show()
39 plt.figure()
40 plt.xlabel("X values",fontsize=24)
41 plt.ylabel("Y values",fontsize=24)
```

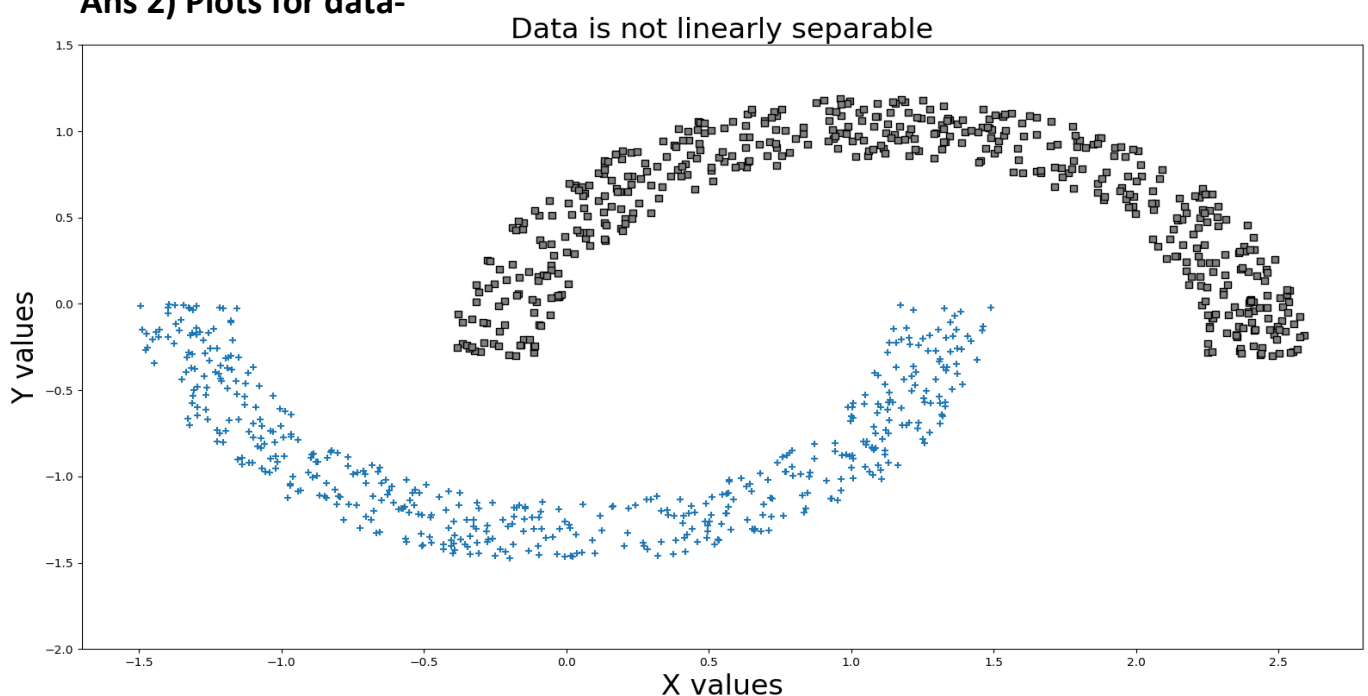
```

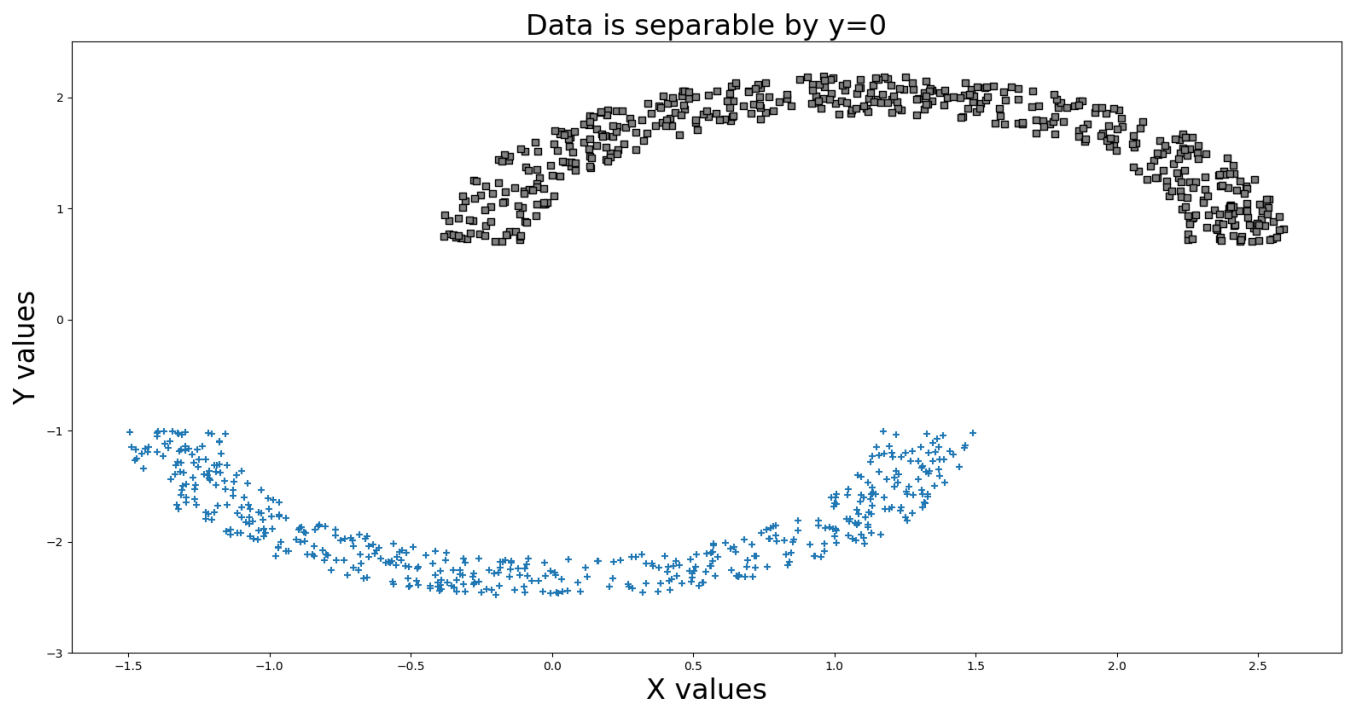
32 y2=[0,2,-3,4,0,1,2,-1,-2,2,1,-2,9,2]
33 y3=[5,-1,0,1,2,2,4,2,1,2,1,2,-4,4]
34 plt.scatter(x1,y1,c="pink",linewidths=2,marker="s",edgecolor="green",s=50)
35 plt.title("Plot for (a)",fontsize=24)
36 plt.xlabel("X values",fontsize=24)
37 plt.ylabel("Y values",fontsize=24)
38 plt.show()
39 plt.figure()
40 plt.xlabel("X values",fontsize=24)
41 plt.ylabel("Y values",fontsize=24)
42 plt.title("Plot for (b)",fontsize=24)
43 plt.scatter(x2,y2,c="yellow",linewidths=2,marker="^",edgecolor="red",s=100)
44 plt.show()
45 plt.figure()
46 plt.scatter(x3,y3,c="blue",linewidths=2,marker="D",edgecolor="yellow",s=50,label="plot for (c)")
47 plt.xlabel("X values",fontsize=24)
48 plt.title("Plot for (c)",fontsize=24)
49 plt.ylabel("Y values",fontsize=24)
50 plt.show()
51
52 #Q1(B) Correlation values
53
54 print("Correlation for (a)",correlation(x1, y1))
55 print("Correlation for (b)",correlation(x2,y2))
56 print("Correlation for (c)",correlation(x3,y3))
57

```

**Q2)** Consider samples belonging to the data of two categories (two different classes) as given below (Each sample is 2D and each class is having 500 samples). Plot this 2D data as a scatter plot, representing the data for each category in a different colour. See if the data for two categories is linearly separated (two classes are linearly separated if a single line can separate those two classes). If not, then perform the following operation on the given samples. Calculate the mean of the individual classes. If the mean of a class lies in the first quadrant (i.e., x and y of the mean are positive) then add +k in y - dimension of each sample. If the mean of a class lies in the third quadrant, then add -k in the y-dimension of each sample. Decide k yourself. Again, plot this 2-class data and comment about the separability of the 2 classes.

**Ans 2) Plots for data-**





Code for Q2 is shown-

```

1  # -*- coding: utf-8 -*-
2  """
3  Created on Thu Apr  8 22:06:33 2021
4  |
5  @author: verma
6  """
7  import pandas as pd
8  import matplotlib.pyplot as plt
9
10
11 # Q2
12 # Reading Data
13 foo=pd.read_csv("Class1.txt",sep=" ",header=None)
14 foo.columns=["x","y"]
15 class2=pd.read_csv("Class2.txt",sep=" ",header=None)
16 class2.columns=["x1","y1"]
17 x1=list(foo.x)
18 y1=list(foo.y)
19 x2=list(class2.x1)
20 y2=list(class2.y1)
21 plt.ylim(-2.0,1.5)
22 plt.scatter(x1,y1,marker="s",c="grey",edgecolor='black')
23 plt.title("Data is not linearly separable",fontsize=24)
24 plt.xlabel("X values",fontsize=24)
25 plt.ylabel("Y values",fontsize=24)
26 plt.scatter(x2,y2,marker="+")
27 mean_x1=sum(x1)/len(x1)
28 mean_y1=sum(y1)/len(y1)
29 mean_x2=sum(x2)/len(x2)
30 mean_y2=sum(y2)/len(y2)
31 k=1
32 if mean_x1>0 and mean_y1>0 :
33     for i in range(0,500):
34         y1[i]+=k
35 elif mean_x1<0 and mean_y1<0 :
36     for j in range(0,500):
37         y1[j]-=k
38 if mean_x2>0 and mean_y2>0 :
39     for l in range(0,500):
40         y2[l]+=k
41 elif mean_x2<0 and mean_y2<0 :

```

```

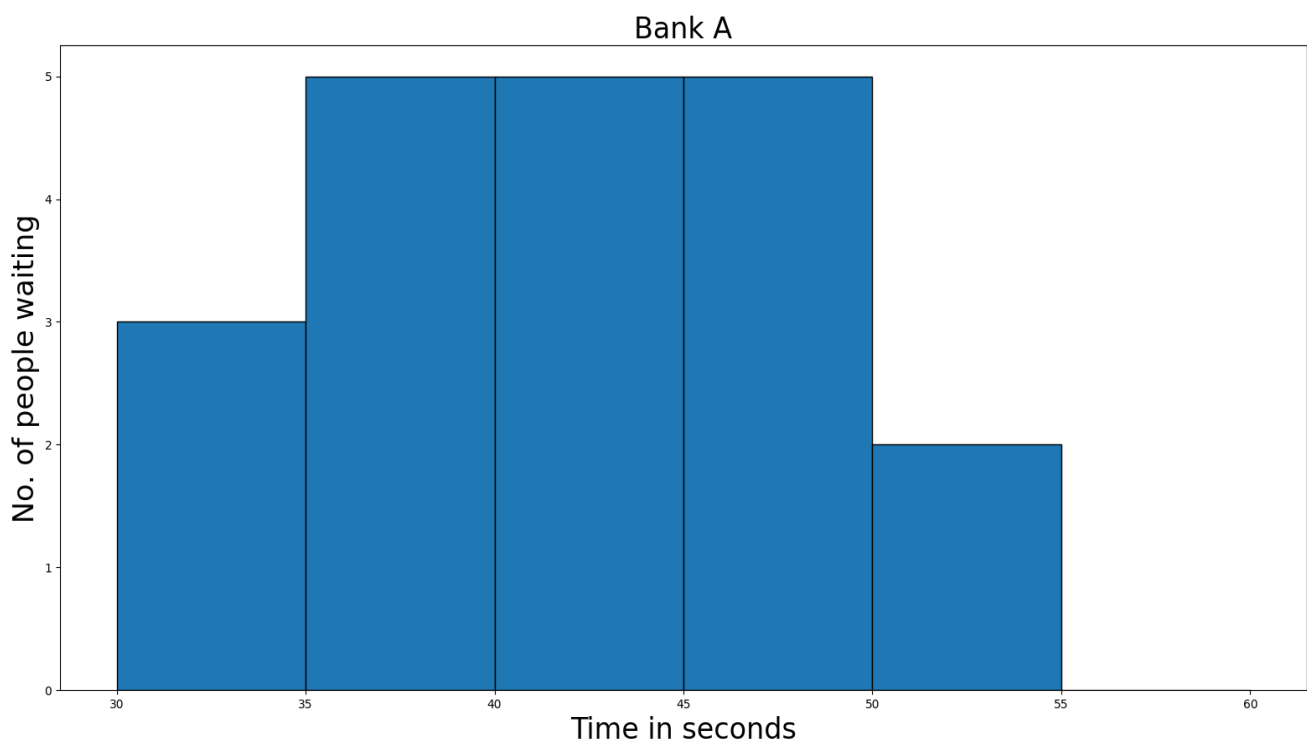
28 mean_y1=sum(y1)/len(y1)
29 mean_x2=sum(x2)/len(x2)
30 mean_y2=sum(y2)/len(y2)
31 k=1
32 if mean_x1>0 and mean_y1>0 :
33     for i in range(0,500):
34         y1[i]+=k
35 elif mean_x1<0 and mean_y1<0 :
36     for j in range(0,500):
37         y1[j]-=k
38 if mean_x2>0 and mean_y2>0 :
39     for l in range(0,500):
40         y2[l]+=k
41 elif mean_x2<0 and mean_y2<0 :
42     for m in range(0,500):
43         y2[m]-=k
44 plt.figure()
45 plt.xlabel("X values",fontsize=24)
46 plt.ylabel("Y values",fontsize=24)
47 plt.ylim(-3.0,2.5)
48 plt.title("Data is separable by y=0",fontsize=24)
49 plt.scatter(x1,y1,marker="s",c="grey",edgecolor='black')
50 plt.scatter(x2,y2,marker="+")
51
52
53

```

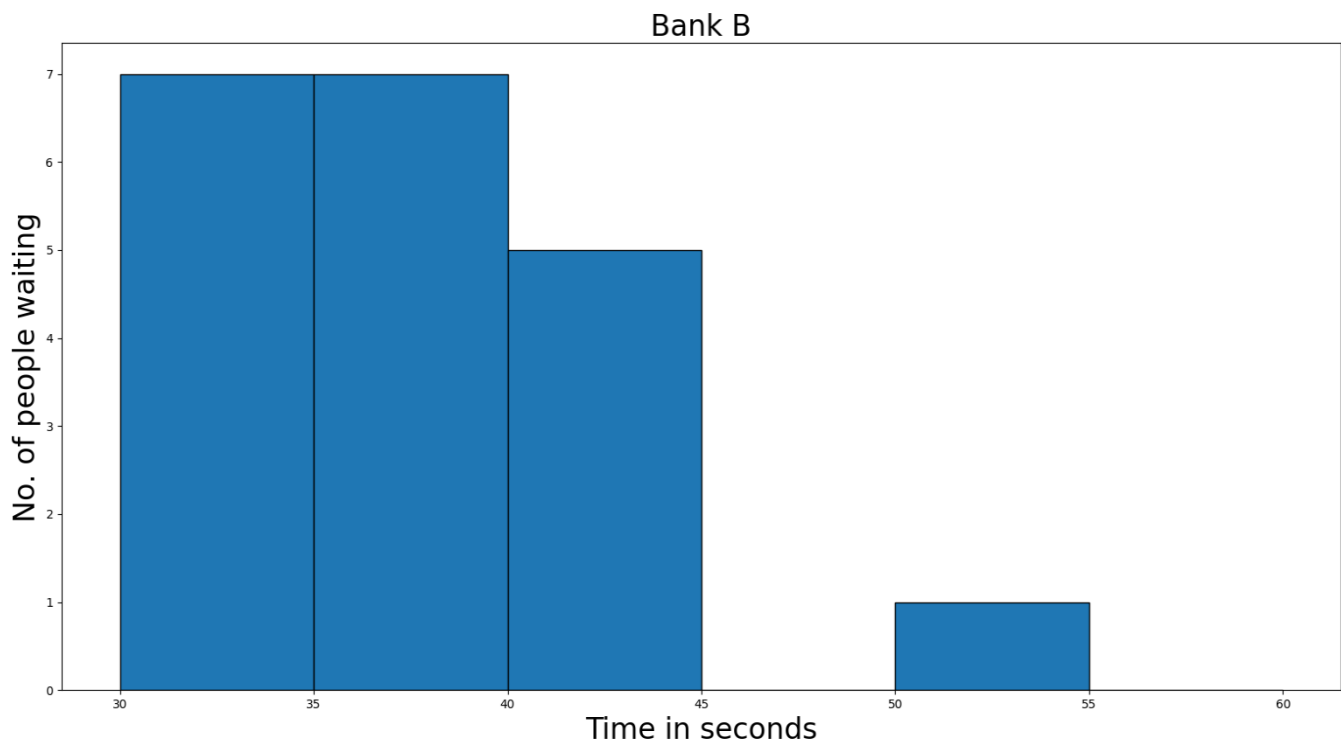
We have taken  $k=1$  to make data separable.

**Q3)** A survey has been conducted by a private firm in between customer care services of two banks. The bank in which the waiting time for a client to be served by a customer service representative is less will be declared as a better bank. The data for 20 customers for the waiting time (in seconds) in bank A & bank B is given. Plot the histogram of the above data with 5-second bins (5-second intervals). Analyse the histogram and argue which bank is giving better services to its customer.

**Ans 3)** Histogram for Bank A is shown -



## Histogram for Bank B-



From both the histogram we can infer that it is the Bank B, which is better than Bank A, since the waiting time of majority of the people is less in comparison to that of Bank A. Hence, Bank B is providing better services to its customer.

## Code for Histograms-

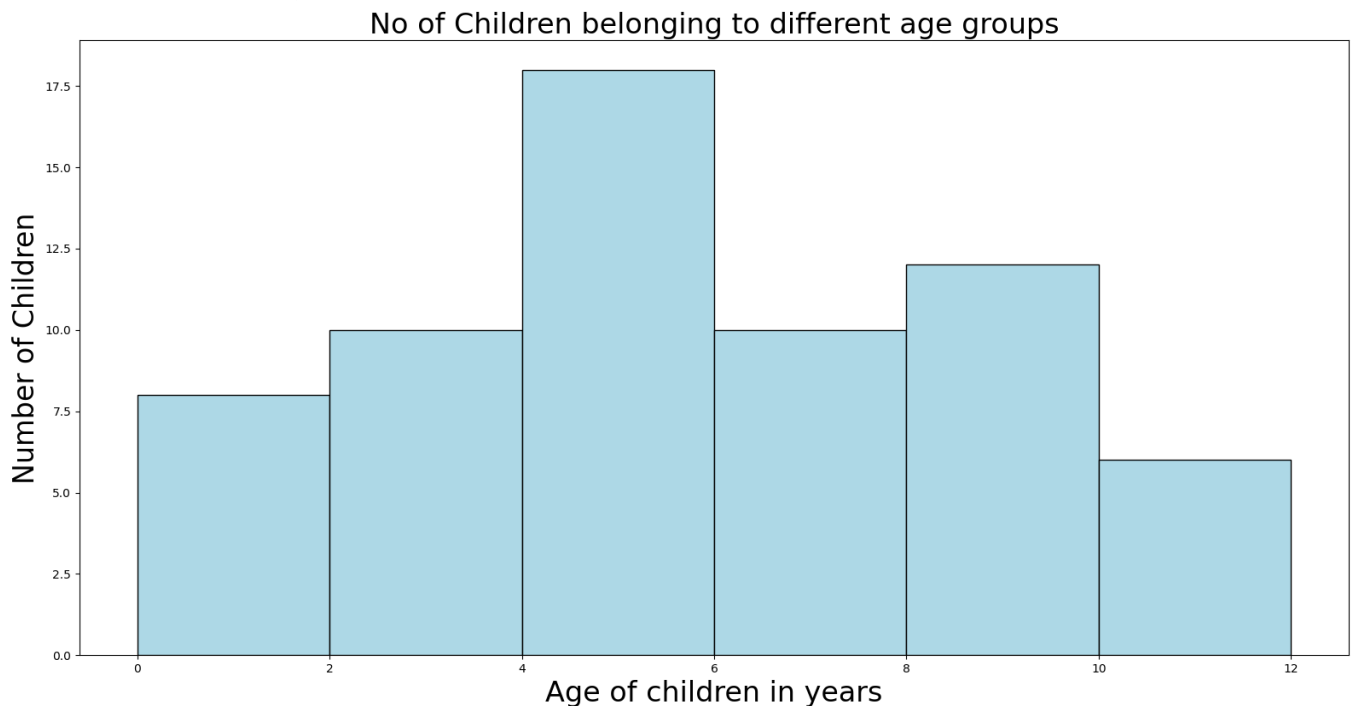
```
1  # -*- coding: utf-8 -*-
2  """
3  Created on Fri Apr  9 09:42:54 2021
4
5  @author: verma
6  """
7  # Q3
8  import matplotlib.pyplot as plt
9
10 Bank_A=[43.1,45.3,47.3,30.3,54.1,
11         35.6,43.5,31.2,31.4,45.6,
12         37.6,40.3,42.2,35.6,36.5,
13         36.5,50.2,45.5,45.2,43.1]
14 Bank_B=[33.1,35.3,37.3,32.3,44.1,
15         31.6,33.5,35.2,30.4,35.6,
16         32.6,42.3,40.2,36.6,38.5,
17         36.5,30.2,42.5,50.2,40.1]
18 plt.figure()
19 plt.hist(Bank_A,bins=[30,35,40,45,50,55,60],edgecolor="black")
20 plt.title("Bank A",fontsize=24)
21 plt.xlabel("Time in seconds",fontsize=24)
22 plt.ylabel("No. of people waiting",fontsize=24)
23 plt.show()
24 plt.figure()
25 plt.title("Bank B",fontsize=24)
26 plt.xlabel("Time in seconds",fontsize=24)
27 plt.ylabel("No. of people waiting",fontsize=24)
28 plt.hist(Bank_B,bins=[30,35,40,45,50,55,60],edgecolor="black")
29 plt.show()
```

**Q4)** The information about the number of children belonging to different age groups tabulated in the table given below.

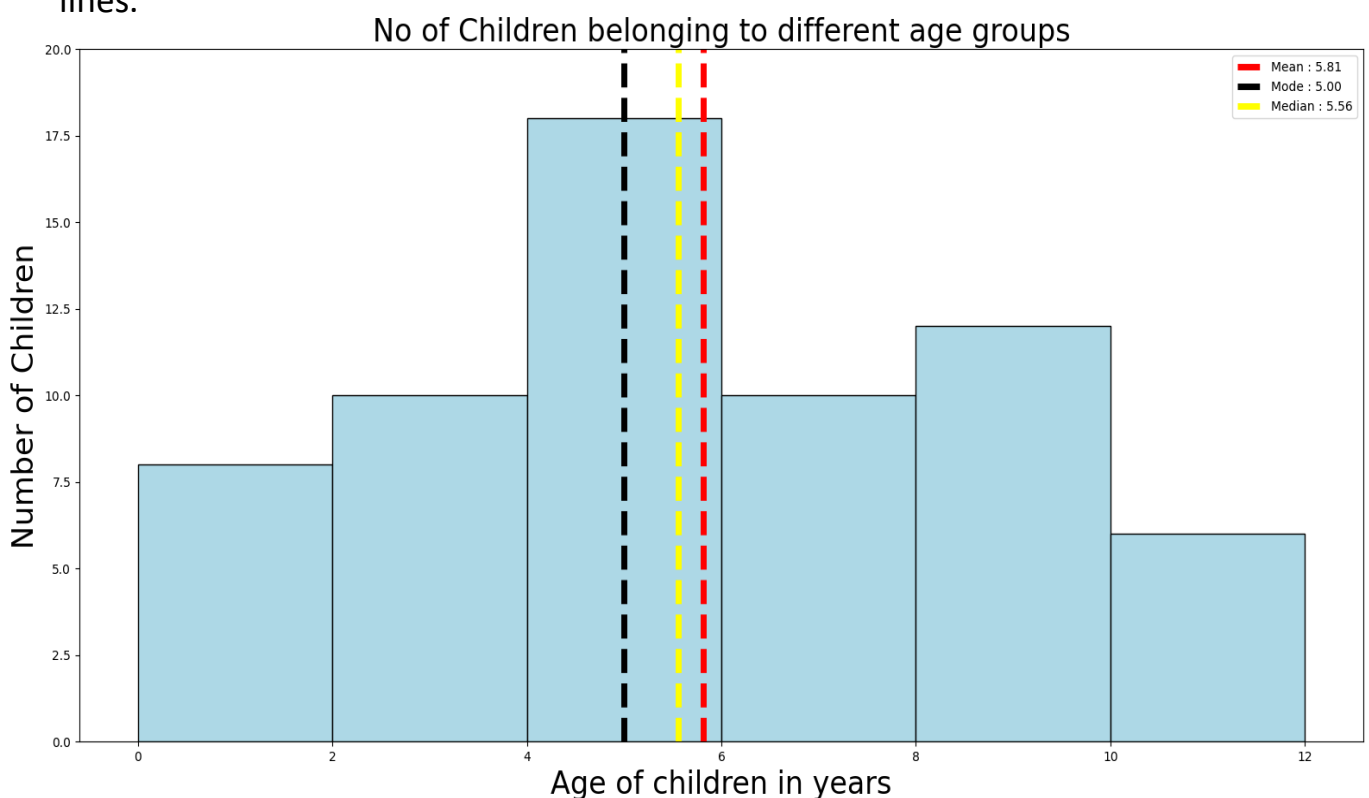
a) Plot this table as a graph, which is essentially the histogram representing the data.

b) Mark the points of mean, mode and median on the histogram, by vertical lines.

**Ans 4) a)** Histogram for table data –



**b)** Histogram marked with mean , mode & median on the histogram by vertical lines.





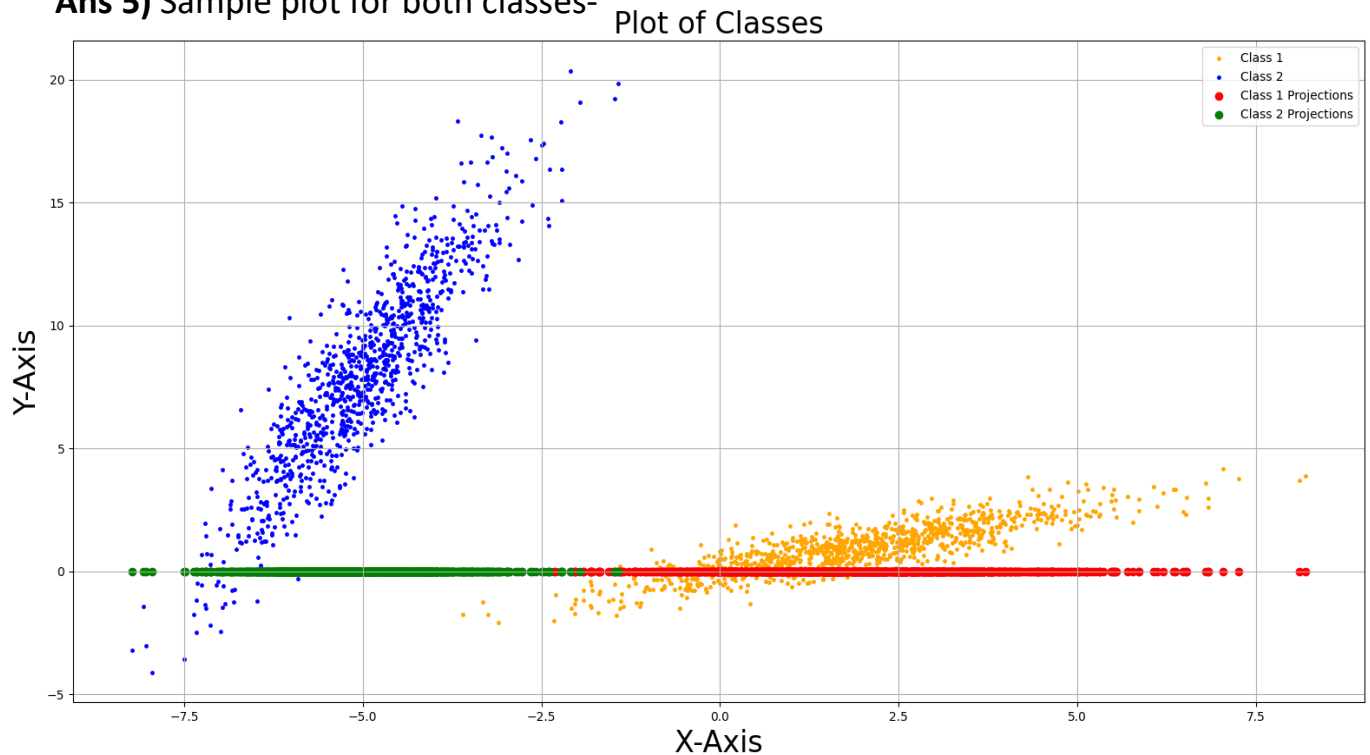
## Code for Histograms-

```
1  # -*- coding: utf-8 -*-
2  """
3  Created on Fri Apr  9 10:07:26 2021
4
5  @author: verma
6  """
7  import matplotlib.pyplot as plt
8
9  def mean(x,y):
10     k=0
11     for i,j in zip(x,y):
12         k+=i*j
13     mn=k/(sum(y)*1.0)
14     return mn
15  def mode(x,y):
16     u=max(y)
17     s=y.index(u)
18     c=x[s]-1
19     f= c+ ((y[s]-y[s-1])/(2*y[s]-y[s-1]-y[s+1]))*2
20     return f
21  def median(x,y):
22     n=sum(y)
23     k=int(n/2)
24     a=0
25     for i in range(0,len(y)):
26         a+=y[i]
27         v=i
28         if a>k:break
29     med=(x[v]-1)+ (((n/2)-(a-y[v]))*2/y[v])
30     return med
31
32  # Q4(a)
33  Age=[1,3,5,7,9,11]
34  no_child=[8,10,18,10,12,6]
35  plt.bar(Age,no_child,width=2,edgecolor="black",color='lightblue',align='center')
36  plt.xlabel("Age of children in years",fontsize=24)
37  plt.ylabel("Number of Children",fontsize=24)
38  plt.title("No of Children belonging to different age groups",fontsize=24)
39  plt.show()
40
```

```
27     v=i
28     if a>k:break
29     med=(x[v]-1)+ (((n/2)-(a-y[v]))*2/y[v])
30     return med
31
32  # Q4(a)
33  Age=[1,3,5,7,9,11]
34  no_child=[8,10,18,10,12,6]
35  plt.bar(Age,no_child,width=2,edgecolor="black",color='lightblue',align='center')
36  plt.xlabel("Age of children in years",fontsize=24)
37  plt.ylabel("Number of Children",fontsize=24)
38  plt.title("No of Children belonging to different age groups",fontsize=24)
39  plt.show()
40
41  # Q4 (b)
42  plt.axvline(mean(Age,no_child),color='red',linestyle='dashed',linewidth='5',label="Mean : 5.81")
43  plt.axvline(mode(Age,no_child),color='black',linestyle='dashed',linewidth='5',label='Mode : 5.00')
44  plt.axvline(median(Age,no_child),color='yellow',linestyle='dashed',linewidth='5',label='Median : 5.56')
45  plt.legend()
46  plt.ylim(0,20)
47  plt.show()
48
```

**Q5)** Consider the given two class linearly separable data (i.e., the two classes can be separated with the help of a line). Project each point of this data on the line  $x = 0$  (i.e., on x-axis). After projection, plot the samples and comment on the separability of the two classes.

**Ans 5)** Sample plot for both classes-



From the plot we can observe that the Class 1 projections and Class 2 projections are not linearly separable since they overlap.

Code for plot-

```

5  @author: verma
6  """
7
8  import matplotlib.pyplot as plt
9
10 # Q5
11 # Reading Data
12 class1=open("C:/Users/verma/OneDrive/Desktop/LAB 6/linearly_seperable_data/Class1.txt",'r+')
13 class2=open("C:/Users/verma/OneDrive/Desktop/LAB 6/linearly_seperable_data/Class2.txt",'r+')
14 x1=[]
15 y1=[]
16 for line in class1.readlines():
17     l,k=line.split()
18     x1.append(float(l))
19     y1.append(float(k))
20
21 x2=[]
22 y2=[]
23 for line in class2.readlines():
24     j,k=line.split()
25     x2.append(float(j))
26     y2.append(float(k))
27

```

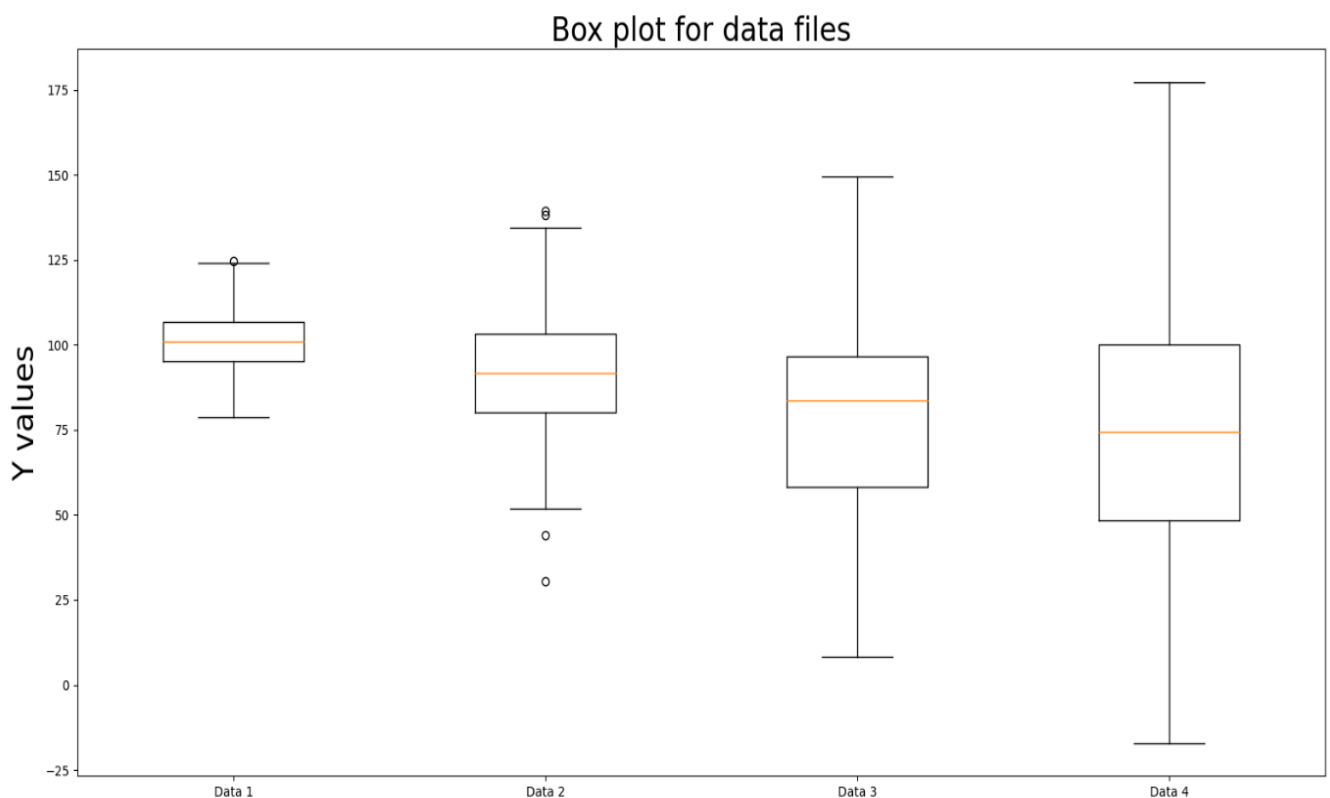
```

23 for line in class2.readlines():
24     j,k=line.split()
25     x2.append(float(j))
26     y2.append(float(k))
27
28 plt.scatter(x1,y1,color="orange",marker='.',label="Class 1",s=25)
29 plt.scatter(x2,y2,color='b',marker='.',label="Class 2",s=25)
30 plt.xlabel("X-Axis",fontsize=24)
31 plt.ylabel("Y-Axis",fontsize=24)
32 plt.title("Plot of Classes",fontsize=24)
33 plt.grid()
34 plt.show()
35
36 #Plotting projections
37
38 for i in range(0,len(x1)):
39     y1[i]=0
40     y2[i]=0
41 plt.scatter(x1,y1,color='red',label="Class 1 Projections")
42 plt.scatter(x2,y2,color='green',label="Class 2 Projections")
43 plt.legend()
44 plt.show()
45

```

**Q6)** Given data (data.csv) is a single column data. Draw the box plot of this single column data to find out the range in which maximum data is lying (Interquartile range).

**Ans)** Box Plot for the data is shown below-



## The interquartile range for the boxplot-

```
In [37]: runfile('C:/Users/verma/OneDrive/Desktop/LAB 6/q6_boxplot/untitled8.py', wdir='C:/Users/verma/OneDrive/Desktop/LAB 6/q6_boxplot')
Reloaded modules: jupyter_client.session, zmq.eventloop, zmq.eventloop.ioloop, tornado.platform, tornado.platform.asyncio, tornado.gen, zmq.eventloop.zmqstream, jupyter_client.jsonutil, jupyter_client.adapter, spyder, spyder.pil_patch, PIL, PIL_version, PIL.Image, PIL.ImageMode, PIL.TiffTags, PIL._binary, PIL._util, PIL._imaging, cffi, cffi.api, cffi.lock, cffi.error, cffi.model, IPython.lib.guisupport, IPython.external, IPython.external.qt_for_kernel, IPython.utils.version, IPython.external.qt_loaders
Interquartile range for data is 31.19

In [38]:
```

## Code for Boxplot-

```
1  # -*- coding: utf-8 -*-
2
3  Created on Sat Apr 10 22:37:08 2021
4
5  @author: verma
6  """
7
8  import pandas as pd
9  import matplotlib.pyplot as plt
10
11  # Q6
12  foo=pd.read_csv("d1.txt",sep=" ",header=None)
13  foo1=pd.read_csv("d2.txt",sep=" ",header=None)
14  foo2=pd.read_csv("d3.txt",sep=" ",header=None)
15  foo3=pd.read_csv("d4.txt",sep=" ",header=None)
16  foo.columns=["x"]
17  foo1.columns=["x"]
18  foo2.columns=["x"]
19  foo3.columns=["x"]
20
21  x1=list(foo.x)
22  x2=list(foo1.x)
23  x3=list(foo2.x)
24  x4=list(foo3.x)
25  data=[x1,x2,x3,x4]
26  plt.title("Box plot for data files",fontsize=24)
27  plt.ylabel("Values",fontsize=24)
28  plt.boxplot(data,labels=["Data 1","Data 2","Data 3","Data 4"])
29  plt.show()
```

## Code for Interquartile range-

```
29  plt.boxplot(data,labels=["Data 1","Data 2","Data 3","Data 4"])
30  plt.show()
31  # for Inter quartile range
32  Q1 = np.percentile(data, 25, interpolation = 'midpoint')
33  Q2= np.percentile(data, 50, interpolation = 'midpoint')
34  Q3 = np.percentile(data, 75, interpolation = 'midpoint')
35  Interquartile= Q3-Q1
36  print("Interquartile range for data is ",round(Interquartile,2))
37
```