

# CS 584-04: Machine Learning

Fall 2019 Assignment 1

## Question 1 (40 points)

Write a Python program to calculate the density estimator of a histogram. Use the field x in the NormalSample.csv file.

- a) (5 points) According to Izenman (1991) method, what is the recommended bin-width for the histogram of x?

**Aa) 0.3998667554864774**

- b) (5 points) What are the minimum and the maximum values of the field x?

**Ab) Minimum(x)=26.3 and Maximum(x)=35.4**

- c) (5 points) Let a be the largest integer less than the minimum value of the field x, and b be the smallest integer greater than the maximum value of the field x. What are the values of a and b?

**Ac) a=26 and b=36**

- d) (5 points) Use  $h = 0.1$ , minimum = a and maximum = b. List the coordinates of the density estimator. Paste the histogram drawn using Python or your favorite graphing tools.

**Ad)**

m	p(m)
26.05	0
26.15	0
	0.00999
26.25	001
	0.00999
26.35	001
26.45	0
26.55	0
26.65	0
26.75	0
26.85	0

26.95	0
27.05	0
	0.00999
27.15	001
	0.00999
27.25	001
27.35	0
27.45	0
27.55	0
	0.01998
27.65	002
	0.01998
27.75	002

	0.02997
27.85	003
	0.03996
27.95	004
	0.01998
28.05	002
	0.05994
28.15	006
	0.07992
28.25	008
	0.04995
28.35	005
	0.05994
28.45	006

	0.07992
28.55	008
	0.08991
28.65	009
	0.12987
28.75	013
	0.12987
28.85	013
	0.09990
28.95	01
	0.08991
29.05	009
	0.14985
29.15	015
	0.25974
29.25	026
	0.22977
29.35	023
	0.21978
29.45	022
	0.32967
29.55	033
	0.27972
29.65	028
	0.18981
29.75	019
	0.30969
29.85	031
	0.36963
29.95	037
	0.30969
30.05	031
	0.36963
30.15	037

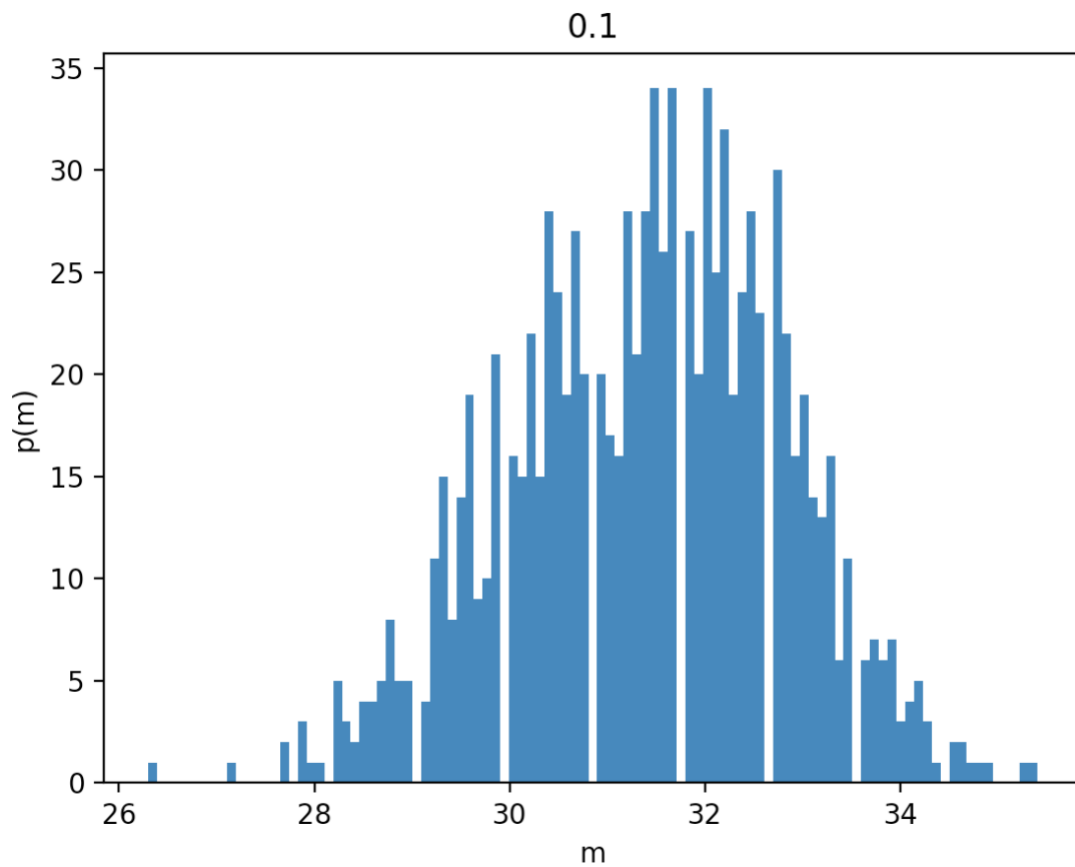
	0.36963
30.25	037
	0.42957
30.35	043
	0.51948
30.45	052
	0.42957
30.55	043
	0.45954
30.65	046
	0.46953
30.75	047
	0.39960
30.85	04
	0.36963
30.95	037
	0.32967
31.05	033
	0.43956
31.15	044
	0.48951
31.25	049
	0.48951
31.35	049
	0.61938
31.45	062
	0.59940
31.55	06
	0.59940
31.65	06
	0.60939
31.75	061
	0.46953
31.85	047

	0.53946
31.95	054
	0.33966
32.05	034
	0.31968
32.15	032
	0.50949
32.25	051
	0.18981
32.35	019
	0.27972
32.45	028
	0.27972
32.55	028
	0.29970
32.65	03
	0.51948
32.75	052
	0.21978
32.85	022
	0.18981
32.95	019
	0.18981
33.05	019
	0.12987
33.15	013
	0.28971
33.25	029
	0.15984
33.35	016
	0.10989
33.45	011
	0.10989
33.55	011

	0.06993
33.65	007
	0.12987
33.75	013
	0.05994
33.85	006
	0.02997
33.95	003
	0.02997
34.05	003
	0.04995
34.15	005
	0.07992
34.25	008

	0.02997
34.35	003
	0.01998
34.45	002
	0.01998
34.55	002
	0.00999
34.65	001
	0.01998
34.75	002
	0.00999
34.85	001
	0
34.95	0
	0
35.05	0

35.15	0
	0.00999
35.25	001
	0.00999
35.35	001
	0
35.45	0
	0
35.55	0
	0
35.65	0
	0
35.75	0
	0
35.85	0
	0
35.95	0



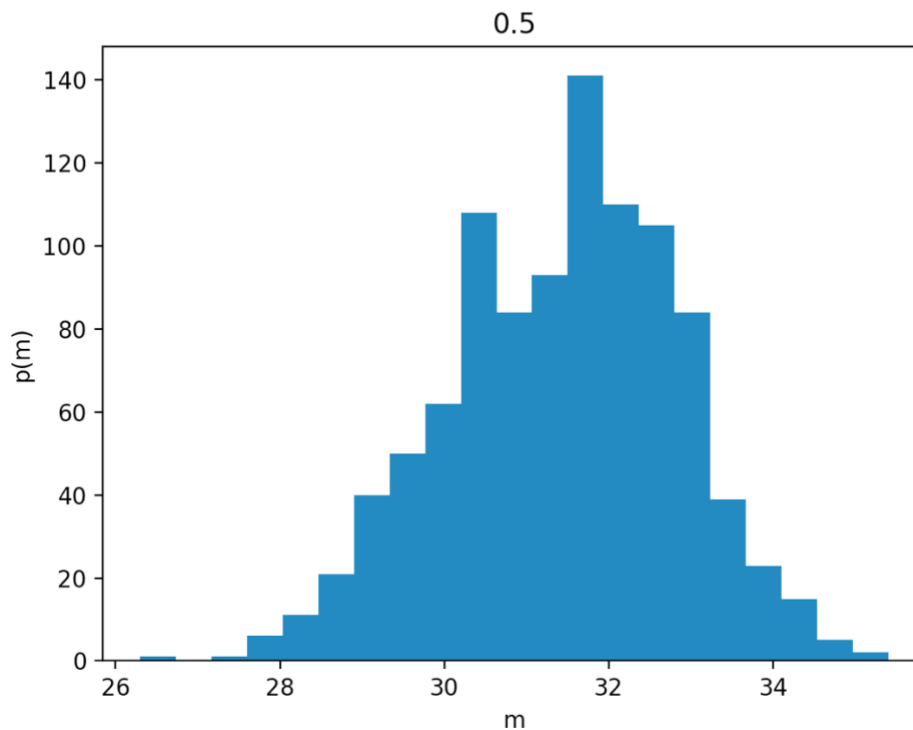
- e) (5 points) Use  $h = 0.5$ , minimum = a and maximum = b. List the coordinates of the density estimator. Paste the histogram drawn using Python or your favorite graphing tools.

Ae)

m	p(m)
	0.00199
26.25	8
26.75	0
	0.00199
27.25	8
	0.01198
27.75	801
	0.03196
28.25	803
	0.06193
28.75	806
	0.11388
29.25	611

	0.17782
29.75	218
	0.23976
30.25	024
	0.25374
30.75	625
	0.28771
31.25	229
	0.34965
31.75	035
	0.32367
32.25	632
	0.27572
32.75	428

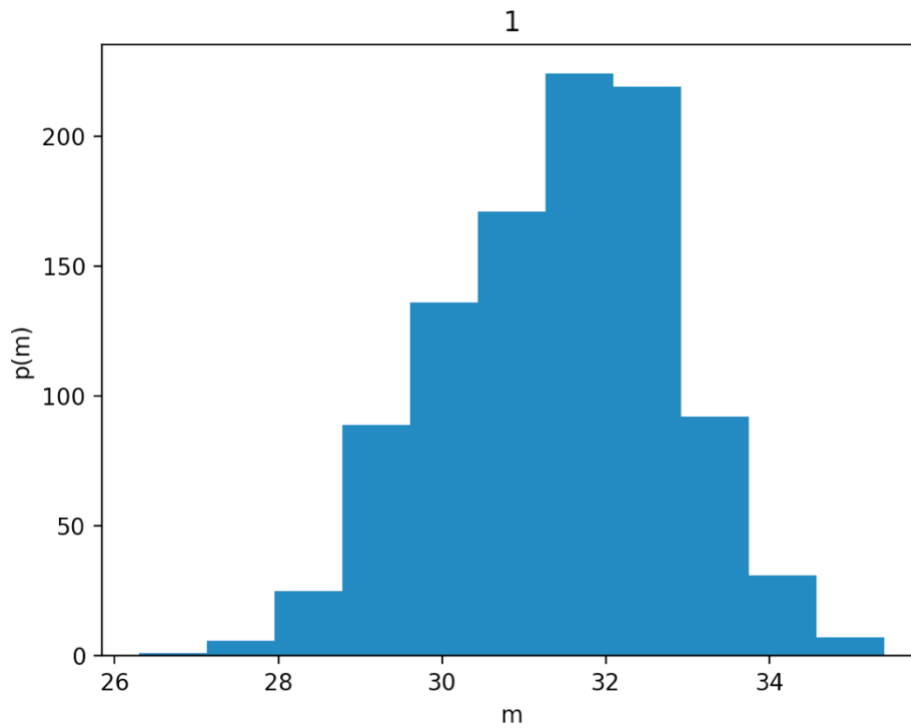
	0.15784
33.25	216
	0.07992
33.75	008
	0.03596
34.25	404
	0.01398
34.75	601
	0.00399
35.25	6
35.75	0



- f) (5 points) Use  $h = 1$ , minimum = a and maximum = b. List the coordinates of the density estimator. Paste the histogram drawn using Python or your favorite graphing tools.

Af)

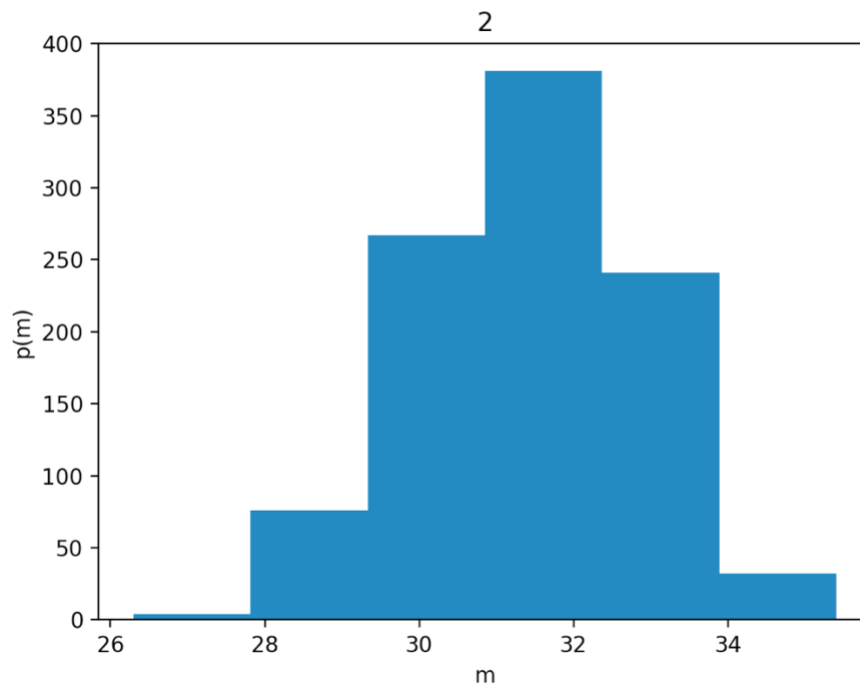
m	p(m)
26.5	0.000999
27.5	0.00699301
28.5	0.04295704
29.5	0.13186813
30.5	0.22277722
31.5	0.28471528
32.5	0.27172827
33.5	0.10789211
34.5	0.02297702
35.5	0.001998



- g) (5 points) Use  $h = 2$ , minimum = a and maximum = b. List the coordinates of the density estimator. Paste the histogram drawn using Python or your favorite graphing tools.

Ag)

m	p(m)
27	0.003996
29	0.08491508
31	0.24525475
33	0.18031968
35	0.01248751



- h) (5 points) Among the four histograms, which one, in your honest opinions, can best provide your insights into the shape and the spread of the distribution of the field  $x$ ? Please state your arguments.

**Ah) I feel that the histogram with the bin-width of 0.5 can provide better insight to the data distribution compared to the others. It doesn't have too many information unlike  $h=0.1$ , which has too many data points and is difficult to form ideas from the graph, while  $h=0.5$  gives a simpler view, but not too simple like the  $h=1$  and  $h=2$ , so that some sort of assessment can be done from the graph.**

## Question 2 (20 points)

Use in the NormalSample.csv to generate box-plots for answering the following questions.

- a) (5 points) What is the five-number summary of x? What are the values of the 1.5 IQR whiskers?

**Aa) median: 31.5**

**Min: 26.3**

**Max: 35.4**

**A: 26**

**B: 36**

**q1: 30.4**

**q2: 32.4**

**IQR 2.0**

**W1 27.4**

**W2 35.4**

- b) (5 points) What is the five-number summary of x for each category of the group? What are the values of the 1.5 IQR whiskers for each category of the group?

**Ab)**

**N:686**

\*\*\*\*\*

**Median: 32.1**

**Min: 29.1**

**Max: 35.4**

**A: 29**

**B: 36**

**q1: 31.4**

**q2: 32.7**

**IQR 1.30000000000000043**

**W1 29.449999999999992**

**W2 34.650000000000006**

**N: 315**

\*\*\*\*\*

**Median: 30.0**

**Min: 26.3**

**Max: 32.2**

**A: 26**

**B: 33**

**q1: 29.4**

**q2: 30.6**

**IQR: 1.20000000000000028**

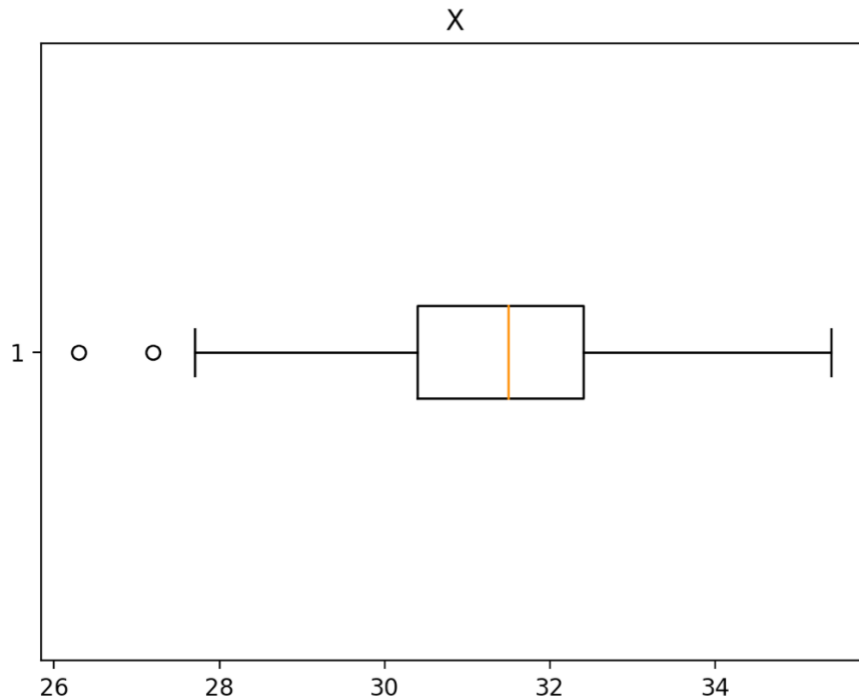
**W1: 27.599999999999994**

**W2: 32.400000000000006**



- c) (5 points) Draw a boxplot of  $x$  (without the group) using the Python boxplot function. Can you tell if the Python's boxplot has displayed the 1.5 IQR whiskers correctly?

Ac)



Yes, it can be clearly seen that there two outliers, i.e, 26.3, and 27.2.

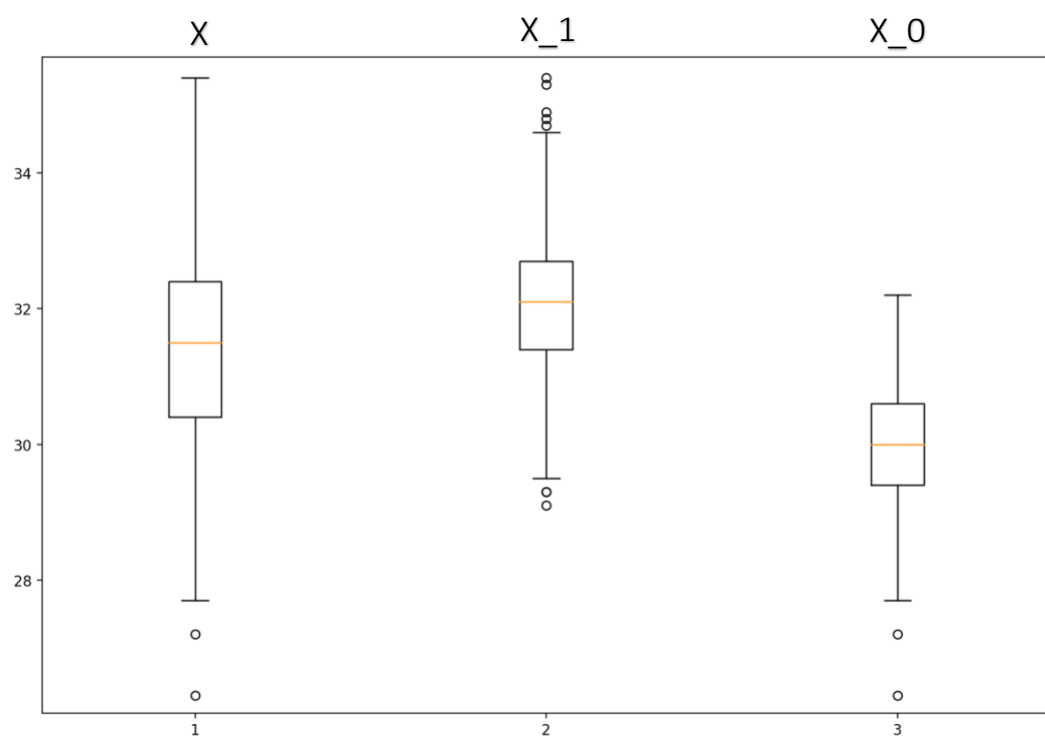
- d) (5 points) Draw a graph where it contains the boxplot of  $x$ , the boxplot of  $x$  for each category of Group (i.e., three boxplots within the same graph frame). Use the 1.5 IQR whiskers, identify the outliers of  $x$ , if any, for the entire data and for each category of the group.

*Hint: Consider using the CONCAT function in the PANDA module to append observations.*

Ad) **The outliers of  $X$  are: (26.3, 27.2)**

**The outliers for  $X_1$  are: (29.1, 29.3, 29.3) and (35.3, 35.4, 34.9, 34.7, 34.8)**

**The outliers for  $X_0$  are: 27.2, 26.3)**



### Question 3 (40 points)

The data, FRAUD.csv, contains results of fraud investigations of 5,960 cases. The binary variable FRAUD indicates the result of a fraud investigation: 1 = Fraudulent, 0 = Otherwise. The other interval variables contain information about the cases.

1. TOTAL\_SPEND: Total amount of claims in dollars
2. DOCTOR\_VISITS: Number of visits to a doctor
3. NUM\_CLAIMS: Number of claims made recently
4. MEMBER\_DURATION: Membership duration in number of months
5. OPTOM\_PRESC: Number of optical examinations
6. NUM\_MEMBERS: Number of members covered

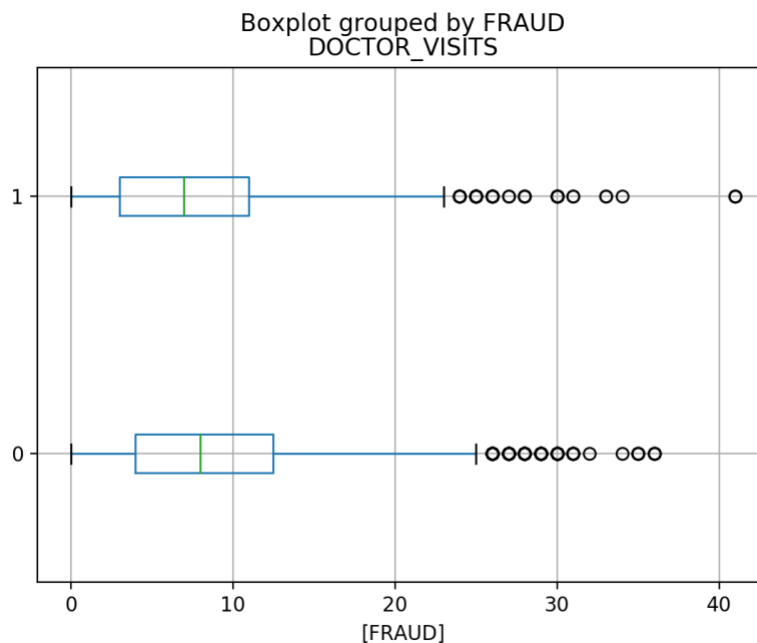
You are asked to use the Nearest Neighbors algorithm to predict the likelihood of fraud.

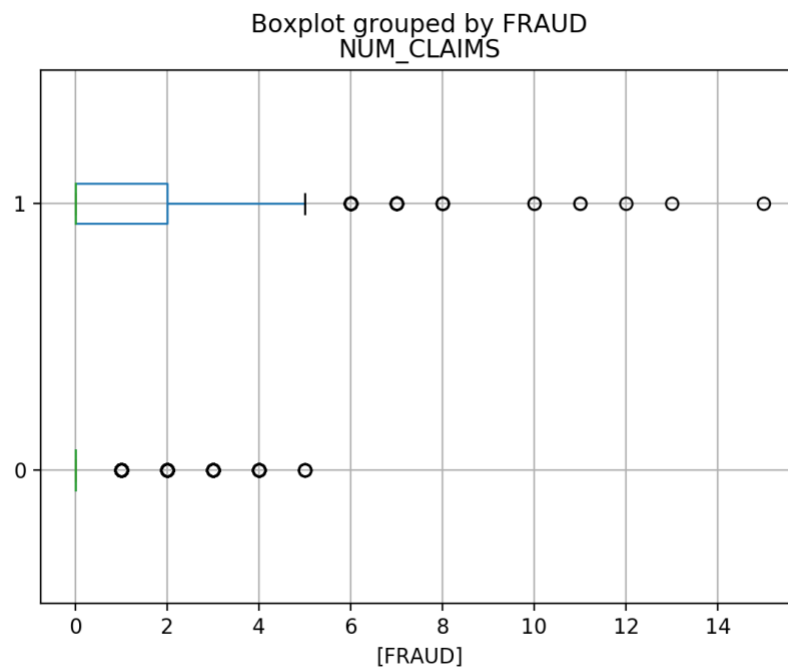
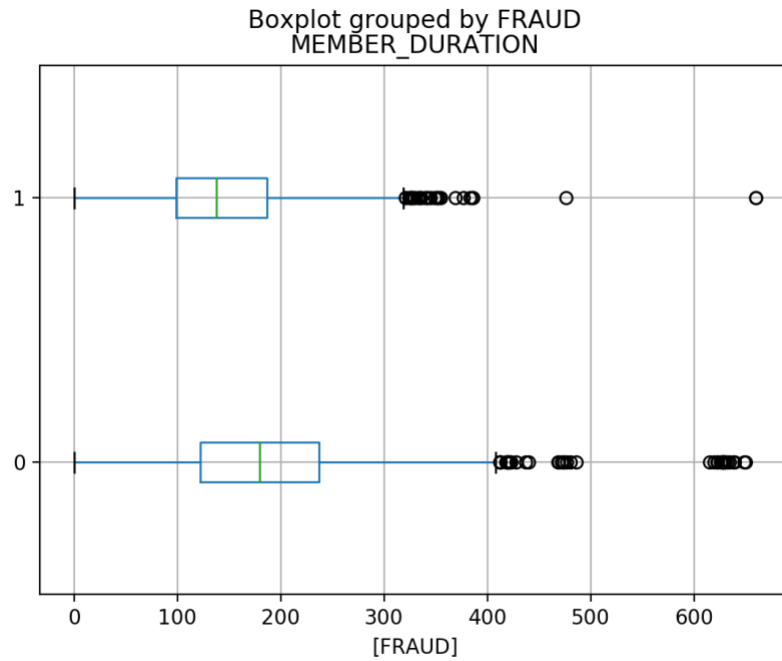
- a) (5 points) What percent of investigations are found to be fraudulent? Please give your answer up to 4 decimal places.

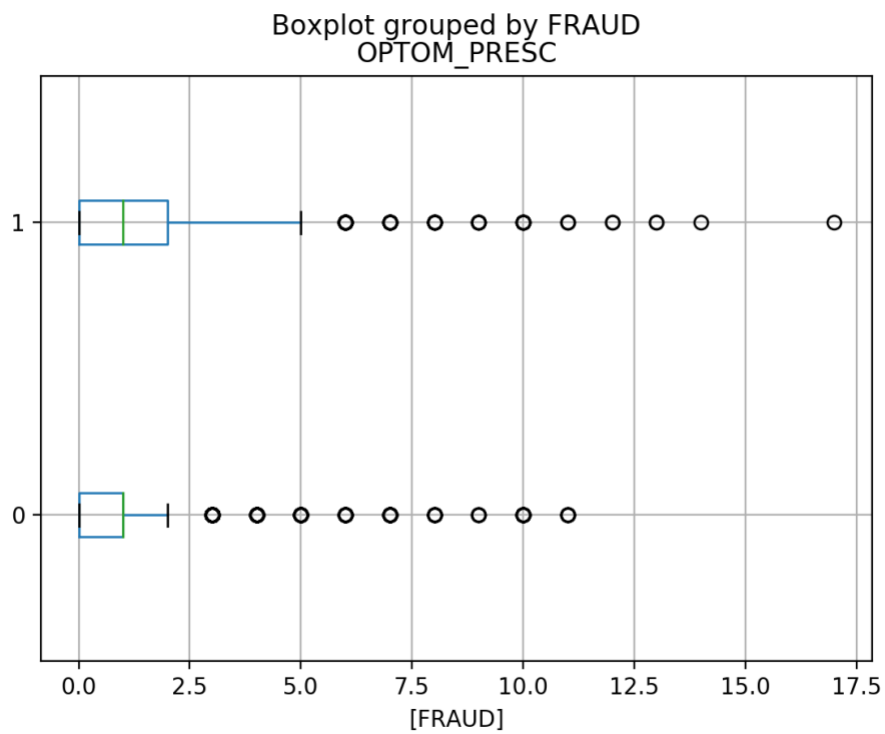
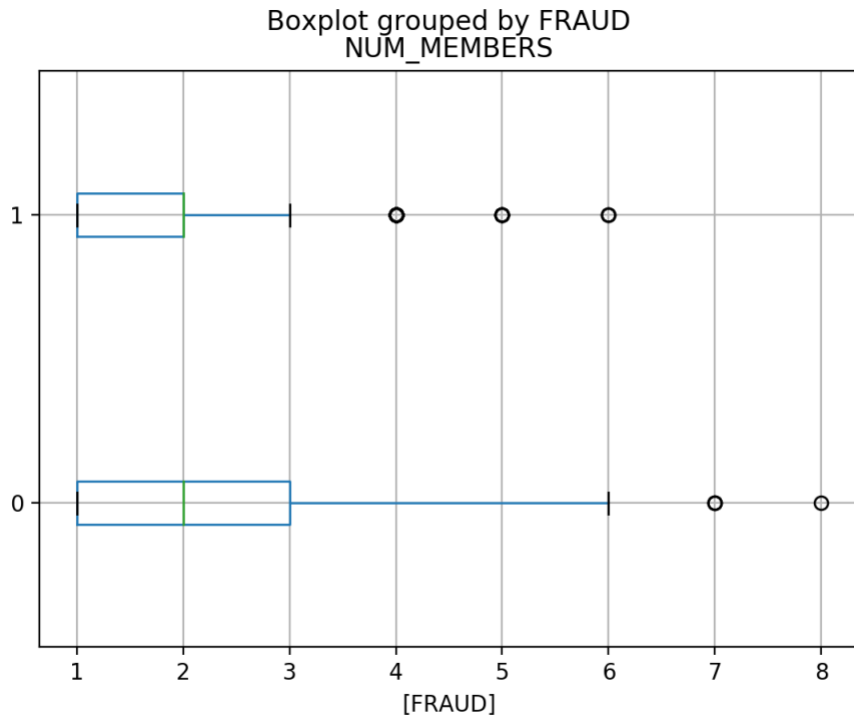
**Aa) Fraudulent percentage: 19.9497%**

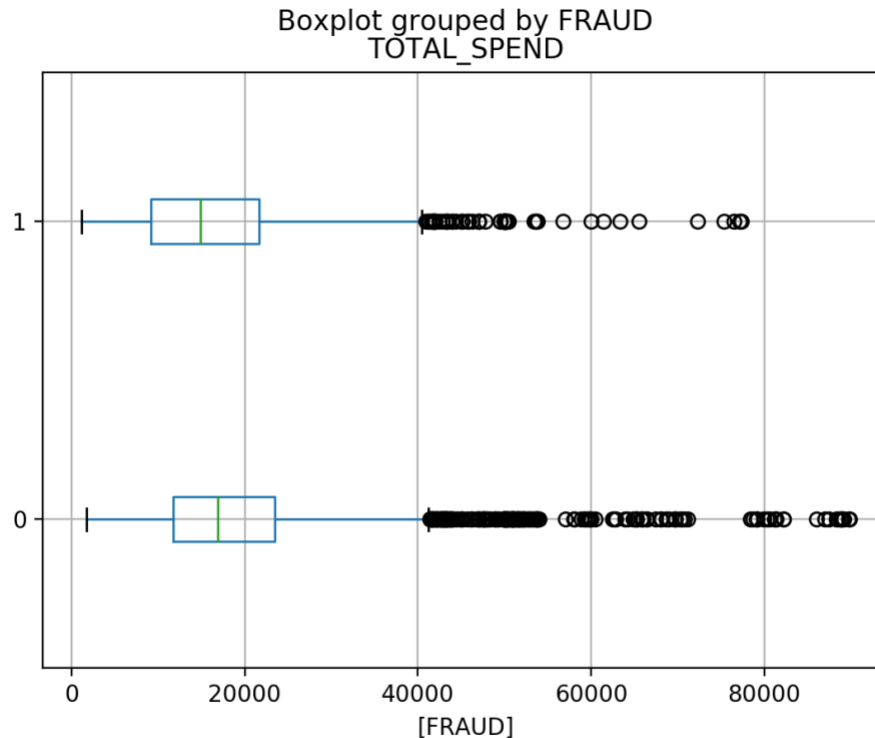
- b) (5 points) Use the BOXPLOT function to produce horizontal box-plots. For each interval variable, one box-plot for the fraudulent observations, and another box-plot for the non-fraudulent observations. These two box-plots must appear in the same graph for each interval variable.

**Ab)**









- c) (10 points) Orthonormalize interval variables and use the resulting variables for the nearest neighbor analysis. Use only the dimensions whose corresponding eigenvalues are greater than one.
- (5 points) How many dimensions are used?
  - (5 points) Please provide the transformation matrix? You must provide proof that the resulting variables are actually orthonormal.

Ac) i. 6 dimensions

## ii. Transformation Matrix

```
[[ -6.49862374e-08  0.00000000e+00  0.00000000e+00 -0.00000000e+00  0.00000000e+00 -0.00000000e+00]
 [ 0.00000000e+00 -2.94741983e-04  0.00000000e+00  0.00000000e+00 -0.00000000e+00 -0.00000000e+00]
 [-0.00000000e+00  0.00000000e+00 -7.68683456e-04  0.00000000e+00 -0.00000000e+00 -0.00000000e+00]
 [ 0.00000000e+00 -0.00000000e+00  0.00000000e+00  5.78327741e-05 -0.00000000e+00 -0.00000000e+00]
 [ 0.00000000e+00 -0.00000000e+00 -0.00000000e+00 -0.00000000e+00 -2.39238772e-07 -0.00000000e+00]
 [ 0.00000000e+00  0.00000000e+00 -0.00000000e+00 -0.00000000e+00 -0.00000000e+00  4.66565230e-11]]
```

If, the transpose of the Transformed\_matrix multiplied by the transformed matrix gives an Identity matrix then the matrix is orthonormal.

i.e,

```
IMat = trans_m.transpose().dot(trans_m)
```

```
IMat = Transpose(Transformed_matrix)*Transformed_matrix:
```

```
IMat =
```

```
[[ 1.00000000e+00 -5.55111512e-17 1.17961196e-16 4.85722573e-17 1.38777878e-17 -3.46944695e-17]
 [-5.55111512e-17 1.00000000e+00 -1.66533454e-16 -2.08166817e-17 6.93889390e-18 -3.46944695e-17]
 [ 1.17961196e-16 -1.66533454e-16 1.00000000e+00 -1.95427442e-17 3.46944695e-18 -3.51281504e-17]
 [ 4.85722573e-17 -2.08166817e-17 -1.95427442e-17 1.00000000e+00 -2.38524478e-16 -1.97758476e-16]
 [ 1.38777878e-17 6.93889390e-18 3.46944695e-18 -2.38524478e-16 1.00000000e+00 -6.85215773e-17]
 [-3.46944695e-17 -3.46944695e-17 -3.51281504e-17 -1.97758476e-16 -6.85215773e-17 1.00000000e+00]]
```

- d) (10 points) Use the NearestNeighbors module to execute the Nearest Neighbors algorithm using exactly five neighbors and the resulting variables you have chosen in c). The KNeighborsClassifier module has a score function.
- (5 points) Run the score function, provide the function return value
  - (5 points) Explain the meaning of the score function return value.

**Ad) i: 0.8825503355704698**

**li: The score gives the accuracy of the predicted data.**

- e) (5 points) For the observation which has these input variable values: TOTAL\_SPEND = 7500, DOCTOR\_VISITS = 15, NUM\_CLAIMS = 3, MEMBER\_DURATION = 127, OPTOM\_PRESC = 2, and NUM\_MEMBERS = 2, find its **five** neighbors. Please list their input variable values and the target values. *Reminder: transform the input observation using the results in c) before finding the neighbors.*

**Ae) The neighbors are: [[ 588 1199 2264 1246 3809]]**

**The distances are:**

```
[[1.64636127e-10 3.78641748e-04 5.36680515e-04 7.69632938e-04 8.81391774e-04]]
```

CASE_ID	589
FRAUD	1
TOTAL_SPEND	7500
DOCTOR_VISITS	15
NUM_CLAIMS	3
MEMBER_DURATION	127
OPTOM_PRESC	2
NUM_MEMBERS	2

CASE\_ID        1200  
 FRAUD           1  
 TOTAL\_SPEND    10000  
 DOCTOR\_VISITS   16  
 NUM\_CLAIMS      3  
 MEMBER\_DURATION 124  
 OPTOM\_PRESC    2  
 NUM\_MEMBERS    1

CASE\_ID        2265  
 FRAUD           1  
 TOTAL\_SPEND    13800  
 DOCTOR\_VISITS   15  
 NUM\_CLAIMS      3  
 MEMBER\_DURATION 121  
 OPTOM\_PRESC    1  
 NUM\_MEMBERS    1

CASE\_ID        1247  
 FRAUD           1  
 TOTAL\_SPEND    10200  
 DOCTOR\_VISITS   13  
 NUM\_CLAIMS      3  
 MEMBER\_DURATION 119  
 OPTOM\_PRESC    2  
 NUM\_MEMBERS    3

CASE\_ID        3810  
 FRAUD           1  
 TOTAL\_SPEND    20000  
 DOCTOR\_VISITS   16  
 NUM\_CLAIMS      3  
 MEMBER\_DURATION 124  
 OPTOM\_PRESC    0  
 NUM\_MEMBERS    2

- f) (5 points) Follow-up with e), what is the predicted probability of fraudulent (i.e., FRAUD = 1)? If your predicted probability is greater than or equal to your answer in a), then the observation will be classified as fraudulent. Otherwise, non-fraudulent. Based on this criterion, will this observation be misclassified?

**Af) The predicted probability is 100% and since it is greater than the fraudulent percentage i.e., 19.9497%, it will be classified as fraudulent. Based on this criterion, the observation is not misclassified.**