

CS 412, HW #1

Pouya Akbarzadeh (pa2)

February 24, 2021

Question 1.

1.1 Answers

A) *Maximum and minimum.*

Mid-Term: 37 , 100

Final: 35 , 100

B) *First quartile Q1, median, and third quartile Q3*

Mid-Term: 68.0 , 77.0 , 87.0

Final: 82.0 , 89.0 , 96.0

C) *Mean*

Mid-Term: 76.715

Final: 87.084

D) *Mode*

Mid-Term: 77 83

Final: 97

E) *Variance*

Mid-Term: 173.10577

Final: 119.11294

1.2 Explanation and Equations

A) Numpy was used to find min and max

B) Numpy's function percentile() was used. 25,50,75 for Q1, Q2, and Q3 relatively.

C) Numpy's function mean() was used.

D) Numpy's function mode()

E) Numpy's function var() was used.

What the code is doing:

$$\text{Mean: } A = \frac{1}{n} \sum_{i=1}^n a_i = \frac{a_1 + a_2 + \cdots + a_n}{n}$$

Median :

$$Q1 : (N + 1) * 1/4$$

$$Q2 : (N + 1) * 2/4$$

$$Q3 : (N + 1) * 3/4$$

Mode : Mostrepeatedvalue

$$\text{Variance : } A = \sigma^2 = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n}$$

1.3 Code (Not including print statements)

```
1 import numpy
2 import pandas
3 from scipy import stats
4
5 # Reading Data
6 data_table = pandas.read_table('data.online.scores.txt', names = ['ID#', 'Mid-Term', 'Final'])
7 finals = data_table['Final']
8 mid = data_table['Mid-Term']
9 # https://numpy.org/doc/stable/reference/generated/numpy.array.html
10 mid_term_score_arr = numpy.array(mid)
11 final_score_arr = numpy.array(finals)
12
13
14 # Part A
15 # Mid-Term Data
16 min_mid = numpy.min(mid_term_score_arr)
17 max_mid = numpy.max(mid_term_score_arr)
18 #Final Data
19 min_val_final = numpy.min(final_score_arr)
20 max_val_final = numpy.max(final_score_arr)
21
22
23 # Part B
24 # https://numpy.org/doc/stable/reference/generated/numpy.around.html
25 #numpy.median(final_score_arr)
26 mq1 = numpy.percentile(mid_term_score_arr,25) # Q1
27 mq2 = numpy.percentile(mid_term_score_arr,50) # median
28 mq3 = numpy.percentile(mid_term_score_arr,75) # Q3
29 q1 = numpy.percentile(final_score_arr,25) # Q1
30 q2 = numpy.percentile(final_score_arr,50) # median
31 q3 = numpy.percentile(final_score_arr,75) # Q3
32
33 # Part C
34 finals_mean = numpy.mean(final_score_arr)
35 # https://numpy.org/doc/stable/reference/generated/numpy.around.html
36 finals_mean_rounded = numpy.around(finals_mean, decimals = 3)
37 mid_mean = numpy.mean(mid_term_score_arr)
38 mid_mean_rounded = numpy.around(mid_mean, decimals =3)
39
40 # Part D
41 # https://stackoverflow.com/questions/16330831/most-efficient-way-to-find-mode-in-numpy-array
42 mid_term_mode = mid.mode()
43 final_mode = finals.mode()
44
45 # Part E
46 # https://www.geeksforgeeks.org/numpy-var-in-python/
47 mid_var = numpy.var(mid_term_score_arr, dtype = numpy.float32)
48 var_final = numpy.var(final_score_arr, dtype = numpy.float32)
```

1.4 Terminal Output

```
A) Maximum and minimum.
=====
Mid-Term 37 , 100
Final: 35 , 100

B) First quartile Q1, median, and third quartile Q3
=====
Mid-Term: 68.0 , 77.0 , 87.0
Final: 82.0 , 89.0 , 96.0

C) Mean
=====
Mid-Term: 76.715
Final: 87.084

D) Mode
=====
Mid-Term: 77 83
Final: 97

E) Variance
=====
Mid-Term: 173.10577
Final: 119.11294
```

Question 2.

2.1 Answers

A) Compute and compare the variance of midterm-original and midterm-normalized, i.e., the midterm scores before and after normalization.

Variance MidTerm: 173.10577

Normalized Variance: 1.0

B) Given an original midterm score of 90, what is the corresponding score after normalization?

Mean: 76.715

Std-Dev: 13.156966785699508

Using $v' = v - \text{Avg} / \text{std-dev}$

Our score of 90 is normalized to 1.009731210573523

C) Compute the Pearson's correlation coefficient between midterm-original and finals-original?

Pearson's correlation coefficient: 0.544424742312412

D) Compute the Pearson's correlation coefficient between midterm-normalized and finals-original.

Pearson's correlation coefficient: 0.544424742312412

E) Compute the covariance between midterm-original and finals-original.

Covariance: 78.25419419

2.2 Explanation and Equations

A) Numpy var() was used. The function performed the following equation:

$$\text{Variance : } A = \sigma^2 = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n}$$

Using the standard dev and mean. I used the same equations as I did in question 1 to find those.

B) Scipy zscore() was used

then using $v' = v - \text{Avg} / \text{std-dev}$, we normalize the score

C) Scipy pearsonr() was used

D) Scipy pearsonr() was used

For both C and D the pearsonr() function utilized the following equation.

Formula

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

r = correlation coefficient

x_i = values of the x-variable in a sample

\bar{x} = mean of the values of the x-variable

y_i = values of the y-variable in a sample

\bar{y} = mean of the values of the y-variable

E) Numpy's stack() was used

This allowed me to find covariance. I could have use the cov() function in Panda, but this worked just as fine. It utilized this equation.

$$\text{cov}_{x,y} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{N - 1}$$

$\text{cov}_{x,y}$ = covariance between variable a and y

x_i = data value of x

y_i = data value of y

\bar{x} = mean of x

\bar{y} = mean of y

N = number of data values

2.3 Code (Not including print statements)

```
1 import numpy
2 import pandas
3 from scipy import stats
4
5 # Reading Data
6 data_table = pandas.read_table('data.online.scores.txt', names = ['ID#', 'Mid-Term', 'Final'])
7 finals = data_table['Final']
8 mid = data_table['Mid-Term']
9 # https://numpy.org/doc/stable/reference/generated/numpy.array.html
10 mid_term_score_arr = numpy.array(mid)
11 final_score_arr = numpy.array(finals)
12
13 # Part A
14 var_final = numpy.var(final_score_arr, dtype = numpy.float32)
15 var_mid = numpy.var(mid_term_score_arr, dtype = numpy.float32)
16
17 mid_norm = stats.zscore(mid_term_score_arr)
18 final_norm = stats.zscore(final_score_arr)
19 var_mid_norm = numpy.var(mid_norm, dtype = numpy.float32)
20 var_final_norm = numpy.var(final_norm, dtype = numpy.float32)
21
22 # Part B
23 given_score = 90
24 mid_std = numpy.std(mid_term_score_arr)
25 mid_mean = numpy.mean(mid_term_score_arr)
26 b_norm = (given_score - mid_mean) / mid_std
27
28 # Part C and D
29 corr_o_mid_o_fin = stats.pearsonr(mid_term_score_arr, final_score_arr)
30
31 # Part E
32 stack = numpy.stack((mid_term_score_arr, final_score_arr), axis = 0)
33 my_cov = numpy.cov(stack)
34
```

2.4 Terminal Output

```
A) Compute and compare the variance of midterm-original and midterm-normalized, i.e., the midterm scores before and after normalization.
=====
Variance MidTerm1: 173.10577
Normalized Variance: 1.0

B) Given an original midterm score of 90, what is the corresponding score after normalization?
=====
Mean: 76.715
Std-Dev: 13.156966785699508
Using v' = v - Avg / std-dev
Our score of 90 is normalized to 1.009731210573523

C) Compute the Pearson's correlation coefficient between midterm-original and finals-original.?
=====
Pearson's correlation coefficient: (0.544424742312412, 3.0298169609726267e-78)

D) Compute the Pearson's correlation coefficient between midterm-normalized and finals-original.
=====
Pearson's correlation coefficient: (0.544424742312412, 3.0298169609726267e-78)

E) Compute the covariance between midterm-original and finals-original.
=====
Covariance: [[173.27985405 78.25419419]
 [ 78.25419419 119.23217618]]
```

Question 3.

3.1 Answers

a) Each library has multiple copies of each book. Based on all the books (treat the counts of the 100 books as a feature vector for each of the libraries), compute the Minkowski distance of the vectors for CML and CBL with regard to different h values:

- i) 6152.0
- ii) 715.3278968417211
- iii) 170.0

b) Compute the cosine similarity between the feature vectors for CML and CBL.
Cosine similarity: 0.8414040256623079

c) Compute the Kullback-Leibler (KL) divergence between CML and CBL
 $KL(CML \parallel CBL)$: 0.149377 \rightarrow 14.94%
 $KL(CBL \parallel CML)$: 0.16565 \rightarrow 16.57%

3.2 Explanation and Equations

a) Using scipy library and used minkowski() function.

$$\left(\sum_{i=1}^n |X_i - Y_i|^p\right)^{1/p}$$

The case where $p = 1$ is equivalent to the [Manhattan distance](#) and the case where $p = 2$ is equivalent to the [Euclidean distance](#).

Based on equation above, this would have been calculated without any specific function. Rather by manual calculation. We could have found difference of CBL and CML data, taken their absolute value and summed these differences together. The $p * 1/p$ cancels out. And thus, we are left with the original equation showed above. The function had the following parameters, (vector 1, vector 2, h value). Thus, for parts i, ii, and iii I simply changed the h value that was given. For part iii, I used float('inf')

b) For this part, I again, used scipy, however, this time I used the cosine distance function.

$$\text{similarity} = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^n A_i \times B_i}{\sqrt{\sum_{i=1}^n (A_i)^2} \times \sqrt{\sum_{i=1}^n (B_i)^2}}$$

I could have simply used this equation, but using the predefined function allowed for a much cleaner code base.

c) KL divergence could have been found using SK-Learns lib, however, I used entropy in numpy. When scipy.stats.entropy is used, it will normalize the probabilities to one. Scipy's entropy function will calculate KL divergence if feed two vectors p and q, each representing a probability distribution. If the two vectors are not pdfs, it will normalize them first.

$$D(p||q) = \sum_{x \in X} p(x) \log \frac{p(x)}{q(x)}.$$

Using the equation shown above, we could 'manually' calculate this.

3.3 Code (Not including print statements)

```
1 import numpy
2 import pandas
3 from scipy import stats
4 from scipy.spatial import distance
5 from scipy.stats import entropy
6
7 # Reading Data
8 # I was changing text file before, I do not think I should do that
9 data_table = "data.libraries.inventories.txt"
10 parse_data = pandas.read_csv(data_table, "\t")
11 cml = numpy.array(parse_data.loc[parse_data['library'] == 'CML'])[0]
12 cbl = numpy.array(parse_data.loc[parse_data['library'] == 'CBL'])[0]
13 #get rid of the book stuff
14 cml = numpy.delete(cml, 0)
15 cbl = numpy.delete(cbl, 0)
16 #Part A
17 q3p1 = distance.minkowski(cml, cbl, 1)
18 q3p2 = distance.minkowski(cml, cbl, 2)
19 q3p3 = distance.minkowski(cml, cbl, float('inf'))
20
21 #Part B
22 q3p4 = 1 - distance.cosine(list(cml), list(cbl))
23
24 #Part C
25 k11 = entropy(list(cml), list(cbl), base = 4)
26 k12 = entropy(list(cbl), list(cml), base = 4)
27 k11perc = k11 * 100
28 k12perc = k12 * 100
29
```

3.4 Terminal Output

a) Each library has multiple copies of each book. Based on all the books (treat the counts of the 100 books as a feature vector for each of the libraries), compute the Minkowski distance of the vectors for CML and CBL with regard to different h values:

```
=====
i)      6152.0
ii)     715.3278968417211
iii)    170.0
```

b) Compute the cosine similarity between the feature vectors for CML and CBL.

```
=====
Cosine similarity:  0.8414040256623079
```

c) Compute the Kullback-Leibler (KL) divergence between CML and CBL

```
=====
KL(CML || CBL):  0.14937732067546497 in '%' it will be  14.9377 %
KL(CBL || CML):  0.16565531653503116 in '%' it will be  16.57 %
```

Question 4.

4.1 Answers

a) Calculate the distance between the binary attributes Buy Beer and Buy Diaper by assuming they are symmetric binary variables.

Distance: 0.015691868758915834

b) Calculate the distance between the binary attributes Buy Beer and Buy Diaper by assuming they are symmetric binary variables.

Jaccard Coefficient : 0.7317073170731707

c) Compute the χ^2 statistic for the contingency table.

χ^2 : 2450.716326822006

d) Consider a hypothesis test based on the χ^2 statistic where the null hypothesis is that Buy Beer and Buy Diaper are independent? Can you reject the null hypothesis? *at a significance level of $\alpha = 0.05$? Explain your answer, and mention the degrees of freedom used for the hypothesis test.*

If the value was more than 0.05 we could not reject the null hypothesis

The value based on info given was: 0.0

We were able to reject

P value: 0.0

Deg of freedom: 1

4.2 Explanation and Equations

A) This was done without using any library. First, I added the items then I divided by the sum of the entire array.

B) The Jaccard Coefficient between buying beer and diaper was done with this equation.

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}.$$

Which translates to

The Jaccard similarity coefficient, J , is given as

$$J = \frac{M_{11}}{M_{01} + M_{10} + M_{11}}.$$

C) Done with the use of `scipy chi2_contingency()`

D) Done with the use of `scipy chi2_contingency()`

4.3 Code (Not including print statements)

```
1 import numpy
2 from scipy import stats
3
4
5 table_q4 = numpy.array([[150,40],[15,3300]])
6 # Part A
7 sum = numpy.sum(table_q4)
8 beer = table_q4[0][1]
9 diaper = table_q4 [1][0]
10 dis = (beer + diaper)/sum;
11
12 # Part B
13 set1 = table_q4[0][0]
14 set2 = table_q4[0][1]
15 set3 = table_q4[1][0]
16 jac = set1 / (set1+set2+set3)
17 # Part C
18 chi = stats.chi2_contingency(table_q4)[0]
19 # Part D
20 chi_2 = stats.chi2_contingency(table_q4)[1]
21 freedom = stats.chi2_contingency(table_q4)[2]
22
```

4.4 Terminal Output

```
a) Calculate the distance between the binary attributes Buy Beer and Buy Diaper
by assuming they are symmetric binary variables.
=====
Distance:  0.015691868758915834

b) Calculate the distance between the binary attributes Buy Beer and Buy Diaper
by assuming they are symmetric binary variables.
=====
Jaccard Coefficient :  0.7317073170731707

c) Compute the  $\chi^2$  statistic for the contingency table.
=====
 $\chi^2$ :  2450.716326822006

d) Consider a hypothesis test based on the  $\chi^2$  statistic where the null hypothesis
is that Buy Beer and Buy Diaper are independent. Can you reject the null hypothesis
at a significance level of  $\alpha = 0.05$ ? Explain your answer, and also mention the degrees
of freedom used for the hypothesis test.
=====
If the value was more than 0.05 we could not reject the null hypothesis
The value based on info given was:  0.0
We were able to reject
P value:  0.0
Deg of freedom:  1
```