

Homework 1: Written Assignment

Handed Out: August 29, 2019

Due: September 16, 2019 11:55pm

Save your homework submission as *NETID-hw1-written.pdf*.

1 Data Description (20 points)

Suppose the population size is $N = 1,000,000$. We sample $n = 9$ examples x_i ($1 \leq i \leq n$) from the data. Suppose the mean value of the sample data is $\mu = 10$ and the variance is $v = 18$. Now we sample one more example $x_{n+1} = 20$ from the data. So the sample size is $n + 1 = 10$. What is the new mean value μ' and the new variance v' ?

Note that the result will be the same no matter what x_i ($1 \leq i \leq n$) are. You are expected to derive functions of calculating $\mu' = f(\mu, n, x_{n+1})$ and $v' = g(v, \mu, n, x_{n+1})$. Making fake numbers to find the answer, not mathematically proving the answer, will have only half of the points.

Solution:

$$\mu = \frac{x_1 + \dots + x_n}{n}.$$

$$v = \frac{(x_1 - \mu)^2 + \dots + (x_n - \mu)^2}{n-1}.$$

$$\mu' = \frac{x_1 + \dots + x_n + x_{n+1}}{n+1} = \frac{n\mu + x_{n+1}}{n+1}. \quad [5']$$

$$v' = \frac{(x_1 - \mu')^2 + \dots + (x_n - \mu')^2 + (x_{n+1} - \mu')^2}{n}.$$

We have

$$\begin{aligned} nv' - (n-1)v &= \{(x_1 - \mu')^2 - (x_1 - \mu)^2\} + \dots + \{(x_n - \mu')^2 - (x_n - \mu)^2\} + (x_{n+1} - \mu')^2 \\ &= (2x_1 - \mu - \mu') \times (\mu - \mu') + \dots + (2x_n - \mu - \mu') \times (\mu - \mu') + (x_{n+1} - \mu')^2 \\ &= \{2 \times (x_1 + \dots + x_n) - n\mu - n\mu'\} \times (\mu - \mu') + (x_{n+1} - \mu')^2 \\ &= (2n\mu - n\mu - n\mu') \times (\mu - \mu') + (x_{n+1} - \mu')^2 \\ &= (n\mu - n\mu') \times (\mu - \mu') + (x_{n+1} - \mu')^2 \\ &= n(\mu - \mu')^2 + (x_{n+1} - \mu')^2 \end{aligned}$$

$$\text{So } v' = v + (\mu - \mu')^2 + \frac{(x_{n+1} - \mu')^2 - v}{n} = \frac{n-1}{n}v + \frac{1}{n+1}(x_{n+1} - \mu)^2. \quad [5']$$

$$\text{The result is } \mu' = \frac{9 \times 10 + 20}{10} = 11 \quad [5'] \text{ and } v' = \frac{8}{9} \times 18 + \frac{1}{10}(20 - 10)^2 = 26. \quad [5']$$

2 Data Integration - Correlation Analysis (10 points)

Suppose two stocks X_1 and X_2 have the following values in one week:

(2, 5), (3, 8), (5, 10), (4, 11), (6, 14)

Are their prices rising/falling together or in different trends?

Solution:

$$E(X_1) = (2 + 3 + 5 + 4 + 6)/5 = 20/5 = 4 \text{ [3']}$$

$$E(X_2) = (5 + 8 + 10 + 11 + 14)/5 = 48/5 = 9.6 \text{ [3']}$$

$$\sigma_{12} = (2 \times 5 + 3 \times 8 + 5 \times 10 + 4 \times 11 + 6 \times 14)/5 - 4 \times 9.6 = 4 \text{ [4']}$$

Thus, X_1 and X_2 rise together since $\sigma_{12} > 0$.

3 Data Reduction - Regression (20 points)

Suppose we have K numerical features and 1 numerical label. A standard linear regression model M makes a coefficient of determination R . Now we add one more feature, that is the average of the original K features, and build a new linear regression model M' . Will the new coefficient of determination R' be bigger than, or small than, or equal to R ? Please mathematically prove your answer.

Solution:

Recall the formula for the coefficient of determination talked in class:

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - y'_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

It's easy to observe that y_i and \bar{y} will not change after the new feature is introduced. Hence the only component that may affect the coefficient of determination is y' , which is the predicted value by the linear regression model.

Now consider how linear regression works, the linear regression model first optimizes the following objective function,

$$\min_{w_1, w_2, \dots, w_k, b} \left(\sum_{i=1}^k w_i x_i + b - y \right)^2$$

and then make prediction with the learnt parameters by

$$y' = \sum_{i=1}^k w_i x_i + b$$

When the new feature $x_{k+1} = \sum_{i=1}^k x_i / k$ is introduced, the objective becomes (denoted

weights as a_i to avoid confusion)

$$\begin{aligned}
& \min_{a_1, a_2, \dots, a_k, a_{k+1}, b} \left(\sum_{i=1}^{k+1} a_i x_i + \hat{b} - y \right)^2 \\
&= \min_{a_1, a_2, \dots, a_k, a_{k+1}, b} \left(a_{k+1} \cdot \frac{\sum_{i=1}^k x_i}{k} + \sum_{i=1}^k a_i x_i + \hat{b} - y \right)^2 \\
&= \min_{a_1, a_2, \dots, a_k, a_{k+1}, b} \left(\sum_{i=1}^k \frac{a_{k+1}}{k} x_i + \sum_{i=1}^k a_i x_i + \hat{b} - y \right)^2 \\
&= \min_{a_1, a_2, \dots, a_k, a_{k+1}, b} \left(\sum_{i=1}^k \left(a_i + \frac{a_{k+1}}{k} \right) x_i + \hat{b} - y \right)^2
\end{aligned}$$

which can be rewritten as

$$\min_{\hat{w}_1, \hat{w}_2, \dots, \hat{w}_k, \hat{b}} \left(\sum_{i=1}^k \hat{w}_i x_i + \hat{b} - y \right)^2$$

the parameters of the new linear regression model $\hat{w}_1, \hat{w}_2, \dots, \hat{w}_k, \hat{b}$ will be learnt to be the same as the parameters of the old model w_1, w_2, \dots, w_k, b , hence the prediction value of the new model will be

$$\hat{y}' = \sum_{i=1}^k \hat{w}_i x_i + \hat{b} = \sum_{i=1}^k w_i x_i + b = y'$$

So after adding the new feature, the predicted value will not change. Therefore, the new coefficient of determination will stay the same.

Homework 1: Programming Assignment

*Handed Out: August 29, 2019**Due: September 16, 2019 11:55pm*

Save your homework submission as *NETID-hw1-programming.zip*. The zip file has one pdf file *NETID-hw1-programming.pdf* and multiple code files.

You are allowed to use any programming language (Python recommended; R, C++, Java, etc.), however, the solutions will be in **Python**. You are allowed to use any public package (including Numpy and Scikit-learn) and any other kind of tools (Excel).

A Film Dataset

File name: data-film.csv

Introduction: Suppose we have 1,000,000 films. We sample 150 films, so our dataset has 150 films (as the data objects). Each film has an ID (" f_{DIGIT} "; $\text{DIGIT} \in \{1, \dots, 150\}$) and a genre (or called category/class) from {"ACTION", "ROMANCE", "COMEDY"}. Each film has four attributes. Each attribute is the average rating of the film given by a specific website. The attribute name is "AVGRATING_WEBSITE_ DIGIT " ($\text{DIGIT} \in \{1, 2, 3, 4\}$). The attribute values are numerical. Note that the rating scales can be different.

The first line is the header of attribute and label names. The following lines (rows) are object ID, attribute values, and label values of the data objects (films). The columns are separated by comma.

Example: The third line is " $f2, 4.9, 3, 1.4, 0.2, \text{ACTION}$ ": The film " $f2$ " is an ACTION film. It is graded as 4.9 on Website 1, 3.0 on Website 2, 1.4 on Website 3, and 0.2 on Website 4.

Data Description (10 points)

The object-feature data matrix, which is denoted as \mathbf{D} , has $m = 150$ objects and $n = 4$ features (and therefore, it has 600 values). Please use *Z-score normalization* to normalize the data matrix by each feature. The normalized data matrix is denoted as \mathbf{A} . What are the maximum/minimum Z score for each feature?

Output: Write down the maximum/minimum Z score in the pdf. Save your code as *NETID-hw1-1.py*.

Solution:

AVGRATING_WEBSITE_1 : 2.49201920212, -1.87002413385

AVGRATING_WEBSITE_2 : 3.0907752483, -2.43394714191

AVGRATING_WEBSITE_3 : 1.78583195363, -1.56757623428

AVGRATING_WEBSITE_4 : 1.71309868604 , -1.44954504204

[If one output is incorrect, you will lose one point. If you code is not runnable, you will lose five points. If you do not submit your code, you will get zero point.]

Data Visualization (30 points)

With the original data matrix D:

(1) generate *boxplot* for the first attribute "AVGRATING_WEBSITE.1;"

Output: Show the figure in the pdf. Save your code as NETID-hw1-2-1.py.

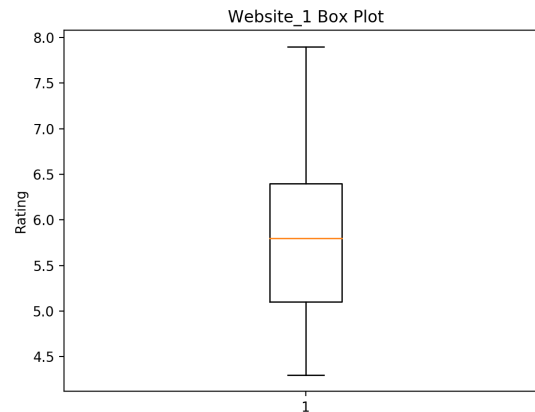


Figure 1: Boxplot

(2) generate *histogram* for the third attribute "AVGRATING_WEBSITE.3". Set the number of histograms/bins as 10;

Output: Show the figure in the pdf. Save your code as NETID-hw1-2-2.py.

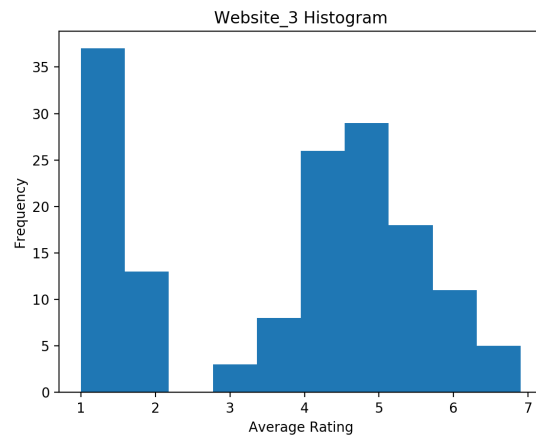


Figure 2: histogram

(3) generate *bar chart* where the X-axis is the genre and the Y-axis is the mean value of average ratings of films given in the first attribute "AVGRATING_WEBSITE.1;"

Output: Show the figure in the pdf. Save your code as NETID-hw1-2-3.py.

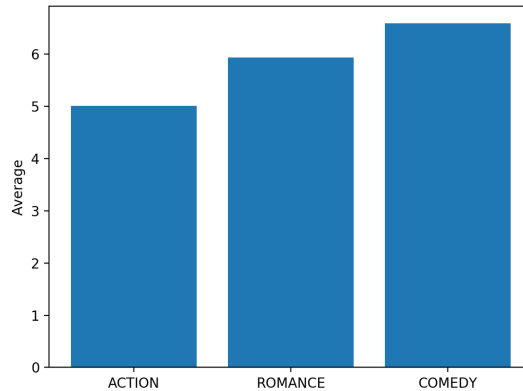


Figure 3: bar chart

(4) generate 2-dimensional *scatter plot* using attributes "AVGRATING_WEBSITE.1" and "AVGRATING_WEBSITE.3", where the marker types and colors should be different for different genres (node labels);

Output: Show the figure in the pdf. Save your code as NETID-hw1-2-4.py.

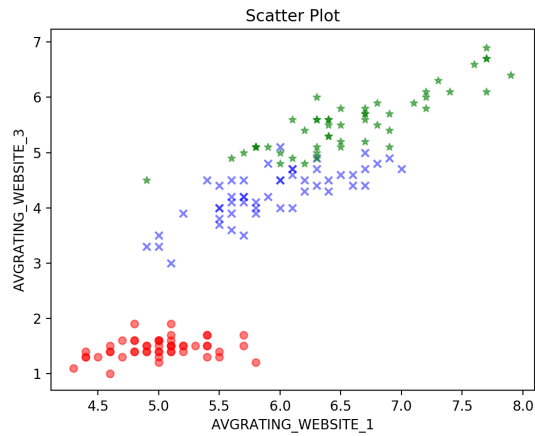


Figure 4: scatter plot

(5) calculate the *KL divergence* between attributes “AVGRATING_WEBSITE_1” and “AVGRATING_WEBSITE_3.” Hint: For each attribute, you need to generate a probability distribution. You can determine bins on the attribute value by yourself. For each bin, you have the frequency (which is the number of films that have an attribute value in the bin) can calculate the probability. For example, if the total number of films is 150, if there are 15 films whose value of attribute “AVGRATING_WEBSITE_1” is in the bin $[1.5, 2)$ (between 1.5 and 2), the probability is 0.1. Then you have two probability distributions: one for each attribute. Just call python libraries to calculate the KL divergence.

Output: Write down the KL divergence score in the pdf. Save your code as NETID-hw1-2-5.py.

Solution:

KL Divergence : 0.0901355700351

[If one output is incorrect, you will lose five points. If one of your code is not runnable, you will lose five points. If you do not submit your code, you will get zero point.]

Data Cleaning and Integration (10 points)

With the original data matrix D:

(1) calculate *correlation coefficients* $\rho_{i,j}$, which is $\rho(\text{“AVGRATING_WEBSITE_i”}, \text{“AVGRATING_WEBSITE_j”})$, between every pair of the attributes using covariance analysis.

Output: Write down the correlation coefficients in the pdf. Save your code as NETID-hw1-3-1.py.

Solution:

$\rho(\text{“AVGRATING_WEBSITE_1”}, \text{“AVGRATING_WEBSITE_2”}) = -0.117569784133$

$$\begin{aligned}\rho(\text{"AVGRATING_WEBSITE_1"}, \text{"AVGRATING_WEBSITE_3"}) &= 0.871753775887 \\ \rho(\text{"AVGRATING_WEBSITE_1"}, \text{"AVGRATING_WEBSITE_4"}) &= 0.817942174858 \\ \rho(\text{"AVGRATING_WEBSITE_2"}, \text{"AVGRATING_WEBSITE_3"}) &= -0.428440104331 \\ \rho(\text{"AVGRATING_WEBSITE_2"}, \text{"AVGRATING_WEBSITE_4"}) &= -0.365430794103 \\ \rho(\text{"AVGRATING_WEBSITE_3"}, \text{"AVGRATING_WEBSITE_4'}) &= 0.962746024624\end{aligned}$$

With the Z-score normalized data matrix **A**:

(2) (the same as (1)) calculate *correlation coefficients* between every pair of the attributes using covariance analysis.

Output: Write down the correlation coefficients in the pdf. Save your code as NETID-hw1-3-2.py.

Solution:

$$\begin{aligned}\rho(\text{"AVGRATING_WEBSITE_1"}, \text{"AVGRATING_WEBSITE_2"}) &= -0.117569784133 \\ \rho(\text{"AVGRATING_WEBSITE_1"}, \text{"AVGRATING_WEBSITE_3"}) &= 0.871753775887 \\ \rho(\text{"AVGRATING_WEBSITE_1"}, \text{"AVGRATING_WEBSITE_4"}) &= 0.817942174858 \\ \rho(\text{"AVGRATING_WEBSITE_2"}, \text{"AVGRATING_WEBSITE_3"}) &= -0.428440104331 \\ \rho(\text{"AVGRATING_WEBSITE_2"}, \text{"AVGRATING_WEBSITE_4"}) &= -0.365430794103 \\ \rho(\text{"AVGRATING_WEBSITE_3"}, \text{"AVGRATING_WEBSITE_4'}) &= 0.962746024624\end{aligned}$$

(3) Are the above results the same? And why?

Output: Write down your reasoning in the pdf.

Solution:

They will be same. Since all the processing we have done are just normalizing, transforming the scales. In other word, the interrelations among them should be kept the same.

[For the first and second question, if only one output is incorrect, you will not lose any point. You will lose one more point with every two more wrong outputs. For the third question, if your answer is wrong, you will lose two points. If one of your code is not runnable, you will lose four points. If you do not submit your code, you will get zero point.]