

112-1 Introduction to Data Science

Midterm Exam

Student ID: 11204604 Name: Tse Wen Hong

44
-5

1. (10%) Please briefly describe the major tasks in data preprocessing.

Ch. 3

Try to find out the relationⁱⁿ between the attributes,

Key word! and the relation in between the data.

Data cleaning and integration so that we can make better prediction or judgement according the finding result.

2. (20%) Given two objects represented by the tuples $X_1 (20, 10, 5, 3)$ and $X_2 (30, 8, 6, 16)$:

- (a) Compute the Euclidean distance between the two objects.
- (b) Compute the Manhattan distance between the two objects.
- (c) Compute the Minkowski distance between the two objects, using $h=3$.
- (d) Compute the Pearson correlation coefficient between the two objects. Please describe the relationship between X_1 and X_2 .

$$(a) i. (20-30)^2 + (10-8)^2 + (5-6)^2 + (3-16)^2 = 100 + 4 + 1 + 169 \\ = 274$$

ii. $\sqrt{274} = 16.5529 \leftarrow$
Euclidean distance is

$$(b) |20-30| + |10-8| + |5-6| + |3-16| = 10 + 2 + 1 + 13 = 26$$

Manhattan distance is: 26

$$(c) i. (20-30)^3 + (10-8)^3 + (5-6)^3 + (3-16)^3 = 1600 + 8 + 1 + 2197 = 3206$$

Minkowski distance is = $\sqrt[3]{3206}$

~~(d)~~ i. $\bar{x}_1 = (20+10+5+3)/4 = 9.5$ (sorry my calculator have no power function)

$$\bar{y}_2 = (30+8+6+16)/4 = 15$$

$$\frac{\text{Cov}(A, B)}{\sigma_A \sigma_B}$$

ii. (turn to the back)

$$(i) \begin{array}{c|ccccc} x-a & | & 10.5 & 0.5 & -4.5 & -6.5 \\ y-a & | & 15 & 22 & 24 & 14 \end{array}$$

$$(ii) r = \frac{(10.5 \times 15) + (0.5 \times 22) + (-4.5 \times 24) + (-6.5 \times 14)}{\sqrt{10.5^2 + 0.5^2 + (-4.5)^2 + (-6.5)^2} \times \sqrt{15^2 + 22^2 + 24^2 + 14^2}}$$

$$= \frac{157.5 + 11 - 108 - 91}{\sqrt{110.25 + 0.25 + 20.25 + 42.25} \times \sqrt{225 + 484 + 576 + 196}}$$

$$= \frac{-30.5}{\sqrt{173} \times \sqrt{1481}}$$

$$= \frac{-30.5}{13.1529 \times 38.4838}$$

$$= \frac{-30.5}{506.1753}$$

$$r = -0.60 - 0.0602$$

$\boxed{0.75}$

參考公式
兩用、 $r = \frac{\text{Var}(A \cap B)}{\delta_A \delta_B}$

3. (20%) Suppose a group of 12 students with the test scores listed as follows:

19, 71, 48, 63, 35, 85, 69, 81, 72, 88, 99, 95.

Partition them into four bins by each of the following methods.

(a) equal-frequency (equi-depth) partitioning

(b) equal-width partitioning

(a) 19, 35, 48, 63, 71, 72, 81, 85, 88, 95, 99

bin 1 = 19, 35, 48

bin 2 = 63, 69, 71

bin 3 = 72, 81, 85

bin 4 = 88, 95, 99

(b) (i) $(99-19)/4 = 20$

(ii) bin 1 = 19, 35

bin 2 = 48,

bin 3 = 63, 69, 71, 72,

bin 4 = 81, 85, 88, 95, 99

4. (30%) For the following group of data with mean=600 and variance=260000

100, 200, 400, 800, 1500

-26

(a) Normalize the above group of data by min-max normalization with min = 0 and max = 10.

(b) In z-score normalization, what value should the third number 400 be transformed to?

(c) Normalize the above values by decimal scaling.

$$(a) 100 \Rightarrow \frac{100-100}{1500-100} \times (10 - 0) = 0$$

$$200 \Rightarrow \frac{200-100}{1500-100} \times 10 = 0.6711 \times$$

$$400 \Rightarrow \frac{400-100}{1500-100} \times 10 = 2.0134 \times$$

$$800 \Rightarrow \frac{800-100}{1500-100} \times 10 = 4.6979 \times$$

$$1500 \Rightarrow 10$$

(b) (c) 本來自己沒寫

這兩天都在想 Arti or i, 現在用到子題 3.

實再回來複習 DB'

3. (20%) Suppose a group of 12 students with the test scores listed as follows:

19, 71, 48, 63, 35, 85, 69, 81, 72, 88, 99, 95.

Partition them into four bins by each of the following methods.

- equal-frequency (equi-depth) partitioning
- equal-width partitioning

(a) $19, 35, 48, \boxed{63}, 71, 72, 81, 85, \boxed{88}, 95, 99$

$$\text{bin 1} = 19, 35, 48$$

$$\text{bin 2} = 63, 69, 71$$

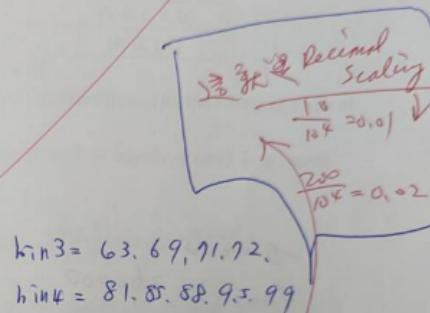
$$\text{bin 3} = 72, 81, 85$$

$$\text{bin 4} = 88, 95, 99$$

$$(b) (i) (99-19)/4 = 20$$

$$(ii) \text{ bin 1} = 19, 35$$

$$\text{bin 2} = 48,$$



$$\text{bin 3} = 63, 69, 71, 72,$$

$$\text{bin 4} = 81, 85, 88, 95, 99$$

4. (30%) For the following group of data with mean=600 and variance=260000

100, 200, 400, 800, 1500 \leftarrow ~~z-score~~ (42%)

-26

- Normalize the above group of data by min-max normalization with min = 0 and max = 10.
- In z-score normalization, what value should the third number 400 be transformed to?
- Normalize the above values by decimal scaling.

$$100 \Rightarrow \frac{100-100}{1500-100} \times (10 - 0) = 0$$

$$X_{\text{norm}} = \frac{X - X_{\text{min}}}{X_{\text{max}} - X_{\text{min}}} \times (Y_{\text{max}} - Y_{\text{min}})$$

$$200 \Rightarrow \frac{200-100}{1500-100} \times 10 = 0.6711 \times \frac{100-0}{1500-100} = 0.6711 \times 10 = 6.711$$

$$400 \Rightarrow \frac{400-100}{1500-100} \times 10 = 2.0134 \times \frac{100-0}{1500-100} = 2.0134 \times 10 = 20.134$$

$$800 \Rightarrow \frac{800-100}{1500-100} \times 10 = 4.6979 \times \frac{100-0}{1500-100} = 4.6979 \times 10 = 46.979$$

$$1500 \Rightarrow 10$$

(b) (c) 本章已經第 3

New value $= (X - \mu) / \sigma$ 這兩天都在搞 Artificial, 現在開始第 3.

$\mu = 600$ 實再回來複習

$$\sigma^2 = \sqrt{\frac{1}{n} \sum (x_i - \mu)^2} = \sqrt{\frac{1}{5} (100-600)^2 + (200-600)^2 + (400-600)^2 + (1500-600)^2} = \sqrt{\frac{1}{5} (160000 + 40000 + 40000 + 810000)} = \sqrt{250000} = 500$$

5. (20%) A database has 5 transactions. Let minimum support be 60% and minimum confidence be 80%.

Customer	Date	Items_boughts
100	10/15	{I, P, A, D, B, C}
200	10/15	{D, A, E, F}
300	10/16	{C, D, B, E}
400	10/18	{B, A, C, K, D}
500	10/19	{A, G, T, C}

Support Count =
 $60\% \times 5 = 3$

(a) List the frequent k-itemset for the largest k. (10%) $\Rightarrow \{B, C, D\}$

(b) List all the strong association rules (with support and confidence) for the following shape of rules.

$$\forall x \in \text{transaction}, \text{buys}(x, \text{item}_1) \wedge \text{buys}(x, \text{item}_2) \Rightarrow \text{buys}(x, \text{item}_3) \quad [\text{sup.}, \text{conf.}]$$

$$(10\%) \quad m=5$$

C	Item	Count	Support
i(i)	A	4	80%
Candidate	B	3	60%
i(i)	C	4	80%
i(i)	D	4	80%
i(i)	E	2	10%
-F	1	1	
-G	1	1	
-H	1	1	
-I	1	1	

(iii)

$$P(D|A) = \frac{60\%}{80\%} = 75\%$$

$$P(A|D) = \frac{60\%}{80\%} = 75\%$$

PFT \Rightarrow T
 $3/3$

$B, C \Rightarrow D$ conf = $3/3$
 $B, D \Rightarrow C$
 $C, D \Rightarrow B$

$= 100\%$
 100%

P			
C ₂			
i(i)	{A, B}	1	20%
V	{A, C}	3	60%
V	{A, D}	3	60%
	{B, C}	2	40%
	{B, D}	3	60%
	{C, D}	2	40%

(ii)

3!

3!

Confidences:

$$\{A, D\} \Rightarrow B \quad P(B|\{\{A, D\}\}) = \frac{60\%}{60\%} = 100\%$$

$$\{A, D\} \Rightarrow C \quad P(C|\{\{A, D\}\}) = \frac{80\%}{60\%} = 133\%$$

$$\{A, D\} \Rightarrow B \rightarrow B$$

C ₃			
A B C	2		
A C D	2		
A B D	2		
V B C D	3	3/5 = 60%	

frequent k-itemset is $\{A, D\}$

Strong association rule is

$\{A, D\} \Rightarrow B$ confidence 100%

$\{A, D\} \Rightarrow C$ confidence 133%