# XI'AN JIAOTONG-LIVERPOOL UNIVERSITY

## 西 交 利 物 浦 大 学

## REMOTE CLOSE BOOK

## EXAMANSWER SUBMISSION

## COVER SHEET

| Name | Zizhe (Surname) | Wang (Given Name) |
|---|---|---|
| Student ID Number | 1929413 | |
| Programme | information and computing science | |
| Module Title | Big Data Analytics | |
| Module Code | INT 303 | |
| Module Examiner | Jia Wang/Pengfei Fan | |

By uploading or submitting the answers of this Remote Close Book Exam, I certify the following:

- I will act fairly to my classmates and teachers by completing all of my academic work with integrity. This means that I will respect the standards and instructions set by the Module Leader and the University, be responsible for the consequences of my choices, honestly represent my knowledge and abilities, and be a community member that others can trust to do the right thing even when no one is watching. I will always put learning before grades, and integrity before performance.
- I have read and understood the definitions of collusion, copying, plagiarism, and dishonest use of data as outlined in the Academic Integrity Policy, and cheating behaviors in the Regulations for the Conduct of Examinations of Xi'an Jiaotong-Liverpool University.
- This work is produced all on my own and can effectively represent my own knowledge and abilities.

I understand collusion, plagiarism, dishonest use of data, submission of procured work, submission of work produced and/or contributed by others are serious academic misconducts. By uploading or submitting the answers with this statement, I acknowledge that I will be subject to disciplinary action if I am found to have committed such acts.

Signature___王玫哲 Zizhe Wang___      Date ___2023/1/5___

# Q1

## (a)

⟨1⟩ missing data: we have two solution: deletion or imputation.

~~When the data without some value~~

When the number of data without some value is overcome the number of data with ~~of~~ full ~~of~~ value, we should use deletion to process data.

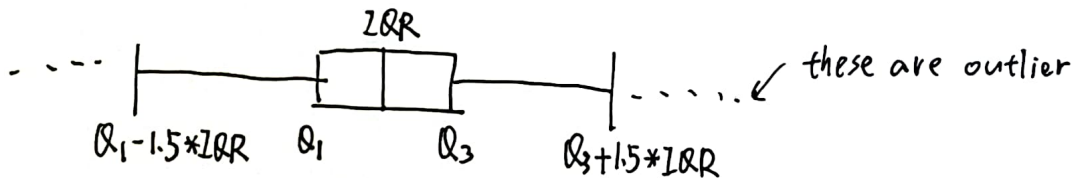~~When~~ In other situations, we consider ~~three~~ ways for ~~imputation. missing.~~ imputation.

~~missing~~ ① regression imputation

② simple imputation

③ KNN imputation

④ piority knoledge of this area

⟨2⟩ ~~data st~~ outlier valnes: we use IQR:



$Q_1 - 1.5*IQR \quad Q_1 \quad Q_3 \quad Q_3 + 1.5*IQR$

✓ these are outlier

we draw box chart to delete them

⟨3⟩ data standardization ~~issea~~ issues

~~we use min-max scale~~, ~~we~~
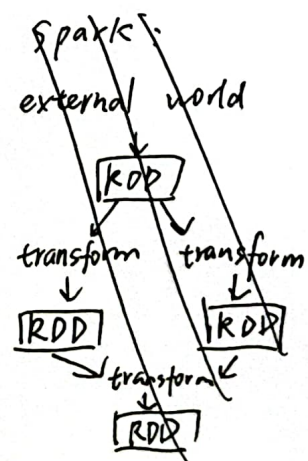
~~$S = \frac{\sqrt{...}}{N}$~~

$$S = \frac{\sqrt{\sum_{i=0}^{\varphi}(x_i - \bar{x})^2}}{n-1}$$

## (b)

Spark is higher level of Hadoop, Hadoop is constructed by Mapreduce and HDFS.
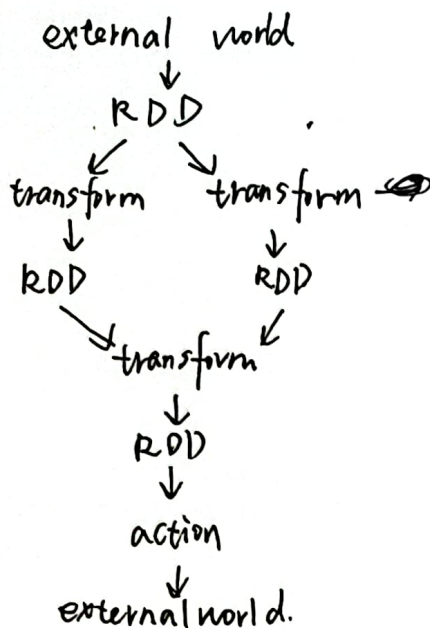
① Mapreduce have three main component:

map, join value by key, reduce

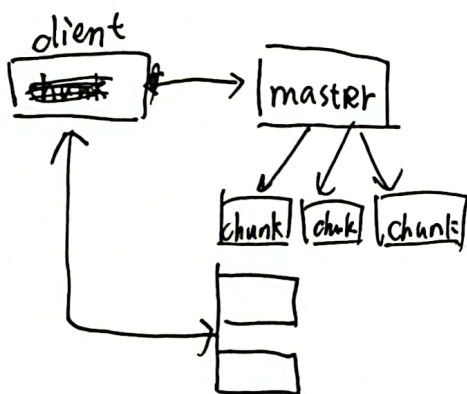② Hadoop have chunk ~~so~~ node, master node and client server ~~while~~.

Spark:

external world

↓

RDD

transform ↙ ↘ transform

RDD     RDD

transform

RDD

① ②

③. spark :

external world
↓
RDD

transform ↘ transform ◁
↓              ↓
RDD           RDD
↘ transform ↙
↓
RDD
↓
action
↓
external world.

spork use memory, but Hadoop use disk.
transform : join, distinct, union and so on.
action : count, save and so on.

Hardoop infrastructure :



client ask master node and find chunk, then chunk return data to client.

(c).

~~ker~~ sklearn, numpy, pandas,

Q2.

(a). 4  2  1

(b).
$$N = \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 4 & 0 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} & 0 \\ 1/\sqrt{2} & -1/\sqrt{2} & 0 \end{bmatrix}$$

$$= \begin{bmatrix} 2\sqrt{2} & \sqrt{2} \\ 2\sqrt{2} & -\sqrt{2} \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} & 0 \\ 1/\sqrt{2} & -1/\sqrt{2} & 0 \end{bmatrix} = \begin{bmatrix} 3 & 1 & 0 \\ 1 & 3 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

②

N singular value decomposition:

$$N = U \Sigma V^T = \begin{bmatrix} \frac{2}{\sqrt{2}} & \frac{\sqrt{2}}{\sqrt{2}} \\ \frac{2}{\sqrt{2}} & -\frac{\sqrt{2}}{\sqrt{2}} \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 4 & 0 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} & 0 \\ 1/\sqrt{2} & -1/\sqrt{2} & 0 \end{bmatrix}$$

7∞

(C)

reconstruction error $\approx \sqrt{\cdots}$

$$= \sqrt{(3-\frac{1}{\sqrt{2}})^2 + (1-\frac{1}{\sqrt{2}})^2 + (1-\frac{1}{\sqrt{2}})^2 + (1 \cdots}$$

$$= \sqrt{(3-3)^2 + (1-1)^2 + 0^2 + (1-1)^2 + (3-3)^2 + (0-0)^2 + (0-0)^2 + (0-0)^2 + (0-1)^2}$$

$$= 1$$

Q3.

(a) Student A is correct. Because it is obrious that if we split
data into training and test sets, the min and max value
should be different in these two sets. if we not do the.
normalize togther, the scale standard will be different on
each sets. it will result in some error.

(b)

(i).  A = 3  B = 4  C = 4

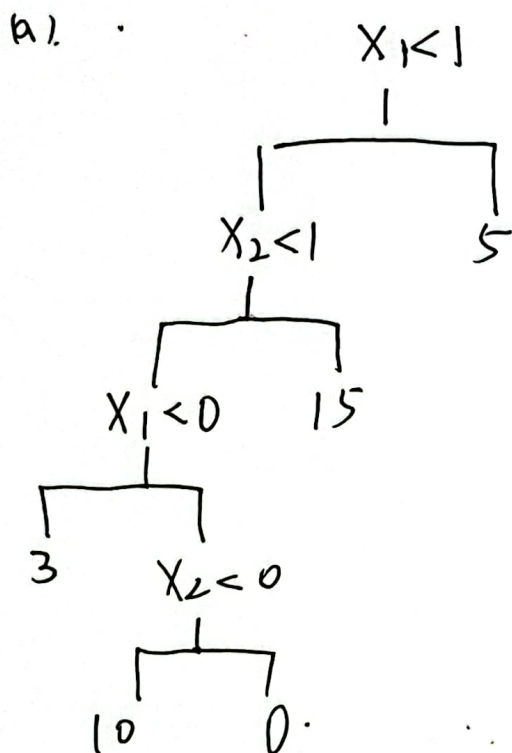(ii)  A = a   B = 12.0   C = 12.0

(iii)  a = 13.0   b = 12.3

(iiii) ~~histogram~~ pie ~~chart~~ chart

Q4

a)

```
                    X₁<1
                     |
          ┌──────────┴──────────┐
          |                     |
        X₂<1                    5
          |
     ┌────┴────┐
     |         |
   X₁<0        15
     |
 ┌───┴───┐
 |       |
 3      X₂<0
          |
      ┌───┴───┐
      |       |
      10      0
```

b)

$$X_1 < 1$$

$$X_2 < 1$$

$$X_1 < 0$$

$$X_2 < 0$$



Grid diagram:

- Top region: 2.49
- X₂ = 2 line
- Left middle region: −1.06
- Right middle region: 0.21
- X₂ = 1 line
- Lower left region: −1.80
- Lower right region: 0.63
- Bottom axis labels: D and X₁

④

## Q5

a). False. We must reduce after all the mappers finished, or we will get wrong number of key-value pairs.

(b) False. Sort and Group should before the reducer.

(c) False. ~~it doesn't matter~~ it doesn't matter.

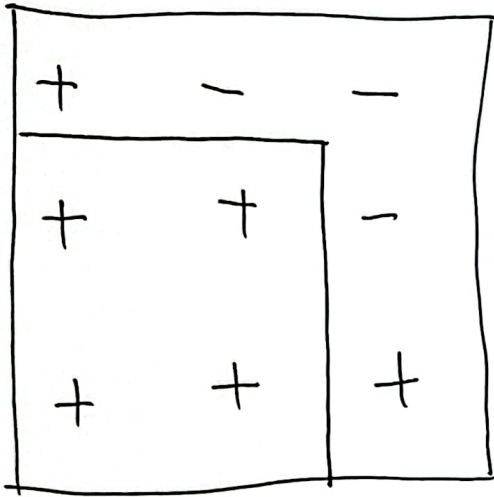(d) False. It up to the request of ~~the~~ reducer.

## Q6.

(a).

i).



ii)



the wrong classification point's weight will increase.

In the first iteration, the "―" in (1, 2) position has increased its weight, so in 2 iteration, it must be classified in right way, the same to "+" in (3,3) position
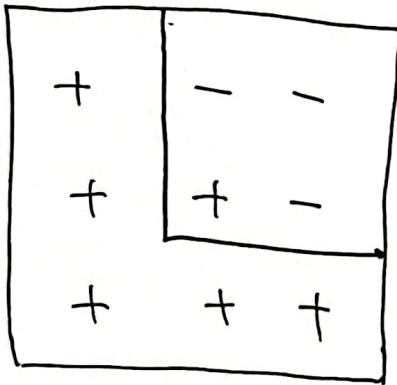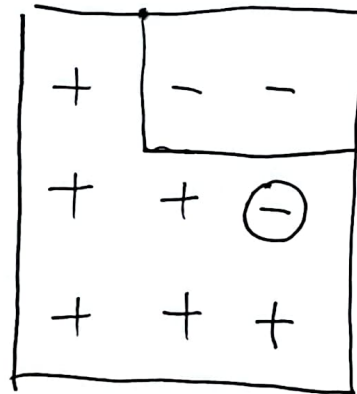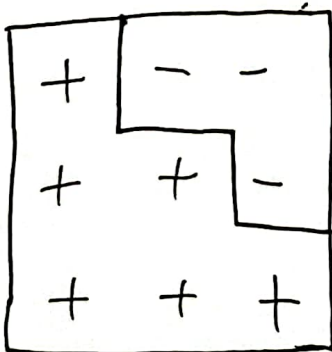
(iii)



(b)

(i)



(ii)



the wrong classification point's weight will increase.

In the first iteration, the "+" in (2,2) position has increased its weight, so in 2 iteration, it must be classified in right way.

iii)



Consider the weight of position (2,2) and position (2,3) has increased, and "−" in (2,3) has increased more, so this "−" must be classified in right way.

**Q7.**

**(a)** Reader 1: mean = 1.5 ~~Similarity~~

Reader 2: mean = $\frac{3+1}{2} = 2$    Similarity $(1,2) = \dfrac{-0.5}{\sqrt{1+1}\cdot\sqrt{\frac{1}{2}+\frac{1}{2}}} = -0.5 = -\sqrt{\frac{5}{10}}$

Reader 3: mean = 1    ~~Similarity (1,3)~~ $= \frac{4}{\sqrt{1}}$ ~~Similarity $= 0$~~ Similarity $= 0$

Reader 4: mean = $\frac{2+1}{2} = 1.5$    Similarity $(1,4) = \dfrac{0.5 \times -0.5}{\sqrt{0.5^2+(-0.5)^2}\cdot\sqrt{0.5^2+(0.5)^2}} = \dfrac{-\frac{1}{4}}{\frac{1}{4}} = -1$

Reader 5: mean = $\frac{0+3}{2} = 1.5$

Similarity $(1,5) = \dfrac{-0.5 \times 1.5}{\sqrt{(0.5)^2+(0.5)^2}\cdot\sqrt{(1.5)^2+(1.5)^2}}$

$= \dfrac{-\frac{3}{4}}{\frac{\sqrt{10}}{4}} = -\dfrac{3}{\sqrt{10}} = -\sqrt{\frac{9}{10}}$

~~similar~~ $\text{sim}(1,3) > \text{sim}(1,2) > \text{sim}(1,5) > \text{sim}(1,4)$

most similar : reader 3 and reader 2.

$\therefore \dfrac{1.0 \times 0 + 3.0 \times (-0.5)}{0 + (-0.5)} = \dfrac{-1.5}{-0.5} = 3$.

$\therefore$ Book1 for reader1 is rating 3.0.

**(b)**,

|        | r1  | r2  | r3  | r4  | r5  |
|--------|-----|-----|-----|-----|-----|
| book1  | ?   | 3.0 | 1.0 | 2.0 | 0.0 |
| book2  | 2.0 | 1.0 |     | 1.0 |     |
| book3  | 1.0 |     |     | 3.0 |     |

|   | r1 | r2 | r3 | r4 | r5 |
|---|----|----|----|----|----|
|   | ~~5~~ | 1.5 | -0.5 | 0.5 | -1.5 |
| $\frac{8}{3}$ |   | $-\frac{1}{3}$ |   | $-\frac{1}{3}$ |   |
| $-1$ |   |   |   | 1 |   |

$\therefore$ book 1 mean = $\dfrac{3+1+2+0}{4} = \dfrac{6}{4} = 1.5$

book 2 mean = $\dfrac{2+1+1}{3} = \dfrac{4}{3}$

book 3 mean = $\dfrac{1+3}{2} = 2$

$\text{Sim}(1,3) = \dfrac{-1.5 \times 1}{\sqrt{1^2+1^2}\cdot\sqrt{(1.5)^2+(0.5)^2+(0.5)^2+(1.5)^2}} = \dfrac{-1.5}{2} = -0.75$

$\text{Sim}(1,2) = \dfrac{1.5 \times (-\frac{1}{3}) + (0.5) \times (-\frac{1}{3})}{\sqrt{\frac{1}{9}+\frac{1}{9}+\frac{64}{9}}\cdot\sqrt{(\frac{3}{2})^2 \cdot (\frac{1}{2})^2}} = \dfrac{(\frac{1}{4})^2 \cdot (\frac{3}{2})^2}{} = 0$

⑧

$sim(1,2) > sim(1,3)$

$\therefore$ rating $D_c 0$