

## 第1题

### (a) 基于 $\alpha, \beta$

三台计算机的向量分别为：

- A: [3.06, 500 $\alpha$ , 6 $\beta$ ]
- B: [2.68, 320 $\alpha$ , 4 $\beta$ ]
- C: [2.92, 640 $\alpha$ , 6 $\beta$ ]

则sim(A,B):

$$\text{sim}(A, B) = \frac{3.06 \times 2.68 + 500\alpha \times 320\alpha + 6\beta \times 4\beta}{\sqrt{3.06^2 + (500\alpha)^2 + (6\beta)^2} \cdot \sqrt{2.68^2 + (320\alpha)^2 + (4\beta)^2}}$$

其余 sim(A,C)、sim(B,C) 同理。

### (b) $\alpha=0.01, \beta=0.5$

三台计算机的向量分别为：

- A: [3.06, 5, 3]
- B: [2.68, 3.2, 2]
- C: [2.92, 6.4, 3]

计算余弦相似度：

$$\text{sim}(A, B) = \frac{3.06 \times 2.68 + 5 \times 3.2 + 3 \times 2}{\sqrt{3.06^2 + 5^2 + 3^2} \cdot \sqrt{2.68^2 + 3.2^2 + 2^2}} \approx \frac{28.9408}{\sqrt{43.3636} \times \sqrt{22.5424}} \approx 0.926$$
$$\text{sim}(A, C) \approx 0.961 \quad \text{sim}(B, C) \approx 0.910$$

## 第2题

效用矩阵：

	a	b	c	d	e	f	g	h
A	4	5		5	1		3	2
B		3	4	3	1	2	1	
C	2		1	3		4	5	3

### (a) Jaccard 距离

转为布尔值（有评分为1，无为0）：

- A: [1 1 0 1 1 0 1 1]
- B: [0 1 1 1 1 1 1 0]
- C: [1 0 1 1 0 1 1 1]
- $A \cap B = 4, A \cup B = 8 \Rightarrow \text{距离} = 1 - 4/8 = 0.5$
- $A \cap C = 4, A \cup C = 8 \Rightarrow \text{距离} = 1 - 4/8 = 0.5$

- $B \cap C = 4, B \cup C = 8 \Rightarrow \text{距离} = 1 - 4/8 = 0.5$

### (b) 布尔余弦距离

用布尔向量计算余弦相似度，再取 1-相似度为距离：

- $A \cdot B = 4, |A| = \sqrt{6}, |B| = \sqrt{6} \Rightarrow \text{sim} \approx 0.667 \Rightarrow \text{dist} \approx 0.333$
- $A \cdot C = 4, |A| = \sqrt{6}, |C| = \sqrt{6} \Rightarrow \text{sim} \approx 0.667 \Rightarrow \text{dist} \approx 0.333$
- $B \cdot C = 4, |B| = \sqrt{6}, |C| = \sqrt{6} \Rightarrow \text{sim} \approx 0.667 \Rightarrow \text{dist} \approx 0.333$

### (c) 评分3-5为1，1-2为0，其余为0

- A: [1 1 0 1 0 0 1 0]
- B: [0 1 1 1 0 0 0 0]
- C: [0 0 0 1 0 1 1 1]

得出余弦距离：

- $A \cdot B = 2, |A| = \sqrt{4}, |B| = \sqrt{3} \Rightarrow \text{sim} \approx 0.577 \Rightarrow \text{dist} \approx 0.423$
- $A \cdot C = 2, |A| = \sqrt{4}, |C| = \sqrt{4} \Rightarrow \text{sim} = 0.5 \Rightarrow \text{dist} = 0.5$
- $B \cdot C = 1, |B| = \sqrt{3}, |C| = \sqrt{4} \Rightarrow \text{sim} \approx 0.289 \Rightarrow \text{dist} \approx 0.711$

### (d) 评分减去各用户均值

用户均值：

- A:  $(4+5+5+1+3+2)/6 = 3.33$
- B:  $(3+4+3+1+2+1)/6 = 2.33$
- C:  $(2+1+3+4+5+3)/6 = 3.0$

归一化：

	a	b	c	d	e	f	g	h
A	0.67	1.67		1.67	-2.33		-0.33	-1.33
B		0.67	1.67	0.67	-1.33	-0.33	-1.33	
C	-1.0		-2.0	0.0		1.0	2.0	0.0

计算余弦距离（找出共同评分项）：

用户对	归一化余弦相似度	归一化余弦距离
A-B	0.823	<b>0.177</b>
A-C	-0.263	<b>1.263</b>
B-C	-0.933	<b>1.933</b>

## 第3题

原始矩阵：

$$M = \begin{bmatrix} 1 & 2 & 3 \\ 3 & 4 & 5 \\ 5 & 4 & 3 \\ 0 & 2 & 4 \\ 1 & 3 & 5 \end{bmatrix}$$

(a)

$$MM^T = \begin{bmatrix} 1 & 2 & 3 \\ 3 & 4 & 5 \\ 5 & 4 & 3 \\ 0 & 2 & 4 \\ 1 & 3 & 5 \end{bmatrix} \cdot \begin{bmatrix} 1 & 3 & 5 & 0 & 1 \\ 2 & 4 & 4 & 2 & 3 \\ 3 & 5 & 3 & 4 & 5 \end{bmatrix} = \begin{bmatrix} 14 & 26 & 22 & 16 & 22 \\ 26 & 50 & 46 & 28 & 40 \\ 22 & 46 & 50 & 20 & 32 \\ 16 & 28 & 20 & 20 & 26 \\ 22 & 40 & 32 & 26 & 35 \end{bmatrix}$$

$$M^T M = \begin{bmatrix} 1 & 3 & 5 & 0 & 1 \\ 2 & 4 & 4 & 2 & 3 \\ 3 & 5 & 3 & 4 & 5 \end{bmatrix} \cdot \begin{bmatrix} 1 & 2 & 3 \\ 3 & 4 & 5 \\ 5 & 4 & 3 \\ 0 & 2 & 4 \\ 1 & 3 & 5 \end{bmatrix} = \begin{bmatrix} 36 & 37 & 38 \\ 37 & 49 & 61 \\ 38 & 61 & 84 \end{bmatrix}$$

(b)

解如下特征多项式：

$$\det(M^T M - \lambda I) = 0$$

$$\lambda_1 \approx 153.57, \lambda_2 \approx 15.43, \lambda_3 \approx 0$$

(c)

这些特征向量已归一化，对应于上面三个特征值（列向量）：

$$V = \begin{bmatrix} -0.409 & -0.816 & 0.408 \\ -0.563 & -0.126 & -0.816 \\ -0.718 & 0.564 & 0.408 \end{bmatrix}$$

(d)

SVD 形式为：

$$M = U \Sigma V^T$$

- 奇异值  $\Sigma$ （对角线）为特征值平方根：

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \sigma_2 & 0 \\ 0 & 0 & \sigma_3 \end{bmatrix} = \begin{bmatrix} 12.392 & 0 & 0 \\ 0 & 3.928 & 0 \\ 0 & 0 & \approx 0 \end{bmatrix}$$

- 矩阵  $V$ （即  $V^T$  的转置）：

$$V = \begin{bmatrix} -0.409 & -0.563 & -0.718 \\ -0.816 & -0.126 & 0.564 \\ 0.408 & -0.816 & 0.408 \end{bmatrix}$$

- 矩阵  $U$ （截取前三列）：

$$U = \begin{bmatrix} -0.298 & 0.159 & 0.941 \\ -0.571 & -0.033 & -0.153 \\ -0.521 & -0.736 & -0.052 \\ -0.323 & 0.510 & -0.196 \\ -0.459 & 0.414 & -0.224 \end{bmatrix}$$

## 第4题

布隆过滤器参数:  $m = 80$  亿 bit,  $n = 10$  亿

$$f = (1 - e^{-kn/m})^k$$

令  $k=3$ :

$$f \approx (1 - e^{-3 \times 10^9 / 8 \times 10^9})^3 \approx (1 - e^{-0.375})^3 \approx (1 - 0.687)^3 \approx 0.313^3 \approx 0.031$$

令  $k=4$  同理, 结果  $\approx 0.017$

## 第5题

(a)  $h(x) = 2x + 1 \bmod 32$

流: 3, 1, 4, 1, 5, 9, 2, 6, 5

哈希结果:

- $3 \rightarrow 7 \rightarrow 111 \rightarrow$  尾长0
- $1 \rightarrow 3 \rightarrow 011 \rightarrow$  尾长0
- ...尾长分别为: [0, 0, 0, 0, 0, 0, 0, 0]

最大尾长=0  $\Rightarrow$  估计数目:  $2^0 / \varphi \approx 1 / 0.773 \approx 1.29$

(b)  $h(x) = 4x \bmod 32$

- $3 \rightarrow 12 \rightarrow 1100 \rightarrow$  尾长2
- $1 \rightarrow 4 \rightarrow 0100 \rightarrow$  尾长2
- ...尾长分别为: [2, 2, 4, 2, 2, 2, 2, 3, 2]

最大尾长=4  $\Rightarrow$  估计数目:  $2^4 / \varphi \approx 16 / 0.773 \approx 20.69$

## 第6题

给定的数据流为:

$x = [3, 1, 4, 1, 3, 4, 2, 1, 2]$

先计算样本的均值:

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\mu = \frac{3 + 1 + 4 + 1 + 3 + 4 + 2 + 1 + 2}{9} = \frac{21}{9} = 2.333 \dots$$

二阶中心矩

$$\mu_2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

计算每一项：

x_i	x_i - μ	(x_i - μ) <sup>2</sup>
3	0.6667	0.4444
1	-1.3333	1.7778
4	1.6667	2.7778
1	-1.3333	1.7778
3	0.6667	0.4444
4	1.6667	2.7778
2	-0.3333	0.1111
1	-1.3333	1.7778
2	-0.3333	0.1111

求和：

$$\frac{0.4444 + 1.7778 + 2.7778 + 1.7778 + 0.4444 + 2.7778 + 0.1111 + 1.7778 + 0.1111}{9} = \frac{12.0}{9} = 1.3333$$

所以：

二阶矩（奇异数）=1.3333

三阶中心矩

$$\mu_3 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^3$$

计算每一项：

x_i	x_i - μ	(x_i - μ) <sup>3</sup>
3	0.6667	0.296
1	-1.3333	-2.370
4	1.6667	4.630
1	-1.3333	-2.370
3	0.6667	0.296
4	1.6667	4.630
2	-0.3333	-0.037
1	-1.3333	-2.370

x_i	x_i - μ	(x_i - μ) <sup>3</sup>
2	-0.3333	-0.037

求和：

$$\frac{0.296 - 2.370 + 4.630 - 2.370 + 0.296 + 4.630 - 0.037 - 2.370 - 0.037}{9} = \frac{2.668}{9} \approx 0.296$$

所以：

三阶矩=0.296

---

## 第7题

---

k 值	真实值	估计值	误差
5	3	3	0
15	9	10	1