**zID: z5224151**

**Name: Zanning Wang**

Q1:



```
Q1
  answer = < >
  while A1 != NULL and B != NULL
      if (docID (P1) = docID (P2)) then
          P1 ← next (P1)
          P2 ← next (P2)
      else if (docID (P2) < docID (P1) then
          Add (answers, docID (P1)),
          P2 ← next (P2)
          P1 ← skipTo (P2)

      else
          P2 ← next (P2)
```

Q2:
(1)



1) To decode the $V$ encoded, we should calculate the $kd$. $kr$.
as we know $kd = \lfloor \log_2 x \rfloor$ ,

length: $|kd| = \lfloor \log_2 x \rfloor + 1$ ; $kr = x - 2^{\lfloor \log_2 x \rfloor}$

$1 \le kr < \frac{x}{2}$ ~~Because.~~

length: ~~length of~~ $|kr| = \lfloor \log_2 kr \rfloor + 1 = \log_2 x$.

The total $r$ code. =. $kd + 0 + kr$ , which the $kd$ no longer than $\lfloor \log_2 x \rfloor + 1$
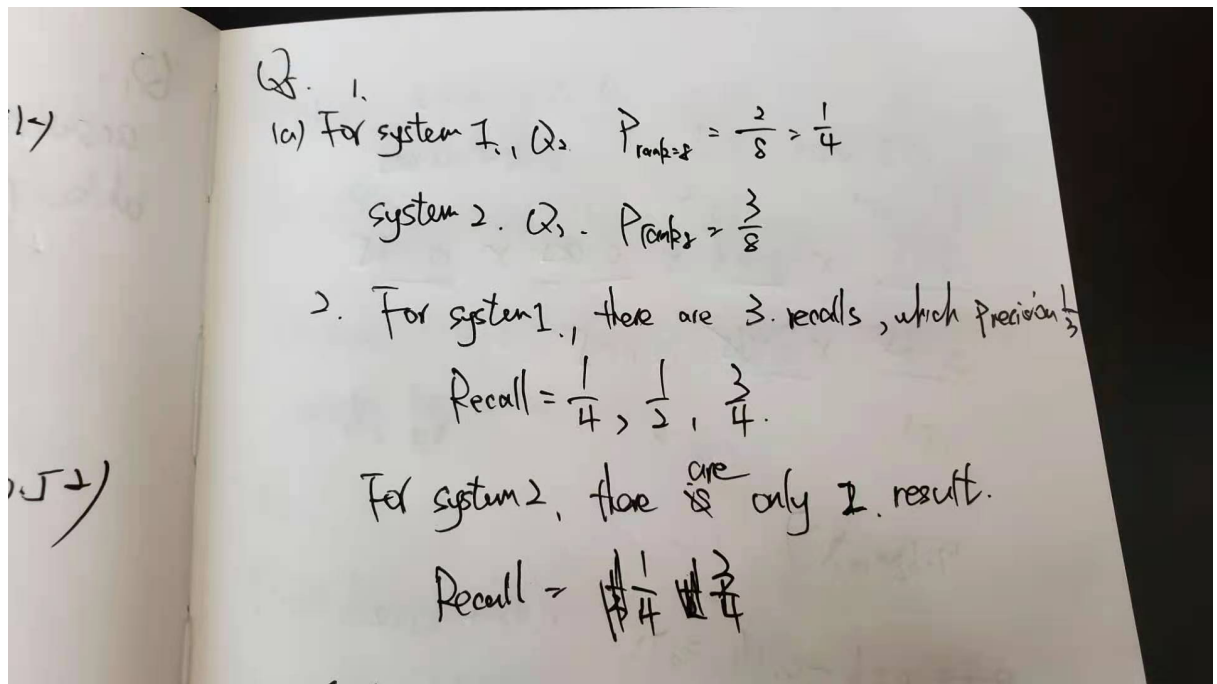the $kr$ is no longer than $\log_2 x$, therefore. $\log_2 x + 1 + \log_2 x < 2\log_2 x + 1$

Q3:

Q4:

Q4

(a) Because, there is already $I_0, I_1, I_2, I_3$, in the disk, we will create a new $I_4$ and merge the above index, therefore, the sub-index is $I_4$.

(b) the number of sub-index is $\dfrac{|C|}{|M|}$

therefore the sub-index that the algorithm create is to take the log of the number of sub-index, which is $\log_2 \dfrac{|C|}{|M|}$

Q5:

Q. 1.

(a) For system $1$, $Q_2$. $P_{rank=8} = \frac{2}{8} = \frac{1}{4}$

System $2$. $Q_3$. $P_{rank=8} = \frac{3}{8}$

2. For system 1, there are 3 recalls, which Precision $\frac{1}{3}$

Recall $= \frac{1}{4}, \frac{1}{2}, \frac{3}{4}$.

For system 2, there are only 2 result.

Recall $= \frac{1}{4} \quad \frac{3}{4}$

$$MAP(\theta) = \frac{1}{|\theta|} \sum_{j=1}^{|\theta|} \frac{1}{m_j} \sum_{k=1}^{m_j} Precision(R_j | y)$$

(b) System 1

$Q_1 : \frac{1}{6}(1+1+1+\frac{4}{9}+\frac{3}{5}) = 0.86$

$Q_2 = \frac{1}{4} \times (1+\frac{1}{3}+\frac{3}{9}+\frac{4}{10}) = 0.52$

Map$(S_1) = \frac{1}{2}(Q_1 + Q_2) = \frac{1}{2}(0.86 + 0.52)$
$= 0.69$

System 2 :

$Q_1 : \frac{1}{6}(1+1+1+\frac{4}{5}+\frac{5}{6}+\frac{2}{3}) = 0.88$

$Q_2 : \frac{1}{4}(1+\frac{2}{4}+\frac{3}{5}) = 0.53$

Map$(S_2) = \frac{1}{2}(Q_1 + Q_2) = 0.71$

c) Q, S₁

$P_1 = 1 \qquad R_1 = \frac{1}{6} = 0.16$
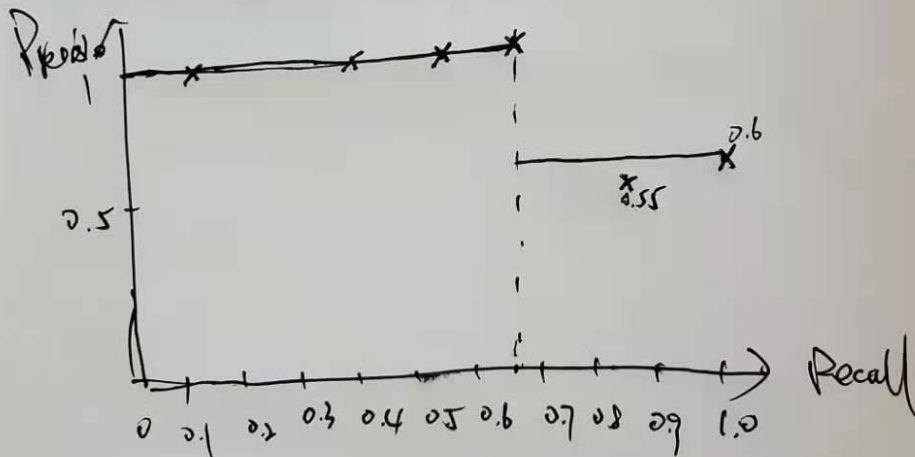
$P_2 = 1 \qquad R_2 = \frac{1}{3} = 0.33$

$P_3 = 1 \qquad R_3 = \frac{1}{2} = 0.50$

$P_4 = 1 \qquad R_4 = \frac{2}{3} = 0.67$

$P_5 = \frac{5}{9} = 0.55 \qquad R_5 = \frac{5}{6} = 0.83$

$P_6 = \frac{6}{10} \qquad R_6 = 1 = 1.0$



therefore when $R = 0.5 \quad P = 1$

when $R = 0.8 \quad P = 0.6$.

Q6:

Q7:

Q7.

(a) $P(Q \mid d_1) = \prod_{x \in Q} P(x \mid d_1) = \frac{2}{10} \cdot \frac{3}{10} \cdot \frac{1}{10} \cdot \frac{2}{10} \cdot \frac{2 \times 0}{10} = \cancel{}$

$= \cancel{2.4 \times 10^{-4}}$ 0

$P(Q \mid d_2) = \prod_{x \in Q} P(x \mid d_2) = \frac{7}{10} \times \frac{1}{10} \times \frac{1}{10} \times \frac{1}{10} \times \frac{0 \times 0}{} = \cancel{7 \times 10} $ 0

$= 0$

Because $P(Q \mid d_2) \overset{\sim}{=} P(Q \mid d_1)$, therefore document id 1

~~would be ranked higher~~

and document rank the same.

(b)

$$P(Q \mid d_1) = \left(0.8 \times \frac{2}{10} + 0.2 \times 0.8\right)^x \left(0.8 \times \frac{3}{10} + 0.2 \times 0.1\right)^x$$

$$\left(0.8 \times \frac{1}{10} + 0.2 \times 0.025\right)^x \left(0.8 \times \frac{2}{10} + 0.2 \times 0.025\right)^x$$

$$\left(0.8 \times \frac{2}{10} + 0.2 \times 0.025\right)^x \left(0.8 \times \frac{0}{10} + 0.2 \times 0.025\right)$$

$$= 9.62 \times 10^{-7}$$

$$P(Q \mid d_2) = \left(0.8 \times \frac{7}{10} + 0.2 \times 0.8\right) \times \left(0.8 \times \frac{1}{10} + 0.2 \times \frac{0.1}{10}\right) \times$$

$$\ast \left(0.8 \times \frac{1}{10} + 0.2 \times 0.025\right) \times \left(0.8 \times \frac{1}{10} + 0.2 \times 0.025\right) \times$$

$$\left(0.8 \times 0 + 0.2 \times 0.025\right) \times \left(0.8 \times 0 + 0.2 \times 0.025\right)$$

$$= 1.30 \times 10^{-8}$$

$P(Q \mid d_1) > P(Q \mid d_2)$    $P(Q \mid d_1)$ rank higher.

Q8:

(a)

In the content seen module, we will check if the current URL has been crawl, in this process, we will store as an index and it will also pass to URL filter to check the URL valid or not, after that the "Dup url elim can check what URL is duplication or not

Qe. (b)

Because green based shingle $\{1, 7, 15, 81\}$.

$$\Rightarrow h_1(x) = \{8, 6, 0, 10\}$$

$$h_2(x) = \{0, 2, 4, 1\}$$

therefore for $h_1(x)$, the third figure is the.

Min Hash ~~for h1~~, which is 15.

for $h_2(x)$, the first figure is the MinHash which is 1.

therefore, the resulting Min Hash signature is $\{1, 0\}$