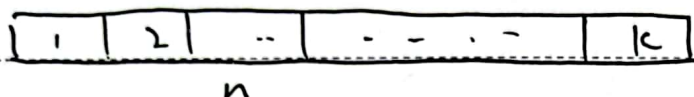


2.  $n$  位内存容量，集合  $S$  中有  $m$  个成员，如果使用  $k$  个哈希函数，则误判率为

$$(1 - e^{-\frac{km}{n}})^k$$

如果将内存分为  $k$  组，每组  $\frac{n}{k}$  位，记  $t = \frac{n}{k}$



则首先假设哈希函数的选择是完全随机的。那么，任意一次哈希选中这一位的概率为  $\frac{1}{t} = \frac{1}{n/k}$ ，因此没有选中这一位的概率为  $1 - \frac{1}{t}$ 。

插入一个元素需要哈希  $k$  次，所以经过  $k$  次哈希之后，某个特定元素未被置为 1 的概率为  $(1 - \frac{1}{t})^k$

在  $m$  个元素都被插入后，该位仍为 0 的概率为

$$P(X_i = 0) = (1 - \frac{1}{t})^{km} \approx e^{-\frac{km}{t}}$$

$\therefore$  在  $m$  个元素都被插入后，该位为 1 的概率是

$$P(X_i = 1) = 1 - P(X_i = 0) \approx 1 - e^{-\frac{km}{t}}$$

只有当某个元素经过  $k$  个组哈希之后，对应  $t$  组中的  $t$  个位直都恰好被置为 1 才会发生误判，则误判率为

$$(1 - e^{-\frac{km}{t}})^k = (1 - e^{-\frac{km^2}{n}})^k$$



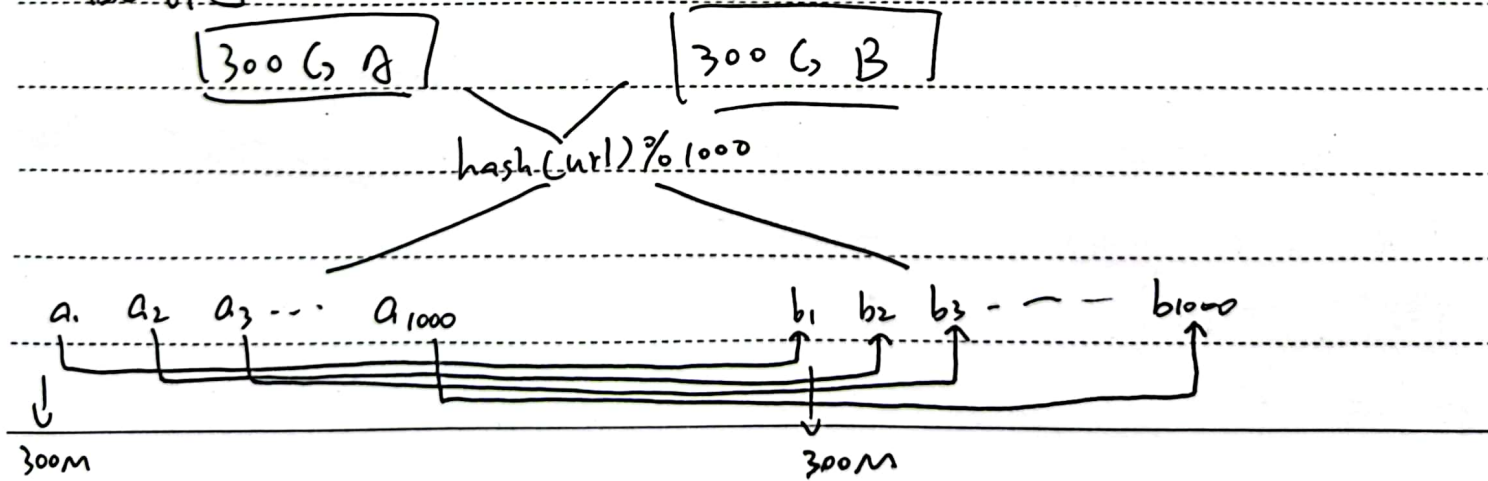
3.  $50 \text{亿} \times 64 \text{B} \div 1024 \div 1024 \div 1024 = 298 \text{G} \approx 300 \text{G}$

而内存限制为 4GB, 所以无法一次读入内存, 所以需要  
采用大文件切割的方法

假设将 a, b 两个大文件都分割成 1000 个小文件, 所以每个小  
文件大小为  $300 \text{G} \div 1000 \times 1024 = 307 \text{M} \approx 300 \text{M}$ , 所以同时加载  
两个文件需要  $300 \text{M} \times 2 = 600 \text{M}$

对于 a, b 两个小文件的集合运用相同的哈希函数, 这个每个  
小文件集合中的小文件都有独立的哈希值

如果两个 url 相同, 则这两个 url 所在的文件序号也一定相  
同, 所以比较的时候只需要在序号相同的两个文件中进行 url 匹  
配即可



5.  $A: \{1, 2, 3, 4\}$   $B: \{2, 3, 5, 7\}$   $C: \{2, 4, 6\}$

$$\text{Jaccard}(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

$$\text{Jaccard}(A, B) = \frac{2}{6} = \frac{1}{3}$$

$$\text{Jaccard}(B, C) = \frac{1}{6}$$

$$\text{Jaccard}(A, C) = \frac{2}{5}$$



## 7. 证明

$$\therefore P(\text{mh}(S_1) = \text{mh}(S_2)) = \text{Jaccard}(S_1, S_2)$$

如果  $\text{Jaccard}(S_1, S_2) = 0$

$$\text{则 } P(\text{mh}(S_1) = \text{mh}(S_2)) = 0$$

即  $\text{mh}(S_1)$  与  $\text{mh}(S_2)$  有 100% 的概率不相等

即用 Min-Hashing 方法可以知道

$S_1$  与  $S_2$  两个集合必不相似

此时 Min-Hashing 一定可以给出正确的估计





6. 题目中的  $X_i$  是独立的伯努利随机向量序列

$$P(X_i=1) = P(h_i(S_1) = h_i(S_2)) = J_S(S_1, S_2)$$

令  $X = \sum_{i=1}^k X_i$  和  $\mu = \sum_{i=1}^k p_i$ , 对任意  $\epsilon \in (0, 1)$ , 则

Chernoff 不等式

$$P(X > (1+\epsilon)\mu) < \exp\left(-\frac{\mu\epsilon^2}{3}\right)$$

$$\hat{J}_S(S_1, S_2) = \frac{1}{k} \sum_{i=1}^k X_i = \frac{X}{k}$$

$$E(\hat{J}_S(S_1, S_2)) = \frac{1}{k} \times k \times E(X_i) = P(h_i(S_1) = h_i(S_2)) = J_S(S_1, S_2) = \frac{\mu}{k}$$

$$P(|\hat{J}_S(S_1, S_2) - J_S(S_1, S_2)| > \epsilon J_S(S_1, S_2))$$

$$= P(\hat{J}_S(S_1, S_2) > J_S(S_1, S_2) + \epsilon J_S(S_1, S_2)) + P(\hat{J}_S(S_1, S_2) < J_S(S_1, S_2) - \epsilon J_S(S_1, S_2)) < \underbrace{P(\hat{J}_S(S_1, S_2) > (1+\epsilon)J_S(S_1, S_2))}$$

$$P(X > (1+\epsilon)\mu) < \exp\left(-\frac{\mu\epsilon^2}{3}\right)$$

$$\Downarrow P\left(\frac{X}{k} < \frac{(1+\epsilon)\mu}{k}\right) < \exp\left(-\frac{\mu\epsilon^2}{3}\right)$$

$$\therefore P(\hat{J}_S(S_1, S_2) > (1+\epsilon)J_S(S_1, S_2)) < \exp\left(-\frac{k J_S(S_1, S_2) \epsilon^2}{3}\right)$$

需比较其与  $\delta$  的大小

$$k = O\left(\frac{-\ln \delta}{J_S \epsilon^2}\right) \therefore \frac{-\ln \delta}{J_S \epsilon^2} \leq k \quad -\ln \delta \leq J_S \cdot \epsilon^2 \cdot k$$

$$\ln \delta \geq -J_S \cdot \epsilon^2 \cdot k$$

$$\delta \geq \exp(-J_S \cdot \epsilon^2 \cdot k) \geq \exp\left(-\frac{k J_S(S_1, S_2) \epsilon^2}{3}\right)$$

$$\therefore P(\hat{J}_S(S_1, S_2) > (1+\epsilon)J_S(S_1, S_2)) < \delta$$

$$\therefore P(|\hat{J}_S(S_1, S_2) - J_S(S_1, S_2)| > \epsilon J_S(S_1, S_2)) < \delta$$

得证



9.

$$(1) \text{Jaccard}(S_1, S_2) = \frac{1}{4}$$

$$\text{Jaccard}(S_2, S_3) = 0$$

$$\text{Jaccard}(S_1, S_3) = \frac{1}{4}$$

$$(7x+1) \bmod 6 \quad (11x+2) \bmod 6 \quad (5x+2) \bmod 6$$

(2)	$S_1$	$S_2$	$S_3$	$h_1(x)$	$h_2(x)$	$h_3(x)$
0	1	1	0	1	2	2
1	0	1	0	2	1	1
2	1	0	0	3	0	0
3	0	0	1	4	5	5
4	1	0	1	5	4	4
5	0	0	0	0	3	3

$$mh_1(S_1) = 1 \quad mh_1(S_2) = 1 \quad mh_1(S_3) = 4$$

$$mh_2(S_1) = 0 \quad mh_2(S_2) = 1 \quad mh_2(S_3) = 4$$

$$mh_3(S_1) = 0 \quad mh_3(S_2) = 1 \quad mh_3(S_3) = 4$$

最小哈希签名矩阵

	$S_1$	$S_2$	$S_3$
$h_1$	1	1	4
$h_2$	0	1	4
$h_3$	0	1	4



11. 每个 signature 向量被分成  $b$  段, 每段  $r$  行

两个 signature 向量的任意一个段所有行都相同  $t^r$

至少有一行不相同  $1 - t^r$

所有  $b$  段都不同  $(1 - t^r)^b$

至少有一个段相同  $1 - (1 - t^r)^b$

由题目可知两个集合成为 candidate 用户 (即被哈希到同一个桶中的概率为  $\frac{1}{2}$ )

$$\therefore 1 - (1 - t^r)^b = \frac{1}{2}$$

$$(1 - t^r)^b = \frac{1}{2}$$

$$1 - t^r = 2^{-\frac{1}{b}}$$

$$t^r = 1 - 2^{-\frac{1}{b}}$$

$$t = [1 - 2^{-\frac{1}{b}}]^{\frac{1}{r}}$$

$$t = \sqrt[r]{1 - 2^{-\frac{1}{b}}}$$

