

温兆和. 10205501432. 数据科学算法作业3

2. 解: 当假设哈希函数均匀, 当一个集合元素被插入过滤器组中某位未置1的概率为 $1 - \frac{k}{n}$ (每组各位)

由于将元素哈希到各组中是分别、独立进行的, 故整个 n 位内存中的任何一位未被置1的概率为 $1 - \frac{k}{n}$.

此处与布隆过滤器不同。布隆过滤器中是 k 个哈希函数将元素哈希进同一集合, 此时某特定位未被置1的概率为 $(1 - \frac{1}{n})^k$.

将所有 m 个元素插入此过滤器后, 整个内存中任一位未被置1的概率为 $(1 - \frac{k}{n})^m$.

故某一位被置1的概率为 $1 - (1 - \frac{k}{n})^m$.

当某元素不在集合中, 误判概率为 $(1 - (1 - \frac{k}{n})^m)^k$.

3. 解: 可将 a 文件中的 50 亿个 url 插入大小为 $4 \times 2^{30} \times 8$ 个比特位的布隆过滤器中, 再分别判断 b 中的各个 url 是否在布隆过滤器中.

为使误判概率更低, 需要 $k = \ln 2 \cdot \frac{4 \times 2^{30} \times 8}{50 \times 10^8} \approx 4.76$ 个哈希函数. 当 $k=5$ 时, 误判率为 0.0369, 比 $k=4$ 时的误判率 0.0579 更低.

所以这个布隆过滤器需要 5 个哈希函数

$$5. \text{解: } J(\{1, 2, 3, 4\}, \{2, 3, 5, 7\}) = \frac{2}{6} = \frac{1}{3}.$$

$$J(\{2, 3, 5, 7\}, \{2, 4, 6\}) = \frac{1}{6}.$$

$$J(\{1, 2, 3, 4\}, \{2, 4, 6\}) = \frac{2}{5}$$

7. 证: 由于 $P(h_{\pi}(C_1) = h_{\pi}(C_2)) = \text{Jaccard}(C_1, C_2) = 0$.
 故对任意随机排列 π , $h_{\pi}(C_1) \neq h_{\pi}(C_2)$ 必定成立.
 若在矩阵分组后, 这两个集合必定被映射到不同的桶中.
 中. 我们认为不映射到同一桶中的集合不相似.
 故: 我们一定可以给出一个正确的估计, 即认为这两个集合不相似.

8. 证: 诸 $X_i \sim b(1, J_S)$, $E(X_i) = J_S$, $\text{Var}(X_i) = J_S(1-J_S)$
 $\hat{J}_S = \frac{1}{k} \sum_{i=1}^k X_i$ 故 $E(\hat{J}_S) = J_S$, $\text{Var}(\hat{J}_S) = \frac{J_S(1-J_S)}{k}$
 由切比雪夫不等式: $P(|\hat{J}_S - J_S| > \epsilon J_S (S_1, S_2))$
 $\leq \frac{\text{Var}(\hat{J}_S)}{\epsilon^2 J_S^2} = \frac{1-J_S}{\epsilon^2 J_S k} \stackrel{②}{=} \frac{1-J_S}{J \ln(1/\delta)}$
 $< k \cdot \frac{1}{\ln(1/\delta)} < \delta$?

9. 证: (1). $J(S_1, S_2) = \frac{1}{4}$, $J(S_1, S_3) = \frac{1}{4}$, $J(S_2, S_3) = 0$.
 (2). 经哈希后 T_3 行号变为:

| $h_1(x)$ | $h_2(x)$ | $h_3(x)$ |
|----------|----------|----------|
| 1 | 2 | 2 |
| 2 | 1 | 1 |
| 3 | 0 | 0 |
| 4 | 5 | 5 |
| 5 | 4 | 4 |
| 0 | 3 | 3 |

故, 每一列的最小哈希值为:

| | s_1 | s_2 | s_3 |
|-------|-------|-------|-------|
| h_1 | 1 | 1 | ✓ |
| h_2 | 0 | 1 | ✓ |
| h_3 | 0 | 1 | ✓ |

11. 证: 当 b, r 较大, 在相似度阈值右侧, 集合有很大可能
能被哈希到同一个桶中; 在左侧, 则不太可能被哈希
到同一个桶中. 故, 我们只要求 s 曲线最陡的地方即
可.

$$f(s) = 1 - (1 - s^r)^b$$

$$f''(s) = -b(b-1)(1-s^r)^{b-2}(rs^{r-1})^2 + b(1-s^r)^{b-1}r(r-1)s^{r-2}$$

$$\text{令 } f''(s) = 0, \text{ 得: } (1-s^r)(r-1) = (b-1)rs^r$$

$$(1-s^r)(r-1) = (b-1)rs^r$$

$$s = \left(\frac{r-1}{br-1} \right)^{\frac{1}{r}} \approx \left(\frac{r}{br} \right)^{\frac{1}{r}} = \left(\frac{1}{b} \right)^{\frac{1}{r}}$$