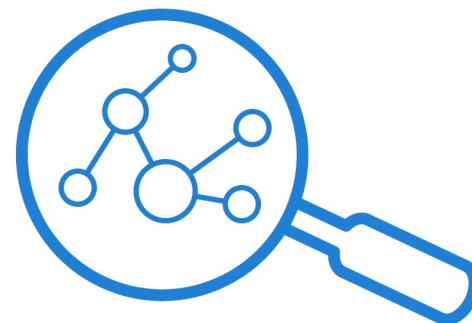


数据科学与大数据技术的 数学基础



第二讲



计算机学院

余皓然

2024/4/25

课程内容

Part1 随机化方法

一致性哈希 布隆过滤器 CM Sketch方法 最小哈希
欧氏距离下的相似搜索 Jaccard相似度下的相似搜索

Part2 谱分析方法

主成分分析 奇异值分解 谱图论

Part3 最优化方法

压缩感知



布隆过滤器

从属问题



内容分发网络

Akamai（阿卡迈）搭建的内容分发网络有部署在136个国家的约30万个边缘服务器，为Apple、Microsoft等公司提供缓存服务

问题：具体应该缓存哪些文件？所有当地用户曾经请求过的文件？



内容分发网络

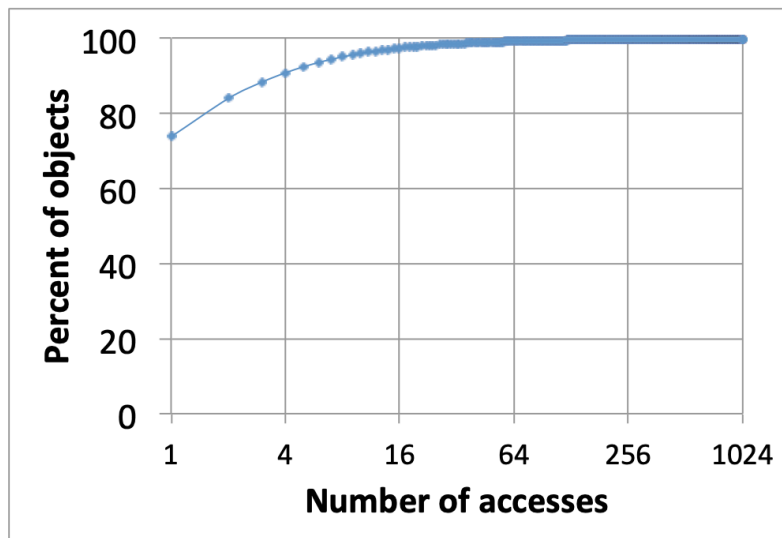


Figure 5: On a typical CDN server cluster serving web traffic over two days, 74% of the roughly 400 million objects in cache were accessed only once and 90% were accessed less than four times.

(2015年的论文) 分析了Akamai的一个有45台服务器的集群，统计在两天时间内4亿份文件被用户访问的频次

说明不需要缓存所有曾经被访问的文件

内容分发网络

一种简单的缓存策略：当用户请求访问某文件，先判断该文件是否是第一次被请求：

- (1) 如果是，不缓存；
- (2) 如果不是，缓存。

* 传输策略：如果用户请求的文件已被缓存，则从本级服务器提取文件进行传输；否则，从上级服务器中查找/提取文件进行传输



内容分发网络

一种简单的缓存策略：当用户请求访问某文件，先判断该文件是否是第一次被请求：

- (1) 如果是，不缓存；
- (2) 如果不是，缓存。

* 传输策略：如果用户请求的文件已被缓存，则从本级服务器提取文件进行传输；否则，从上级服务器中查找/提取文件进行传输

可以**减少**存储空间的浪费（即缓存的文件从未被用户再次请求）



如何判断文件是否是第一次被请求？

本质上需要解决“从属问题”



内容分发网络

从属判断问题 (membership query)

如何存储关于集合 S 的信息从而可以准确判断关于“元素 x 是否属于集 S ”的问题?

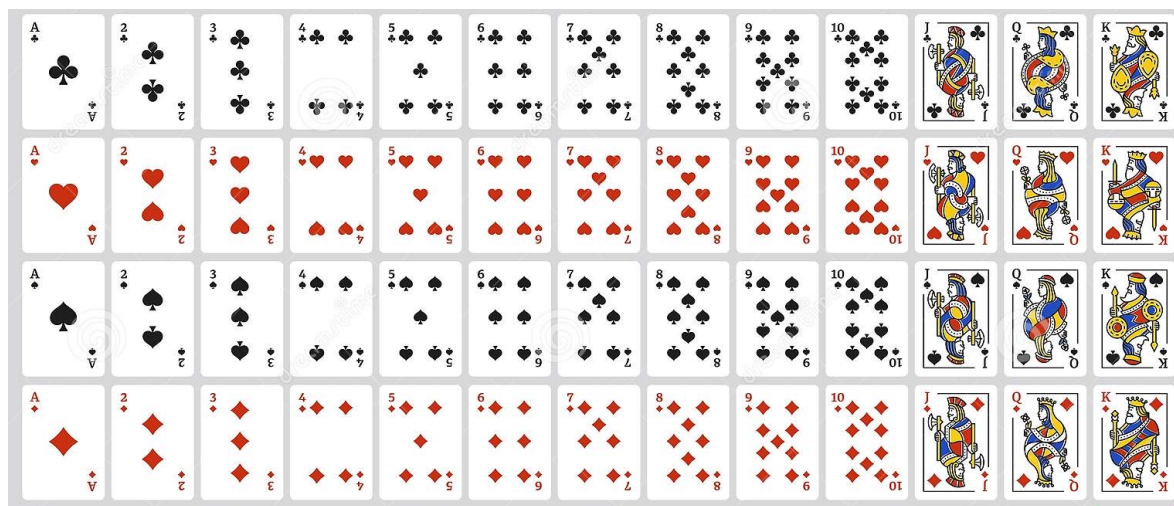
需要尽量减少用于存储的空间

先分析所需存储空间的下限



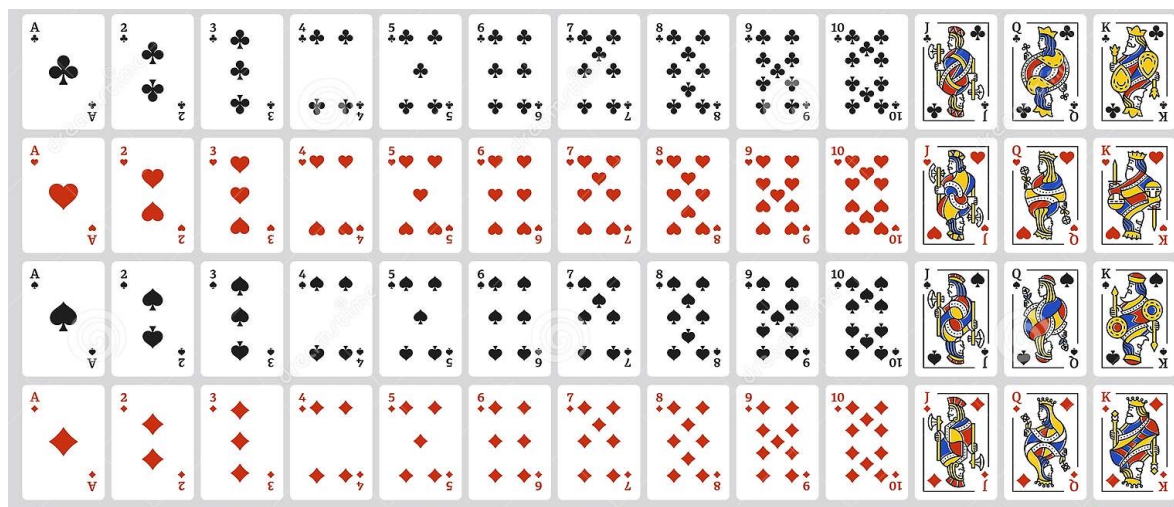
内容分发网络

已知52张牌，（不重复地）选取4张展示。至少需要多少空间（多少比特）完整存储该信息（即哪4张牌被展示了）？



内容分发网络

已知52张牌，（不重复地）选取4张展示。至少需要多少空间（多少比特）完整存储该信息（即哪4张牌被展示了）？



$$\left\lceil \log_2 \binom{52}{4} \right\rceil$$

内容分发网络

已知有 $|u|$ 个不同的物品，选取 n 件展示。至少需要多少空间（多少比特）完整存储该信息（即哪 n 件被展示了）？

$$\left\lceil \log_2 \binom{|u|}{n} \right\rceil$$



内容分发网络

已知有 $|U|$ 个不同的物品，选取 n 件展示。至少需要多少空间（多少比特）完整存储该信息（即哪 n 件被展示了）？

$$\begin{aligned}\left\lceil \log_2 \binom{|U|}{n} \right\rceil &= \lg \binom{|U|}{n} \\&= \lg \left(\frac{|U|!}{n! (|U|-n)!} \right) \\&\geq \lg \left(\frac{(|U|-n)^n}{n^n} \right) \\&= n \lg \left(\frac{|U|-n}{n} \right) \\&= n \lg \left(\frac{|U|}{n} - 1 \right) \\&\geq n \lg \left(\frac{|U|}{n} \right) - n \\&\geq n \lg |U| - n \lg n - n \\&= \mathbf{\Omega(n \lg |U| - n \lg n)}\end{aligned}$$

$$f(x) = \Omega(g(x)) \text{ as } x \rightarrow \infty \text{ if } \limsup_{x \rightarrow \infty} \left| \frac{f(x)}{g(x)} \right| > 0.$$

内容分发网络

已知有 $|U|$ 个不同的物品，选取 n 件展示。至少需要多少空间（多少比特）完整存储该信息（即哪 n 件被展示了）？

当 $|U| \gg n$ 时，需要的存储空间至少是 $\Omega(n \log_2 |U|)$ 比特



内容分发网络

已知有 $|U|$ 个不同的物品，选取 n 件展示。至少需要多少空间（多少比特）完整存储该信息（即哪 n 件被展示了）？

当 $|U| \gg n$ 时，需要的存储空间至少是 $\Omega(n \log_2 |U|)$ 比特

如果 $|U|$ 代表所有可能的网站URL、 n 是当地用户已经请求访问的网站URL，该数字非常大

准确记忆用户已经请求过哪些网站/文件将消耗大量的存储资源



如何以较小的存储空间解决从属问题？



内容分发网络

已知有 $|U|$ 个不同的物品，选取 n 件展示。至少需要多少空间（多少比特）完整存储该信息（即哪 n 件被展示了）？

当 $|U| \gg n$ 时，需要的存储空间至少是 $\Omega(n \log_2 |U|)$ 比特

如果 $|U|$ 代表所有可能的网站URL、 n 是当地用户已经请求访问的网站URL，该数字非常大

准确记忆用户已经请求过哪些网站/文件将消耗大量的存储资源



如何以较小的存储空间解决从属问题？

无论采用何种办法，如果要实现准确记忆，必然占用大量存储空间（已证理论下界）

采用“以精确度换存储空间”的思想

内容分发网络

近似从属判断问题 (approximate membership query)

如何存储关于集合 S 的信息从而可以近似判断关于“元素 x 是否属于集 S ”的问题?

当 $x \in S$ 时，要100%正确判断；

当 $x \notin S$ 时，要大概率正确判断：即以 $1 - \varepsilon$ 的概率返回“不属于”的答案



内容分发网络

近似从属判断问题 (approximate membership query)

如何存储关于集合 S 的信息从而可以近似判断关于“元素 x 是否属于集 S ”的问题?

当 $x \in S$ 时，要100%正确判断；

当 $x \notin S$ 时，要大概率正确判断：即以 $1 - \varepsilon$ 的概率返回“不属于”的答案

允许以小概率出现误报 (false positive) ，不允许出现漏报 (false negative)
满足缓存的实际应用



内容分发网络

近似从属判断问题 (approximate membership query)

如何存储关于集合 S 的信息从而可以近似判断关于“元素 x 是否属于集 S ”的问题?

当 $x \in S$ 时，要100%正确判断；

当 $x \notin S$ 时，要大概率正确判断：即以 $1 - \varepsilon$ 的概率返回“不属于”的答案

允许以小概率出现误报 (false positive) ，不允许出现漏报 (false negative)
满足缓存的实际应用



设计的存储方式的空间复杂度不能随 $|U|$ 增长

$|U|$: 所有可能的网站/文件的个数

$|S|$: 用户访问过的网站/文件的个数 (同前例子中的 n)

布隆过滤器

布隆过滤器的方法：思想一



布隆过滤器

布隆过滤器 (Bloom Filters) 由Burton Howard Bloom在1970年提出

基于以下两个思想解决近似从属判断问题：

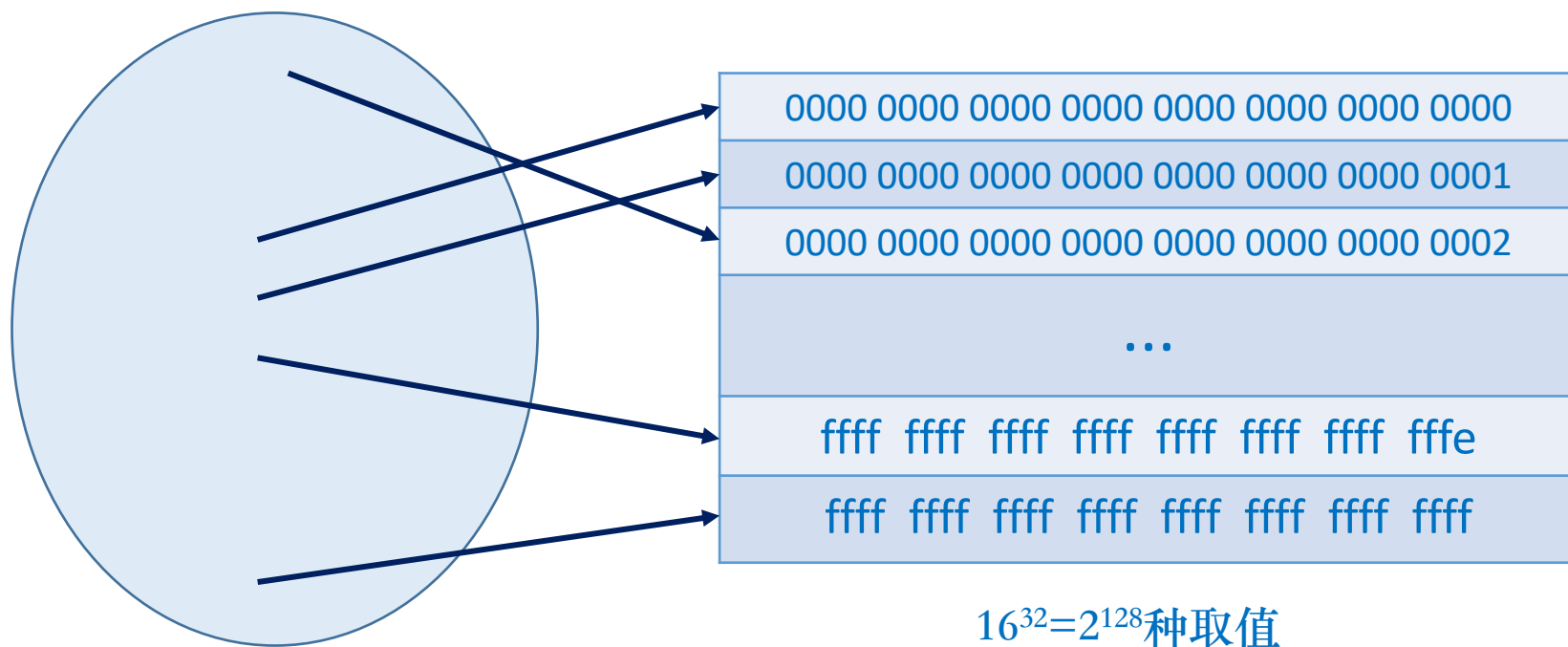
(1) 通过哈希函数将网站/文件名进行映射



布隆过滤器

(1) 通过哈希函数将网站/文件名进行映射

回顾哈希函数

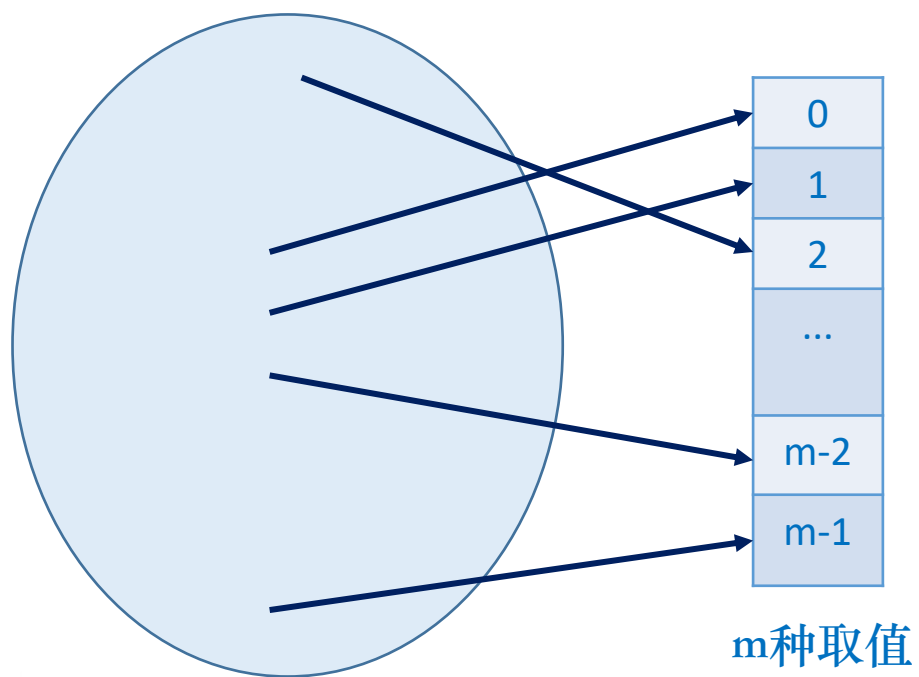


比如URL的集合

上一讲中关于MD5的例子

布隆过滤器

(1) 通过哈希函数将网站/文件名进行映射



比如URL的集合

布隆过滤器

(1) 通过哈希函数将网站/文件名进行映射

As an example, let's have
 $S = \{103, 137, 166, 271, 314\}$

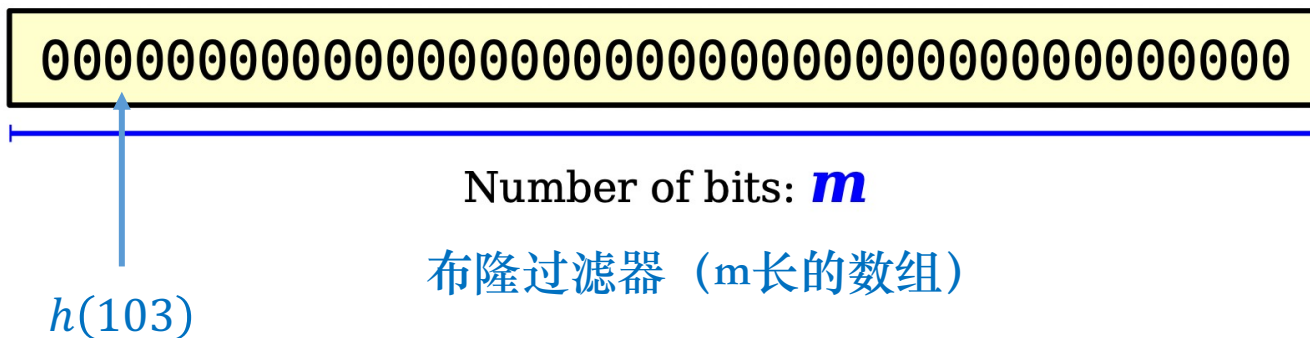


布隆过滤器

(1) 通过哈希函数将网站/文件名进行映射

As an example, let's have
 $S = \{103, 137, 166, 271, 314\}$

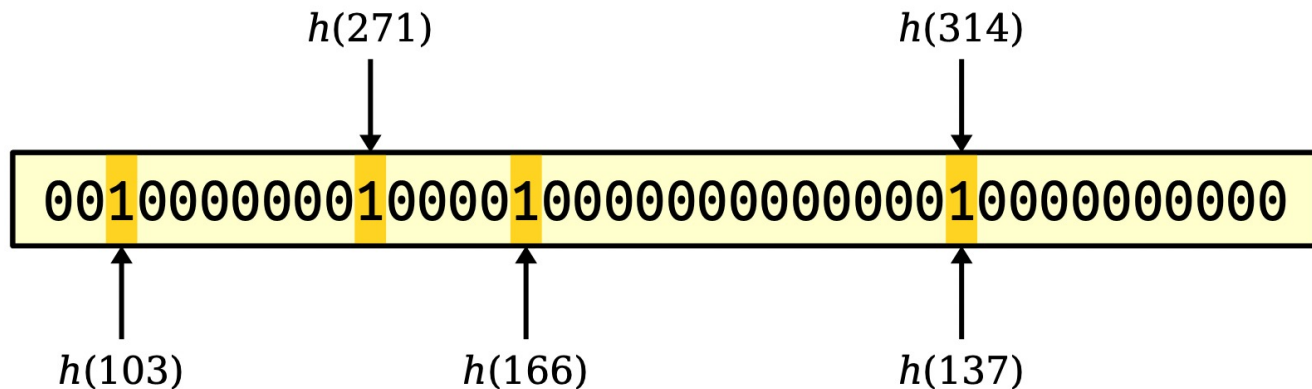
例如 $h(103) = 2$



布隆过滤器

(1) 通过哈希函数将网站/文件名进行映射

As an example, let's have
 $S = \{103, 137, 166, 271, 314\}$



可能出现两个不同元素哈希值相同的情况

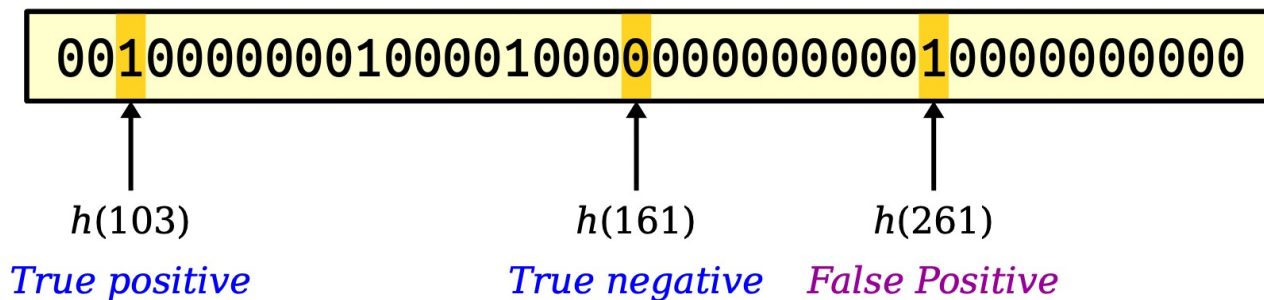
布隆过滤器

(1) 通过哈希函数将网站/文件名进行映射

As an example, let's have

$$S = \{103, 137, 166, 271, 314\}$$

用布隆过滤器判断103、161、261是否存在于集合S中



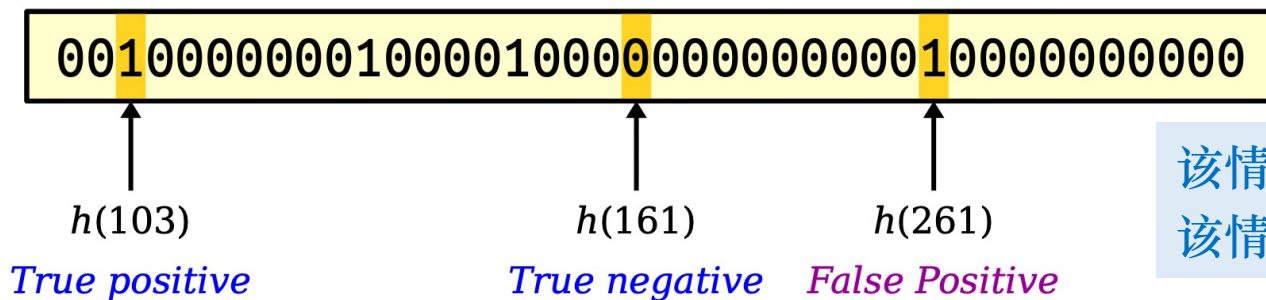
即便261不在集合中，布隆过滤器仍会判断“在集合中”
原因是261的哈希值恰好等于137和314的哈希值

布隆过滤器

(1) 通过哈希函数将网站/文件名进行映射

As an example, let's have
 $S = \{103, 137, 166, 271, 314\}$

用布隆过滤器判断103、161、261是否存在于集合 S 中



该情况不可避免，当 m 增大时
该情况出现的可能性如何变化？

即便261不在集合中，布隆过滤器仍会判断“在集合中”
原因是261的哈希值恰好等于137和314的哈希值

布隆过滤器

(1) 通过哈希函数将网站/文件名进行映射

当布隆过滤器（数组）的长度 m 增大时，哈希值的范围增大，出现“误报”（false positive）的可能性降低



布隆过滤器

(1) 通过哈希函数将网站/文件名进行映射

当布隆过滤器（数组）的长度 m 增大时，哈希值的范围增大，出现“误报”（false positive）的可能性降低



如果目标是当 $x \notin S$ 时，出现误报的概率小于等于 ε ，应该如何设置 m ?

* 已知集合 S 的大小为 n （即 $|S| = n$ ）



布隆过滤器

(1) 通过哈希函数将网站/文件名进行映射

当布隆过滤器（数组）的长度 m 增大时，哈希值的范围增大，出现“误报”（false positive）的可能性降低



如果目标是当 $x \notin S$ 时，出现误报的概率小于等于 ε ，应该如何设置 m ?

* 已知集合 S 的大小为 n （即 $|S| = n$ ）

元素 x 对应的哈希值 $h(x)$ 与集合 S 中第一个元素的哈希值相同的概率?



布隆过滤器

(1) 通过哈希函数将网站/文件名进行映射

当布隆过滤器（数组）的长度 m 增大时，哈希值的范围增大，出现“误报”（false positive）的可能性降低



如果目标是当 $x \notin S$ 时，出现误报的概率小于等于 ε ，应该如何设置 m ?

* 已知集合 S 的大小为 n （即 $|S| = n$ ）

元素 x 对应的哈希值 $h(x)$ 与集合 S 中第一个元素的哈希值相同的概率为 $\frac{1}{m}$

元素 x 对应的哈希值 $h(x)$ 与集合 S 中至少一个元素的哈希值相同的概率为



布隆过滤器

(1) 通过哈希函数将网站/文件名进行映射

当布隆过滤器（数组）的长度 m 增大时，哈希值的范围增大，出现“误报”（false positive）的可能性降低



如果目标是当 $x \notin S$ 时，出现误报的概率小于等于 ε ，应该如何设置 m ?

* 已知集合 S 的大小为 n （即 $|S| = n$ ）

元素 x 对应的哈希值 $h(x)$ 与集合 S 中第一个元素的哈希值相同的概率为 $\frac{1}{m}$

元素 x 对应的哈希值 $h(x)$ 与集合 S 中至少一个元素的哈希值相同的概率为 $1 - \left(1 - \frac{1}{m}\right)^n$

即误报的概率为 $1 - \left(1 - \frac{1}{m}\right)^n \leq \varepsilon$

布隆过滤器

(1) 通过哈希函数将网站/文件名进行映射

当布隆过滤器（数组）的长度 m 增大时，哈希值的范围增大，出现“误报”（false positive）的可能性降低



如果目标是当 $x \notin S$ 时，出现误报的概率小于等于 ε ，应该如何设置 m ?

* 已知集合 S 的大小为 n （即 $|S| = n$ ）

元素 x 对应的哈希值 $h(x)$ 与集合 S 中第一个元素的哈希值相同的概率为 $\frac{1}{m}$

元素 x 对应的哈希值 $h(x)$ 与集合 S 中至少一个元素的哈希值相同的概率为 $1 - \left(1 - \frac{1}{m}\right)^n$

即误报的概率为 $1 - \left(1 - \frac{1}{m}\right)^n \approx 1 - e^{-\frac{n}{m}} \leq \varepsilon$ 经整理得 $m \geq \frac{n}{\ln\left(\frac{1}{1-\varepsilon}\right)}$

$$\lim_{m \rightarrow \infty} \left(1 - \frac{1}{m}\right)^m = \frac{1}{e}$$

布隆过滤器

(1) 通过哈希函数将网站/文件名进行映射

当布隆过滤器（数组）的长度 m 增大时，哈希值的范围增大，出现“误报”（false positive）的可能性降低



如果目标是当 $x \notin S$ 时，出现误报的概率小于等于 ε ，应该如何设置 m ?

* 已知集合 S 的大小为 n （即 $|S| = n$ ）

元素 x 对应的哈希值 $h(x)$ 与集合 S 中第一个元素的哈希值相同的概率为 $\frac{1}{m}$

元素 x 对应的哈希值 $h(x)$ 与集合 S 中至少一个元素的哈希值相同的概率为 $1 - \left(1 - \frac{1}{m}\right)^n$

即误报的概率为 $1 - \left(1 - \frac{1}{m}\right)^n \approx 1 - e^{-\frac{n}{m}} \leq \varepsilon$ 经整理得 $m \geq \frac{n}{\ln\left(\frac{1}{1-\varepsilon}\right)}$

需要约 $\frac{n}{\ln\left(\frac{1}{1-\varepsilon}\right)}$ 比特的存储空间

例如，当 ε 为0.01时，该值约 $99.5n$ 能否改进从而降低该值?

布隆过滤器

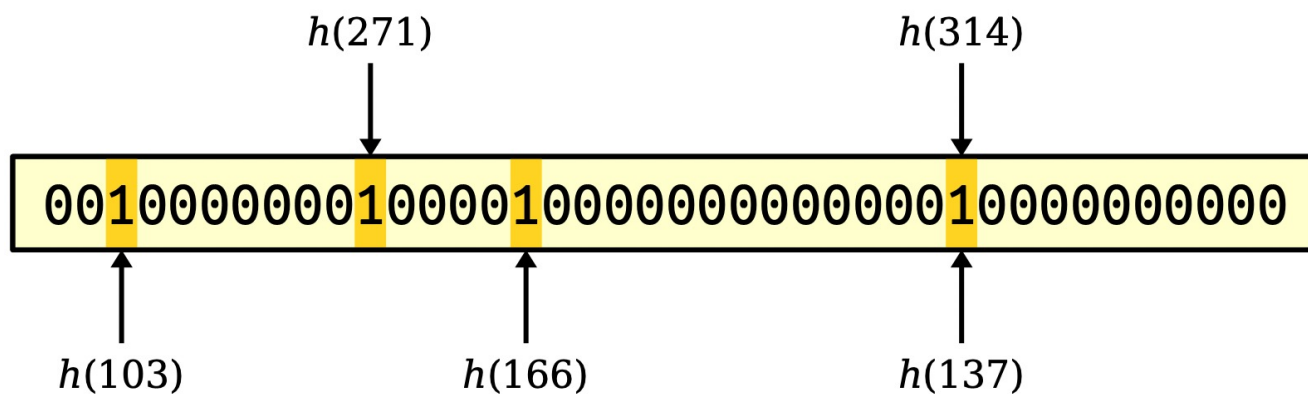
布隆过滤器的方法：思想二



布隆过滤器

- (1) 通过哈希函数将网站/文件名进行映射
- (2) 利用多个 (k 个) 哈希函数

只考虑一个哈希函数:



- (1) 通过哈希函数将网站/文件名进行映射
- (2) 利用多个 (k 个) 哈希函数

The diagram shows a horizontal bar representing a 100-bit array. The bits are displayed as a sequence of '0's and '1's. Three specific bits are highlighted with yellow vertical bars: the 103rd bit (a '1'), the 166th bit (a '1'), and the 314th bit (a '1'). Arrows point from labels below the bar to these highlighted bits: $h(103)$ points to the first highlighted bit, $h(166)$ points to the second, and $h(314)$ points to the third. Above the bar, two labels $h(271)$ and $h(314)$ have arrows pointing down to the bar, but they do not point to any highlighted bits.

The diagram illustrates a 32-bit hash function. A 32-bit input string, `0010010100100000100000010000000100000001000000`, is shown in a yellow box. Above the box, two 16-bit outputs are indicated: $h_2(271) h_1(271)$ and $h_1(314) h_2(314)$. Below the box, two 16-bit outputs are indicated: $h_1(103) h_2(137)$ and $h_1(166) h_2(166)$. Arrows point from the input string to the output labels. A red label $h_2(103) h_1(137)$ is also present, pointing to the input string.

$$h_1(103) \text{ } h_2(137) \quad h_1(166) \quad h_2(166) \quad h_1(137)h_2(103)$$

- (1) 通过哈希函数将网站/文件名进行映射
- (2) 利用多个 (k个) 哈希函数

Diagram illustrating a 32-bit hash function structure. The input sequence is 001001010010000100000010000000100001000000. The inputs are grouped into pairs: $h_2(271)$ and $h_1(271)$ for the first two segments, $h_1(314)$ and $h_2(314)$ for the seventh and eighth segments, $h_1(103)$ and $h_2(137)$ for the first segment, $h_1(166)$ and $h_2(166)$ for the third segment, and $h_1(137)$ and $h_2(103)$ for the seventh segment. The label h 为 m is shown on the left.

注意过滤器长度不变，仍为 m

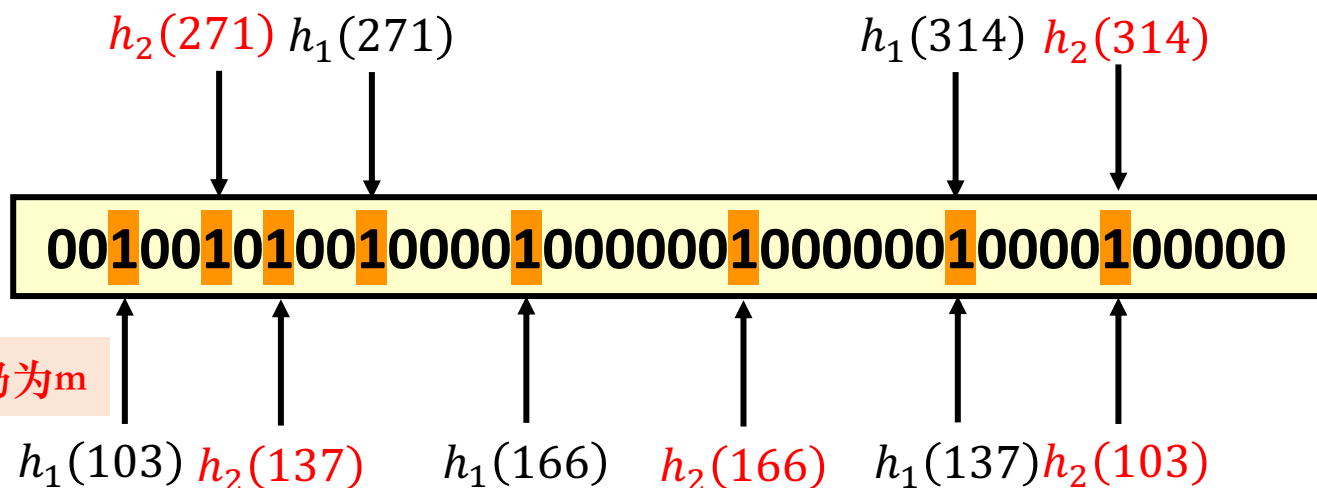
如何利用布隆过滤器判断元素是否在集合内? 比如查询261是否在集合内?

布隆过滤器

- (1) 通过哈希函数将网站/文件名进行映射
- (2) 利用多个 (k个) 哈希函数

考虑两个哈希函数:

注意过滤器长度不变, 仍为m

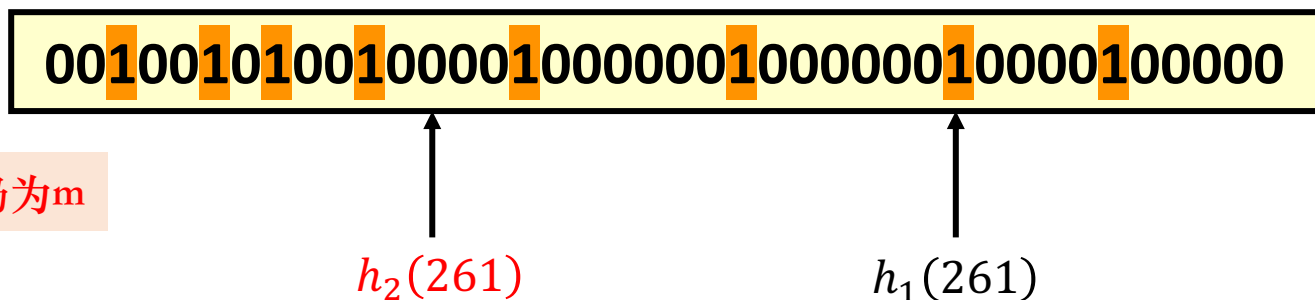


如何查询261是否在集合内? 计算 $h_1(261)$ 和 $h_2(261)$, 检查过滤器相应位置是否都是1

布隆过滤器

- (1) 通过哈希函数将网站/文件名进行映射
- (2) 利用多个 (k个) 哈希函数

考虑两个哈希函数:



注意过滤器长度不变，仍为m

如何查询261是否在集合内? 计算 $h_1(261)$ 和 $h_2(261)$, 检查过滤器相应位置是否都是1

因为只有当相应位置都是1的时候才会判断“在集合中”，可以减小“误报”的概率

如何计算“误报”的概率?

布隆过滤器

- (1) 通过哈希函数将网站/文件名进行映射
- (2) 利用多个 (k个) 哈希函数



如果目标是当 $x \notin S$ 时，出现误报的概率小于等于 ε ，应该如何设置 m ?

回顾：若仅用一个哈希函数，为使误报率小于等于 ε ，所需存储空间（即 m 大小）约 $\frac{n}{\ln\left(\frac{1}{1-\varepsilon}\right)}$ 比特

元素 x 对应的哈希值 $h_1(x)$ 在过滤器的位置取值为1的概率?



布隆过滤器

- (1) 通过哈希函数将网站/文件名进行映射
- (2) 利用多个 (k 个) 哈希函数



如果目标是当 $x \notin S$ 时，出现误报的概率小于等于 ε ，应该如何设置 m ?

回顾：若仅用一个哈希函数，为使误报率小于等于 ε ，所需存储空间（即 m 大小）约 $\frac{n}{\ln\left(\frac{1}{1-\varepsilon}\right)}$ 比特

元素 x 对应的哈希值 $h_1(x)$ 在过滤器的位置取值为 1 的概率?

集合中每个元素分别对应 k 个哈希值，即共有 nk 个哈希值，其中至少有一个值对应该位置的概率为?



布隆过滤器

- (1) 通过哈希函数将网站/文件名进行映射
- (2) 利用多个 (k个) 哈希函数



如果目标是当 $x \notin S$ 时，出现误报的概率小于等于 ε ，应该如何设置 m ?

回顾：若仅用一个哈希函数，为使误报率小于等于 ε ，所需存储空间（即 m 大小）约 $\frac{n}{\ln\left(\frac{1}{1-\varepsilon}\right)}$ 比特

元素 x 对应的哈希值 $h_1(x)$ 在过滤器的位置取值为1的概率?

集合中每个元素分别对应 k 个哈希值，即共有 nk 个哈希值，其中至少有一个值对应该位置的概率为？ $1 - \left(1 - \frac{1}{m}\right)^{kn} \approx 1 - e^{-\frac{kn}{m}}$



布隆过滤器

- (1) 通过哈希函数将网站/文件名进行映射
- (2) 利用多个 (k个) 哈希函数



如果目标是当 $x \notin S$ 时，出现误报的概率小于等于 ε ，应该如何设置 m ?

回顾：若仅用一个哈希函数，为使误报率小于等于 ε ，所需存储空间（即 m 大小）约 $\frac{n}{\ln\left(\frac{1}{1-\varepsilon}\right)}$ 比特

元素 x 对应的哈希值 $h_1(x)$ 在过滤器的位置取值为1的概率为 $1 - e^{-\frac{kn}{m}}$

元素 x 对应的哈希值 $h_2(x)$ 在过滤器的位置取值为1的概率为 $1 - e^{-\frac{kn}{m}}$

元素 x 对应的哈希值 $h_3(x)$ 在过滤器的位置取值为1的概率为 $1 - e^{-\frac{kn}{m}}$

...

元素 x 对应的哈希值 $h_k(x)$ 在过滤器的位置取值为1的概率为 $1 - e^{-\frac{kn}{m}}$

元素 x 对应的 k 个哈希值在过滤器的相应位置取值都为1的概率为? $\left(1 - e^{-\frac{kn}{m}}\right)^k$?

布隆过滤器

- (1) 通过哈希函数将网站/文件名进行映射
- (2) 利用多个 (k个) 哈希函数



如果目标是当 $x \notin S$ 时，出现误报的概率小于等于 ε ，应该如何设置 m ?

回顾：若仅用一个哈希函数，为使误报率小于等于 ε ，所需存储空间（即 m 大小）约 $\frac{n}{\ln\left(\frac{1}{1-\varepsilon}\right)}$ 比特

元素 x 对应的哈希值 $h_1(x)$ 在过滤器的位置取值为1的概率为 $1 - e^{-\frac{kn}{m}}$

元素 x 对应的哈希值 $h_2(x)$ 在过滤器的位置取值为1的概率为 $1 - e^{-\frac{kn}{m}}$

元素 x 对应的哈希值 $h_3(x)$ 在过滤器的位置取值为1的概率为 $1 - e^{-\frac{kn}{m}}$

...

元素 x 对应的哈希值 $h_k(x)$ 在过滤器的位置取值为1的概率为 $1 - e^{-\frac{kn}{m}}$

元素 x 对应的 k 个哈希值在过滤器的相应位置取值都为1的概率为? $\left(1 - e^{-\frac{kn}{m}}\right)^k$?

其实不等于 $\left(1 - e^{-\frac{kn}{m}}\right)^k$ ，但实际概率与 $\left(1 - e^{-\frac{kn}{m}}\right)^k$ 相近。为便于分析，将 $\left(1 - e^{-\frac{kn}{m}}\right)^k$ 作为近似值

布隆过滤器

- (1) 通过哈希函数将网站/文件名进行映射
- (2) 利用多个 (k个) 哈希函数



如果目标是当 $x \notin S$ 时，出现误报的概率小于等于 ε ，应该如何设置 m ?

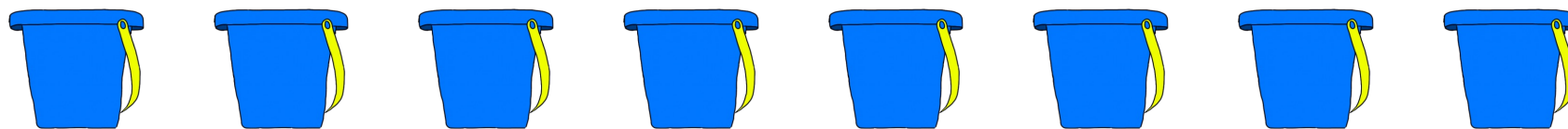
回顾：若仅用一个哈希函数，为使误报率小于等于 ε ，所需存储空间（即 m 大小）约 $\frac{n}{\ln\left(\frac{1}{1-\varepsilon}\right)}$ 比特

$\Pr[\text{元素 } x \text{ 对应的 } k \text{ 个哈希值在过滤器的相应位置取值都为 } 1]$

$= \Pr[h_1(x) \text{ 位置取值为 } 1; h_2(x) \text{ 位置取值为 } 1; \dots; h_k(x) \text{ 位置取值为 } 1]$

是否等于 $\Pr[h_1(x) \text{ 位置取值为 } 1] \times \Pr[h_2(x) \text{ 位置取值为 } 1] \dots \times \Pr[h_k(x) \text{ 位置取值为 } 1]$

布隆过滤器



有8个桶子，每个桶子里可能有水也可能没水，具体有多少个桶里有水未知

先后依次让20个人来随机选择1个桶，随机事件“第 i 个人选择的桶中有水”与随机事件“第 j 个人选择的桶中有水”是否是独立事件？



布隆过滤器

- (1) 通过哈希函数将网站/文件名进行映射
- (2) 利用多个 (k个) 哈希函数



如果目标是当 $x \notin S$ 时，出现误报的概率小于等于 ε ，应该如何设置 m ?

回顾：若仅用一个哈希函数，为使误报率小于等于 ε ，所需存储空间（即 m 大小）约 $\frac{n}{\ln\left(\frac{1}{1-\varepsilon}\right)}$ 比特

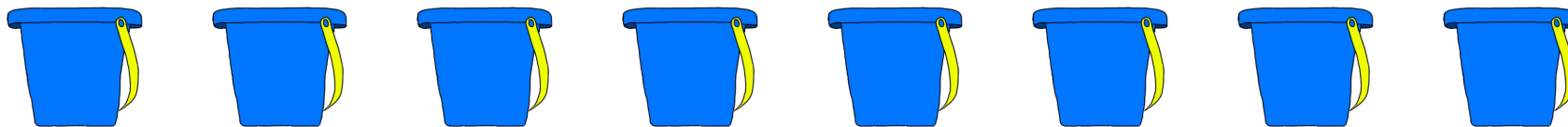
$\Pr[\text{元素 } x \text{ 对应的 } k \text{ 个哈希值在过滤器的相应位置取值都为 } 1]$

$= \Pr[h_1(x) \text{ 位置取值为 } 1; h_2(x) \text{ 位置取值为 } 1; \dots; h_k(x) \text{ 位置取值为 } 1]$

不等于! $\Pr[h_1(x) \text{ 位置取值为 } 1] \times \Pr[h_2(x) \text{ 位置取值为 } 1] \dots \times \Pr[h_k(x) \text{ 位置取值为 } 1]$

元素 x 对应的 k 个哈希值在过滤器的相应位置取值都为 1 的概率不等于 $\left(1 - e^{-\frac{kn}{m}}\right)^k$

布隆过滤器



有8个桶子，每个桶子里可能有水也可能没水，**已知有4个桶子里面有水**

先后依次让20个人来随机选择1个桶，随机事件“第 i 个人选择的桶中有水”与随机事件“第 j 个人选择的桶中有水”是否是独立事件？



布隆过滤器

- (1) 通过哈希函数将网站/文件名进行映射
- (2) 利用多个 (k个) 哈希函数



如果目标是当 $x \notin S$ 时，出现误报的概率小于等于 ε ，应该如何设置 m ?

回顾：若仅用一个哈希函数，为使误报率小于等于 ε ，所需存储空间（即 m 大小）约 $\frac{n}{\ln\left(\frac{1}{1-\varepsilon}\right)}$ 比特

$\Pr[\text{元素 } x \text{ 对应的 } k \text{ 个哈希值在过滤器的相应位置取值都为 } 1]$

$= \Pr[h_1(x) \text{ 位置取值为 } 1; h_2(x) \text{ 位置取值为 } 1; \dots; h_k(x) \text{ 位置取值为 } 1]$

$= \sum_{t=0}^m \Pr[m \text{ 位置中有 } t \text{ 个为 } 0] \Pr[h_1(x) \text{ 位置取值为 } 1; \dots; h_k(x) \text{ 位置取值为 } 1 \mid m \text{ 个位置中有 } t \text{ 个为 } 0]$

$= \sum_{t=0}^m \Pr[m \text{ 位置中有 } t \text{ 个为 } 0] (\Pr[h_1(x) \text{ 位置取值为 } 1 \mid m \text{ 位置中有 } t \text{ 个为 } 0] \times \dots \times \Pr[h_k(x)$

$\text{位置取值为 } 1 \mid m \text{ 位置中有 } t \text{ 个为 } 0])$

布隆过滤器

- (1) 通过哈希函数将网站/文件名进行映射
- (2) 利用多个 (k个) 哈希函数



如果目标是当 $x \notin S$ 时，出现误报的概率小于等于 ε ，应该如何设置 m ?

回顾：若仅用一个哈希函数，为使误报率小于等于 ε ，所需存储空间（即 m 大小）约 $\frac{n}{\ln\left(\frac{1}{1-\varepsilon}\right)}$ 比特

$\Pr[\text{元素 } x \text{ 对应的 } k \text{ 个哈希值在过滤器的相应位置取值都为 } 1]$

$= \Pr[h_1(x) \text{ 位置取值为 } 1; h_2(x) \text{ 位置取值为 } 1; \dots; h_k(x) \text{ 位置取值为 } 1]$

$= \sum_{t=0}^m \Pr[m \text{ 位置中有 } t \text{ 个为 } 0] \Pr[h_1(x) \text{ 位置取值为 } 1; \dots; h_k(x) \text{ 位置取值为 } 1 \mid m \text{ 个位置中有 } t \text{ 个为 } 0]$

$= \sum_{t=0}^m \Pr[m \text{ 位置中有 } t \text{ 个为 } 0] (\Pr[h_1(x) \text{ 位置取值为 } 1 \mid m \text{ 位置中有 } t \text{ 个为 } 0] \times \dots \times \Pr[h_k(x) \text{ 位置取值为 } 1 \mid m \text{ 位置中有 } t \text{ 个为 } 0])$

$= \sum_{t=0}^m \Pr[m \text{ 位置中有 } t \text{ 个为 } 0] \left(1 - \frac{t}{m}\right)^k$

$\approx \sum_{t=0}^m \Pr[m \text{ 位置中有 } t \text{ 个为 } 0] e^{-\frac{kt}{m}}$

后续计算比较复杂，最终结果并不易于分析

布隆过滤器

- (1) 通过哈希函数将网站/文件名进行映射
- (2) 利用多个 (k个) 哈希函数



如果目标是当 $x \notin S$ 时，出现误报的概率小于等于 ε ，应该如何设置m?

回顾：若仅用一个哈希函数，为使误报率小于等于 ε ，所需存储空间（即m大小）约 $\frac{n}{\ln\left(\frac{1}{1-\varepsilon}\right)}$ 比特

元素 x 对应的 k 个哈希值在过滤器的相应位置取值都为1的概率约为 $\left(1 - e^{-\frac{kn}{m}}\right)^k$ ，令其小于等于 ε

布隆过滤器

- (1) 通过哈希函数将网站/文件名进行映射
- (2) 利用多个 (k个) 哈希函数



如果目标是当 $x \notin S$ 时，出现误报的概率小于等于 ε ，应该如何设置 m ?

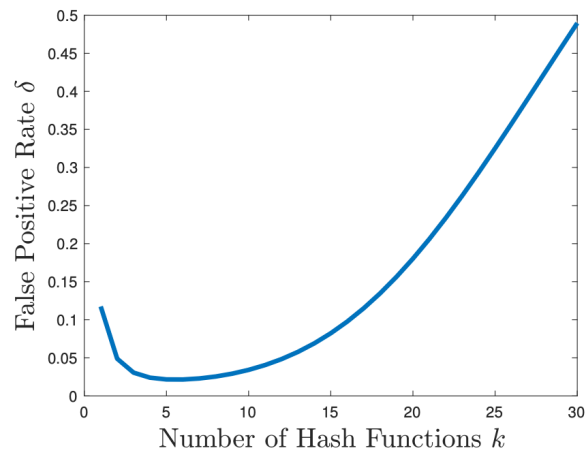
回顾：若仅用一个哈希函数，为使误报率小于等于 ε ，所需存储空间（即 m 大小）约 $\frac{n}{\ln\left(\frac{1}{1-\varepsilon}\right)}$ 比特

元素 x 对应的 k 个哈希值在过滤器的相应位置取值都为1的概率约为 $\left(1 - e^{-\frac{kn}{m}}\right)^k$ ，令其小于等于 ε

先选择 k ，最小化 $\left(1 - e^{-\frac{kn}{m}}\right)^k$

根据先减后增的趋势，导数为0处取最优的 k 值

为什么先减后增?



布隆过滤器

- (1) 通过哈希函数将网站/文件名进行映射
- (2) 利用多个 (k个) 哈希函数



如果目标是当 $x \notin S$ 时，出现误报的概率小于等于 ε ，应该如何设置 m ?

回顾：若仅用一个哈希函数，为使误报率小于等于 ε ，所需存储空间（即 m 大小）约 $\frac{n}{\ln\left(\frac{1}{1-\varepsilon}\right)}$ 比特

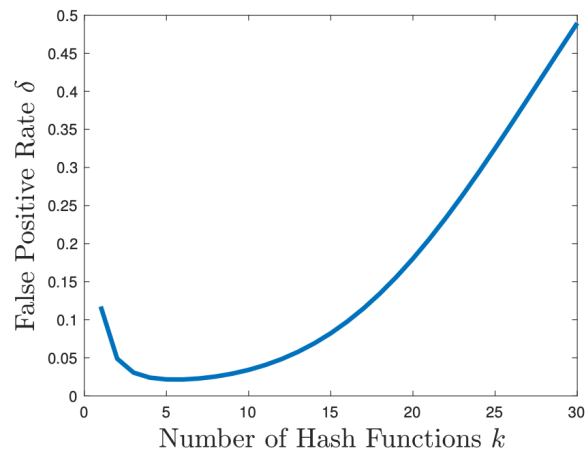
元素 x 对应的 k 个哈希值在过滤器的相应位置取值都为1的概率约为 $\left(1 - e^{-\frac{kn}{m}}\right)^k$ ，令其小于等于 ε

先选择 k ，最小化 $\left(1 - e^{-\frac{kn}{m}}\right)^k$

根据先减后增的趋势，导数为0处取最优的 k 值

方法一：求导找导数为零的条件

$$\frac{d}{dx} \left[a^{f(x)} \right] = \underbrace{a^{f(x)}}_{\text{rewrite}} \cdot \underbrace{\ln a}_{\text{natural log of base}} \cdot \underbrace{f'(x)}_{\text{derivative of exponent}}$$



布隆过滤器

- (1) 通过哈希函数将网站/文件名进行映射
- (2) 利用多个 (k个) 哈希函数



如果目标是当 $x \notin S$ 时，出现误报的概率小于等于 ε ，应该如何设置 m ?

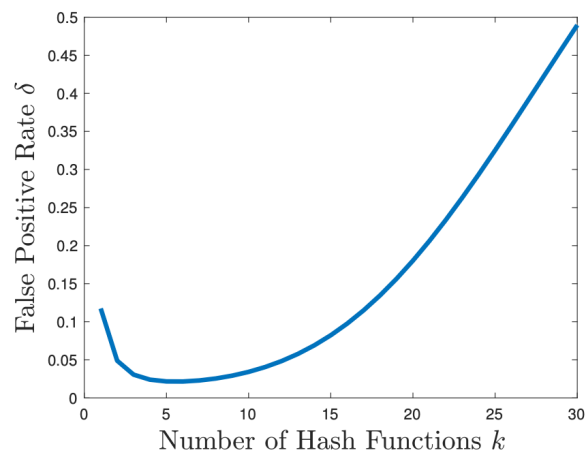
回顾：若仅用一个哈希函数，为使误报率小于等于 ε ，所需存储空间（即 m 大小）约 $\frac{n}{\ln\left(\frac{1}{1-\varepsilon}\right)}$ 比特

元素 x 对应的 k 个哈希值在过滤器的相应位置取值都为1的概率约为 $\left(1 - e^{-\frac{kn}{m}}\right)^k$ ，令其小于等于 ε

先选择 k ，最小化 $\left(1 - e^{-\frac{kn}{m}}\right)^k$

根据先减后增的趋势，导数为0处取最优的 k 值

方法二：利用关系 $k \ln\left(1 - e^{-\frac{kn}{m}}\right) = -\frac{m}{n} \ln e^{-\frac{kn}{m}} \ln\left(1 - e^{-\frac{kn}{m}}\right)$



布隆过滤器

- (1) 通过哈希函数将网站/文件名进行映射
- (2) 利用多个 (k个) 哈希函数



如果目标是当 $x \notin S$ 时，出现误报的概率小于等于 ε ，应该如何设置 m ?

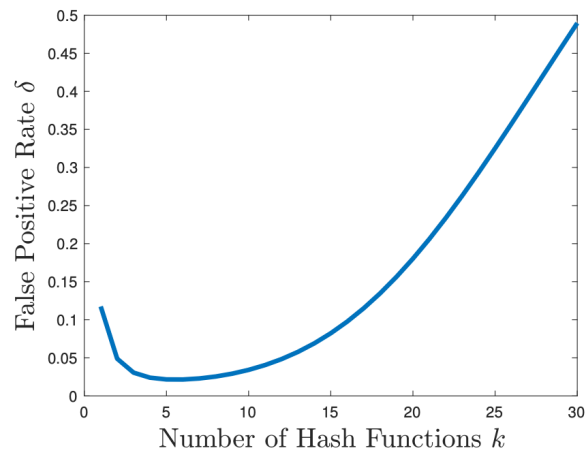
回顾：若仅用一个哈希函数，为使误报率小于等于 ε ，所需存储空间（即 m 大小）约 $\frac{n}{\ln\left(\frac{1}{1-\varepsilon}\right)}$ 比特

元素 x 对应的 k 个哈希值在过滤器的相应位置取值都为1的概率约为 $\left(1 - e^{-\frac{kn}{m}}\right)^k$ ，令其小于等于 ε

先选择 k ，最小化 $\left(1 - e^{-\frac{kn}{m}}\right)^k$

根据先减后增的趋势，导数为0处取最优的 k 值

最优 k 值为 $\frac{m}{n} \ln 2$ ，相应概率为 $2^{-\frac{m}{n} \ln 2}$



布隆过滤器

- (1) 通过哈希函数将网站/文件名进行映射
- (2) 利用多个 (k个) 哈希函数



如果目标是当 $x \notin S$ 时，出现误报的概率小于等于 ε ，应该如何设置 m ?

回顾：若仅用一个哈希函数，为使误报率小于等于 ε ，所需存储空间（即 m 大小）约 $\frac{n}{\ln\left(\frac{1}{1-\varepsilon}\right)}$ 比特

元素 x 对应的 k 个哈希值在过滤器的相应位置取值都为 1 的概率约为 $2^{-\frac{m}{n} \ln 2}$ ，令其小于等于 ε

m 的取值应该满足： $m \geq \frac{\log_2 \frac{1}{\varepsilon}}{\ln 2} n$



布隆过滤器

- (1) 通过哈希函数将网站/文件名进行映射
- (2) 利用多个 (k个) 哈希函数



如果目标是当 $x \notin S$ 时，出现误报的概率小于等于 ε ，应该如何设置 m ?

回顾：若仅用一个哈希函数，为使误报率小于等于 ε ，所需存储空间（即 m 大小）约 $\frac{n}{\ln\left(\frac{1}{1-\varepsilon}\right)}$ 比特

元素 x 对应的 k 个哈希值在过滤器的相应位置取值都为 1 的概率约为 $2^{-\frac{m}{n} \ln 2}$ ，令其小于等于 ε

m 的取值应该满足： $m \geq \frac{\log_2 \frac{1}{\varepsilon}}{\ln 2} n$

若用 k 个哈希函数，为使误报率小于等于 ε ，所需存储空间（即 m 大小）约 $\frac{\log_2 \frac{1}{\varepsilon}}{\ln 2} n$ 比特

布隆过滤器

- (1) 通过哈希函数将网站/文件名进行映射
- (2) 利用多个 (k个) 哈希函数



如果目标是当 $x \notin S$ 时，出现误报的概率小于等于 ε ，应该如何设置 m ?

回顾：若仅用一个哈希函数，为使误报率小于等于 ε ，所需存储空间（即 m 大小）约 $\frac{n}{\ln\left(\frac{1}{1-\varepsilon}\right)}$ 比特

当 ε 为 0.01 时，该值约 $99.5n$

若用 k 个哈希函数，为使误报率小于等于 ε ，所需存储空间（即 m 大小）约 $\frac{\log_2 \frac{1}{\varepsilon}}{\ln 2} n$ 比特

当 ε 为 0.01 时，该值约 $9.6n$

m 值的计算：<https://hur.st/bloomfilter/?n=4000&p=1.0E-7&m=&k=>

布隆过滤器

实验效果



缓存策略回顾

一种简单的缓存策略：当用户请求访问某文件，先判断该文件是否是第一次被请求：

- (1) 如果是，不缓存；
- (2) 如果不是，缓存。

* 传输策略：如果用户请求的文件已被缓存，则从本级服务器提取文件进行传输；否则，从上级服务器中查找/提取文件进行传输



内容分发网络

回顾

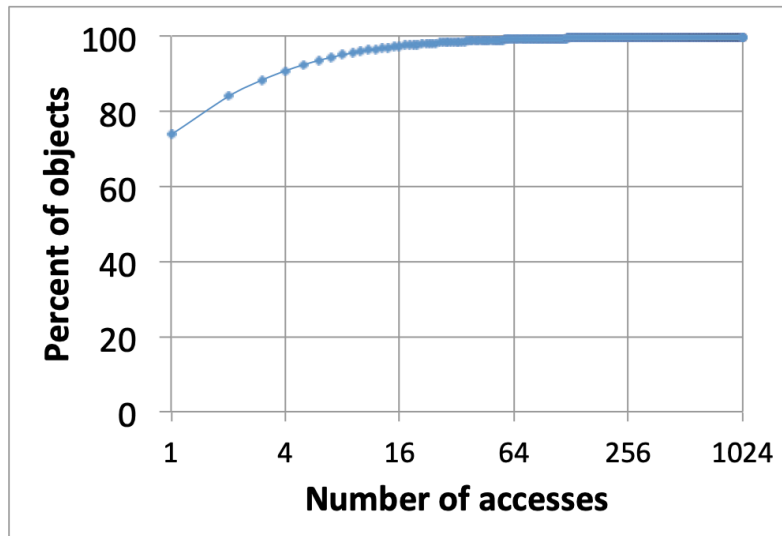


Figure 5: On a typical CDN server cluster serving web traffic over two days, 74% of the roughly 400 million objects in cache were accessed only once and 90% were accessed less than four times.

(2015年的论文) 分析了Akamai的一个有45台服务器的集群，统计在两天时间内4亿份文件被用户访问的频次

内容分发网络

在一个有47台服务器（47×8个硬盘）的集群上测试布隆过滤器的效果（关闭 — 打开 — 关闭）

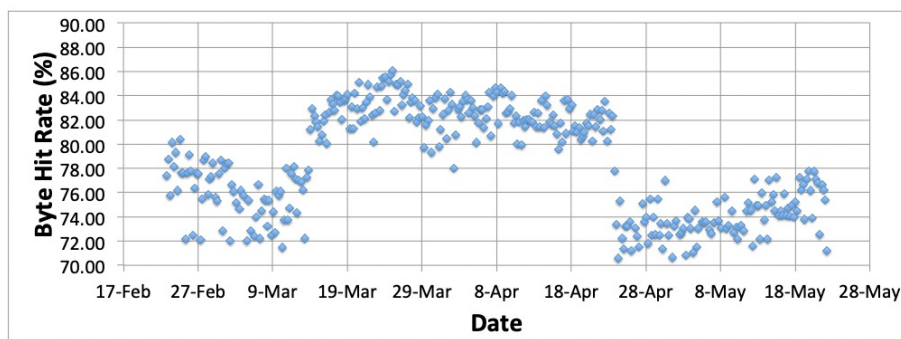


Figure 6: Byte hit rates increased when cache filtering was turned on between March 14th and April 24th because not caching objects that are accessed only once leaves more disk space to store more popular objects.

Byte hit rate: 用户请求文件时文件已被缓存的比例

比如75%说明对每100个比特的网络流量，其中有75个比特被缓存/25个比特通过从上级服务器获取

内容分发网络

在一个有47台服务器（ 47×8 个硬盘）的集群上测试布隆过滤器的效果（关闭 — 打开 — 关闭）

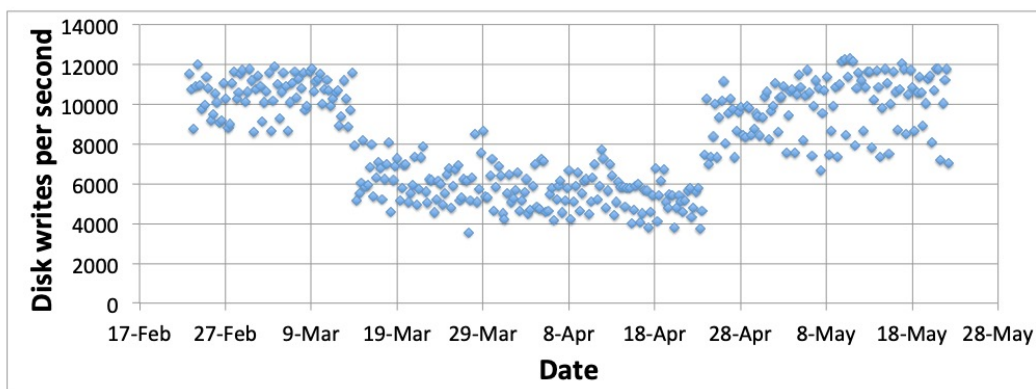


Figure 7: Turning on cache filtering decreases the rate of disk writes by nearly one half because objects accessed only once are not written to disk.

硬盘写入频率的变化

内容分发网络

在一个有47台服务器（47×8个硬盘）的集群上测试布隆过滤器的效果（关闭 — 打开 — 关闭）

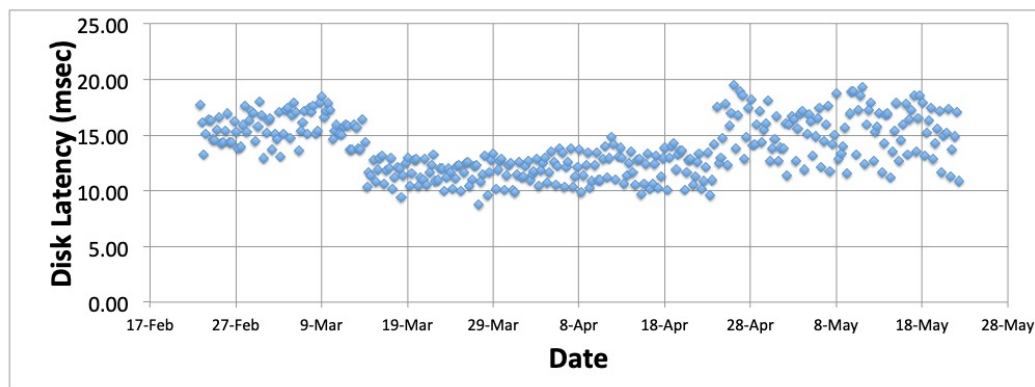


Figure 8: The latency of reading content from disk decreases when cache filtering is turned on, since there are fewer competing disk writes.

文件读取时延的变化

本讲小结



从属问题



布隆过滤器的方法

主要参考资料

Tim Roughgarden and Gregory Valiant <CS 168 - The Modern Algorithmic Toolbox> Lecture Notes

Cameron Musco <COMPSCI 514 - Algorithms for Data Science> Slides

B.M. Maggs and R.K. Sitaraman <Algorithmic Nuggets in Content Delivery> Paper

Keith Schwarz <Stanford CS166 - Data Structures> Slides

谢谢!

