







static const unsigned int

以垃圾邮件过滤中黑白名单为例:现有1亿个email的黑名单,每个都拥有8 bytes的指纹信息,则可能的

 $2^{64}=8*10^{18}~bits=10^9~GB$  ,对于bit array来说是根本不可能的范围,而且元素的数量(即emai

若采用哈希表,由于大多数采用open addressing来解决collision,而此时的search时间复杂度为:

即若哈希表半满(n/m = 1/2),则每次search需要probe 2次,因此在保证效率的情况下哈希表的存储

若采用Perfect hashing(这里可以采用Perfect hashing是因为主要操作是search/query,而并不是。

remove),虽然保证worst-case也只有一次probe,但是空间利用率更低,一般情况下为50%,wor

若采用布隆过滤器,取k=8。因为n为1亿,所以总共需要  $8*10^8~bits$  被置位为1,又因为在保证误判

MAX\_HASH\_FUNCS = 50;

36000/1024 = 3.5k

,相比于元素范围过于稀疏,而且还没有考虑到哈希表中的collision问题。

过50%。此时每个元素占8 bytes,总空间为:

取合适时,空间利用率为50%(后面会解释),所以总空间为:

所需空间比上述哈希结构小得多,并且误判率在万分之一以下。

 $m = \frac{10^8 * 8 \ bits}{50\%} = 1.6 * 10^9 bits = 200 MB$ 

 $\frac{10^8 * 8 \ bytes}{50\%} = 1.6 \ GB$ 

到一半的概率为25%。

三。举例说明

m bits k hashs

n elements

0 1 0 0 1 0 1 0 1 0 1 0