

## 前言

一般而言，标题含有“秒杀”，“99%”，“史上最全/最强”等词汇的往往都脱不了哗众取宠之嫌，但进一步来讲，如果读者读罢此文，却无任何收获，那么，我也甘愿背负这样的罪名，:-)，同时，此文可以看做是对这篇文章：[十道海量数据处理面试题与十个方法大总结](#)的一般抽象性总结。

毕竟受文章和理论之限，本文将摒弃绝大部分的细节，只谈方法/模式论，且注重用最通俗最直白的语言阐述相关问题。最后，有一点必须强调的是，全文行文是基于面试题的分析基础之上的，具体实践过程中，还是得具体情况具体分析，且各个场景下需要考虑的细节也远比本文所描述的任何一种解决方法复杂得多。

OK，若有任何问题，欢迎随时不吝赐教。谢谢。

## 何谓海量数据处理？

所谓海量数据处理，无非就是基于海量数据上的存储、处理、操作。何谓海量，就是数据量太大，所以导致要么是无法在较短时间内迅速解决，要么是数据太大，导致无法一次性装入内存。

那解决办法呢？针对时间，我们可以采用巧妙的算法搭配合适的数据结构，如Bloom filter/Hash/bit-map/堆/数据库或倒排索引/trie树，针对空间，无非就一个办法：大而化小，分而治之（hash映射），你不是说规模太大嘛，那简单啊，就把规模大化为规模小的，各个击破不就完了嘛。

至于所谓的单机及集群问题，通俗点来讲，单机就是处理装载数据的机器有限(只要考虑cpu，内存，硬盘的数据交互)，而集群，机器有多辆，适合分

布式处理，并行计算(更多考虑节点和节点间的数据交互)。

再者，通过本blog内的有关海量数据处理的文章：[Big Data Processing](#)，我们已经大致知道，处理海量数据问题，无非就是：

1. 分而治之/hash映射 + hash统计 + 堆/快速/归并排序；
2. 双层桶划分
3. Bloom filter/Bitmap；
4. Trie树/数据库/倒排索引；
5. 外排序；
6. 分布式处理之Hadoop/Mapreduce。

下面，本文第一部分、从set/map谈到hashtable/hash\_map/hash\_set，简要介绍下set/map/multiset/multimap，及hash\_set/hash\_map/hash\_multiset/hash\_multimap之区别(万丈高楼平地起，基础最重要)，而本文第二部分，则针对上述那6种方法模式结合对应的海量数据处理面试题分别具体阐述。

## 第一部分、从set/map谈到hashtable/hash\_map/hash\_set

稍后本文第二部分中将多次提到hash\_map/hash\_set，下面稍稍介绍下这些容器，以作为基础准备。一般来说，STL容器分两种，

- 序列式容器(vector/list/deque/stack/queue/heap)，
- 关联式容器。关联式容器又分为set(集合)和map(映射表)两大类，以及这两大类的衍生体multiset(多键集合)和multimap(多键映射表)，这些容

器均以RB-tree完成。此外，还有第3类关联式容器，如hashtable(散列表)，以及以hashtable为底层机制完成的hash\_set(散列集合)/hash\_map(散列映射表)/hash\_multiset(散列多键集合)/hash\_multimap(散列多键映射表)。也就是说，set/map/multiset/multimap都内含一个RB-tree，而hash\_set/hash\_map/hash\_multiset/hash\_multimap都内含一个hashtable。

所谓关联式容器，类似关联式数据库，每笔数据或每个元素都有一个键值(key)和一个实值(value)，即所谓的Key-Value(键-值对)。当元素被插入到关联式容器中时，容器内部结构(RB-tree/hashtable)便依照其键值大小，以某种特定规则将这个元素放置于适当位置。

包括在非关联式数据库中，比如，在MongoDB内，文档(document)是最基本的数据组织形式，每个文档也是以Key-Value（键-值对）的方式组织起来。一个文档可以有多个Key-Value组合，每个Value可以是不同的类型，比如String、Integer、List等等。

```
{ "name" : "July",  
  "sex" : "male",  
  "age" : 23 }
```

### **set/map/multiset/multimap**

set，同map一样，所有元素都会根据元素的键值自动被排序，因为set/map两者的所有各种操作，都只是转而调用RB-tree的操作行为，不过，值得注意的是，两者都不允许两个元素有相同的键值。

不同的是：set的元素不像map那样可以同时拥有实值(value)和键值

(key), set元素的键值就是实值, 实值就是键值, 而map的所有元素都是pair, 同时拥有实值(value)和键值(key), pair的第一个元素被视为键值, 第二个元素被视为实值。

至于multiset/multimap, 他们的特性及用法和set/map完全相同, 唯一的差别就在于它们允许键值重复, 即所有的插入操作基于RB-tree的insert\_equal()而非insert\_unique()。

### hash\_set/hash\_map/hash\_multiset/hash\_multimap

hash\_set/hash\_map, 两者的一切操作都是基于hashtable之上。不同的是, hash\_set同set一样, 同时拥有实值和键值, 且实质就是键值, 键值就是实值, 而hash\_map同map一样, 每一个元素同时拥有一个实值(value)和一个键值(key), 所以其使用方式, 和上面的map基本相同。但由于hash\_set/hash\_map都是基于hashtable之上, 所以不具备自动排序功能。为什么? 因为hashtable没有自动排序功能。

至于hash\_multiset/hash\_multimap的特性与上面的multiset/multimap完全相同, 唯一的差别就是它们hash\_multiset/hash\_multimap的底层实现机制是hashtable (而multiset/multimap, 上面说了, 底层实现机制是RB-tree), 所以它们的元素都不会被自动排序, 不过也都允许键值重复。

所以, 综上, 说白了, 什么样的结构决定其什么样的性质, 因为set/map/multiset/multimap都是基于RB-tree之上, 所以有自动排序功能, 而hash\_set/hash\_map/hash\_multiset/hash\_multimap都是基于hashtable之上, 所以不含有自动排序功能, 至于加个前缀multi\_无非就是允许键值重复而已。

此外,

- 关于什么hash, 请看blog内此篇[文章](#);

- 关于红黑树，请参看blog内系列[文章](#)，
- 关于hash\_map的具体应用：请看[这里](#)，关于hash\_set：请看[此文](#)。

OK，接下来，请看本文第二部分、处理海量数据问题之六把密匙。

## 第二部分、处理海量数据问题之六把密匙

### 密匙一、分而治之/Hash映射 + Hash\_map统计 + 堆/快速/归并排序

#### 1、海量日志数据，提取出某日访问百度次数最多的那个IP。

既然是海量数据处理，那么可想而知，给我们的数据那就一定是海量的。针对这个数据的海量，我们如何着手呢？对的，无非就是分而治之/hash映射 + hash统计 + 堆/快速/归并排序，说白了，就是先映射，而后统计，最后排序：

1. 分而治之/hash映射：针对数据太大，内存受限，只能是：把大文件化成(取模映射)小文件，即16字方针：大而化小，各个击破，缩小规模，逐个解决
2. hash\_map统计：当大文件转化了小文件，那么我们便可以采用常规的hash\_map(ip, value)来进行频率统计。
3. 堆/快速排序：统计完了之后，便进行排序(可采取堆排序)，得到次数最多的IP。

具体而论，则是：“首先是这一天，并且是访问百度的日志中的IP取出来，逐个写入到一个大文件中。注意到IP是32位的，最多有个 $2^{32}$ 个IP。同样可以采用映射的方法，比如%1000，把整个大文件映射为1000个小文件，再找出每个小文中出现频率最大的IP（可以采用hash\_map对那1000个文件中的所有IP进行频率统计，然后依次找出各个文件中频率最大的那个IP）及相应的频率。然后再在这1000个最大的IP中，找出那个频率最大的IP，即为所求。”--[十道海量数据处理面试题与十个方法大总结](#)。

关于本题，还有几个问题，如下：

1、Hash取模是一种等价映射，不会存在同一个元素分散到不同小文件中的情况，即这里采用的是mod1000算法，那么相同的IP在hash取模后，只可能落在同一个文件中，不可能被分散的。因为如果两个IP相等，那么经过Hash(IP)之后的哈希值是相同的，将此哈希值取模（如模1000），必定仍然相等。

2、那到底什么是hash映射呢？简单来说，就是为了便于计算机在有限的内存中处理big数据，从而通过一种映射散列的方式让数据均匀分布在对应的内存位置(如大数据通过取余的方式映射成小树存放在内存中，或大文件映射成多个小文件)，而这个映射散列方式便是我们通常所说的hash函数，设计的好的hash函数能让数据均匀分布而减少冲突。尽管数据映射到了另外一些不同的位置，但数据还是原来的数据，只是代替和表示这些原始数据的形式发生了变化而已。

OK，有兴趣的，还可以再了解下一致性hash算法，见blog内此文第五部分：[http://blog.csdn.net/v\\_july\\_v/article/details/6879101](http://blog.csdn.net/v_july_v/article/details/6879101)。

## 2、寻找热门查询，300万个查询字符串中统计最热门的10个查询

原题：搜索引擎会通过日志文件把用户每次检索使用的所有检索串都记录下来，每个查询串的长度为1-255字节。假设目前有一千万个记录（这些查询串的重复度比较高，虽然总数是1千万，但如果除去重复后，不超过3百万个。一个查询串的重复度越高，说明查询它的用户越多，也就是越热门），请你统计最热门的10个查询串，要求使用的内存不能超过1G。

解答：由上面第1题，我们知道，数据大则划为小的，如如一亿个Ip求Top

10, 可先%1000将ip分到1000个小文件中去, 并保证一种ip只出现在一个文件中, 再对每个小文件中的ip进行hashmap计数统计并按数量排序, 最后归并或者最小堆依次处理每个小文件的top10以得到最后的结。

但如果数据规模比较小, 能一次性装入内存呢? 比如这第2题, 虽然有一千万个Query, 但是由于重复度比较高, 因此事实上只有300万的Query, 每个Query255Byte, 因此我们可以考虑把他们都放进内存中去 (300万个字符串假设没有重复, 都是最大长度, 那么最多占用内存 $3M * 1K / 4 = 0.75G$ 。所以可以将所有字符串都存放在内存中进行处理), 而现在只是需要一个合适的数据结构, 在这里, HashTable绝对是我们优先的选择。

所以我们放弃分而治之/hash映射的步骤, 直接上hash统计, 然后排序。So, 针对此类典型的TOP K问题, 采取的对策往往是: hashmap + 堆。如下所示:

1. hash\_map统计: 先对这批海量数据预处理。具体方法是: 维护一个Key为Query字串, Value为该Query出现次数的HashTable, 即hash\_map(Query, Value), 每次读取一个Query, 如果该字串不在Table中, 那么加入该字串, 并且将Value值设为1; 如果该字串在Table中, 那么将该字串的计数加一即可。最终我们在 $O(N)$ 的时间复杂度内用Hash表完成了统计;
2. 堆排序: 第二步、借助堆这个数据结构, 找出Top K, 时间复杂度为 $N' \log K$ 。即借助堆结构, 我们可以在log量级的时间内查找和调整/移动。因此, 维护一个K(该题目中是10)大小的小根堆, 然后遍历300万的Query, 分别和根元素进行对比。所以, 我们最终的时间复杂度是:  $O(N) + N' * O(\log K)$ , (N为1000万, N'为300万)。

别忘了这篇文章中所述的堆排序思路：“维护k个元素的最小堆，即用容量为k的最小堆存储最先遍历到的k个数，并假设它们即是最大的k个数，建堆费时 $O(k)$ ，并调整堆(费时 $O(\log k)$ )后，有 $k_1 > k_2 > \dots > k_{\min}$  ( $k_{\min}$ 设为小顶堆中最小元素)。继续遍历数列，每次遍历一个元素x，与堆顶元素比较，若 $x > k_{\min}$ ，则更新堆 (x入堆，用时 $\log k$ )，否则不更新堆。这样下来，总费时 $O(k \cdot \log k + (n-k) \cdot \log k) = O(n \cdot \log k)$ 。此方法得益于在堆中，查找等各项操作时间复杂度均为 $\log k$ 。”--第三章续、Top K算法问题的实现。

当然，你也可以采用trie树，关键字域存该查询串出现的次数，没有出现为0。最后用10个元素的最小堆来对出现频率进行排序。

**3、有一个1G大小的一个文件，里面每一行是一个词，词的大小不超过16字节，内存限制大小是1M。返回频数最高的100个词。**

由上面那两个例题，分而治之 + hash统计 + 堆/快速排序这个套路，我们已经开始有了屡试不爽的感觉。下面，再拿几道再多多验证下。请看此第3题：又是文件很大，又是内存受限，咋办？还能怎么办呢？无非还是：

1. 分而治之/hash映射：顺序读文件中，对于每个词x，取 $\text{hash}(x) \% 5000$ ，然后按照该值存到5000个小文件（记为 $x_0, x_1, \dots, x_{4999}$ ）中。这样每个文件大概是200k左右。如果其中的有的文件超过了1M大小，还可以按照类似的方法继续往下分，直到分解得到的小文件的大小都不超过1M。
2. hash\_map统计：对每个小文件，采用trie树/hash\_map等统计每个文件中出现的词以及相应的频率。
3. 堆/归并排序：取出出现频率最大的100个词（可以用含100个结点的最小堆）后，再把100个词及相应的频率存入文件，这样又得到了5000个文件。最后就是把这5000个文件进行归并（类似于归并排序）的过



程了。

#### 4、海量数据分布在**100**台电脑中，想个办法高效统计出这批数据的**TOP10**。

如果每个数据元素只出现一次，而且只出现在某一台机器中，那么可以采取以下步骤统计出现次数**TOP10**的数据元素：

1. 堆排序：在每台电脑上求出**TOP10**，可以采用包含**10**个元素的堆完成（**TOP10**小，用最大堆，**TOP10**大，用最小堆，比如求**TOP10**大，我们首先取前**10**个元素调整成最小堆，如果发现，然后扫描后面的数据，并与堆顶元素比较，如果比堆顶元素大，那么用该元素替换堆顶，然后再调整为最小堆。最后堆中的元素就是**TOP10**大）。
2. 求出每台电脑上的**TOP10**后，然后把这**100**台电脑上的**TOP10**组合起来，共**1000**个数据，再利用上面类似的方法求出**TOP10**就可以了。

但如果同一个元素重复出现在不同的电脑中呢，如下例子所述：

这个时候，你可以有两种方法：

- 遍历一遍所有数据，重新hash取摸，如此使得同一个元素只出现在单独的一台电脑中，然后采用上面所说的方法，统计每台电脑中各个元素的出现次数找出**TOP10**，继而组合**100**台电脑上的**TOP10**，找出最终的**TOP10**。
- 或者，暴力求解：直接统计统计每台电脑中各个元素的出现次数，然后把同一个元素在不同机器中的出现次数相加，最终从所有数据中找出**TOP10**。

#### 5、有**10**个文件，每个文件**1G**，每个文件的每一行存放的都是用户的

**query**，每个文件的**query**都可能重复。要求你按照**query**的频度排序。

方案1：直接上：

1. hash映射：顺序读取10个文件，按照 $\text{hash}(\text{query})\%10$ 的结果将**query**写入到另外10个文件（记为a0,a1,...a9）中。这样新生成的文件每个的大小大约也1G（假设hash函数是随机的）。
2. hash\_map统计：找一台内存在2G左右的机器，依次对用`hash_map(query, query_count)`来统计每个**query**出现的次数。注：`hash_map(query,query_count)`是用来统计每个**query**的出现次数，不是存储他们的值，出现一次，则count+1。
3. 堆/快速/归并排序：利用快速/堆/归并排序按照出现次数进行排序，将排序好的**query**和对应的**query\_cout**输出到文件中，这样得到了10个排好序的文件（记为）。最后，对这10个文件进行归并排序（内排序与外排序相结合）。根据此方案1，这里有一份实现：<https://github.com/ooooola/sortquery/blob/master/querysort.py>。

除此之外，此题还有以下两个方法：

方案2：一般**query**的总量是有限的，只是重复的次数比较多而已，可能对于所有的**query**，一次性就可以加入到内存了。这样，我们就可以采用trie树/hash\_map等直接来统计每个**query**出现的次数，然后按出现次数做快速/堆/归并排序就可以了。

方案3：与方案1类似，但在做完hash，分成多个文件后，可以交给多个文件来处理，采用分布式的架构来处理（比如MapReduce），最后再进行合并。

**6、** 给定a、b两个文件，各存放50亿个url，每个url各占64字节，内存限制

是4G，让你找出a、b文件共同的url？

可以估计每个文件安的大小为 $5G \times 64 = 320G$ ，远远大于内存限制的4G。所以不可能将其完全加载到内存中处理。考虑采取分而治之的方法。

1. 分而治之/hash映射：遍历文件a，对每个url求取，然后根据所取得的值将url分别存储到1000个小文件（记为，这里漏写个了a1）中。这样每个小文件的大约为300M。遍历文件b，采取和a相同的方式将url分别存储到1000小文件中（记为）。这样处理后，所有可能相同的url都在对应的小文件（）中，不对应的小文件不可能有相同的url。然后我们只要求出1000对小文件中相同的url即可。
2. hash\_set统计：求每对小文件中相同的url时，可以把其中一个小文件的url存储到hash\_set中。然后遍历另一个小文件的每个url，看其是否在刚才构建的hash\_set中，如果是，那么就是共同的url，存到文件里面就可以了。

OK，此第一种方法：分而治之/hash映射 + hash统计 + 堆/快速/归并排序，再看最后4道题，如下：

## 7、怎么在海量数据中找出重复次数最多的一个？

方案：先做hash，然后求模映射为小文件，求出每个小文件中重复次数最多的一个，并记录重复次数。然后找出上一步求出的数据中重复次数最多的一个就是所求（具体参考前面的题）。

## 8、上千万或上亿数据（有重复），统计其中出现次数最多的前N个数据。

方案：上千万或上亿的数据，现在的机器的内存应该能存下。所以考虑采用hash\_map/搜索二叉树/红黑树等来进行统计次数。然后利用堆取出前N个

出现次数最多的数据。

**9、**一个文本文件，大约有一万行，每行一个词，要求统计出其中最频繁出现的前**10**个词，请给出思想，给出时间复杂度分析。

方案1：如果文件比较大，无法一次性读入内存，可以采用hash取模的方法，将大文件分解为多个小文件，对于单个小文件利用hash\_map统计出每个小文件中10个最常出现的词，然后再进行归并处理，找出最终的10个最常出现的词。

方案2：通过hash取模将大文件分解为多个小文件后，除了可以用hash\_map统计出每个小文件中10个最常出现的词，也可以用trie树统计每个词出现的次数，时间复杂度是 $O(n \cdot l_e)$ （ $l_e$ 表示单词的平均长度），最终同样找出出现最频繁的前10个词（可用堆来实现），时间复杂度是 $O(n \cdot \lg 10)$ 。

**10. 1000万字符串**，其中有些是重复的，需要把重复的全部去掉，保留没有重复的字符串。请怎么设计和实现？

- 方案1：这题用trie树比较合适，hash\_map也行。
- 方案2：from xjbzju:，1000w的数据规模插入操作完全不现实，以前试过在stl下100w元素插入set中已经慢得不能忍受，觉得基于hash的实现不会比红黑树好太多，使用vector+sort+unique都要可行许多，建议还是先hash成小文件分开处理再综合。

上述方案2中读者xbzju的方法让我想到了一些问题，即是set/map，与hash\_set/hash\_map的性能比较？共计3个问题，如下：

- 1、hash\_set在千万级数据下，insert操作优于set？这位blog: <http://t.cn/zOibP7t> 给的实践数据可靠不？

- 2、那map和hash\_map的性能比较呢？谁做过相关实验？
- 3、那查询操作呢，如下段文字所述？

或者小数据量时用map，构造快，大数据量时用hash\_map？

### **rbtree PK hashtable**

据朋友No邦卡猫No的做的红黑树和hash table的性能测试中发现：当数据量基本上int型key时，hash table是rbtree的3-4倍，但hash table一般会浪费大概一半内存。

因为hash table所做的运算就是个%，而rbtree要比较很多，比如rbtree要看value的数据，每个节点要多出3个指针（或者偏移量）如果需要其他功能，比如，统计某个范围内的key的数量，就需要加一个计数成员。

且1s rbtree能进行大概50w+次插入，hash table大概是差不多200w次。不过很多的时候，其速度可以忍了，例如倒排索引差不多也是这个速度，而且单线程，且倒排表的拉链长度不会太大。正因为基于树的实现其实不比hashtable慢到哪里去，所以数据库的索引一般都是用的B/B+树，而且B+树还对磁盘友好(B树能有效降低它的高度，所以减少磁盘交互次数)。比如现在非常流行的NoSQL数据库，像MongoDB也是采用的B树索引。关于B树系

列，请参考本blog内此篇文章：从B树、B+树、B\*树谈到R 树。更多请待后续实验论证。

**11. 一个文本文件，找出前10个经常出现的词，但这次文件比较长，说是上亿行或十亿行，总之无法一次读入内存，问最优解。**

方案1：首先根据用hash并求模，将文件分解为多个小文件，对于单个文件利用上题的方法求出每个文件中10个最常出现的词。然后再进行归并处理，找出最终的10个最常出现的词。

**12. 100w个数中找出最大的100个数。**

方案1：采用局部淘汰法。选取前100个元素，并排序，记为序列L。然后一次扫描剩余的元素x，与排好序的100个元素中最小的元素比，如果比这个最小的要大，那么把这个最小的元素删除，并把x利用插入排序的思想，插入到序列L中。依次循环，知道扫描了所有的元素。复杂度为 $O(100w*100)$ 。

方案2：采用快速排序的思想，每次分割之后只考虑比轴大的一部分，知道比轴大的一部分在比100多的时候，采用传统排序算法排序，取前100个。复杂度为 $O(100w*100)$ 。

方案3：在前面的题中，我们已经提到了，用一个含100个元素的最小堆完成。复杂度为 $O(100w*\lg 100)$ 。

接下来，咱们来看第二种方法，双层桶划分。

## 密匙二、多层划分

多层划分----其实本质上还是分而治之的思想，重在“分”的技巧上！

适用范围：第k大，中位数，不重复或重复的数字

基本原理及要点：因为元素范围很大，不能利用直接寻址表，所以通过多次划分，逐步确定范围，然后最后在一个可以接受的范围内进行。

问题实例：

**13、2.5亿个整数中找出不重复的整数的个数，内存空间不足以容纳这2.5亿个整数。**

有点像鸽巢原理，整数个数为 $2^{32}$ ，也就是，我们可以将这 $2^{32}$ 个数，划分为 $2^8$ 个区域(比如用单个文件代表一个区域)，然后将数据分离到不同的区域，然后不同的区域在利用bitmap就可以直接解决了。也就是说只要有足够的磁盘空间，就可以很方便的解决。

**14、5亿个int找它们的中位数。**

1. 思路一：这个例子比上面那个更明显。首先我们将int划分为 $2^{16}$ 个区域，然后读取数据统计落到各个区域里的数的个数，之后我们根据统计结果就可以判断中位数落到那个区域，同时知道这个区域中的第几大数刚好是中位数。然后第二次扫描我们只统计落在这个区域中的那些数就可以了。

实际上，如果不是int是int64，我们可以经过3次这样的划分即可降低到可以接受的程度。即可以先将int64分成 $2^{24}$ 个区域，然后确定区域的第几大数，在将该区域分成 $2^{20}$ 个子区域，然后确定是子区域的第几大数，然后子区域里的数的个数只有 $2^{20}$ ，就可以直接利用 direct addr table进行统计了。

2. 思路二@绿色夹克衫：同样需要做两遍统计，如果数据存在硬盘上，就需要读取2次。

方法同基数排序有些像，开一个大小为65536的Int数组，第一遍读取，统计Int32的高16位的情况，也就是0-65535，都算作0,65536 - 131071都算作1。就相当于用该数除以65536。Int32 除以 65536的结果不会超过65536种情况，因此开一个长度为65536的数组计数就可以。每读取一个数，数组中对应的计数+1，考虑有负数的情况，需要将结果加32768后，记录在相应的数组内。

第一遍统计之后，遍历数组，逐个累加统计，看中位数处于哪个区间，比如处于区间k，那么0- k-1的区间里数字的数量sum应该 $< n/2$ （2.5亿）。而k+1 - 65535的计数和也 $< n/2$ ，第二遍统计同上面的方法类似，但这次只统计处于区间k的情况，也就是说 $(x / 65536) + 32768 = k$ 。统计只统计低16位的情况。并且利用刚才统计的sum，比如sum = 2.49亿，那么现在就是要在低16位里面找100万个数(2.5亿-2.49亿)。这次计数之后，再统计一下，看中位数所处的区间，最后将高位和低位组合一下就是结果了。

### 密钥三：Bloom filter/Bitmap

#### Bloom filter

关于什么是Bloom filter，请参看blog内此文：

- [海量数据处理之Bloom Filter详解](#)

适用范围：可以用来实现数据字典，进行数据的判重，或者集合求交集  
基本原理及要点：



对于原理来说很简单，位数组+k个独立hash函数。将hash函数对应的值的位数组置1，查找时如果发现所有hash函数对应位都是1说明存在，很明显这个过程并不保证查找的结果是100%正确的。同时也不支持删除一个已经插入的关键字，因为该关键字对应的位会牵动到其他的关键字。所以一个简单的改进就是 counting Bloom filter，用一个counter数组代替位数组，就可以支持删除了。

还有一个比较重要的问题，如何根据输入元素个数 $n$ ，确定位数组 $m$ 的大小及hash函数个数。当hash函数个数 $k=(\ln 2)*(m/n)$ 时错误率最小。在错误率不大于 $E$ 的情况下， $m$ 至少要等于 $n*\lg(1/E)$ 才能表示任意 $n$ 个元素的集合。但 $m$ 还应该更大些，因为还要保证bit数组里至少一半为0，则 $m$ 应该 $\geq n\lg(1/E)*\lg e$  大概就是 $n\lg(1/E)1.44$ 倍( $\lg$ 表示以2为底的对数)。

举个例子我们假设错误率为0.01，则此时 $m$ 应大概是 $n$ 的13倍。这样 $k$ 大概是8个。

注意这里 $m$ 与 $n$ 的单位不同， $m$ 是bit为单位，而 $n$ 则是以元素个数为单位(准确的说是不同元素的个数)。通常单个元素的长度都是有很多bit的。所以使用bloom filter内存上通常都是节省的。

扩展：

Bloom filter将集合中的元素映射到位数组中，用 $k$  ( $k$ 为哈希函数个数)个映射位是否全1表示元素在不在这个集合中。Counting bloom filter (CBF) 将位数组中的每一位扩展为一个counter，从而支持了元素的删除操作。Spectral Bloom Filter (SBF) 将其与集合元素的出现次数关联。SBF采用counter中的最小值来近似表示元素的出现频率。

可以看下上文中的第6题：

**“6、给你A,B两个文件，各存放50亿条URL，每条URL占用64字节，内存限制是4G，让你找出A,B文件共同的URL。如果是三个乃至n个文件呢？”**

根据这个问题我们来计算下内存的占用， $4G=2^{32}$ 大概是40亿\*8大概是340亿， $n=50$ 亿，如果按出错率0.01算需要的大概是650亿个bit。现在可用的是340亿，相差并不多，这样可能会使出错率上升些。另外如果这些urlip是一一对应的，就可以转换成ip，则大大简单了。

同时，上文的第5题：给定a、b两个文件，各存放50亿个url，每个url各占64字节，内存限制是4G，让你找出a、b文件共同的url？如果允许有一定的错误率，可以使用Bloom filter，4G内存大概可以表示340亿bit。将其中一个文件中的url使用Bloom filter映射为这340亿bit，然后挨个读取另外一个文件的url，检查是否与Bloom filter，如果是，那么该url应该是共同的url（注意会有一定的错误率）。”

## Bitmap

- 关于什么是Bitmap，请看blog内此文第二部分：[http://blog.csdn.net/v\\_july\\_v/article/details/6685962](http://blog.csdn.net/v_july_v/article/details/6685962)。

下面关于Bitmap的应用，可以看下上文中的第13题，以及另外一道新题：

**“13、在2.5亿个整数中找出不重复的整数，注，内存不足以容纳这2.5亿个整数。**

方案1：采用2-Bitmap（每个数分配2bit，00表示不存在，01表示出现一次，10表示多次，11无意义）进行，共需内存 $2^{32} * 2 \text{ bit} = 1 \text{ GB}$ 内存，还可以接受。然后扫描这2.5亿个整数，查看Bitmap中相对应位，如果是00变

01, 01变10, 10保持不变。所描完事后, 查看bitmap, 把对应位是01的整数输出即可。

方案2: 也可采用与第1题类似的方法, 进行划分小文件的方法。然后在小文件中找出不重复的整数, 并排序。然后再进行归并, 注意去除重复的元素。”

**15、给40亿个不重复的unsigned int的整数, 没排过序的, 然后再给一个数, 如何快速判断这个数是否在那40亿个数当中?**

方案1: frome oo, 用位图/Bitmap的方法, 申请512M的内存, 一个bit位代表一个unsigned int值。读入40亿个数, 设置相应的bit位, 读入要查询的数, 查看相应bit位是否为1, 为1表示存在, 为0表示不存在。

#### 密匙四、Trie树/数据库/倒排索引

##### Trie树

适用范围: 数据量大, 重复多, 但是数据种类小可以放入内存

基本原理及要点: 实现方式, 节点孩子的表示方式

扩展: 压缩实现。

问题实例:

1. 上面的第2题: 寻找热门查询: 查询串的重复度比较高, 虽然总数是1千万, 但如果除去重复后, 不超过3百万个, 每个不超过255字节。
2. 上面的第5题: 有10个文件, 每个文件1G, 每个文件的每一行都存放的是用户的query, 每个文件的query都可能重复。要你按照query的频率排序。

3. 1000万字符串，其中有些是相同的(重复),需要把重复的全部去掉，保留没有重复的字符串。请问怎么设计和实现？
4. 上面的第8题：一个文本文件，大约有一万行，每行一个词，要求统计出其中最频繁出现的前10个词。其解决方法是：用trie树统计每个词出现的次数，时间复杂度是 $O(n \cdot l_e)$  ( $l_e$ 表示单词的平均长度)，然后是找出出现最频繁的前10个词。

更多有关Trie树的介绍，请参见此文：[从Trie树（字典树）谈到后缀树](#)。

## 数据库索引

适用范围：大数据量的增删改查

基本原理及要点：利用数据的设计实现方法，对海量数据的增删改查进行处理。

- 关于数据库索引及其优化，更多可参见此文：<http://www.cnblogs.com/pkuoliver/archive/2011/08/17/mass-data-topic-7-index-and-optimize.html>;
- 关于MySQL索引背后的数据结构及算法原理，这里还有一篇很好的文章：<http://blog.codinglabs.org/articles/theory-of-mysql-index.html>;
- 关于B 树、B+ 树、B\* 树及R 树，本blog内有篇绝佳文章：[http://blog.csdn.net/v\\_JULY\\_v/article/details/6530142](http://blog.csdn.net/v_JULY_v/article/details/6530142)。

## 倒排索引(Inverted index)

适用范围：搜索引擎，关键字查询

基本原理及要点：为何叫倒排索引？一种索引方法，被用来存储在全文搜索下某个单词在一个文档或者一组文档中的存储位置的映射。

以英文为例，下面是要被索引的文本：

T0 = "it is what it is"

T1 = "what is it"

T2 = "it is a banana"

我们就能得到下面的反向文件索引：

"a": {2}

"banana": {2}

"is": {0, 1, 2}

"it": {0, 1, 2}

"what": {0, 1}

检索的条件"what","is"和"it"将对应集合的交集。

正向索引开发出来用来存储每个文档的单词的列表。正向索引的查询往往满足每个文档有序频繁的全文查询和每个单词在校验文档中的验证这样的查询。在正向索引中，文档占据了中心的位置，每个文档指向了一个它所包含的索引项的序列。也就是说文档指向了它包含的那些单词，而反向索引则是单词指向了包含它的文档，很容易看到这个反向的关系。

扩展：

问题实例：文档检索系统，查询那些文件包含了某单词，比如常见的学术论文的关键字搜索。

关于倒排索引的应用，更多请参见：

- 第二十三、四章：杨氏矩阵查找，倒排索引关键词Hash不重复编码实践，
- 第二十六章：基于给定的文档生成倒排索引的编码与实践。

## 密钥五、外排序

适用范围：大数据的排序，去重

基本原理及要点：外排序的归并方法，置换选择败者树原理，最优归并树

问题实例：

1).有一个1G大小的一个文件，里面每一行是一个词，词的大小不超过16个字节，内存限制大小是1M。返回频数最高的100个词。

这个数据具有很明显的特点，词的大小为16个字节，但是内存只有1M做hash明显不够，所以可以用来排序。内存可以当输入缓冲区使用。

关于多路归并算法及外排序的具体应用场景，请参见blog内此文：

- 第十章、如何给 $10^7$ 个数据量的磁盘文件排序

## 密钥六、分布式处理之Mapreduce

MapReduce是一种计算模型，简单的说就是将大批量的工作（数据）分解（MAP）执行，然后再将结果合并成最终结果（REDUCE）。这样做的好处是可以在任务被分解后，可以通过大量机器进行并行计算，减少整个操作的时间。但如果你要我再通俗点介绍，那么，说白了，Mapreduce的原理就是一个归并排序。

适用范围：数据量大，但是数据种类小可以放入内存

基本原理及要点：将数据交给不同的机器去处理，数据划分，结果归约。

问题实例：

1. The canonical example application of MapReduce is a process to count the appearances of each different word in a set of documents:
2. 海量数据分布在100台电脑中，想个办法高效统计出这批数据的TOP10。
3. 一共有N个机器，每个机器上有N个数。每个机器最多存O(N)个数并对它们操作。如何找到 $N^2$ 个数的中数(median)?

更多具体阐述请参见blog内：

- 从Hadoop框架与MapReduce模式中谈海量数据处理，
- 及MapReduce技术的初步了解与学习。

其它模式/方法论，结合操作系统知识

至此，六种处理海量数据问题的模式/方法已经阐述完毕。据观察，这方面的面试题无外乎以上一种或其变形，然题目为何取为是：秒杀99%的海量数据处理面试题，而不是100%呢。OK，给读者看最后一道题，如下：

非常大的文件，装不进内存。每行一个int类型数据，现在要你随机取100个数。

我们发现上述这道题，无论是以上任何一种模式/方法都不好做，那有什么好的别的方法呢？我们可以看看：操作系统内存分页系统设计(说白了，就是映射+建索引)。

Windows 2000使用基于分页机制的虚拟内存。每个进程有4GB的虚拟地址空间。基于分页机制，这4GB地址空间的一些部分被映射了物理内存，一些部分映射硬盘上的交换文件，一些部分什么也没有映射。程序中使用的都是4GB地址空间中的虚拟地址。而访问物理内存，需要使用物理地址。关于什么是物理地址和虚拟地址，请看：

- 物理地址 (physical address): 放在寻址总线上的地址。放在寻址总线

上，如果是读，电路根据这个地址每位的值就将相应地址的物理内存中的数据放到数据总线中传输。如果是写，电路根据这个地址每位的值就将相应地址的物理内存中放入数据总线上的内容。物理内存是以字节(8位)为单位编址的。

- 虚拟地址 (virtual address): 4G虚拟地址空间中的地址，程序中使用的都是虚拟地址。使用了分页机制之后，4G的地址空间被分成了固定大小的页，每一页或者被映射到物理内存，或者被映射到硬盘上的交换文件中，或者没有映射任何东西。对于一般程序来说，4G的地址空间，只有一小部分映射了物理内存，大片大片的部分是没有映射任何东西。物理内存也被分页，来映射地址空间。对于32bit的 Win2k，页的大小是4K字节。CPU用来把虚拟地址转换成物理地址的信息存放在叫做页目录和页表的结构里。

物理内存分页，一个物理页的大小为4K字节，第0个物理页从物理地址 0x00000000 处开始。由于页的大小为4KB，就是0x1000字节，所以第1页从物理地址 0x00001000 处开始。第2页从物理地址 0x00002000 处开始。可以看到由于页的大小是4KB，所以只需要32bit的地址中高20bit来寻址物理页。

返回上面我们的题目：非常大的文件，装不进内存。每行一个int类型数据，现在要你随机取100个数。针对此题，我们可以借鉴上述操作系统中内存分页的设计方法，做出如下解决方案：

操作系统中的方法，先生成4G的地址表，在把这个表划分为小的4M的小文件做个索引，二级索引。30位前十位表示第几个4M文件，后20位表示在这个4M文件的第几个，等等，基于key value来设计存储，用key来建索引。

但如果现在只有10000个数，然后怎么去随机从这一万个数里面随机取



100个数？请读者思考。更多海量数据处理面试题，请参见此文第一部分：  
[http://blog.csdn.net/v\\_july\\_v/article/details/6685962](http://blog.csdn.net/v_july_v/article/details/6685962)。

## 参考文献

1. 十道海量数据处理面试题与十个方法大总结；
2. 海量数据处理面试题集锦与Bit-map详解；
3. 十一、从头到尾彻底解析Hash表算法；
4. [海量数据处理之Bloom Filter详解](#)；
5. 从Trie树（字典树）谈到后缀树；
6. 第三章续、Top K算法问题的实现；
7. 第十章、如何给 $10^7$ 个数据量的磁盘文件排序；
8. 从B树、B+树、B\*树谈到R 树；
9. 第二十三、四章：杨氏矩阵查找，倒排索引关键词Hash不重复编码实践；
10. [第二十六章：基于给定的文档生成倒排索引的编码与实践](#)；
11. 从Hadoop框架与MapReduce模式中谈海量数据处理；
12. 第十六~第二十章：全排列，跳台阶，奇偶排序，第一个只出现一次等问题；
13. [http://blog.csdn.net/v\\_JULY\\_v/article/category/774945](http://blog.csdn.net/v_JULY_v/article/category/774945)；
14. STL源码剖析第五章，侯捷著；
15. 2012百度实习生招聘笔试题：<http://blog.csdn.net/hackbuteer1/article/details/7542774>。

## 后记

经过上面这么多海量数据处理面试题的轰炸，我们依然可以看出这类问题是有一定的解决方案/模式的，所以，不必将其神化。然这类面试题所包含的问题还是比较简单的，若您在这方面有更多实践经验，欢迎随时来信与我不吝分享：[zhoulei0907@yahoo.cn](mailto:zhoulei0907@yahoo.cn)。当然，自会注明分享者及来源。

不过，相信你也早就意识到，若单纯论海量数据处理面试题，本blog内的有关海量数据处理面试题的文章已涵盖了你能在网上所找到的70~80%。但有点，必须负责任的敬告大家：无论是这些海量数据处理面试题也好，还是算法也好，**面试时**，70~80%的人不是倒在这两方面，而是倒在基础之上(诸如语言，数据库，操作系统，网络协议等等)，所以，**无论任何时候，基础最重要**，没了基础，便什么都不是。

最后，推荐国外一面试题网站：<http://www.careercup.com/>，以及个人正在读的Redis/MongoDB绝佳站点：<http://blog.nosqlfan.com/>。

OK，本文若有任何问题，欢迎随时不吝留言，评论，赐教，谢谢。完。