

44、热点网站统计，考点 or 实现——字符串，数组，集合操作

题目描述

企业路由器的统计页面，有一个功能需要动态统计公司访问最多的网页URL top N。请设计一个算法，可以高效动态统计Top N的页面。

输入描述

每一行都是一个URL或一个数字，如果是URL，代表一段时间内的网页访问；  
如果是一个数字N，代表本次需要输出的Top N个URL。

输入约束：

- 1、总访问网页数量小于5000个，单网页访问次数小于65535次；
- 2、网页URL仅由字母，数字和点分隔符组成，且长度小于等于127字节；
- 3、数字是正整数，小于等于10且小于当前总访问网页数；

输出描述

每行输入要对应一行输出，输出按访问次数排序的前N个URL，用逗号分隔。

输出要求：

- 1、每次输出要统计之前所有输入，不仅是本次输入；
- 2、如果有访问次数相等的URL，按URL的字符串字典序升序排列，输出排序靠前的URL；

用例

输入	news.qq.com news.sina.com.cn news.qq.com news.qq.com game.163.com game.163.com www.huawei.com www.cctv.com 3 www.huawei.com www.cctv.com www.huawei.com www.cctv.com www.huawei.com www.cctv.com www.huawei.com www.cctv.com www.huawei.com 3
输出	news.qq.com,game.163.com,news.sina.com.cn www.huawei.com,www.cctv.com,news.qq.com
说明	无

输入	news.qq.com www.cctv.com 1 www.huawei.com www.huawei.com 2 3
输出	news.qq.com www.huawei.com,news.qq.com www.huawei.com,news.qq.com,www.cctv.com
说明	无

## 题目解析

本题需要注意：“每次输出要统计之前所有输入，不仅是本次输入”。

这句话，可以看第一个用例的第二行输出，其中news.qq.com的输出就取决于上面这句话。

另外，还有这句话：“总访问网页数量小于5000个，单网页访问次数小于65535次；”

也就说，最多有不超过5000 \* 65535条访问URL记录，这个规模，我们需要尽可能地优化代码时间复杂度到O(1)左右，特别是“每次输出要统计之前所有输入，不仅是本次输入”，我们最好创建缓存表，而不是重新统计。

其余地就感觉没什么需要注意的了，都是逻辑实现了。

## JavaScript算法源码

```

1  /* JavaScript Node ACM模式 控制台输入获取 */
2  const readline = require("readline");
3
4  const rl = readline.createInterface({
5    input: process.stdin,
6    output: process.stdout,
7  });
8
9  const lines = [];
10 const cache = {};
11 rl.on("line", (line) => {
12   lines.push(line);
13
14   if (/^\d+$/.test(line)) {
15     console.log(sortURL(lines, cache));
16     lines.length = 0;
17   }
18 });
19
20 function sortURL(lines, cache) {
21   let n = parseInt(lines.pop());
22
23   lines.forEach((url) => {
24     cache[url] ? cache[url]++ : (cache[url] = 1);
25   });
26
27   let arr = [];
28   for (let key in cache) {

```

```

29     arr.push({
30         count: cache[key],
31         url: key,
32     });
33 }
34
35 return arr
36     .sort((a, b) => {
37         let res = b.count - a.count;
38         return res === 0 ? (a.url > b.url ? 1 : -1) : res;
39     })
40     .slice(0, n)
41     .map(ele => ele.url)
42     .join(",");
43 }

```

#### Java算法源码

```

1  import java.util.*;
2  import java.util.regex.Pattern;
3
4  public class Main {
5      public static void main(String[] args) {
6          Scanner sc = new Scanner(System.in);
7          // 改正正则用于判断一个字符串是否为正整数
8          Pattern reg = Pattern.compile("^\\d+$");
9          // urls 保存每次输入的url
10         ArrayList<String> urls = new ArrayList<>();
11         // cache 记录每个url出现次数
12         HashMap<String, Integer> cache = new HashMap<>();
13         // ans 保存输出
14         ArrayList<String> ans = new ArrayList<>();
15
16         while (sc.hasNextLine()) {
17             String tmp = sc.nextLine();
18
19             if ("".equals(tmp)) { // 由于本题没有说明输入结束条件，我这里将输入空行作为输入结束的条件，此时打印前面统计的结果
20                 for (String an : ans) System.out.println(an);
21                 sc.close();
22                 break;
23             }

```

```
23     }
24
25     // 如果输入的是正整数
26     if (reg.matcher(tmp).find()) {
27         ans.add(sortURL(urls, Integer.parseInt(tmp), cache)); // 将统计结果记录到ans
28         urls = new ArrayList<>(); // 清空urls, 继续下次输入记录
29     } else { // 输入的是url
30         urls.add(tmp);
31     }
32 }
33 }
34
35 public static String sortURL(ArrayList<String> urls, int n, HashMap<String, Integer> cache) {
36     for (String url : urls) {
37         cache.put(url, cache.getOrDefault(url, 0) + 1);
38     }
39
40     StringJoiner sj = new StringJoiner(",");
41
42     cache.entrySet().stream()
43         .sorted(
44             (a, b) ->
45                 Objects.equals(a.getValue(), b.getValue())
46                     ? a.getKey().compareTo(b.getKey())
47                     : b.getValue() - a.getValue()
48         )
49         .limit(n)
50         .map(ele -> ele.getKey())
51         .forEach(url -> sj.add(url));
```

```
51
52     return sj.toString();
53 }
54 }
```

## Python算法源码

```
1 import functools
2
3
4 def cmp(a, b): # a,b是dict_item类型数据, a[0]是url, a[1]是url出现次数
5     if a[1] != b[1]:
6         return b[1] - a[1]
7     elif a[0] != b[0]:
8         return 1 if a[0] > b[0] else -1
9     else:
10        return 0
11
12
13 def sortURL(urls, cache, n):
14     for url in urls:
15         if cache.get(url) is None:
16             cache[url] = 1
17         else:
18             cache[url] += 1
19
20     urlCount = list(cache.items())
21     urlCount.sort(key=functools.cmp_to_key(cmp))
22
23     return ",".join(map(lambda x: x[0], urlCount[:n]))
24
25
26 # 输入获取
27 urls = [] # 记录输入的url
28 cache = {} # 记录每个url出现的次数
29 ans = [] # 记录题解
30
31 while True:
32     tmp = input()
33
34     if tmp == "": # 由于本题没有说明输入结束的条件, 因此, 我将输入空行视为输入结束
35         for an in ans:
36             print(an) # 打印结果
37         break
38
39     if tmp.isnumeric(): # 如果输入的是数字, 则需要对前面输入的url进行排序统计
40         ans.append(sortURL(urls, cache, int(tmp)))
41         urls = []
42     else:
43         urls.append(tmp)
```