



目录



1

简单的倒排索引实现

2

网站KPI数据统计

3

典型运营商基站用户停留数据统计

4

PeopleRank算法

5

6

7



实现简单的倒排索引

倒排索引简单的可以理解为全文检索某个词

例如：在a.txt 和b.txt两篇文章分别中查找统计hello这个单词出现的次数，出现次数越多，和关键词的吻合度就越高

现有a.txt内容如下：

hello tom

hello jerry

hello kitty

hello world

hello tom

b.txt内容如下：

hello jerry

hello tom

hello world

在hadoop平台上编写mr代码分析统计各个单词在两个文本中出现的次数

其实也只是WordCount程序的改版而已~

MapReduce应用案例-网站kpi数据统计



- 1.browser: 用户使用的浏览器统计
- 2.ips: 页面用户独立ip数统计
- 3.pv: 网站pv量统计
- 4.source: 用户来源网址统计
- 5.time: 时间段用户访问量统计

	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7
1	ip地址	客户端用户名	请求时间	请求方法	请求页面	http协议信息	返回的状态码	发送的页面字节数	从什么页面跳转进来	用户使用的客户端							
2	120.197.87.216	-	[04/Jan/2012:00:00:02 +0800]	"GET	/home.php?mod=space&uid=563413&mobile=yes	HTTP/1.1"	200	3388	"-"	"-"							
3	123.126.50.73	-	[04/Jan/2012:00:00:02 +0800]	"GET	/thread-679411-1-1.html	HTTP/1.1"	200	5251	"-"	"Sogou web spider/4.0(+ http://www.sogou.com/docs/help/webmasters.htm							
4	116.205.130.2	-	[04/Jan/2012:00:00:02 +0800]	"GET	/popwin_js.php?fid=6	HTTP/1.1"	200	32	http://www.itpub.net/forum-6-1.html?ts=28	"Mozilla/4.0 (compatible; MSIE 8.							
5	218.186.15.10	-	[04/Jan/2012:00:00:16 +0800]	"GET	/forum.php?mod=ajax&action=forumchecknew&fid=61&time=1325606260&inajax=yes	HTTP/1.1"	200	92	http://www.itpub.net/f								

MapReduce应用案例-电信运营商用户基站停留数据统计



原始数据分为位置和网络两种

位置数据格式为:

用户标识 设备标识 开关机信息 基站位置 通讯的时间

example:

0000009999 0054785806 3 00000089 2016-02-21 21:55:37

网络数据格式为:

用户标识 设备标识 基站位置 通讯的时间 访问的URL

example:

0000000999 0054776806 00000109 2016-02-21 23:35:18 www.baidu.com

需要得到的数据格式为:

用户标识 时段 基站位置 停留时间

example:

00001 09-18 00003 15

用户00001在09-18点这个时间段在基站00003停留了15分钟

两个reducer:

- 1.统计每个用户在不同时段中各个基站的停留时间
- 2.在1的结果上只保留停留时间最长的基站位置信息

MapReduce应用案例-电信运营商用户基站停留数据统计



位置
数据

IMSI	IMEI	UPDATETYPE	LOC	TIME
...				
A	001	0	X基站	2013-09-12 09:00:00
A	001	2	Y基站	2013-09-12 09:45:00
...				

数据文件名
以POS开头

两种数据最大的差别在于文件名

上网
数据

IMSI	IMEI	LOC	TIME	URL
...				
A	001	X基站	2013-09-12 09:15:00	www.baidu.com
A	001	Y基站	2013-09-12 09:30:00	www.google.com
...				

数据文件名
以NET开头

MapReduce应用案例-电信运营商用户基站停留数据统计



认为用户在任何时间的停留位置都取决于之前一次位置更新的基站位置

时间间隔超过超过60分钟的判定为关机

IMSI	IMEI	UPDATETYPE	LOC	TIME
...				
A	001	0	X基站	2013-09-12 09:00:00
A	001	2	Y基站	2013-09-12 09:45:00
...				

IMSI	IMEI	LOC	TIME	URL
...				
A	001	X基站	2013-09-12 09:15:00	www.baidu.co m
A	001	Y基站	2013-09-12 09:30:00	www.google.co m
...				



用户A在X基站
停留了30分钟

MapReduce应用案例-电信运营商用户基站停留数据统计



输入数据

IMSI	IMEI	UPDATETYPE	LOC	TIME
...				
A	001	0	X基站	2013-09-12 09:00:00
A	001	2	Y基站	2013-09-12 09:45:00
...				

IMSI	IMEI	LOC	TIME	URL
...				
A	001	X基站	2013-09-12 09:15:00	www.baidu.com
A	001	Y基站	2013-09-12 09:30:00	www.google.com
...				



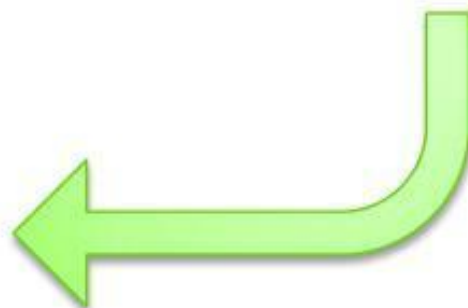
根据文件名
提取字段

IMSI	LOC	TIME
...		
A	X基站	2013-09-12 09:00:00
A	Y基站	2013-09-12 09:45:00
...		

IMSI	LOC	TIME
...		
A	X基站	2013-09-12 09:15:00
A	Y基站	2013-09-12 09:30:00
...		

<http://blog.csdn.net/>

IMSI	LOC	TimeFlag	TIME
...			
A	X基站	09-17	1386579600
A	Y基站	09-17	1386582300
A	X基站	09-17	1386580500
A	Y基站	09-17	1386581400
...			



计算时间所属时间段
把日期转换为UNIX格式

MapReduce应用案例-电信运营商用户基站停留数据统计



IMSI	TimeFlag
...	
A	09-17
A	09-17
A	09-17
A	09-17
...	

LOC	TIME
X基站	1386579600
Y基站	1386582300
X基站	1386580500
Y基站	1386581400

Map
输出

IMSI	LOC	TIME
...		
A	X基站	2013-09-12 09:15:00
A	Y基站	2013-09-12 09:45:00
...		

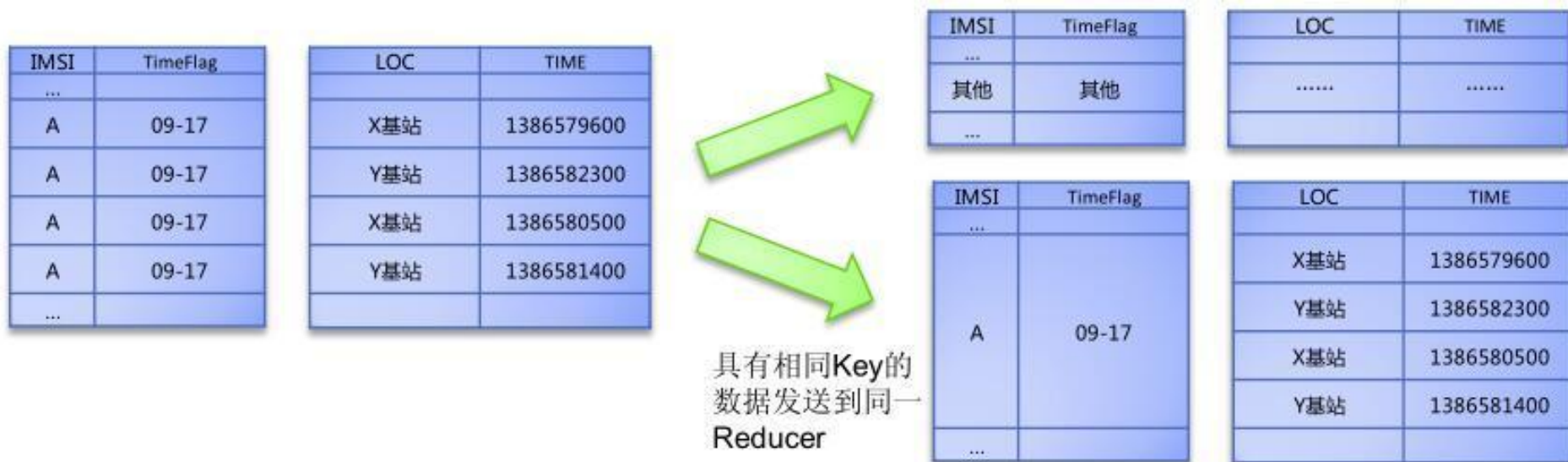
IMSI	LOC	TIME
...		
A	X基站	2013-09-12 11:15:00
A	Y基站	2013-09-12 09:30:00
...		

以IMSI和TimeFlag作为Key
以LOC和TIME作为VALUE

IMSI	LOC	TimeFlag	TIME
...			
A	X基站	09-17	1386579600
A	Y基站	09-17	1386582300
A	X基站	09-17	1386580500
A	Y基站	09-17	1386581400
...			

1. 计算时间所属时间段
2. 把日期转换为UNIX格式

MapReduce应用案例-电信运营商用户基站停留数据统计



<http://blog.csdn.net/>



MapReduce应用案例-电信运营商用户基站停留数据统计



IMSI	TimeFlag
...	
A	09-17
...	

LOC	TIME
X基站	1386579600
X基站	1386580500
Y基站	1386581400
Y基站	1386582300
OFF基站	1386608400

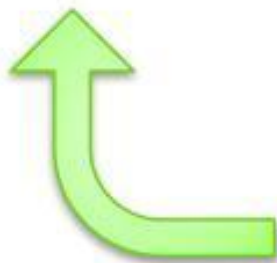


计算停留时间

IMSI	TimeFlag
...	
A	09-17
...	

LOC	STAY_TIME
X基站	15分钟
X基站	15分钟
Y基站	15分钟
Y基站	435分钟

<http://blog.csdn.net/>

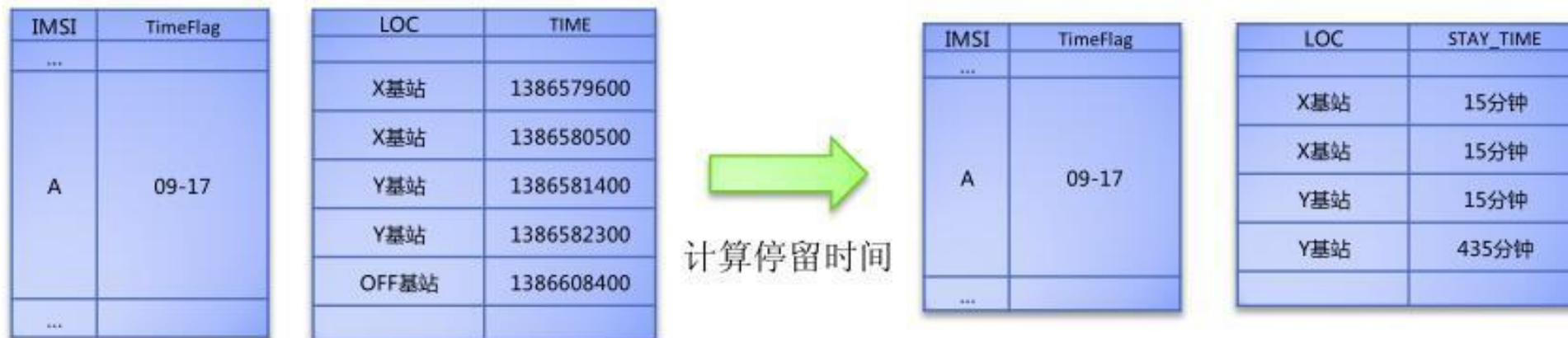


在最后添加特殊基站
时间是这个时段的结束时间

IMSI	TimeFlag
...	
A	09-17
...	

LOC	TIME
X基站	1386579600
X基站	1386580500
Y基站	1386581400
Y基站	1386582300

MapReduce应用案例-电信运营商用户基站停留数据统计



输出数据



- PageRank是Google专有的算法，用于衡量特定网页相对于搜索引擎索引中的其他网页而言的重要程度。它由Larry Page 和 Sergey Brin在20世纪90年代后期发明。
PageRank实现了将链接价值概念PageRank是Google的核心算法，用于给每个网页做评分，是google在“垃圾中找黄金”的关键算法，这个算法成就了今天的google。
- 作为排名因素。
- PageRank有两大特性：
 - ✓ PR值的传递性：网页A指向网页B时，A的PR值也部分传递给B
 - ✓ 重要性的传递性：一个重要网页比一个不重要网页传递的权重要多

➤ 计算公式:

$$PR(p_i) = \frac{1-d}{n} + d \sum_{p_j \in M(i)} \frac{PR(p_j)}{L(j)}$$

PR(pi): pi页面的PageRank值

n: 所有页面的数量

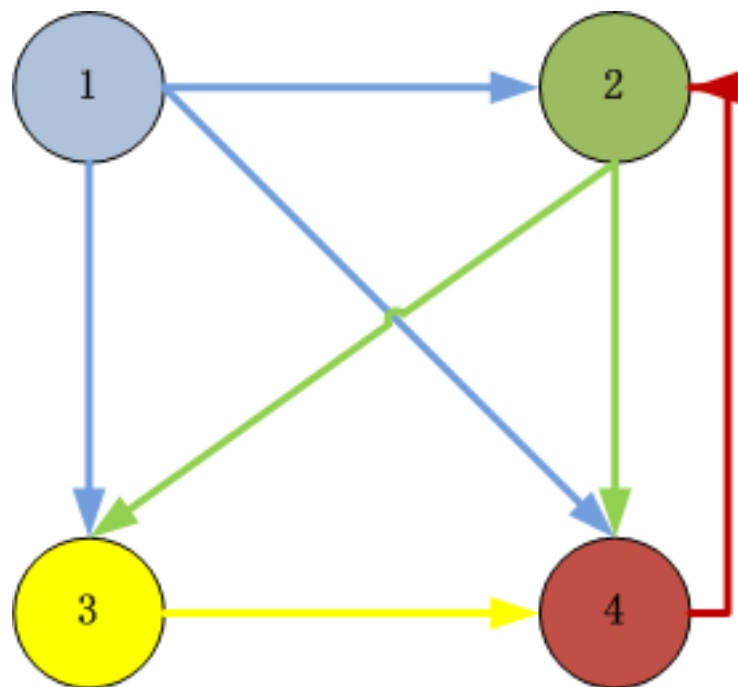
pi: 不同的网页p1,p2,p3

M(i): pi链入网页的集合

L(j): pj链出网页的数量

d: 阻尼系数, 任意时刻, 用户到达某页面后并继续向后浏览的概率。(1-d=0.15): 表示用户停止点击, 随机跳到新URL的概率取值范围: $0 < d \leq 1$, Google设为0.85

MapReduce应用案例-PageRank算法



链接源页面 链接目标页面

1	2,3,4
2	3,4
3	4
4	2

邻接矩阵：列=源链接，行=目标链接

$$\begin{bmatrix} 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 \end{bmatrix}$$

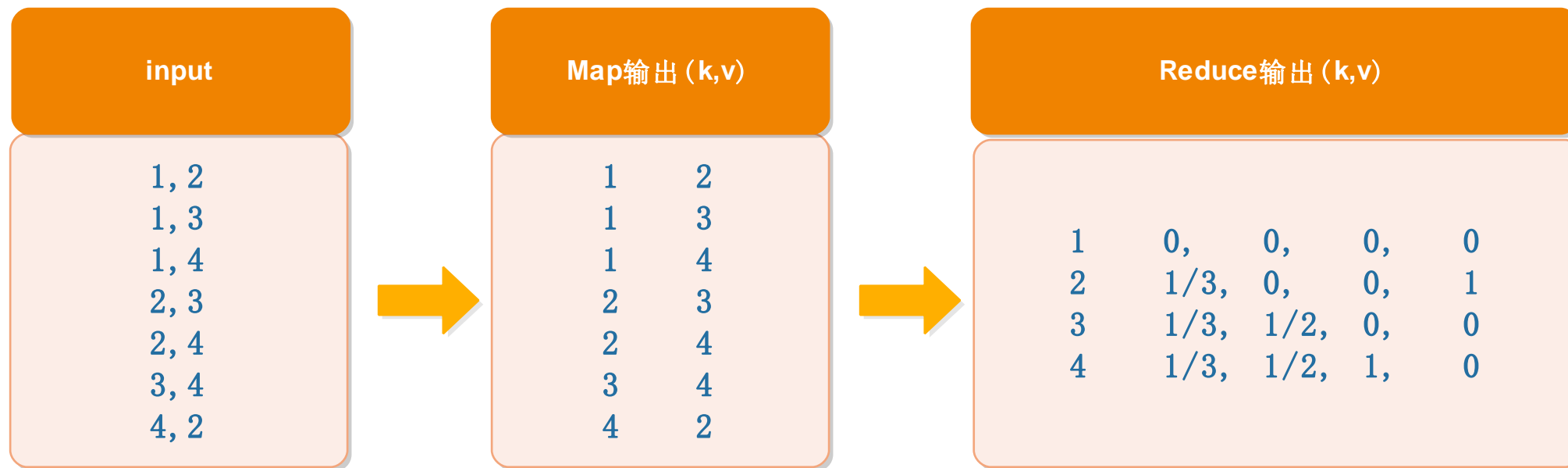
$$\begin{bmatrix} 0 & 0 & 0 & 0 \\ 1/3 & 0 & 0 & 1 \\ 1/3 & 1/2 & 0 & 0 \\ 1/3 & 1/2 & 1 & 0 \end{bmatrix}$$

概率矩阵：列=源链接，行=目标链接

MapReduce应用案例-PageRank算法



➤ 生成邻接概率矩阵。



MapReduce应用案例-PageRank算法



假设每个链接初始PageRank值为1。

$$\begin{bmatrix} 0 & 0 & 0 & 0 \\ 1/3 & 0 & 0 & 1 \\ 1/3 & 1/2 & 0 & 0 \\ 1/3 & 1/2 & 1 & 0 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 * 1 + 0 * 1 + 0 * 1 \\ \frac{1}{3} * 1 + 0 * 1 + 0 * 1 + 0 * 1 \\ \frac{1}{3} * 1 + \frac{1}{2} * 1 + 0 * 1 + 0 * 1 \\ \frac{1}{3} * 1 + \frac{1}{2} * 1 + 1 * 1 + 0 * 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1/3 \\ 5/6 \\ 11/6 \end{bmatrix}$$

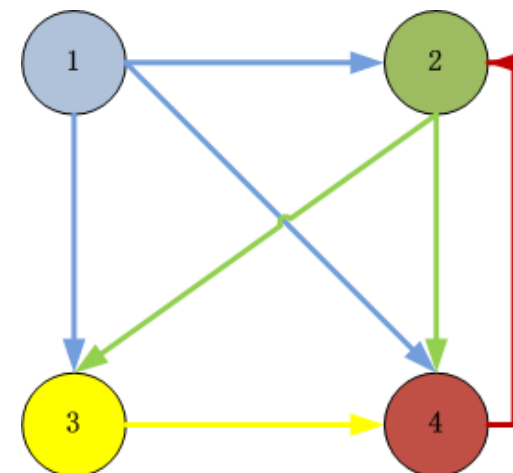
概率矩阵

初始PR矩阵

第一轮PR值

$$\sum_{p_j \in M(i)} \frac{PR(p_j)}{L(j)}$$

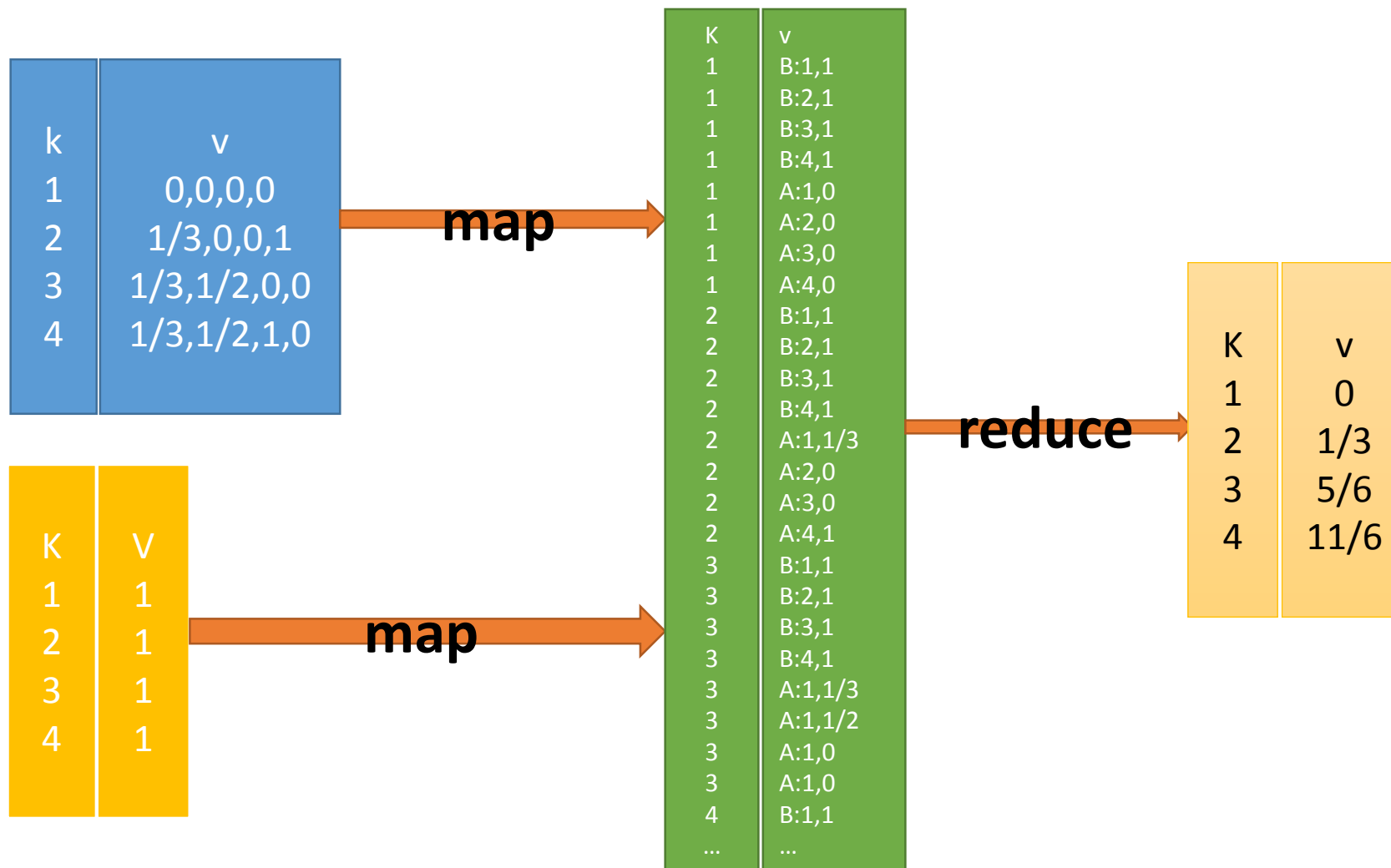
$$PR(p_i) = \frac{1-d}{n} + d \sum_{p_j \in M(i)} \frac{PR(p_j)}{L(j)}$$



MapReduce应用案例-PageRank算法



➤ 实现邻接与PR矩阵的乘法

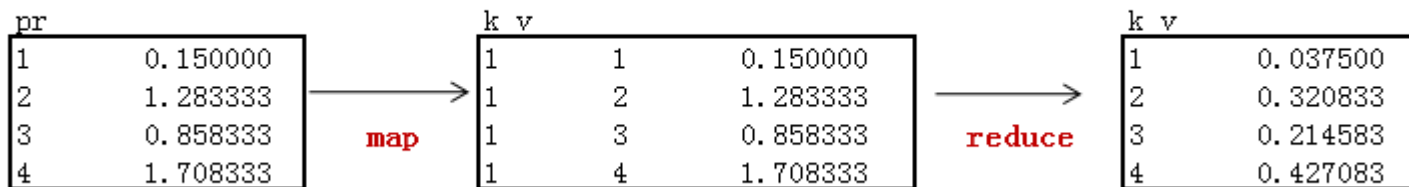


MapReduce应用案例-PageRank算法



➤ 对PR的计算结果标准化，让所以PR值落在(0,1)区间

Normal



- 1.说明Mapreduce中，如何实现PageRank算法矩阵相乘的原理，可以使用自己认为合理的任何表达方式，包括伪代码，图表，文字。
- 2.用MapReduce实现PageRank算法（提供算法程序和运行结果的截图）。
- 3.查找Mapreduce关于计数器相关的资料，mapreduce输出结果中这些计数器都代表什么含义？如何自定义计数器？
- 4.自己动手实现今天所讲的案例代码，并运行出结果。
- 5.描述一下Mapreduce的shuffle过程。
- 6.描述几条你所知道的关于Mapreduce的优化策略。

THANKS

