

MapReduce 数据挖掘基础算法(I)

K-Means 聚类算法



摘要

2

- 为什么选择数据挖掘作为并行计算的研究点
- **K-Means** 聚类算法介绍
- **K-Means** 算法为什么适合使用并行方法
- 基于 **MapReduce** 的 **K-Means** 并行算法
- 问题讨论



数据挖掘

3

- 定义：数据挖掘是通过对大规模观测数据集的分析，寻找确信的关系，并将数据以一种可理解的、且利于使用的新颖方式概括数据的方法。
- 数据挖掘的特征之一：海量数据
 - ▣ Small data does not require data mining, large data causes problems.
- 可见，数据挖掘是并行计算中值得研究的一个领域。

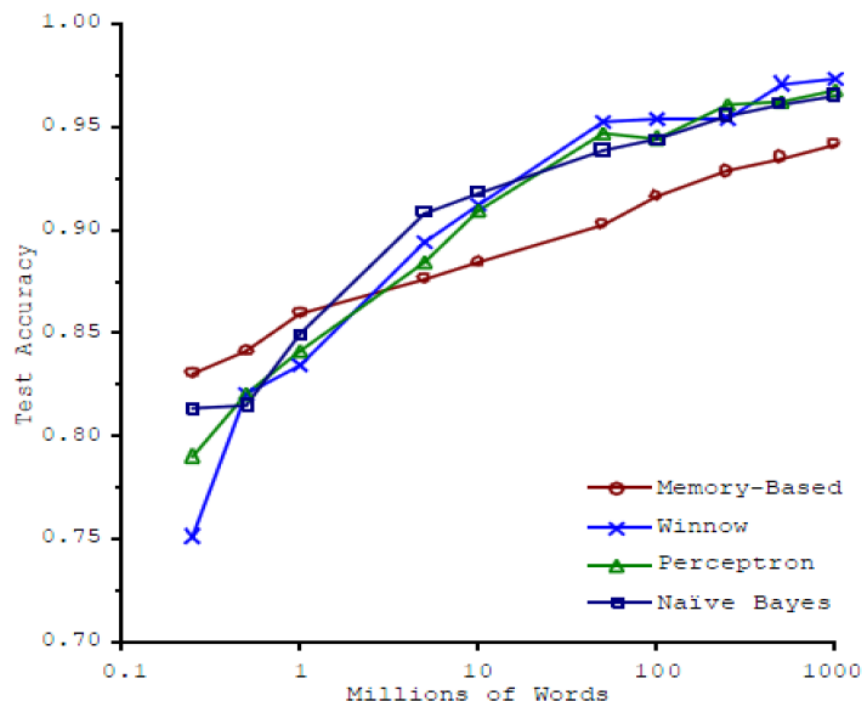


数据挖掘

4

□ 研究发现大数据隐含着更为准确的事实

2001年微软研究院的 Banko and Brill 等研究发现数据越大，机器学习的精度越高；当数据不断增长时，不同算法的分类精度趋向于相同！



M. Banko and E. Brill (2001). **Scaling to very very large corpora for natural language disambiguation.** ACL 2001.

Figure 1. Learning Curves for Confusion Set Disambiguation



数据挖掘

5

□ 研究发现大数据隐含着更为准确的事实

2007年Google公司的Brants等基于MapReduce研究了一个2万亿单词训练数据集的语言模型，发现大数据集上的简单算法能比小数据集上的复杂算法产生更好的结构！

T. Brants, A.C. Popat, et al. (2007).
Large language models in machine translation. *In EMNLP-CoNLL 2007.*

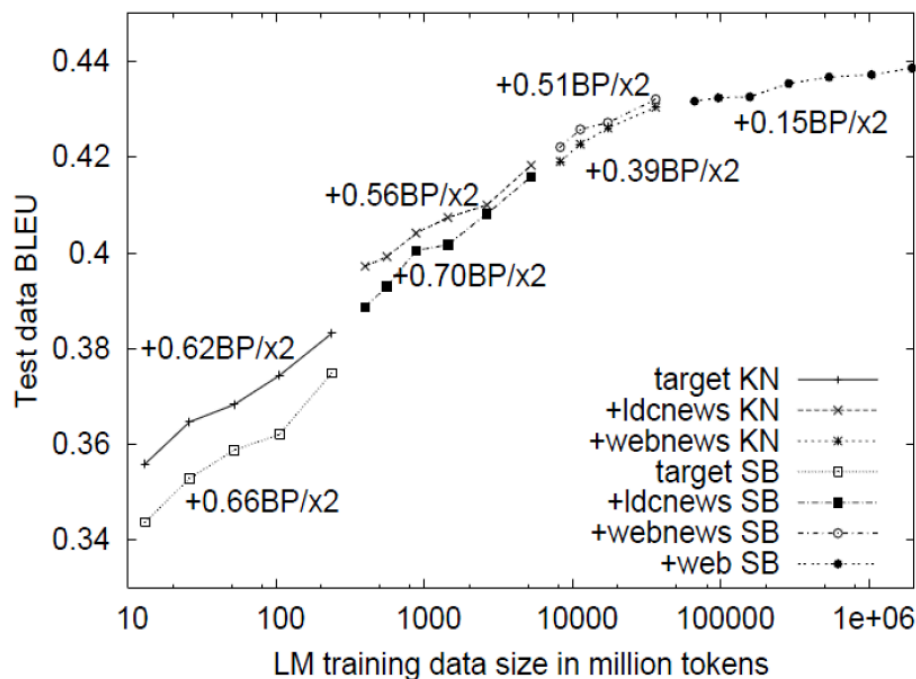


Figure 5: BLEU scores for varying amounts of data using Kneser-Ney (KN) and Stupid Backoff (SB).



聚类的应用领域

6

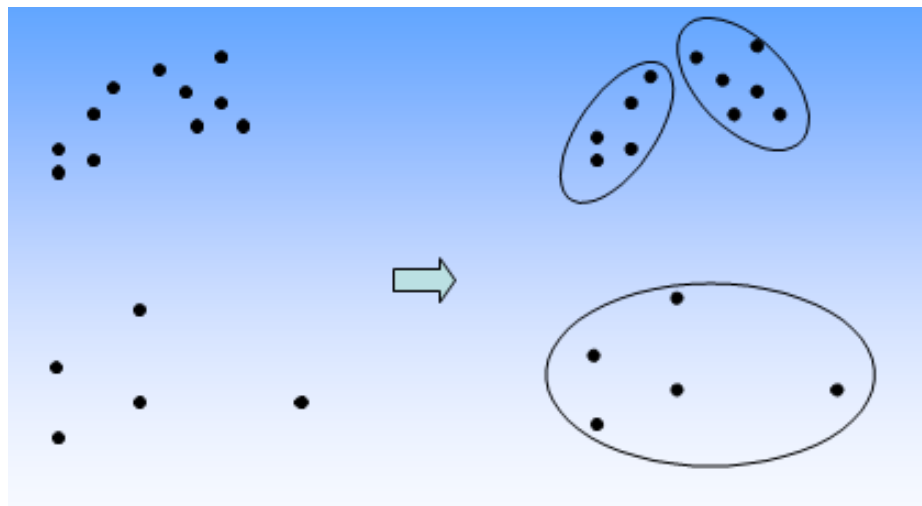
- 市场营销
 - ▣ 给定一个很大的顾客交易集，找出有类似购买行为的顾客分组
- 文档分类
 - ▣ 对**Web**日志数据聚类，发现有类似访问模式的分组
- 保险
 - ▣ 通过识别可能的欺诈行为找出平均索赔支出很高的车险投保人群。
- ...



聚类过程

7

- 定义：将给定的多个对象分成若干组，组内的各个对象是相似的，组间的对象是不相似的。进行划分的过程就是聚类过程，划分后的组称为簇（cluster）。
- 几种聚类方法：
 - 基于划分的方法；
 - 基于层次的方法；
 - 基于密度的方法；
 -





点、空间和距离

8

- 点集是一种适合于聚类的数据集，每个点都是某空间下的对象。
 - ▣ 欧式空间下的点就是实数向量。
 - ▣ 向量的长度是空间的维数。
 - ▣ 向量的分量通常称为所表示点的坐标。
- 能够进行聚类的所有空间下都有一个距离测度，即给出空间下任意两点的距离。
 - ▣ 距离永远非负，只有点到自身的距离为0。
 - ▣ 距离具有对称性。
 - ▣ 距离遵守三角不等式。



数值类型

9

- 数据点的类型可分为：
 - 欧氏 (**Euclidean**) 空间：空间中的点的平均总是存在，并且也是空间中的一个点。一般用欧几里得距离来衡量两个点之间的距离。
 - 非欧空间：**Jaccard**距离，**Cosine**距离，**Edit**编辑距离等多种距离衡量方法。
- 这二者在数据的表示以及处理上有较大的不同：
 - 怎样来表示**cluster**?
 - 怎样来计算相似度?



Cluster 的表示

10

- 欧氏空间：
 - 取各个数据点的平均值 (**centroid**)

- 非欧空间：
 - 取某个处于最中间的点
 - 取若干个最具代表性的点 (**clustroid**)
 -



相似度（距离）的计算

11

- 欧氏空间：可以有较为简单的方法

$$d([x_1, x_2, \dots, x_n], [y_1, y_2, \dots, y_n]) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

- 非欧氏空间：通常不能直接进行简单的数字计算

- Jaccard 距离：两个集合中不同元素占所有元素的比例

$$J_s = 1 - J(A, B) = \frac{|A \cup B| - |A \cap B|}{|A \cup B|}$$

- Cosine 距离：两个向量的夹角大小

$$\cos A = \frac{\langle b, c \rangle}{|b| |c|}, \quad \text{sim}(X, Y) = \cos \theta = \frac{\vec{x} \cdot \vec{y}}{\|\vec{x}\| \cdot \|\vec{y}\|}$$

- Edit 距离：适合于string类型的数据
- Hamming 距离：两个向量中不同分量的个数



基于划分 (partitioning) 的聚类方法

12

- 给定 N 个对象，构造 K 个分组，每个分组就代表一个聚类。
- 这 K 个分组满足以下条件：
 - 每个分组至少包含一个对象；
 - 每个对象属于且仅属于一个分组。
- **K-Means算法**是最常见和典型的基于划分的聚类方法

注：本节课只讨论**欧氏空间**里的K-Means聚类方法



K-Means 算法

13

输入：待聚类的N个数据点，期望生成的聚类的个数K

输出：K个聚类

-----算法描述-----

选出K个点作为初始的cluster center

Loop:

对输入中的每一个点p: {
 计算p到各个cluster center的距离;
 将p归入最近的cluster;
}

重新计算各个cluster的中心

如果不满足停止条件, goto Loop; 否则, 停止



过程示例 (1)

14

初始数据

$K = 2$



选择初始中心



第1次聚类：计算距离

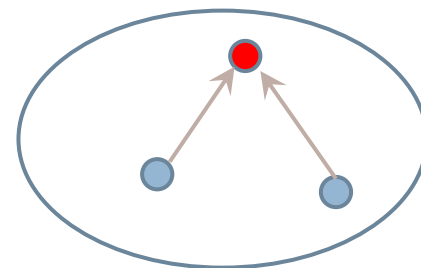
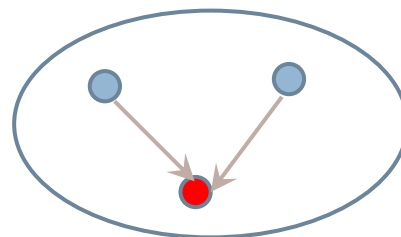




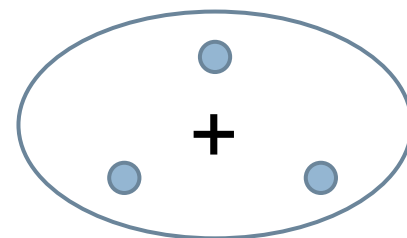
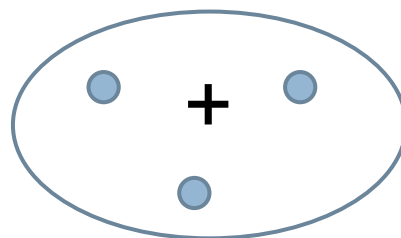
过程示例 (2)

15

第1次聚类：归类各点



重新计算聚类中心





过程示例 (3)

16

第2次聚类：计算距离



第2次聚类：归类各点



聚类无变化，迭代终止



K-Means是个不断迭代的过程

17

- 第 i 轮迭代：
 - 生成新的clusters，并计算cluster centers
- 第 $i+1$ 轮迭代：
 - 根据第 i 轮迭代中生成的clusters和计算出的cluster centers，进行新一轮的聚类
- ◇ 如此不断迭代直到满足终止条件



K-Means 算法的局限性

18

- 对初始 **cluster centers** 的选取会影响到最终的聚类结果
- 由此带来的结果是：能得到局部最优解，不保证得到全局最优解
- 相似度计算和比较时的计算量较大



K-Means计算性能的瓶颈

19

- 如果样本数据有 n 个，预期生成 k 个cluster，则K-Means算法 t 次迭代过程的时间复杂度为 $O(n*k*t)$ ，需要计算 $n*t*k$ 次相似度
- 如果能够将各个点到cluster center相似度的计算工作分摊到不同的机器上并行地计算，则能够减少计算时间
- 利用MapReduce将K-Means聚类过程并行化



考虑数据相关度

20

- 在进行K-Means聚类中，在处理每一个数据点时
- 只需要知道：
 - ▣ 各个cluster的中心信息
- 不需要知道：
 - ▣ 关于其他数据点的任何信息
- ◆ 所以，如果涉及到全局信息，只需要知道关于各个cluster center的信息即可



并行化改造的出发点

21

- 将所有数据分布到不同的**node**上，每个**node**只对自己的数据进行计算
- 每个**node**能够读取上一次迭代生成的**cluster centers**，并判断自己的各个数据点应该属于哪一个**cluster**
- 每个**node**在每次迭代中根据自己的数据点计算出相关数据
- 综合每个节点计算出的相关数据，计算出最终的**实际 cluster centers**



需要全局共享的数据

22

- 每一个节点需要访问如下的全局文件

- ☞ 当前的迭代计数

- ☞ K 个 如下结构

cluster id

cluster center

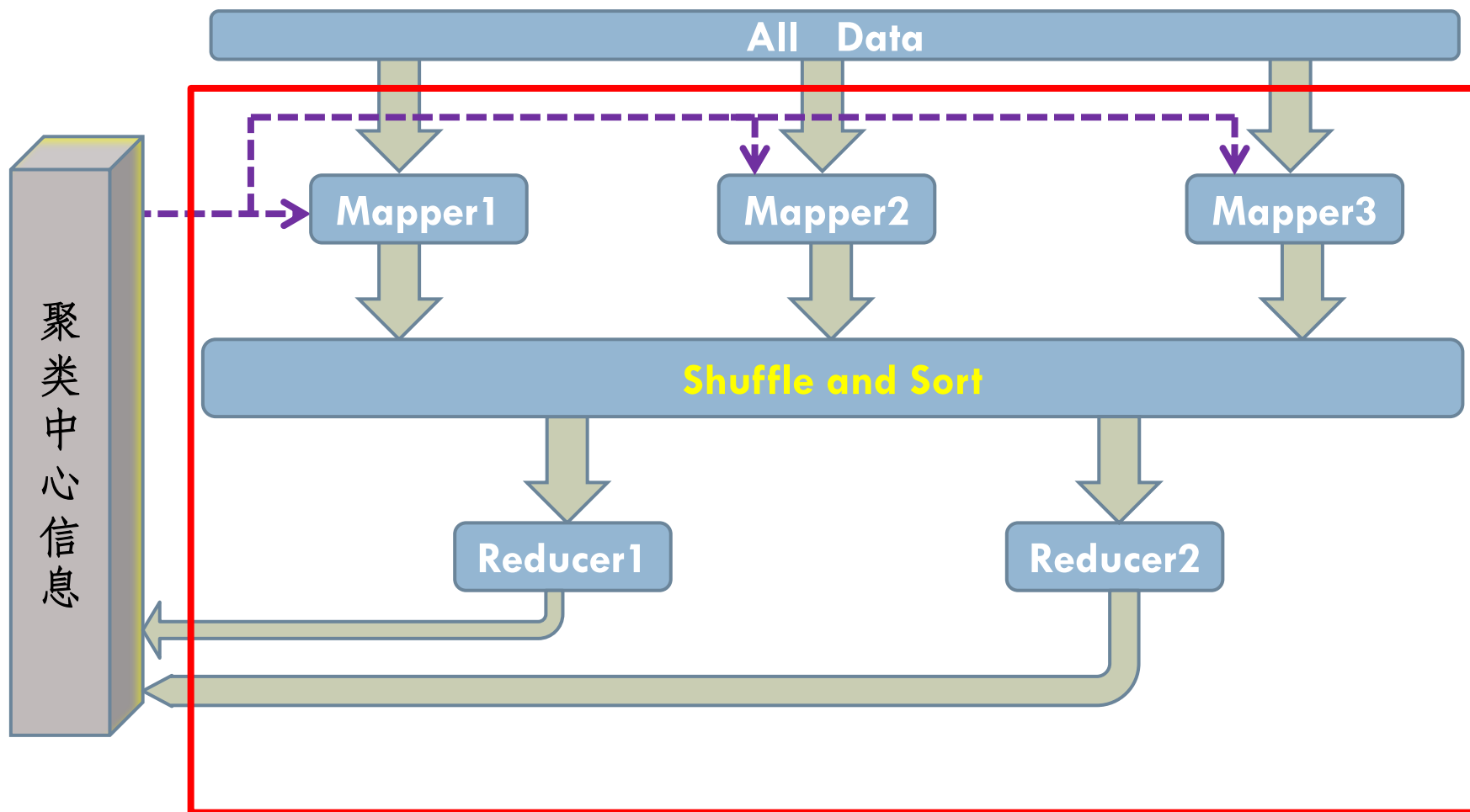
属于该**cluster center**的数据点的个数

- 这是**唯一**的全局文件



架构图

23





数据举例说明

24

Object	Attribute 1 (x)	Attribute 2(y)
A	1	1
B	2	1
C	4	3
D	5	4

- 令 $k = 2$ ，欲生成 cluster-0 和 cluster-1
- 随机选取 $A(1,1)$ 作为 cluster-0 的中心， $C(4,3)$ 作为 cluster-1 的中心
- 假定将所有数据分布到 2 个节点 node-0 和 node-1 上,即
node-0 : $A(1,1)$ 和 $C(4,3)$
node-1 : $B(2,1)$ 和 $D(5,4)$



准备开始K-Means聚类

25

- 在开始之前，首先创建一个如前所述的全局文件

Iteration No.	0	
Cluster id	cluster 中心	# of data points assigned
cluster -0	$A(1, 1)$	0
cluster-1	$C(4, 3)$	0



Map阶段

26

- 每个节点读取全局文件，以获得上一轮迭代生成的 **cluster centers** 等信息
- 计算本节点上的每一个点到各个 **cluster center** 的距离
- 为每一个数据点
Emit <cluster assigned to, 数据点自身>



Map: Iteration 1

27

- 计算各个数据点到各个cluster的距离，假定是如下的结果

	数据点	到各个 cluster 的距离比较		Assigned to
		cluster-0	cluster-1	
Node-0	A(1, 1)	近		cluster-0
	C(4, 3)		近	cluster-1
Node-1	B(2, 1)		近	cluster-1
	D(5, 4)		近	cluster-1

node-0 输出：

<cluster-0, A(1,1)>

<cluster-1, C(4,3)>

node-1 输出：

<cluster-1, B(2,1)>

<cluster-1, D(5,4)>



Combine 阶段

28

- 利用 **combiner** 减少 **map** 阶段 **emit** 的大量数据
- **Combiner** 计算要 **emit** 的所有数据点的均值，以及这些数据点的个数
- 然后，为每一个 **cluster** 发射 **key-value pair**
 - key** - cluster id,
 - value** - [# of data points of this cluster, mean]



Combine: Iteration 1

29

- Map的输出（即Combiner的输入）：

node-0 发射：

<cluster-0, A(1,1)>

<cluster-1, C(4,3)>

node-1 发射：

<cluster-1, B(2,1)>

<cluster-1, D(5,4)>

- Combiner的输出

- ▣ key - cluster id

- ▣ value - [# of data points of this cluster, mean]

node-0 发射：

<cluster-0, [1,(1,1)]>

<cluster-1, [1,(4,3)]>

node-1 发射：

<cluster-1, [2,(3.5,2.5)]>



Reduce 阶段

30

- 由于map阶段emit的key是cluster-id，所以每个cluster的全部数据将被发送同一个reducer，包括：
 - ▣ 该cluster 的id
 - ▣ 该cluster的数据点的均值，及对应于该均值的数据点的个数
- 然后经计算后输出
 - ☞ 当前的迭代计数
 - ☞ cluster id
 - ☞ cluster center
 - ☞ 属于该cluster center的数据点的个数



Reduce: Iteration 1

31

- 在上一阶段，Combiner的输出

node-0 发射：

<cluster-0, [1,(1,1)]>

<cluster-1, [1,(4,3)]>

node-1 发射：

<cluster-1, [2,(3.5,2.5)]>

- 两个reducer分别收到

Reducer-0：

<cluster-0, [1,(1,1)]>

Reducer-1：

<cluster-1, [1,(4,3)]>

<cluster-1, [2,(3.5,2.5)]>



Reduce: Iteration 1

32

□ 计算可得

□ cluster-0 的中心 = (1, 1)

□ cluster-1 的中心 = $\left(\frac{4+2 \times 3.5}{1+2}, \frac{3+2 \times 2.5}{1+2} \right)$
 $= (3.67, 2.67)$



Reduce: Iteration 1

33

□ Reducer 输出

Iteration No.	1	
Cluster id	cluster 中心	# of data points assigned
cluster-0	(1, 1)	1
cluster-1	(3.67, 2.67)	3



第1轮迭代结束

34

- 下面开始第2轮迭代，此时，全局文件已经更新为：

Iteration No.	1	
Cluster id	cluster 中心	# of data points assigned
cluster-0	(1, 1)	1
cluster-1	(3.67, 2.67)	3



终止迭代

35

- 在第 i 次迭代后，已经生成了 K 个聚类。如果满足了终止条件，即可停止迭代，输出 K 个聚类

- 终止条件：
 - 设定迭代次数;
 - 均方差的变化（非充分条件）
 - 每个点固定地属于某个聚类
 - 其他设定条件
- 与具体的应用高度相关



Test Case

36

□ 参数 $K=3$, $T=3$

□ Job List

ID	User	Name	Application Type	Queue	StartTime	FinishTime	State	FinalStatus	Progress	Tracking UI	Blacklisted Nodes
application_1510553148930_0004	yuping	KMeansClusterJob	MAPREDUCE	default	Mon Nov 13 14:10:44 +0800 2017	Mon Nov 13 14:10:57 +0800 2017	FINISHED	SUCCEEDED	<div></div>	History	0
application_1510553148930_0003	yuping	clusterCenterJob2	MAPREDUCE	default	Mon Nov 13 14:10:23 +0800 2017	Mon Nov 13 14:10:42 +0800 2017	FINISHED	SUCCEEDED	<div></div>	History	N/A
application_1510553148930_0002	yuping	clusterCenterJob1	MAPREDUCE	default	Mon Nov 13 14:10:03 +0800 2017	Mon Nov 13 14:10:21 +0800 2017	FINISHED	SUCCEEDED	<div></div>	History	N/A
application_1510553148930_0001	yuping	clusterCenterJob0	MAPREDUCE	default	Mon Nov 13 14:09:47 +0800 2017	Mon Nov 13 14:10:01 +0800 2017	FINISHED	SUCCEEDED	<div></div>	History	N/A

□ output

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
drwxr-xr-x	yuping	supergroup	0 B	2017/11/13 下午2:09:46	0	0 B	cluster-0
drwxr-xr-x	yuping	supergroup	0 B	2017/11/13 下午2:10:01	0	0 B	cluster-1
drwxr-xr-x	yuping	supergroup	0 B	2017/11/13 下午2:10:21	0	0 B	cluster-2
drwxr-xr-x	yuping	supergroup	0 B	2017/11/13 下午2:10:41	0	0 B	cluster-3
drwxr-xr-x	yuping	supergroup	0 B	2017/11/13 下午2:10:57	0	0 B	clusteredInstances



小结

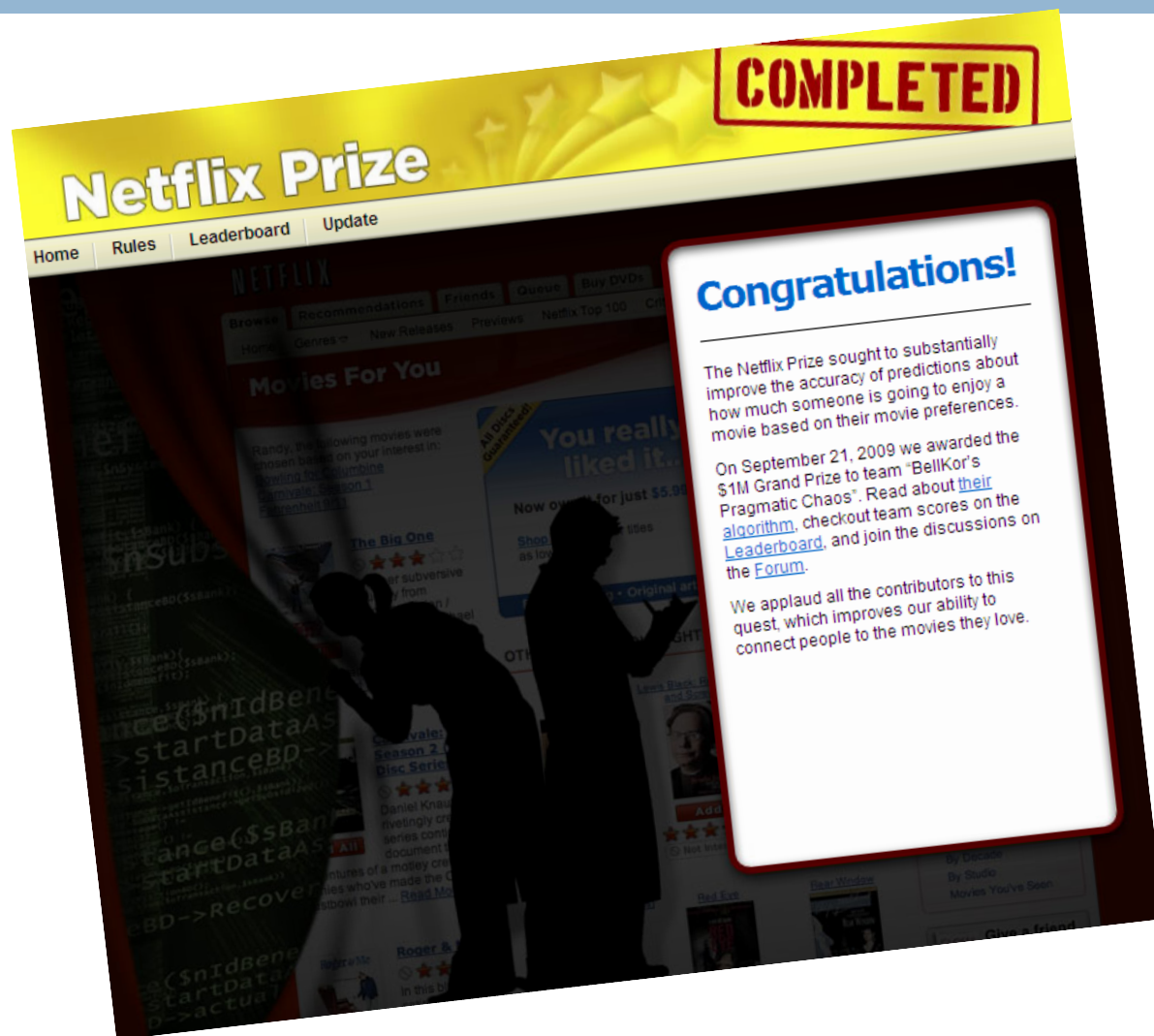
37

- 利用 **MapReduce** 来并行化 **K-Means** 聚类过程是可行的
- 每个节点计算一部分数据的归属，从而实现并行
- 数据间是无关的，但是数据和聚类中心是相关的，因此需要全局文件，但不构成性能瓶颈
- 没有因为并行而降低了算法的精确度（每一个点均保证与每一个 **cluster center** 进行了比较）



Open Problem

NetFlix 百万美元 大奖赛





Netflix公司

- 美国的一家电影在线租赁公司
- 拥有大量用户的影评记录
- 影片推荐系统基于这些影评记录
- 能够将推荐正确率提高10%者，将获得100万美元的奖励



竞赛数据

- 有 **17770** 个影片的影评文件，每个文件代表一部影片，描述了观众对于该影片的评价
- 评价的指数从 ★ 到 ★★★★★ ，每个影评文件的格式如 mv_0000001.txt
- 两个影评文件中的观众id可能有相同的，也可能不同



Netflix 中的聚类问题

- **目的：** 根据观众的评价对这**17770**部影片进行数据挖掘，输出约**400**个**聚类**，使得每个聚类中的影片是相似的

- **考虑：**
 - 1. 怎样定义及计算这类聚类问题中的相似度
 - 2. 怎样表示一个聚类（或聚类中心）
 - 3. 对于高维数据怎样预处理
 - 4. 怎样减少高维数据的计算量



参考文献

- *Web Mining – II: Parallelizing K-Means Clustering with MapReduce.* By Tushar Deshpande , Tejas Vora .
- *Mining of Massive Datasets.* By Anand Rajaraman, Jeffrey D. Ullman.
- *Data-Intensive Text Processing with MapReduce.* By Jimmy Lin and Chris Dyer
- *Data Algorithms – Recipes for Scaling Up with Hadoop and Spark.* By Mahmoud Parsian