

个人资料



fuck_prometheus

访问：59790次

积分：1133

等级：BLOG > 4

排名：千里之外

原创：58篇 转载：13篇

译文：0篇 评论：15条

文章搜索

文章存档

2017年07月 (2)

2017年06月 (3)

2017年05月 (2)

2017年03月 (2)

2017年01月 (5)

展开

阅读排行

QQ邮箱模拟登录 (2582)

sparkmlib朴素贝叶斯分 (2346)

sparkmlib逻辑回归源码 (2282)

sparkmlib协同过滤推荐 (2089)

sparkmlib关联规则算法 (2003)

sparkmlib决策树算法 (1965)

网易163邮箱模拟登录 (1895)

sparkmlib聚类算法：k-m (1739)

sparkmlib矩阵向量 (1590)

sparkmlib线性回归源码 (1574)

评论排行

QQ邮箱模拟登录 (5)

网易163邮箱模拟登录 (5)

新浪邮箱模拟登录java (3)

sparkmlib聚类算法：k-m (2)

赠书 | AI专栏 (AI圣经！《深度学习》中文版)

评论送书 | 机器学习、Java虚拟机、微信开发

sparkmlib协同过滤推荐算法

标签：sparkmlib 协同过滤算法

2016-11-20 15:39

2111人阅读

评论(0)

收藏

举报

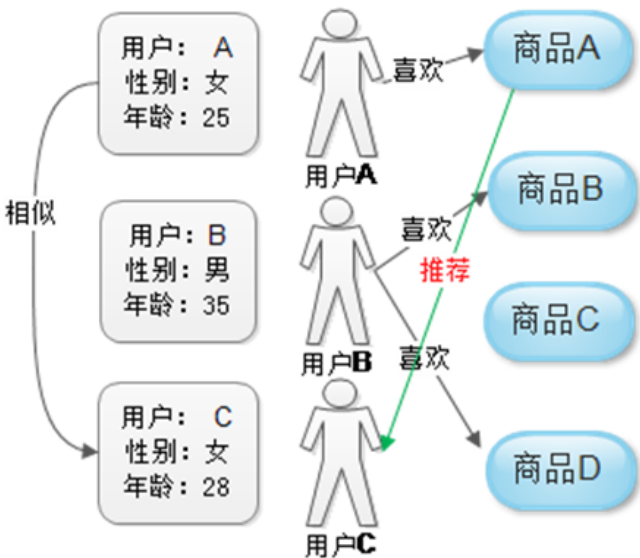
版权声明：本文为博主原创文章，未经博主允许不得转载。

目录(?)

[+]

协同过滤推荐算法，是最经典、最常用的推荐算法。通过分析用户兴趣，在用户群中找到相似用户，综合这些相似用户对某一信息的评价，形成系统关于该指定用户对此信息的喜好。要实现协同过滤，需要以下几个步骤：

- 1
- 2
- 3
- 1) 收集用户偏好；
- 2) 找到相似的用户或物品；
- 3) 计算推荐。



用户评分

从用户的行为和偏好中发现规律，并基于此进行推荐，所以收集用户的偏好信息成为系统推荐效果最基础的决定因素。用户有很多种方式向系统提供自己的偏好信息，比如：评分、投票、转发、保存书签、购买、点击流、页面停留时间等。

1. 将不同的行为分组

一般可以分为查看和购买，然后基于不同的用户行为，计算不同用户或者物品的相似度。

2. 对不同行为进行加权

对不同行为产生的用户喜好进行加权，然后求出用户对物品的总体喜好。当我们收集好用户的行为数据后，还要对数据进行预处理，最核心的工作就是减噪和归一化。

表 14-1 用户评分表

用户/物品	物品A	物品B	物品C
用户A	0.1	0.8	1
用户B	0.1	0	0.02
用户C	0.5	0.3	0.1

对用户的行为分析得到用户的偏好后，可以根据用户的偏好计算相似用户和物品，然后可以基于相似用

sparkmllib算法实例	(2)
Spark相关	(0)
scala学习	(0)
CDH安装	(0)
Elasticsearch线程池配置	(0)
Java执行js文件	(0)

推荐文章

* CSDN日报20170725——《新的开始，从研究生到入职亚马逊》

* 深入剖析基于并发AQS的重入锁(ReentrantLock)及其Condition实现原理

* Android版本的"Wannacry"文件加密病毒样本分析(附带锁机)

* 工作与生活真的可以平衡吗？

* 《Real-Time Rendering 3rd》 提炼总结——高级着色：BRDF 及相关技术

* 《三体》读后思考-泰勒展开/维度打击/黑暗森林

最新评论

sparkmllib算法实例
fuck_prometheus: @chuncun:网
上有很多类libsvm的数据,我这
个好像是官网自带的

sparkmllib算法实例
silencewinter: 楼主，这份libsvm
data file 哪里下载？

网易163邮箱模拟登录
fuck_prometheus:
@Henryone_L:谢啦，找机会试
下

网易163邮箱模拟登录
Henryone_L: 在进入页面的时候
会请求这个
<https://dl.reg.163.com/ini>，写一个cookie：

网易163邮箱模拟登录
冷若寒冰: 虽然我有其他的方法可以以更简单的方式搞定。但是您能免费公布出来，还是必须给您个赞。如果可以交个朋友加我...

QQ邮箱模拟登录
caishengkai: 最后一步显示
ptuiCB('4','0','0','您输入的验证
码不正确, 请重新输入。': ...

sparkmllib聚类算法：k-means算
fuck_prometheus:
@chenchaofuck1:首先它是随机
选取一点作为一个簇中心，然后
选取离这个中心点最远的距离作
为...

sparkmllib聚类算法：k-means算法
记忆力不好：
initKMeansParallel(data)方法是怎么做的？

QQ邮箱模拟登录
fuck_prometheus:
@guowujun321:这个如果你在js
里面找不到它的定义，你需要自
己去定义它 像pt={}这样

QQ邮箱模拟登录
炮孩子: @illbehere:是的。我引入
qqmail.js(加密js), 在
pt.setHeader ...

户或物品进行推荐。这就是协同过滤中的两个分支了，即基于用户的协同过滤和基于物品的协同过滤。

	物品1	物品2	物品3	物品4	物品5	物品6
用户A	1	1	0	0	0	1
用户B	1	0	1	1	0	0
用户C	1	1	1	0	0	0
用户D	0	0	1	1	1	0
用户E	0	0	0	1	1	1
用户F	1	0	1	0	1	0

1. 同现相似度

物品 i 和物品 j 的同现相似度公式定义:

$$w_{i,j} = \frac{|N(i) \cap N(j)|}{|N(i)|}$$

其中,分母 $|N(i)|$ 是喜欢物品 i 的用户数,而分子 $|N(i) \cap N(j)|$ 是同时喜欢物品 i 和物品 j 的用户数据。因此,上述公式可以理解成喜欢物品 i 的用户中有多少比例的用户也喜欢物品 j 。

上述公式存在一个问题,如果物品 j 是热门物品,很多人都喜欢,那么 $w_{i,j}$ 就会很大,因此,该公式会造成任何物品都会和热门物品有很大的相似度。为了避免推荐出热物品,可以用如下公式:

$$w_{i,j} = \frac{|N(i) \cap N(j)|}{\sqrt{|N(i)| |N(j)|}}$$

这个公式惩罚了物品 j 的权重，因此减轻了热门物品与很多物品相似的可能性。

2. 欧氏距离 (Euclidean Distance)

最初用于计算欧几里得空间中两个点的距离，假设 x 、 y 是 n 维空间的两个点，它们之间的欧几里得距离是：

$$d(x, y) = \sqrt{\sum (x_i - y_i)^2}$$

可以看出, 当 $n=2$ 时, 欧几里得距离就是平面上两个点的距离。

当用欧几里得距离表示相似度时，一般采用以下公式进行转换：距离越小，相似度越大。

$$\text{sim}(x, y) = \frac{1}{1 + d(x, y)}$$

3. 皮尔逊相关系数 (Pearson Correlation Coefficient)

4. Cosine 相似度 (Cosine Similarity)

5. Tanimoto 系数 (Tanimoto Coefficient)

推荐计算

1. 基于用户的CF (User CF)

基于用户的 CF 的基本思想相当简单：基于用户对物品的偏好找到相邻的邻居用户，然后将邻居用户喜欢的推荐给当前用户。在计算上，就是将一个用户对所有物品的偏好作为一个向量来计算用户之间的相似度，找到K 邻居后，根据邻居的相似度权重及其对物品的偏好，预测当前用户没有偏好的未涉及物品，计算得到一个排序的物品列表作为推荐。图14-1 给出了一个例子，对于用户A，根据用户的历史偏好，这里只计算得到一个邻居-用户C，然后将用户C 喜欢的物品D 推荐给用户

上海
¥260起

上海
¥120起

上海
¥80起

立秀宝Little Sociu
m韩国顶...

《魔逗先生》—亚
洲亲子魔术喜剧秀

首部“纽伯瑞金奖
“凯迪克银奖”

A。

用户/物品	物品A	物品B	物品C	物品D
用户A	✓		✓	推荐
用户B		✓		
用户C	✓		✓	✓

2. 基于物品的CF (Item CF)

基于物品的CF 的原理和基于用户的CF 类似，只是在计算邻居时采用物品本身，而不是从用户的角度。即基于用户对物品的偏好找到相似的物品，然后根据用户的历史偏好，推荐相似的物品给他。从计算的角度看，就是将所有用户对某个物品的偏好作为一个向量来计算物品之间的相似度，得到物品的相似物品后，根据用户历史的偏好预测当前用户还没有表示偏好的物品，计算的物品列表作为推荐。

用户/物品	物品A	物品B	物品C
用户A	✓		✓
用户B	✓	✓	✓
用户C	✓		推荐

图 14-2 基于物品的 CF 的基本原理

根据用户评分矩阵采用同现相似度计算物品相似度矩阵。

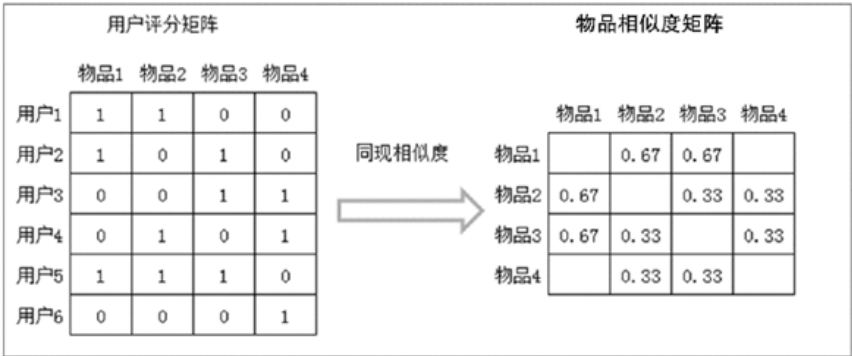


图 14-3 物品相似度矩阵

其相似度计算实现了分布式计算，实现过程如下：



图 14-4 分布式同现相似度矩阵计算过程

对于欧氏相似度的计算，采用离散计算公式 $d(x, y) = \sqrt{\sum((x(i)-y(i)) * (x(i)- y(i)))}$ 。其中，i 只取 x、y 同现的点，未同现的点不参与相似度计算； $sim(x, y) = m / (1 + d(x, y))$ ，m 为 x、y 重叠数，同现次数



图 14-5 分布式欧氏距离相似度矩阵计算过程

根据物品相似度矩阵和用户评分计算用户推荐列表，计算公式是 $R=W*A$ ，取推荐评分过的物品，并且按照计算结果倒序推荐给用户。

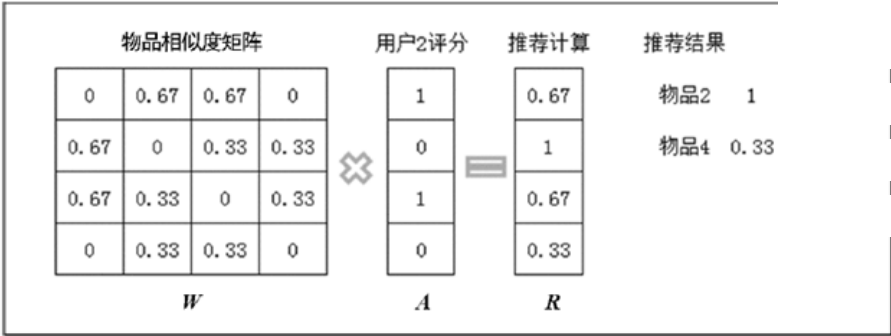


图 14-6 协同推荐计算

其推荐计算实现了分布式计算。

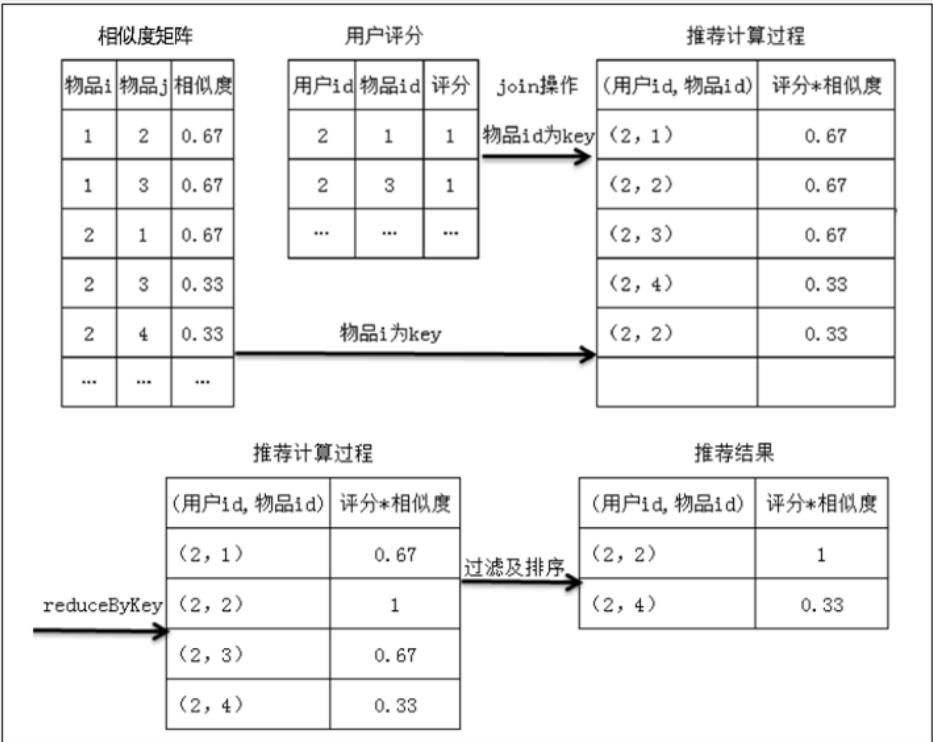


图 14-7 分布式协同推荐计算过程

源码分析

表 14-2 ItemCF协同过滤源码分解

源码分解说明		
1. 物品相似度类	ItemSimilarity	ItemSimilarity类
1.1 相似度计算	Similarity	ItemSimilarity类的Similarity方法。该方法计算物品与物品之间的相似度
2. 基于物品协同推荐类	RecommendedItem	RecommendedItem类
2.1 推荐计算	Recommend	RecommendedItem类的Recommend方法。该方法根据物品与物品的相似度模型和用户评分，计算用户的推荐物品列表

ALS算法:

```
1 def train(  
2     ratings: RDD[Rating],  
3     rank: Int,  
4     iterations: Int,  
5     lambda: Double,  
6     blocks: Int,  
7     seed: Long  
8 ): MatrixFactorizationModel = {  
9     new ALS(blocks, blocks, rank, iterations, lambda, false, 1.0, seed).run(  
10 )  
  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45
```

```
1 private def recommend(  
2     recommendToFeatures: Array[Double],  
3     recommendableFeatures: RDD[(Int, Array[Double])],  
4     num: Int): Array[(Int, Double)] = {  
5     val scored = recommendableFeatures.map { case (id, features) =>  
6         (id, blas.ddot(features.length, recommendToFeatures, 1, features, 1))  
7     }  
8     scored.top(num) (Ordering.by(_._2))  
9 }
```

实例：

```
1 import org.apache.spark.mllib.recommendation.ALS  
2 import org.apache.spark.mllib.recommendation.MatrixFactorizationModel  
3 import org.apache.spark.mllib.recommendation.Rating  
4  
5 // Load and parse the data  
6 val data = sc.textFile("data/mllib/als/test.data")  
7 val ratings = data.map(_ .split(',')) match { case Array(user, item, rate) =>  
8     Rating(user.toInt, item.toInt, rate.toDouble)  
9 }  
10  
11 // Build the recommendation model using ALS  
12 val rank = 10  
13 val numIterations = 10  
14 val model = ALS.train(ratings, rank, numIterations, 0.01)  
15  
16 // Evaluate the model on rating data  
17 val usersProducts = ratings.map { case Rating(user, product, rate) =>  
18     (user, product)  
19 }  
20 val predictions =  
21     model.predict(usersProducts).map { case Rating(user, product, rate) =>  
22         ((user, product), rate)  
23     }  
24 val ratesAndPreds = ratings.map { case Rating(user, product, rate) =>  
25     ((user, product), rate)  
26 }.join(predictions)  
27 val MSE = ratesAndPreds.map { case ((user, product), (r1, r2)) =>  
28     val err = (r1 - r2)  
29     err * err  
30 }.mean()  
31 println("Mean Squared Error = " + MSE)  
32  
33 // Save and load model  
34 model.save(sc, "target/tmp/myCollaborativeFilter")  
35 val sameModel = MatrixFactorizationModel.load(sc, "target/tmp/myCollaborativeFilter")
```

顶 踩
0 0

上一篇 [sparkmllib关联规则算法 \(FPGrowth,Apriori \)](#)

下一篇 [QQ邮箱模拟登录](#)

相关文章推荐

- [Spark MLlib系列\(二\):基于协同过滤的电影推荐系统](#)
- [spark基于用户的协同过滤算法与坑点，提交job](#)
- [基于Spark构建推荐引擎之一：基于物品的协同过...](#)
- [spark中协同过滤算法分析](#)

- 《Spark MLlib 机器学习》第二章代码
 - spark/MLlib 协同过滤算法
 - 基于Spark MLlib平台的协同过滤算法---电影推荐...
- 协同过滤itembase增量计算Spark实现(一)
 - Spark学习笔记 - 推荐系统 (协同过滤算法为用户...
 - 基于Spark MLlib平台的协同过滤算法---电影推荐...



整牙年龄



整牙的危害



猎头公司收费



oa系统



牙套多少钱



补

猜你在找

- 【直播】机器学习&深度学习系统实战（唐宇迪）
 - 【直播回放】深度学习基础与TensorFlow实践（王琛）
 - 【直播】机器学习之凸优化（马博士）
 - 【直播】机器学习之概率与统计推断（冒教授）
 - 【直播】TensorFlow实战进阶（智亮）
- 【直播】Kaggle 神器：XGBoost 从基础到实战（冒教授）
 - 【直播】计算机视觉原理及实战（屈教授）
 - 【直播】机器学习之矩阵（黄博士）
 - 【直播】机器学习之数学基础
 - 【直播】深度学习30天系统实训（唐宇迪）

查看评论

暂无评论

您还没有登录,请[\[登录\]](#)或[\[注册\]](#)

* 以上用户言论只代表其个人观点，不代表CSDN网站的观点或立场

公司简介 | 招贤纳士 | 广告服务 | 联系方式 | 版权声明 | 法律顾问 | 问题报告 | 合作伙伴 | 论坛反馈

网站客服 杂志客服 微博客服 webmaster@csdn.net 400-660-0108 | 北京创新乐知信息技术有限公司 版权所有 | 江苏知之为计算机有限公司 |

江苏乐知网络技术有限公司

京 ICP 证 09002463 号 | Copyright © 1999-2017, CSDN.NET, All Rights Reserved

