

文章编号: 1007-5321(2016)增-0037-05

DOI: 10.13190/j.jbupt.2016.s.009

基于 Spark 的 UCSLIM 推荐算法研究及实现

杨 娟, 张鹏业

(北京邮电大学 智能通信软件与多媒体北京市重点实验室, 北京 100876)

摘要: 稀疏线性(SLIM)推荐算法侧重于通过挖掘物品与物品之间的关系进而产生推荐结果. 为了提高推荐质量, 借鉴了 SLIM 算法和协同过滤算法的思想, 将用户划分为用户集合, 进一步挖掘用户与用户集合之间的隐含关系, 并综合考虑用户与用户相关性、用户与用户集合相关性这两个因素, 提出了融合用户集合关系的稀疏线性(UCSLIM)推荐算法. 实验结果表明, UCSLIM 算法能够提高推荐结果质量. 同时为了提高算法的执行效率, 分别在 Spark 和 Hadoop 云计算平台上实现了 UCSLIM 并行推荐算法, 并通过实验表明, UCSLIM 的 Spark 版本具有更高的计算效率.

关 键 词: TopN 推荐; 稀疏线性; 用户集合; Spark

中图分类号: TP391.3

文献标志码: A

Design and Implementation of Parallel UCSLIM Algorithm Based on Spark

YANG Juan, ZHANG Peng-ye

(Key Laboratory of Intelligent Telecommunication Software and Multimedia, Beijing University of Posts and Telecommunications, Beijing 100876, China)

Abstract: A top-N recommendation method—user class sparse linear methods (UCSLIM) based on sparse linear method (SLIM) was proposed. In order to improve the quality of recommendation, we learn from the idea of SLIM algorithm and collaborative filtering algorithm. The users are divided into different sets. So the correlation was analyzed between the user and the set of users and the correlation between user and user. Based on these two factors, UCSLIM was proposed. Experiments show that, compared with SLIM, UCSLIM can improve the quality of results. Furthermore, in order to improve the computational efficiency, the UCSLIM in Hadoop and Spark was implemented. Experiments show that the implementation by Spark has higher efficiency than that of Hadoop.

Key words: TopN recommender system; user class sparse linear methods; user class; Spark

目前, TopN 推荐系统已经得到广泛的应用与发展, 稀疏线性(SLIM, sparse linear method)^[1]推荐算法是一种基于模型的 TopN 推荐算法. 此算法先计算出物品之间的相关系数, 然后将用户物品评分与物品相关系数相乘, 从而预测出用户对其他物品的

评分.

在借鉴 SLIM 算法思想的基础上, 提出一种新的推荐算法——UCSLIM 算法. UCSLIM 算法综合考虑了用户-用户之间以及用户-用户集合之间的相关关系. 同时为了解决单机程序运行速度慢的问

收稿日期: 2015-09-17

基金项目: 国家高技术研究发展计划(国家 863 计划)项目(2015AA050204)

作者简介: 杨 娟(1972—), 女, 博士, 副教授, E-mail: yangjuan@bupt.edu.cn.

题,同时在 Spark^[2]和 Hadoop 平台上实现了 UCSLIM 推荐算法的并行化版本,最后给出实验结果以及结论。

1 UCSLIM 算法

1.1 符号定义

符号 u 和 t 分别代表用户和物品, μ_i 和 t_j 分别代表用户 i 和物品 j . 用户全集和物品全集分别用 U 和 T 表示,用户数量为 $n(|U| = n)$,物品数量为 $m(|T| = m)$. 所有矩阵均用大写字母表示,如 A 表示用户-物品评分矩阵, a_{ij} 表示用户 i 对物品 j 的评分. 列向量用加粗小写字母表示,如 a_j , 行向量用 a_i^T 表示.

1.2 SLIM 算法简介

SLIM 算法通过式(1)预测用户 i 对物品 j 的评分有

$$\bar{a}_{ij} = a_i^T w_j \quad (1)$$

其中: a_i^T 为用户-物品评分矩阵 A 的第 i 行,即用户 i 对所有物品的评分记录; w_j 为物品 j 和其余物品之间的相关系数. 从式(1)可以看出,此算法是通过分析物品与物品之间的相关系数进而预测用户对物品的评分.

1.3 UCSLIM 算法

1.3.1 算法思想

在借鉴 SLIM 算法和协同过滤算法^[3-4]思想的基础上,提出了一种综合考虑个体(用户与用户)维度间的相关性以及个体和集合(用户与用户集合)之间的隐含相关性的推荐算法——UCSLIM 算法.

UCSLIM 算法首先会将用户划分为用户集合,然后计算出用户之间的相关系数以及用户与用户集合的相关系数. 在预测用户 i 对物品 j 的评分时,会同时将用户 i 和其他用户的相关性以及用户 i 和其他用户集合的相关性都考虑在内.

1.3.2 划分用户集合

在划分用户集合时,会将最相似的 C 个用户划分到同一个集合中. C 的取值不宜过大, C 的值越大,产生的用户集合数量越少,反而会影响对用户集合有效信息的挖掘,故采取 $C=2$. 采用余弦相似度来衡量用户之间的相似度,在计算出所有用户之间的相似度后,对于用户 i ,选取其最相似的用户 q ,将用户 i 和用户 q 划分到同一个用户集 k 中. 以此类推,直到所有用户都找到了属于自己的用户集合,从而形成了一个用户集-物品矩阵 A' . 在 A' 中,如果

用户集合 k 中的用户 i 或用户 q 对物品 j 评分过,则 $a'_{kj} = 1$, 否则为 0. 经过以上计算即可得出完整的用户集-物品评分矩阵 A' . 因为采取的是每 2 个用户划分为一个用户集合,故矩阵 A' 的大小为 $(n/C) \times m$.

1.3.3 UCSLIM 计算模型

UCSLIM 算法通过式(2)预测用户 i 对物品 j 的评分有

$$\bar{a}_{ij} = r_i^T a_j + r_i^T a'_j \quad (2)$$

其中: r_i^T 为一个维度为 n 的行向量,表示用户 i 和其他用户的相关系数; a_j 为用户评分矩阵 A 的第 j 列,即所有用户对于物品 j 的评分; r_i^T 为一个维度为 n/C 的行向量,表示用户 i 和其他用户集的相关系数; a'_j 为用户集-物品评分矩阵 A' 的第 j 列,即所有用户集合对于物品 j 的评分. 将式(2)扩展到矩阵相乘的形式为

$$\bar{A} = RA + R'A' \quad (3)$$

其中: R 为大小为 $n \times n$ 的用户-用户相关系数矩阵, R' 为大小为 $n \times (n/C)$ 的用户-用户集相关系数矩阵, A' 是大小为 $(n/C) \times m$ 的用户集-物品评分矩阵,即根据原始的用户-物品矩阵计算出的新的用户集-物品矩阵. 通过如下公式可以计算出 R 和 R' :

$$\min \frac{1}{2} \|A - (RA + R'A')\|_F^2 + \frac{\beta_1}{2} \|R\|_F^2 +$$

$$\theta_1 \|R\|_1 + \frac{\beta_2}{2} \|R'\|_F^2 + \theta_2 \|R'\|_1$$

$$\text{约束} \quad R > 0 \text{ and } \text{Diag}(R) = 0$$

$$R' > 0 \text{ and } \text{Diag}(R') = 0 \quad (4)$$

其中: $\frac{1}{2} \|A - (RA + R'A')\|_F^2$ 评估了预测值 $(RA + R'A')$ 和实际值 (A) 之间的偏离程度; $\|R\|_1 =$

$\sum_{i=1}^m \sum_{j=1}^m |r_{ij}|$ 是 R 的 L1 范式,即 L1 正则化约束,

$\|R\|_F$ 为矩阵 R 的 F 范式,即 L2 正则化约束. L1 和 L2 正则化约束的作用是防止计算得到的模型过拟合^[5]. 常量 β_1 和 θ_1 分别是 L1 正则化和 L2 正则化约束参数,一般来说,当常量 β_1 和 θ_1 越大,其正则化程度越高;约束条件 $R > 0$ 和 $R' > 0$ 保证了用户集合和用户集合以及用户和用户之间的正相关性, $\text{Diag}(R) = 0$ 和 $\text{Diag}(R') = 0$ 保证了其对角线元素都为 0,即保证用户或用户集合对自身的推荐程度为 0.

求解式(4)本质上是一个最优化问题的求解,采用随机梯度下降法^[6]求解。在求解梯度过程中,向量的L1正则化约束在0附近不可导,采用Tsu-ruoka等^[7]提出的方法处理0附近不可导问题。

1.3.4 并行求解 UCSLIM 算法

由线性代数知识可知,当使用式(4)求解矩阵 R 和矩阵 R 时,矩阵 R 和矩阵 R 内的每一行之间不存在直接相关性,故可以将式(4)转化为一组规模更小的相同的求解问题,其优化之后为式(5)。

$$\min \frac{1}{2} \|a_i^T - (r_i^T A + r_i^T A)\|_F^2 + \frac{\beta_1}{2} \|r_i^T\|_F^2 + \theta_1 \|r_i^T\|_1 + \frac{\beta_2}{2} \|r_i^T\|_F^2 + \theta_2 \|r_i^T\|_1$$

$$\text{约束} \quad r_i^T > 0 \text{ and } r_{ii}^T = 0$$

$$r_i^T > 0 \text{ and } r_{ii}^T = 0 \quad (5)$$

其中: a_i^T 为矩阵 A 的第 i 行, r_i^T 为矩阵 R 的第 i 行, r_i^T 为矩阵 R 的第 i 行。通过此种变换,只需要并行求解矩阵 R 和矩阵 R 的每一行,然后将所有行的结果合并成一个完整矩阵即可。

1.3.5 算法复杂度

对于UCSLIM算法,在划分用户集合时,要计算所有用户之间的相似度,时间复杂度为 $O(n^2)$ 。在计算式(5)时,假设进行了 k 次迭代,每次迭代计算目标函数值时,复杂度最高的部分为将维度为 n 的 r_i^T 与矩阵 $A(n \times m)$ 相乘,其复杂度为 $O(nm)$ 。综上,UCSLIM算法总的时间复杂度为 $O(knm + n^2)$ 。空间上主要是矩阵 $R(m \times m)$ 和矩阵 $A(n \times m)$ 占用空间,故其空间复杂度为 $O(nm + m^2)$ 。

2 实验与分析

2.1 实验平台介绍

实验环境是由5台配置相同的服务器组成。Hadoop版本为2.60,Spark版本为1.3.0。每台服务器的配置情况为:双CPU 24核,内存64GB;每个核的型号为:Intel(R) Xeon(R) CPU E5-2620 v2 @ 2.10 GHz。

2.2 数据集以及预处理

为了验证算法的有效性,采用了5个数据集:movieLens100K、movieLens1M电影评分数据集,lastFM音乐收听数据集,news新闻阅读数据集,Book-Crossing(BX)图书阅读数据集。5个数据集的详细信息如表1所示。

表1 数据集详细信息

数据集	用户数	物品数	评分项	稠密度/%
mL100K	943	1 682	55 375	3.49
lastFM	1 865	2 000	55 963	1.50
news	1 149	5 035	41 057	0.71
ML1M	6 040	3 952	71 6307	3.00
BX	5 698	6 305	86 237	0.24

2.3 评价标准

2.3.1 推荐质量评价标准

实验采用TopN准确率作为评价指标。准确率是通过计算用户在TopN列表中实际打分的物品数量的比例来衡量推荐系统的预测质量。准确率越高,意味着其推荐质量越高。假设算法对于用户 X 推荐的 N 个物品为 $T\{t_1, t_2, \dots, t_N\}$,对应的用户 X 的实际购买物品为 $F\{f_1, f_2, \dots, f_N\}$,则该推荐算法对于用户 X 的准确率为

$$P = \frac{|T \cap F|}{|T|} \quad (6)$$

2.3.2 运行效率评价标准

采用加速比来评价UCSLIM算法在不同的云计算平台(Spark和Hadoop)上的计算效率。加速比是指同一个任务在单机系统和并行系统中运行消耗的时间的比率,一般用来衡量系统并行化的性能和效果。其计算公式为

$$S_p = \frac{U_1}{U_n} \quad (7)$$

其中: U_1 为任务在单机系统中的运行时间, U_n 为任务在并行系统中的运行时间。 S_p 的值越大,表示其加速性能越好,即并行效率越高。

2.4 实验结果及分析

2.4.1 3种算法在不同TopN时的推荐效果

为了验证算法的有效性,采用了SLIM、itemKNN两种算法作为UCSLIM算法的对比算法,每种算法设定的参数取值均为多次实验后确定的最优值。为了排除实验结果的偶然性,实验采用5折交叉验证的方式,其实验结果如图1~图5所示。

从实验结果可以看出,UCSLIM算法在各个数据集上的推荐效果都要优于SLIM算法。3种算法在news数据集上表现最好,因为相对于电影、图书等领域,用户偏好的新闻分类较为固定,故推荐算法能够较为准确地分析出用户喜好。在其余的数据集中,mL100K数据集的稠密度最高。稠密度越高,说

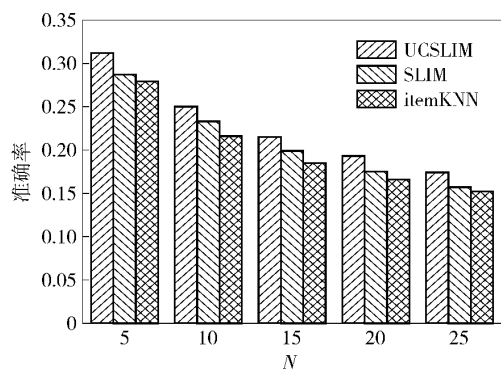


图 1 在 mLI100K 数据集中各推荐算法效果

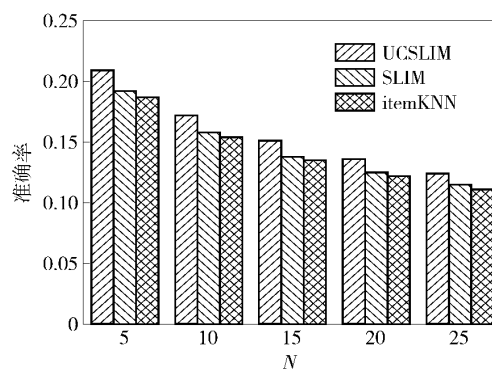


图 4 在 mLI1M 数据集中各推荐算法效果

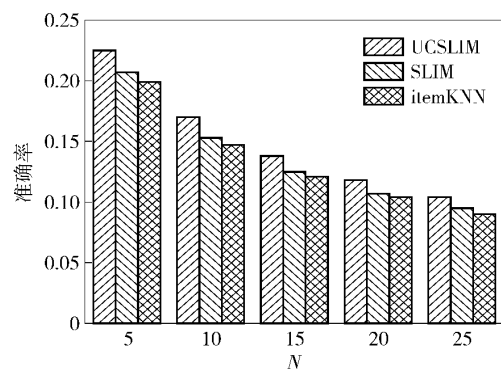


图 2 在 lastFM 数据集中各推荐算法效果

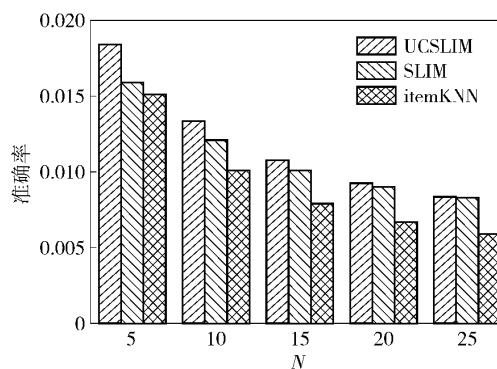


图 5 在 BX 数据集中各推荐算法效果

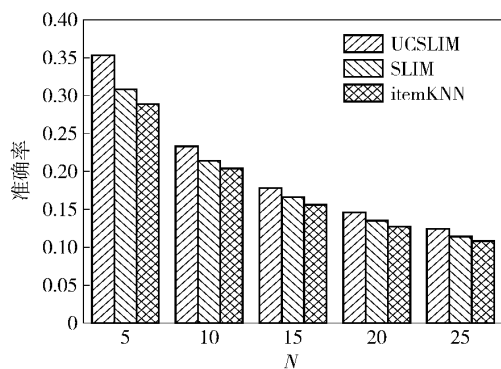


图 3 在 news 数据集中各推荐算法效果

明其中包含有效的信息越多,故其推荐效果要优于其他几个数据集。反之,BX 数据集的稠密度最低,其推荐效果是相对较差的。

2.4.2 运行效率比较

分别在 Spark 和 Hadoop 两种云计算平台上实现了 UCSLIM 算法,其在两个平台上的运行时间以及加速比如表 2 所示。

从表 2 可以看出,Spark 的加速效果要优于 Hadoop。一方面是因为 Hadoop 的 map 阶段需要将文件写入到硬盘中,并且会按照 key 排序,但是排序操

表 2 运行时间以及加速比

时间/s	mLI100K	lastFM	news	mLI1M	BX
Spark	4	9	6	87	75
Hadoop	37	65	28	300	192
单机	50	3 415	905	41 009	12 037
Sp(Spark)	12.5	397.4	150.8	471.4	160.5
Sp(Hadoop)	9.25	52.5	32.3	136.7	62.7

作在计算用户相关系数矩阵时是没有必要的。Spark 则是将数据保存在内存中,且没有排序操作,从而节省了时间。另一方面则是并行程度不同,Hadoop 是进程级并行,Spark 则是线程级并行,并行程度更高,故在计算用户相关系数矩阵时,能够更高效地并行计算。

3 结束语

在借鉴 SLIM 推荐算法思想的基础上,提出了基于用户集合的 UCSLIM 算法。UCSLIM 算法将用户划分为用户集合,并将用户-用户、用户-用户集合这两个隐含关系加入到算法当中,更加充分地挖掘了原始数据中的隐含信息,并通过实验证明 UCS-

LIM 算法比 SLIM 算法推荐质量更高。同时,在 Spark 和 Hadoop 两个云计算平台上实现了 UCSLIM 算法,并通过实验表明 UCSLIM 算法的 Spark 版具有更高的计算效率。提出的 UCSLIM 算法也存在可改进之处,如在划分用户集合时,可使用标准的频繁项集划分方式,这样会增加预处理时的时间复杂度,但是相应的集合划分可能更为精确。

参考文献:

- [1] Ning X, Karypis G. SLIM: sparse linear methods for Top-N recommender systems [C]//2011 11th IEEE International Conference on Data Mining. Vancouver, Canada: IEEE ICDM, 2011: 497-506.
- [2] Zaharia M, Chowdhury M, Das T, et al. Fast and interactive analytics over Hadoop data with Spark [J]. USENIX, 2012, 37(4): 45-51.
- [3] Marlin B M. Modeling user rating profiles for collaborative filtering [C]//Advances in Neural Information Processing Systems 16. Vancouver, Canada: NIPS, 2003: 627-634.
- [4] Sarwar B, Karypis G, Konstan J, et al. Item-based collaborative filtering recommendation algorithms [C]//Proceedings of the 10th International Conference on World Wide Web. Hong Kong: ACM, 2001: 285-295.
- [5] Zou H, Hastie T. Regularization and variable selection via the elastic net [J]. Journal of the Royal Statistical Society: Series B (Statistical Methodology), 2005, 67(2): 301-320.
- [6] Zhang T. Solving large scale linear prediction problems using stochastic gradient descent algorithms [C]//Proceedings of the Twenty-first International Conference on Machine Learning. New York, USA: ACM, 2004: 116.
- [7] Tsuruoka Y, Tsujii J, Ananiadou S. Stochastic gradient descent training for ℓ_1 -regularized log-linear models with cumulative penalty [C]//Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP: Volume 1-Volume 1. Suntec, Singapore: Association for Computational Linguistics, 2009: 477-485.