

BC-AW 协同过滤推荐算法研究^①

张志强^{1,2}, 李 改³

¹(顺德职业技术学院 电子与信息工程学院, 顺德 528300)

²(华中科技大学 计算机科学与技术学院, 武汉 430074)

³(中山大学 信息科学与技术学院, 广州 510006)

通讯作者: 张志强, E-mail: gd229@126.com

摘 要: 针对现有协同过滤算法普遍存在数据稀疏、可扩展性低、计算量大的缺点, 提出一种基于 BC-AW 的协同过滤推荐算法, 引入联合聚类 (BlockClust, BC) 和正则化迭代最小二乘法 (Alternating least squares with Weighted regularization, AW), 首先对原评分矩阵进行用户—项目双维度的联合聚类, 接着产生具有相同模式评分块的多个子矩阵, 通过分析得出这些子矩阵规模远小于原评分矩阵, 从而有效降低预测阶段的计算量. 然后分别对每个子矩阵应用正则化迭代最小二乘法来预测子矩阵的未知评分, 进而实现推荐. 经仿真实验表明, 本文算法与传统的协同过滤算法比较, 能有效改善稀疏性、可扩展性和计算量的问题.

关键词: 协同; 过滤; 联合聚类; 正则化迭代最小二乘法

引用格式: 张志强, 李改. BC-AW 协同过滤推荐算法研究. 计算机系统应用, 2018, 27(5): 198-202. <http://www.c-s-a.org.cn/1003-3254/6338.html>

Study on BC-AW Collaborative Filtering Recommendation Algorithm

ZHANG Zhi-Qiang^{1,2}, LI Gai³

¹(School of Electronics and Information Engineering, Shunde Polytechnic, Shunde 528300, China)

²(School of Computer Science and Technology, Huazhong University of Science & Technology, Wuhan 430074, China)

³(School of Information Science & Technology, Sun Yat-Sen University, Guangzhou 510006, China)

Abstract: Aiming at the weaknesses of sparse data, low scalability and large computing existing in the current collaborative filtering algorithm, a BlockClust-Alternating least squares with Weighted regularization (BC-AW) collaborative filtering recommendation algorithm is proposed. Firstly, the user and the item of the original scoring matrix are jointly clustered and several submatrixes with the same scoring mode are generated. According to the research, the scale of these submatrixes is far less than the original scoring matrix which effectively decreases the computational complexity in the prediction process. Then, the regularized iterative least-square method is applied to each submatrix to predict its score. Hence recommendation is realized. The simulation results reveal that the proposed algorithm can effectively improve sparsity, expand scalability, and reduce computing compared with the traditional one.

Key words: collaborative; filtering; joint cluster; regularized iterative least-square method

当今社会, 基于互联网技术的电子商务不断普及, 大数据分析和大数据挖掘已成为迫切要解决的问题. 推荐系统是指根据用户的兴趣特点和购买行为, 向用户推荐用户感兴趣的信息和商品的系统. 其中协同过

滤推荐 (collaborative filtering recommendation) 是推荐系统中应用最早和最为成功的技术之一, 它一般采用最近邻技术, 利用用户的历史喜好信息计算用户之间的距离, 然后利用目标用户的最近邻居用户对商品评

① 基金项目: 国家自然科学基金(41072247); 广东省自然科学基金(2016A030310018); 广东省哲学社会科学项目(GD16XJY36); 顺德职业技术学院重点教研项目 (2014-SZJGM06)

收稿时间: 2017-06-21; 修改时间: 2017-07-17; 采用时间: 2017-09-25; csa 在线出版时间: 2018-04-23

价的加权评价来预测目标用户对特定商品的喜好程度,系统从而根据这一喜好程度来对目标用户进行推荐.协同过滤最大优点是对推荐对象没有特殊的要求,能处理非结构化的复杂对象,如音乐、电影等媒体数据^[1].其中的潜在因素模型(latent factor models)的核心思想是:通过分析历史数据来预测事物的发展趋势,矩阵分解算法(matrix factorization)是其中最为成功的算法之一^[2].但随着互联网数据的不断扩大,应用传统的矩阵分解算法构造大数据潜在因素模型,算法性能会大大下降,难以实现有效的协同过滤,其关键问题体现在数据稀疏性和算法可扩展性的改进上.

针对上述问题,国内外学者专家们提出了一些改进想法: Pilsazy I 最先提出基于 ALS 的协同过滤算法,相对传统算法而言,能有效提高过滤质量和推荐速度,但该算法只考虑用户的相似性^[3];李改等对前者进行改进,提出了 ALS-WR 协同过滤算法,该算法考虑到可扩展性方面,但需要专门设计迭代式算法^[4];王辉等把用户聚类思想引入到协同过滤算法中,改进了传统算法的数据稀疏性和可扩展性,但精确度较低^[5].

本文提出一种基于 BC-AW 的协同过滤推荐算法,在传统算法的基础上,引入联合聚类 BC(BlockClust)和正则化迭代最小二乘法(Alternating least squares with Weighted regularization, AW).首先分解原始评分矩阵的用户项目双维度,然后使用联合聚类计算出多个同类评分块的子矩阵,最后使用正则化迭代最小二乘法预测各个子矩阵的未知评分,从而计算出推荐结果.使用联合聚类所求得的子矩阵远比初此矩阵规模小,可以大大减少过滤推荐的计算量.通过与 BaseMF 算法、RSTE 算法、TidaTrust 算法、MoleTrust 算法进行对比实验,分析结果表明,本文算法可以有效缓解传统算法并行化、可扩展性差的问题.

1 AW 算法

设定矩阵 R 代表 m 个用户、 n 个对象的评分矩阵、矩阵 U 代表用户特征矩阵、矩阵 V 代表推荐对象特征矩阵. $R_{i,j}$ 代表 R 中的某个元素, R_i 代表 R 的第 i 行, R_j 代表 R 的第 j 列, R^T 代表 R 的转置. R^{-1} 代表 R 的逆. 传算法一般使用 SVD 算法来分解 R ^[6], 但本文算法是通过计算低秩矩阵 X 来逼近 R , 使得 $X = UV^T$, $U \in C^{m \times d}$, $V \in C^{n \times d}$ (其中 d 代表 R 的特征数), 一般状态下, $d \leq r$, r 代表 R 的秩, $r \leq \min(m, n)$.

Step1: 我们通过最小化 Frobenius 的损失函数来算出一个低秩矩阵 X , 令 X 的得尽量接近 R , 如公式 (1) 所示:

$$L(X) = \sum_{ij} (R_{ij} - X_{ij})^2 \quad (1)$$

上述公式 (1) 中, $(R_{ij} - X_{ij})^2$ 为低秩逼近中一般平方误差项.

Step 2: 求解 $L(X)$ 的最优化问题 $\arg \min_x L(X)$, 于是改写公式 (1) 如下:

$$L(rmU, V) = \sum_{ij} (R_{ij} - U_i V_j^T)^2 \quad (2)$$

上述式 (2) 会产生过于拟合的问题, 因此需要通过正则化项来解决这个问题, 则改写如下:

$$L(U, V) = \sum_{ij} (R_{ij} - U_i V_j^T)^2 + \lambda (U_i^2 + V_j^2) \quad (3)$$

Step3: 对上面公式 (3), 设定 V 不变, 对 U_i 求导 $\partial L(U, V) / \partial U_i = 0$, 求解 U_i , 如公式 (4) 所示:

$$U_i = R_i V_{ui} (V_{ui}^T V_{ui} + \lambda n_{ui} I)^{-1}, i \in [1, m] \quad (4)$$

公式 (4) 中, R_i 为用户所评过影片的向量矩阵, 由评分组组成, V_{ui} 为用户 i 所评过影片的特征矩阵, 由特征向量组成, n_{ui} 为用户 i 所评过的影片数, I 为 $d \times d$ 阶的单位矩阵.

Step4: 上述对公式 (4) 中, 设定 U 不变, 可以通过公式 (5) 求解 V_j :

$$V_j = R_j^T U_{mj} (U_{mj}^T U_{mj} + \lambda n_{mj} I)^{-1}, j \in [1, m] \quad (5)$$

上述公式 (5) 中, R_j 为用户 i 所评过影片 j 的评分所形成的向量, U_{mj} 为用户 i 评过影片 j 的特征矩阵, 它由用户群中的特征向量组成, n_{mj} 为用户 i 评过影片 j 的用户数量.

2 BC 算法

BC 算法能够将用多个较小且高相似度的评分矩阵来替代原始数据^[7]. 通过扫描该矩阵中的评分, 使用聚类方法计算行和列, 计算一次后, 重新评估每个子矩阵中用户、项目的概率, 如此类推, 直到产生收敛为止^[8]. 然后将评分匹配到最大概率的子矩阵中. 其算法思想如下:

首先扫描评分矩阵, 通过计算各评分的所属概率 $p(k|u, v, r)$ 来判断该评分是否属于某个子矩阵, 其参数包括: 用户—项目属于某子矩阵的概率 $p(k|u)$ 和 $p(k|v)$ 、该评分值属于某子矩阵中的概率 $p(r|k)$. 其算法公式如下 (其中 $\alpha = 0.000\ 0001$, $\beta = 0.000\ 0001$, $\theta = 0.000\ 0001$):

$$p(k|u, v, r) = \frac{(p(k|u) + \alpha) \times (p(k|v) + \beta) \times (p(r|k) + \theta)}{\sum_{k' \in k} (p(k'|u) + \alpha) \times (p(k'|v) + \beta) \times (p(r|k') + \theta)} \quad (6)$$

$$p(k|u) = \frac{\sum_{v \in V(u)} p(k|u, v, r)}{\sum_{k' \in k} \sum_{v \in V(u)} p(k'|u, v, r)} \quad (7)$$

$$p(k|v) = \frac{\sum_{u \in U(V)} p(k|u, v, r)}{\sum_{k' \in k} \sum_{u \in U(V)} p(k'|u, v, r)} \quad (8)$$

$$p(r|k) = \frac{\sum_p (k|u, v, r)}{\sum_{r'} \sum_{v \in U(V)} p(k'|u, v, r)} \quad (9)$$

以上公式中, u 表示当前用户, v 表示当前项目, r 表示评分值, r' 表示1到5的整数, k 表示当前聚类, k' 表示当前累加聚类, $U(v)$ 表示对 v 评了分的所有用户集合, $V(u)$ 表示 u 已评过分的的所有项目集合. BC算法把原矩阵分为 k 个具有高密度、高相似度的子矩因此,该算法具有一定的降维作用^[9].

3 BC-AW 算法

以下通过联合聚类和 AW 协同过滤两阶段来对评分矩阵的未知项进行预测.

1) 联合聚类

输入: 子矩阵个数 k , 用户—项目评分矩阵 R .

输出: 子矩阵数 k .

Step1: u 表示用户量, v 表示项目数, $p(k|u, v, r)$ 表示评分 r 属于聚类 k 的概率, 分别初始化 u 、 v 、 $p(k|u, v, r)$, 使得 $\sum' k p(k'|u, v, r) = 1$.

Step2: 应用公式(7)计算用户 u 属于聚类 k 的概率 $p(k|u)$; 应用公式(8)计算用户 v 属于聚类 k 的概率 $p(k|v)$; 应用公式(9)算出分值概率 $p(r|k)$; 应用公式(6)算出 $p(k|u, v, r)$.

Step3: 最后选择概率 k 的最大值作为该评分的子矩阵.

Step4: 跳回步骤2, 直到出现收敛, 结束循环.

2) AW 协同过滤

首先用偏差小于等于的高斯随机数初始化矩阵 V , 接着分别用公式(4)和公式(5)更新 U 和 V , 重复进行多次迭代, 直到得出的 $RMSE$ 值出现收敛为止, 结束迭代.

输出: 用户评分矩阵 R , 特征个数 d .

输出: R 的逼近矩阵 X .

Step1: 用偏差小于等于0.01的高斯随机数初始化 V ;

Step2: 分别用公式(4)和公式(5)更新 U 和 V ;

Step3: 重复迭代公式(4)、公式(5), 判断得出的 $RMSE$ 值是否收敛, 如果是, 则迭代结束;

Step4: 令 $X = UV^T$, 返回矩阵 X .

下面我们分析该算法的时间复杂度: 设 n_r 为在矩阵 R 里的评分点个数, $n_r \leq m \times n$ (一般情况下, 矩阵 R 稀疏). n_f 为特征个数, n_i 为该算法的迭代数; 设 n_u 为用户数; 设 n_m 为对象推荐数. 变量 U 每次更新的时间为 $O(n_f^2(n_r + n_f n_m))$, 变量 V 每次更新的时间为 $O(n_f^2(n_r + n_f n_m))$, 如此类推, 在迭代 n_i 次后的时间则为 $O(n_f^2(n_r + n_f n_u + n_f n_m) n_i)$. 由此可见, 在 n_f 、 n_r 、 n_u 以及 n_i 恒定的情况下, $O(n_f^2(n_r + n_f n_u + n_f n_m) n_i)$ 取决于 n_r , 也就是说, 算法的总时间复杂度与 n_r 成正比, 可证该算法具有一定的可扩展性. 在AW协同过滤中, Step 2是最关键的步骤, 在这个步骤里, 每次应用公式(4)和公式(5)只能更新 U 和 V 中的一行值, 效率较低. 因此可以先将 U 和 V 分解为 n 个列相等的子矩阵, 然后, 应用公式(4)和公式(5)对每个子矩阵并行更新, 这样可以大大提高运算效率.

综上所述, 本文基于矩阵分解的BC-AW协同过滤推荐算法改进传统算法的串行运算方式, 实现并行运算, 可以有效地解决传统算法中存在的可扩展性差问题.

4 实验分析

4.1 数据准备

本实验目的是通过分析比较本文算法和传统算法的性能, 硬件方面选用基于云计算的分布式硬件实验平台, 该平台由20台电脑组成20个运算节点, 每个节点配有4.0 Hz的Intel处理器和32 G的内存, Linux操作系统. 实验数据来自美国Minnesota大学GroupLens项目组的MovieLens数据集. MovieLens数据集是GroupLens项目组开发的一个互联网研究型推荐系统数据集. MovieLens数据集中含有943个用户对1682部电影的100 000条评分数据, 平均每个用户评分的电影大于等于20部^[10].

4.2 评价指标

本次实验采用均方根误差 $RMSE$ 和评分覆盖率 COV_R (Rating Coverage)作为算法推荐质量的评价指标.

$$RMSE = \sqrt{\sum_{(u,i) \in R_{test}} (r_{ui} - \hat{r}_{ui})^2 / R_{test}} \quad (10)$$

其中 r_{ui} 代表用户的实际评分, \hat{r}_{ui} 代表用户对影片的预测评分, R_{test} 代表影片的评分数. $RMSE$ 的值越小表示算法的精度越高. Rating Coverage 是指在测试集中, 推荐系统可预测的评分数量占总评分的百分比^[11].

$$COV_R = \frac{\text{可预测的评分数}}{\text{整个测试集的评分个数}} \quad (11)$$

4.3 实验结果及分析

本文算法分别与 BaseMF (基于矩阵分解算法)^[12], RSTE (概率矩阵分解算法)^[13], TidaTrust (基于 TidaTrust 模型推荐算法)^[14], MoleTrust (基于 MoleTrust 模型推荐算法)^[15] 这 4 种推荐方法通过分析 $RMSE$ 和 COV_R 的结果值进行分析比较. 实验开始, 首先随机选取两组训练集 GROUP1 和 GROUP2: GROUP1 为 80% 的 Epinions 数据集, GROUP2 为 90% 的 Epinions 数据集. 然后分别把本文算法、BaseMF 算法、RSTE 算法、TidaTrust 算法、MoleTrust 算法所建立的模型应用到训练集上进行学习, 指定 γ 、 λ 和 f 为潜在维度 ($\gamma = 0.005$, $\lambda = 0.003$, $f = 0.001$), k 为实验基础 ($k = 4$). 反复实验 5 次, 求平均值作为最终实验结果. 图 1 为这 5 种算法得出的 $RMSE$ 值比较图.

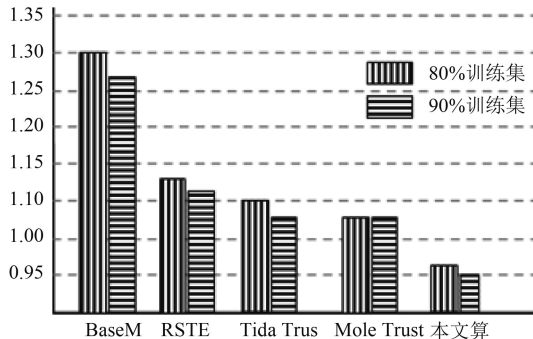
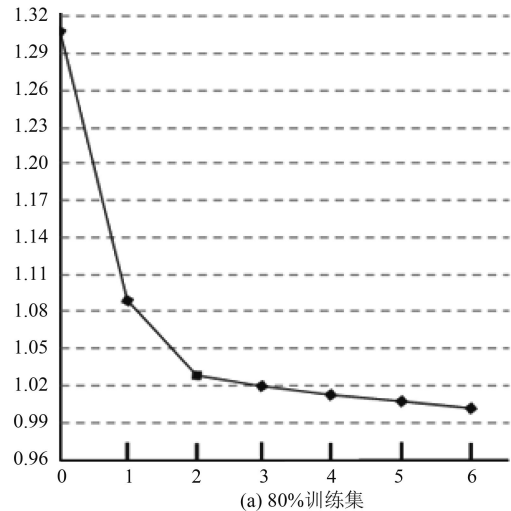


图 1 5 种算法的 $RMSE$ 值比较图

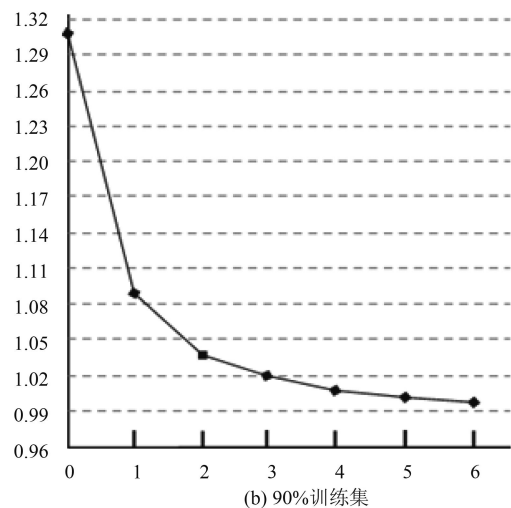
从上图可以看出, 针对 $RMSE$ 值, 当训练集为 80% 时, 本文算法比 BaseMF 算法降低了 22%, 比 RSTE 算法降低了 11%, 比 TidaTrust 算法降低了 8%, 比 MoleTrust 算法降低了 7%; 当训练集为 90%, 本文算法比 BaseMF 算法降低了 21%, 比 RSTE 算法降低了 10%, 比 TidaTrust 算法降低了 7%, 比 MoleTrust 算法降低了 7%; 由此可知, 本文算法的误差值明显低于传统的 4 种推荐算法.

图 2(a)、图 2(b) 分别表示本文算法在 GROUP1、GROUP2 这两组训练集中, $RMSE$ 在不同 k 值下的变化情况. k 值越大, $RMSE$ 值越小, 证明推荐精度越高. 当

$k \in [0, 4]$ 时, $RMSE$ 减小的幅度最大, 迭代时间的增幅最小; 当 $k \geq 4$ 时, $RMSE$ 的变化幅度趋于 0. 从而得出当 $k = 4$ 时, 算法的推荐性能最佳.



(a) 80% 训练集



(b) 90% 训练集

图 2 $RMSE$ 随参数的变化曲线

下面分析 5 种算法的 COV_R . 我们知道, COV_R 直接影响用户对影片的可选择范围, COV_R 越高, 用户对影片的可选择范围越广, 从而用户满意度越高. 采用 GROUP1 数据集 (即 80% 的训练集), 实验结果如表 1 所示.

表 1 5 种算法的 COV_R

方法	$COV_R(\%)$
BaseMF	20.19
RSTE	51.62
TidaTrust	84.87
MoleTrust	82.13
本文算法	98.71

分析表1可以看出,当训练集为80%时,本文算法的 cov_R 比BaseMF算法提高了68.52%,比RSTE算法提高了38.09%,比TidaTrust算法提高了4.84%,比MoleTrust算法提高了6.58%。由此可见,本文算法比传统4种算法具有更广的覆盖范围。

4 结论与展望

本文针对现有协同过滤算法中普遍存在的数据稀疏、可扩展性低这两个核心问题,提出BC-AW(基于联合聚类和迭代最小二乘法)两阶段协同过滤算法。首先由该算法通过对原评分矩阵进行用户—项目双维度联合聚类后得到的若干个子矩阵(这些子矩阵均为模式相同的评分块),其规模比原矩阵小得多,因此可以有效减少预测工作量;再者,采用正则化迭代最小二乘法预测子矩阵的未知评分可以优化推荐效率。该算法在模拟大数据实验中(美国Minnesota大学GroupLens项目组的MovieLens数据集),通过与几个经典的协同过滤算法(BaseMF算法、RSTE算法、TidaTrust、MoleTrust算法)作比较。实验结果表明,本文算法能有效改进推荐系统的稀疏性、可扩展性问题,系统预测评分与用户实际评分接近,并能为用户提供良好的使用体验。

参考文献

- 1 Shuo LX, Chai BF, Zhang XD. Collaborative filtering algorithm based on improved nearest neighbors. *Computer Engineering and Applications*, 2015, 51(5): 137–141.
- 2 Zhang HJ. The research and application of distributed matrix factorization algorithm in recommend system. *Bulletin of Science and Technology*, 2013, 29(12): 151–153.
- 3 Wang QM, Miao Y, He M, *et al.* Parallelized research on collaborative filtering algorithm based on matrix factorization. *Computer Technology and Development*, 2015, 25(2): 55–59.
- 4 Zhu X, Song AB, Dong F, *et al.* A collaborative filtering recommendation mechanism for cloud computing. *Journal of Computer Research and Development*, 2014, 51(10): 2255–2269.
- 5 Bi XR. Collaboration filtering recommendation algorithm of sub-similarity integration between items. *Computer Systems & Applications*, 2015, 24(1): 147–150.
- 6 Wang L, Fu XF, Wang XM. Hybrid dynamic collaborative filtering algorithm based on big data sets. *Journal of Guangdong University of Technology*, 2014, 31(3): 44–48.
- 7 Chen W, Shi QL. An adaptive algorithm for collaborative filtering recommendation. *Journal of Xianyang Normal University*, 2014, 29(6): 47–49.
- 8 Li G, Li L. One-class collaborative filtering based on matrix factorization. *Application Research of Computers*, 2012, 29(5): 1662–1665.
- 9 Ke LW, Wang J. Collaborative filtering recommendation based on user feature transfer. *Computer Engineering*, 2015, 41(1): 37–43.
- 10 Yi ZA, Mu CM. Collaborative filtering algorithm based on subtractive clustering and genetic fuzzy. *Computer & Digital Engineering*, 2014, 42(8): 1363–1367.
- 11 Xu W, Duan F. Combining clustering and collaborative filtering for implicit recommender system. *Computer Engineering and Design*, 2014, 35(12): 4181–4185.
- 12 Liu L. A collaborative filtering recommender algorithm based on iterative kernel method. *Information Technology & Informatization*, 2014, (12): 76–81.
- 13 Zha J, Li ZB, Xu GQ. An optimised collaborative filtering algorithm based on combined similarity. *Computer Applications and Software*, 2014, 31(12): 323–328.
- 14 Wu HC, Wang XJ, Cheng Y, *et al.* Advanced recommendation based on collaborative filtering and partition clustering. *Journal of Computer Research and Development*, 2011, 48(S3): 205–212.
- 15 Luo Q, Miao XJ, Wei Q. Further research on collaborative filtering algorithm for sparse data. *Computer Science*, 2014, 41(6): 264–268.