

doi:10.3772/j.issn.2095-915x.2016.01.011

基于二分网络的推荐算法研究进展

周波, 杨朝峰

(中国科学技术信息研究所 北京 100038)

摘要: 全面总结和介绍基于二分网络的推荐算法研究现状, 旨在介绍基于二分网络推荐算法的思想和各种优化方法, 帮助读者了解这个研究领域。先介绍了二分网络推荐算法研究的背景和基于物质扩散和热传导的两种基本二分网络推荐算法, 然后总结了8大类的优化算法, 最后指出了当前还未研究的但还值得进一步研究的地方, 并对大数据环境下基于二分网络的推荐算法进行了展望。

关键词: 二分网络, 推荐算法, 热传导, 物质扩散

The Research Progress of Recommendation Algorithm Based on Bipartite Network

ZHOU Bo, YANG Chaofeng

(Institute of Scientific and Technical Information of China, Beijing 100038)

Abstract: this paper makes a comprehensive research of recommendation algorithm based on bipartite network aiming to make readers understand this field well. Firstly, it introduced the background of the bipartite network recommendation algorithm and two basic algorithm which is called material diffusion and heat conduction, secondly it summarizes the eight major method used to improve algorithm performance, Finally it points out the field and problem which is worth studying while the current research has not been studied and discuss the further direction of recommendation algorithm in the big data age.

Keywords: Bipartite network, recommendation algorithm, heat conduction, material diffusion

作者简介: 周波 (1991-), 中国科学技术信息研究所情报学硕士研究生, 研究方向: 数据挖掘与知识组织, E-mail: hubeizhoubo@163.com; 杨朝峰 (1975-), 男, 中国科学技术信息研究所副研究员, 同济大学技术经济及管理专业博士, 研究方向: 科技政策、宏观经济政策。

1 引言

推荐系统在信息大爆炸时代的重要性不言而喻,推荐系统在社交网络中可以增加社交网络黏性,在电子商务中可以提高销售收入,科研合作网络中可以发现学术团体和学科研究现状,在蛋白质等化合物网络中起到了发现新的科学知识的作用。在未来,推荐系统将在人类生活的各个方面发挥重大作用,学术界对这一领域进行了大量的研究。复杂网络是新兴的交叉学科,涵盖了数学、物理、计算机,人工智能等领域的知识,同时又对这些领域的研究起到推动作用。复杂网络产生的理论和方法对数据挖掘,机器学习以及推荐算法等产生了较大的影响。基于二分网络的推荐算法由周涛^[1]2007基于复杂网络理论而提出,该算法提出后,受到广泛关注,基于二分网络的推荐算法较传统的协同推荐算法有着更高的推荐精度,此后衍生出一系列的基于二分网络的推荐算法以及各种优化算法。这些优化的推荐算法以物质扩散算法和热传导算法为基础,以后的研究主要针对具体的情况做出优化,并且都取得了不错的效果,随着人人互联,人机互联和物物互联时代和大数据时代的来临,更多的基于复杂网络的新的推荐算法将会被提出。本文对过去基于二分网络的各种推荐算法进行总结和展望,方便后来的学者的进一步的研究。

2 基于二分网络的两种基本推荐算法

2.1 基于物质扩散的二分网络推荐算法

该算法假设每个对象均有一定的初始资源,通过对象的度将资源平均地分配给相邻的用户,然后每个用户又将自己所有分到的资源再次平均

地分配给所选择的对象,通过汇总对象的所有相邻用户分配的资源,得到该对象获得的资源。现已证明通过初始资源在网络上的扩散原理进行链路预测,比协同过滤算法具有更高的精度。^[1]基于物质扩散的算法,由于不考虑用户和对象的内容特征和思想特征,只是把他们抽象成二分网络的节点,算法仅利用隐藏在用户和产品的选择关系之中的信息。文献^[2]中,在二分图上引入了扩散动力学,实现了基于物质扩散(MD)算法,证明了这些算法的结果明显好于经典的协同过滤。

用 X, Y 分别表示二分网络两类节点用户和产品的集合,其中 X 类节点有 m 个, Y 类节点有 n 个, A 表示 X 同 Y 的连接关系, B 是 A 的转置矩阵,如果 X_i 同 Y_j 相连,那么 $a(i, j)$

等于1, $b(j, i)$ 等于1,否则为零, $d(X_i)$ 和 $d(Y_j)$ 分别表示的是用户 X_i 和物品 Y_j 的度。(下文

相关的公式符号含义同此)设用户 X_i 对物品 Y_j 的初始资源为 $h(j, i)$,物品将资源平均地分配给每个用户,每个用户又将所得到资源平均分配给他选择过的物品。用物质扩散算法对物品进行推荐主要可以分为两个步骤:

第一步:资源由物品 Y_j 向用户 X_l 扩散,设用户 X_l 的分配的资源为 $f(X_l)$,则

$$f(X_l) = \sum_{j=1}^n \frac{b(j, l) * h(j, i)}{d(Y_j)}$$

第二步:资源由用户 X_l 向物品 Y_k 扩散,物品 Y_k 所得的资源为 $g(Y_k)$

$$g(Y_k) = \sum_{l=1}^m \frac{f(X_l) * b(k, l)}{d(X_l)}$$

$$= \sum_{l=1}^m \frac{b(k, l) \sum_{j=1}^n \frac{b(j, l) * h(j, i)}{d(Y_j)}}{d(X_l)}$$

矩阵 $W(k, j)$ 记为资源分配矩阵,表示产品 j 愿意分配给产品 k 的资源配额,则 $W(k, j)$ 的一般表达式为:

$$W(kj) = \frac{1}{d(Yj)} \sum_{l=1}^m \frac{b(k,l) * b(j,l)}{d(Xl)}$$

在矩阵 B 的第 i 列是一个 n 维的 0/1 矢量，用户选择的商品其对应的值为 1，其他设为 0，可用来表示用户 Xi 的选择信息，一种常用的初始资源分配方法是假设所有的用户初始资源均为 1，则用户 i 对物品 j 的初始资源分配为 b(j,i)/d(Xi)，即 h(ji) = h'(ji) = b(ji)/d(Xi)。则物质扩散过程中用户 Xi 对各个物品最终资源分配量为：

$$g(Yk) = \sum_{j=1}^n W(kj) * h'(ji)$$

对用户未选择过的产品按资源值大小进行排序，选择 top-N 作为用户 Xi 的推荐列表。

2.2 基于热传导的二分网络推荐算法

基于热传导的链路预测方法其思想来自于物理学的热传导链路预测过程，但是其将物理学中的热传导方程进行离散化，物理学中的热平衡方程是 $-k \nabla^2 T(r) = \nabla \cdot j(r)$ ，k 代表导热性， $\nabla^2 T(r)$ 代表温度梯度。^[3] 用物理学中的基本理论来进行链路预测存在两个问题，一是算法的时间复杂度和空间复杂度都很高，影响大规模的数据处理和运算，二是理解起来较为复杂。鉴于此，周涛根据热传导的思想提出了一种更加形象的热传导算法，基本形式同热传导算法相似，但是有区别，热传导算法链路预测也可分为两步。^[4] 借用在物质分配算法中的符号，可表示如下：

第一步：热量由物品 Yj 向用户 Xl 扩散，设用户 Xl 的分配的资源为 f(Xl)，则

$$f(Xl) = \sum_{j=1}^n \frac{b(j,l) * h(ji)}{d(Xl)}$$

第二步：热量由用户 Xl 向物品 Yi 扩散，物品 Yi 所得的资源为 g(Yi)。

$$g(Yk) = \sum_{l=1}^m \frac{f(Xl) * b(k,l)}{d(Yk)}$$

$$= \sum_{l=1}^m \frac{b(k,l) \sum_{j=1}^n \frac{b(j,i) * h(ji)}{d(Xl)}}{d(Yk)}$$

矩阵 W(kj) 记为资源分配矩阵，表示产品 j 愿意分配给产品 k 的热量，则 W(kj) 的一般表达式为：

$$W(kj) = \frac{1}{d(Yk)} \sum_{l=1}^m \frac{b(k,l) * b(j,l)}{d(Xl)}$$

同物质扩散算法类似，热传导算法的初始资源分配为 h(ji) = h'(ji) = b(j,i)/d(Yj)，则热传导算法中用户 Xi 对各个物品最终资源分配量为

$$g(Yk) = \sum_{j=1}^n W(kj) * h'(ji)$$

对用户未选择过的产品按资源值大小，进行排序，选择 top-N 作为用户 Xi 的推荐列表。

2.3 基于物质扩散和热传导的二分网络推荐算法比较

1) 物质传播矩阵中每个元素的意义是物品 yi 可以分配给用户 yj 的资源量，其体现的源物品对其他物品的资源分配能力，源物品的度越多其资源分配能力就越强，热量扩散矩阵中每个元素的意义是，物品 yj 可以接受用户 yi 的资源量，反映出的物品 yj 的资源接受能力，yj 的度越大其资源接收能力就越弱。^[5]

2) 针对具体的用户来说，在物质扩散的过程中，资源总量总是守恒；而基于热传导算法中，资源是不守恒的。而从整体上来看，物质扩散最终生成的资源总量恒等于 n，热传导算法最终产生的资源总量恒等于 m。^[5]

3) 物质扩散算法的资源传播矩阵各列元素之和等于 1，热量传播矩阵各行元素和等于 1。

4) 已有实验证明，热传导算法倾向于为用户推荐冷门产品，而物质扩散算法倾向于推荐热门产品。^[6]

5) 同传统的协同过滤算法比较

文献 [7] 证明了在不同的推荐列表长度时, 基于资源分配的二分网络推荐算法的命中率和排序分均优于协同推荐算法。

3 二分网络推荐优化算法

推荐算法是各种推荐系统的核心, 算法的优劣决定了推荐系统的好坏。由于基于二分网络的推荐算法的优势, 因此引起很大的关注, 目前很多学者从不同的角度对基于二分网络的推荐算法进行了研究, 各种优化算法被提出。这些优化算法可以分为以下几类, 考虑初始资源分配差异的推荐算法, 考虑用户的度和物品的度的优化算法, 考虑用户的个性化偏好的优化算法, 几种算法的组合算法, 考虑冗余属性的优化算法, 基于内容的优化算法, 基于社团的优化算法和基于三部曲的优化算法, 下面将介绍各个方法。

3.1 基于初始资源分配的优化算法

基于标准的热传导和物质分配算法, 一般假设各个用户的初始资源为 1, 然而在实际中这种初始资源分配不能完全适应情况, 例如, 经济条件好的人购物的次数比经济条件差的人要多, 因此推荐系统给这两种用户的推荐力度是不同的。基于此, 文献 [8] 研究了商品的度即商品的流行性同推荐能力的关系, 周等提出用户的初始资源分配算法, 将用户的初始资源分配同用户的度关联起来。用户 i 对物品 j 的初始资源分配为:

$$h'(ji) = \frac{d(Xi)^\lambda}{d(Xj)}$$

采用基于物质分配算法进行实验, 得出 λ 等于 -0.8 时推荐精度最高。说明降低度大的节点的推荐能力, 能有相对较好的推荐精度。另外文章通过海明距离衡量出该算法有比较好的推荐列表

的多样性。这就表明用户在同样喜好程度下, 推荐非流行的商品比推荐流行的商品更有意义。

3.2 考虑流行性的优化算法

考虑流行度的优化方法是基于以下情景: 夏天某款产品裙子非常流行, 有很多人购买, 那么该款裙子将会很快被推荐给所有的人, 这样一来推荐结果的多样性就会降低, 与此同时那些冷门的产品推荐的力度就会减小, 这会强化推荐系统的马太效应, 如何降低马太效应同时提高推荐精度就是考虑流行性的优化方法主要出发点。

文献 [9] 考虑了这种情况, 并以热传导算法为基础, 提出了偏热传导算法, 基本原理同热传导算法相似, 但是考虑大度产品的影响力, 在最后一步传导过程中考虑产品度的影响, 产品最后得到资源不是除以产品的度, 而是除以产品度的 λ 次方, λ 是一个可调节的参数, 这样热传导算法的转移矩阵为

$$W(kj) = \frac{1}{d(Yk)^\lambda} \sum_{l=1}^m \frac{b(k,l) * b(j,l)}{d(Xl)}$$

将标准热传导算法和偏热传导算法进行对比实验, 结果表明 $\lambda = 0.85$ 时, 基于偏热传导算法的推荐结果比标准的热传导算法提高了很多。证明降低大度节点权重有利于提高推荐结果。

3.3 基于用户的个性化优化算法

3.3.1 考虑不同用户接收能力不同的优化方法

文献 [7] 还研究了资源接受者的接收能力的差异对推荐精度的影响。例如, 将一滴墨水滴入水杯中将一滴水滴入大海中, 其影响效果就截然不同, 将一滴水滴入水杯中, 水杯中的水的颜色就改变了, 而滴入大海中却没有任何变化。考虑这种情况对推荐算法的影响, 基于资源接收能力不同的个性化算法被提出。资源传播矩阵 $W(kj)$ 表示物品 Y_k 从物品 Y_j 那里得到的资源数量, 在

资源扩散过程中，考虑接收者 Y_k 对所收到资源的响应，假定响应反比于它们自身度的 λ 次方，则得到资源传播矩阵为：

$$W(kj) = \frac{1}{d(Y_k)^\lambda d(Y_j)} \sum_{l=1}^m \frac{b(k,l)b(j,l)}{d(Xl)}$$

利用 MovieLens 的数据集进行了实验，得出结论是当 $\lambda=0.91$ ，推荐结果最优。在推荐算法中适当降低流行的产品的推荐力度将有利于提高推荐结果的精度，但是不宜过高，否则将会降低推荐结果的精度，这说明不同的资源接受者的能力将会影响到推荐结果。

3.3.2 考虑用户偏好的优化方法

用户对所购买的产品的偏好程度是不一样的，有些产品是极度的喜欢，有些产品是一般喜欢，有些产品是用户所讨厌的，因此把用户所有选择过的产品一视同仁是不合理的，在基于标准的物质分配和热传导的二分网络推荐算法中，每一条连边都是一样的，其权重都是一，然而在实际中，用户对各个物品的偏爱是不一样的，例如在淘宝上，用户 X_i 购买了 Y_1, Y_2, Y_3 三种产品，用户对 Y_1 是好评， Y_2 是中评， Y_3 是差评，因此用户更倾向于购买与 Y_1 联系密切的产品，不倾向于购买同 Y_3 联系密切的产品，因此推荐系统应该根据用户的这种倾向进行推荐，客户的满意度就会更高。^[9]

在很多推荐系统中，都有评分机制，让用户对所选择的物品进行评分，那么就可以用评分来表示用户偏好。文献^[8]提出了评分加权的二分网络推荐算法，设 B 表示用户评分矩阵，对各个用户的评分进行归一化处理，得到用户 - 评分权重矩阵 R ， $r(j,i)$ 表示用户 X_i 对物品 Y_j 的评分，则资源传播矩阵为：

$$W(kj) = \frac{1}{d(Y_j)} \sum_{l=1}^m \frac{r(k,l)r(j,l)}{d(Xl)}$$

文献 [10] 采用 MovieLens 数据集进行了实验，

证明采用评分加权的二分网络的推荐算法的精度要优于标准的二分网推荐算法。

3.3.3 考虑用户所选择产品特性的优化算法

在现实生活中我们经常可以看到一些情况是：一些富豪喜欢购买奢侈品；还有一些人喜欢的物品其他的人很少喜欢；有些产品是大众喜欢老少皆宜的，因此，如果在推荐系统发现某用户所购买的商品都是大众购买的，那么继续向其推荐大众产品，发现某用户所购买的商品都是很少人购买的奢侈品，那么应继续向其推荐奢侈品，而不应该向其推荐大众商品，因此文献^[5]提出了一种基于用户选择产品特性的个性推荐方法。

其资源传播矩阵如下

$$W(kj) = \frac{1}{d(Y_k)^\lambda d(Y_j)^{1-\lambda_i}} \sum_{l=1}^m \frac{b(k,l)b(j,l)}{d(Xl)} \lambda_i = \left(\frac{d(X_i)}{dmax}\right)^\gamma$$

其中， $d(X_i)$ 表示用户 i 所选择物品的平均度， $dmax$ 表示所有物品的最大度， $0 \leq \gamma \leq 3$ ，由于 $d(X_i) < dmax$ ， $0 \leq \frac{d(X_i)}{dmax} \leq 1$ ，通过 λ_i 来控制用户的选择偏好，即用户是喜欢选择流行性产品还是喜欢选择冷门产品，其与用户选择产品的平均度呈正相关。

两个不同用户，所选择的产品不同，其所选择的产品的平均度存在差异，而推荐算法将这种差异表达出来了，用户的个性化混合参数正比于所选择产品的平均度，因此可以提高预测的精确度。文献^[5]采用 movielens 数据集进行实验，结果证明推荐算法的准确率和召回率都有很大的提升。

3.4 消除冗余属性的优化算法

在基于二分网络的推荐算法中，推荐结果中存在大量的冗余属性^[11]。例如：用户喜欢作者 A 写的书和出版社 B 出版的书。假设该用户只看过两本书籍，这两本书籍一本是由 A 创作的，另一本是由 B 出版的。如果恰有一本书籍 C 由 A 创作

B出版,那么这两本书籍将分别对书籍产生作用,推荐的总强度为2,因为C和A关联,C和B关联,关联的强度都是1。考虑另一种情况,该用户看过的两本书籍都是由出版B的,但是作者都不是A,那么对于另外一本也是由B出版的书籍D,那么书籍D推荐的总强度也是2。显然地,这里来自书籍的推荐包含了重复的属性(出版社B),所以虽然具有同样的强度,用户应该更喜欢书籍既是B出版的,又是A创作的,而不是仅仅是B出版的书^[11]。

因此,文献^[12]使用一种简单的算法来减少这种重复属性的影响。考虑两个书籍在对另一本书籍的推荐中包含了重复的属性,自然这个属性会导致这两本书籍自身的关联,从原来的关联矩阵中合理调节二阶关联,可以提高算法的精确度。

考虑消除冗余属性的算法是在物质扩散算法基础上考虑二阶资源分配,其资源转移矩阵为:

$$W'(kj)=W(kj)+\alpha W(kj)*W(kj)$$

其中 α 为可调节参数。

采用 movielens 数据集进行实验,实验结果证明 $a=-0.75$ 时算法的准确率最高,相比没有消除二阶冗余属性的算法,命中率提高了12%,提高程度很大。当将此算法进一步推广到三阶时,命中率提高1%左右,且时间复杂性增加很大,故在实际中运用到二阶运算就可以了。

3.5 考虑内容的优化算法

3.5.1 基于热传导的三部图优化算法

基于二分网络的三部图推荐算法有多种,基本原理类似,本文只是介绍较为简单的一种。基于标签的推荐算法也是学术界研究的热点,其基本思想是基于用户选择的物品标签,向用户推荐具有相同标签的物品。因此有学者将基于标签推荐算法同二分网络相结合,文献^[13]提出了用户-商品-标签的三部图,将物品的属性(标签)融入二分网络中,提出了新的优化方法。如下图所示,Z表示标签,3幅图展示在三部图上的资源扩散过程。

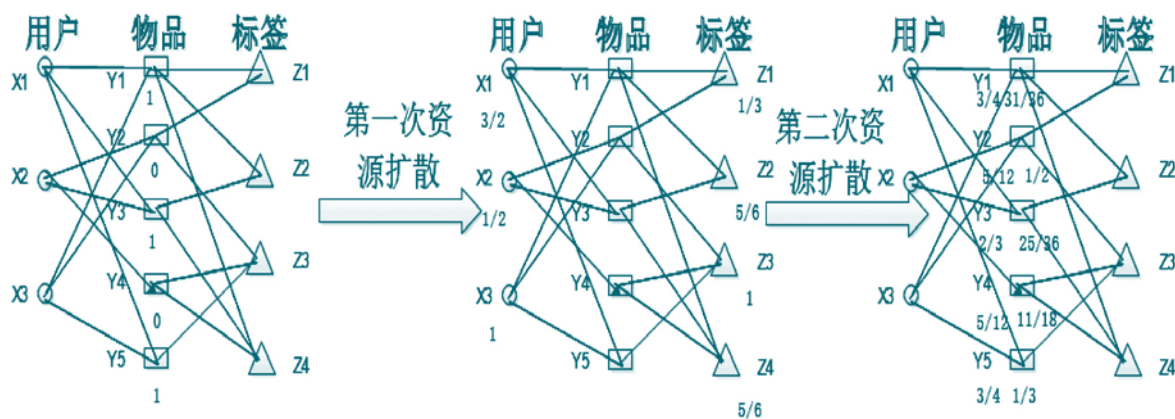


图1 基于三部图资源扩散算法

基于三部图物质扩散方法是将这个三部图看做是由两个二分网络所构成,左边是用户-物品的二分网络,右边是物品-标签的二分网络,对于用户 x_i 来说,分别对这两个二分网络进行物质扩散的推荐算法。

对于左侧的二分网络,同标准的物质扩散算法相同,其最后的资源分配结果为:

$$g(Yk) = \sum_{j=1}^n Wkj * h'(ji)$$

对于右侧的二分网络,也采用标准的物

质扩散算法，其资源过程为：

$$W'(kj) = \frac{1}{d(Yj)} \sum_{l=1}^k \frac{c(i,l) * c(j,l)}{d(Zl)}$$

$$g'(Yk) = \sum_{j=1}^n W'(kj) * h'(ji)$$

其中， W'_{ij} 表示在物品 - 标签二分网络上的资源分配矩阵， $c(i,l)$ 表示物品 i 同标签 l 的连接关系， $d(Zl)$ 表示标签 l 的度。

将上述两个转移矩阵进行线性组合，就可以得到基于物质扩散的用户 - 物品 - 标签的三部曲的推荐算法。

$$g''(Yk) = \lambda g'(Yk) + (1 - \lambda)g(Yk), \text{ 其中 } 0 \leq \lambda \leq 1$$

实验证明，在 λ 最优解的情况下，使用 *movielens*, *delicious* 两类不同的数据集，基于三部曲的推荐算法，精确度比基于二分网络的物质扩散方法算法大约提高了 9.86%，此外，优化的算法在推荐的多样性和新颖性两项指标上优于没有优化之前的三部曲推荐算法。

3.5.2 基于文本内容的优化算法

文献^[14]提出一种将文本内容的推荐同二分网络相结合的算法，在许多推荐系统中有大量的文本，例如豆瓣的推荐系统，就包含大量的文本内容，其中就包含着用户的偏好，将这些文本内容同二分网络推荐算法相结合可以优化推荐的精度。

T 表示所有物品的描述文本集， T_i 表示用户 X_i 购买的物品所形成的文本，然后用 TF-IDF 算法计算出用户对各个主题词的偏好权重，得出用户的主题偏好向量为 $p(X_i)$ ，物品 Y_k 的主题分布是 $p(Y_k)$ ，则用户 X_i 对物品 Y_k 的主题偏好权重为 $p(X_i)$ 和 $p(Y_k)$ 的余弦值，记为：

$$p(ik) = \cos(p(X_i), p(Y_k))$$

则形成了一种基于文本的推荐方法，该算法的模型为

$$g'(Yk) = g(Yk) * p(ik)^\lambda$$

其中， $g(Yk)$ 为任意一种基于网络的推荐模型所得出的物品 Yk 从用户 X_i 所得的资源值， $p(ik)$ 是基于 TF-IDF 为用户 X_i 对物品 Yk 的偏好值，引入可调参数 λ 调整用户偏好的权重。

采用偏热传导算法作为对比实验，实验结果表明加入用户偏好以后能够提高算法的精确度，同时推荐结果的个性化和多样性也有所提升。

3.6 基于社团划分的优化算法

在现实我们可以看到，许多连接关系有明显的社区区别，例如，在科学家合作网络中，同学科的科学家合作密度大于不同学科之间的合作密度，在社交关系中，这种社区就更加明显。基于社区之间更容易发生联系的假设，基于社区划分的二分网络推荐方法就被提出。其中文献 [15] 首次开始对二分网络社区结构进行研究，目前对二分网络社区进行划分大多是采用基于边结构进行划分。

当前，对于基于二分网络社团的推荐算法的划分有两种思路，一种思路是对二分网络进行社团划分后，在社区内部采用基于物质分配或热传导算法的推荐算法。另一种思路是对二分网络进行社团划分后，采用协同过滤算法的思想进行推荐。下面将介绍这两种思路。

文献^[15]采用第一种思路，用基于仿射传播聚类的二分重叠社区发现方法，对 *movielens* 数据集进行了社区划分，然后在社区内部采用基于二分网络的推荐算法，实验结果证明，采用社区的划分的推荐方法比不采用社区划分的推荐方法，推荐精确度提高 5.53%，排序准确度提高 2.57%，推荐结果多样性提高 5%，流行性降低 7%。

文献^[16]是采用第二种思路，用基于边传播社区划分算法对 *movielens* 数据集形成的二分网络进行社区划分，之后，可以得到用户 X_i 的社区近

邻, 设 X_j 是 X_i 的社区近邻, 用户 X_i 与 X_j 的社区近邻相关性为 $CR(i,j)$, 用户 X_j 对 Y_l 的评分为 $r(l,j)$, 则用户 X_i 对未评分物品 Y_l 的评分为

$$\text{score}(i,l) = \sum_{j=1}^n CR(i,j) * r(j,l)$$

实验结果: 该算法相对于其他的算法在提高命中率和降低流行性方面有着明显的优势。该实验还证明, 较高的重叠程度和明显的嵌套结构有助于从整体上改善推荐效果。社区划分的重叠程度越高, 节点的社区隶属度信息越丰富, 就可以更加准确刻画近邻群的规模和近邻之间的相关性。

3.7 考虑时间效应优化算法

时间因素在推荐系统也很重要, 主要体现在三个方面, 第一用户的兴趣是随着时间的变化而改变; 第二, 物品的种类和特性也会随着时间而改变; 第三, 季节效应, 例如人们夏天吃雪糕, 冬天吃火锅, 夏天穿裙子, 冬天穿羽绒服等; 第四, 时间效应可以使得某些联系失效, 例如在论文合作网络中, 80年以前的合作关系, 可能由于作者去世而失效。

对时间建模有很多种方法, 有的推荐系统倾向于推荐最近最热门的物品, 有的系统倾向于根据用户购物时间习惯来进行推荐, 本文就简单列举两种算法。

3.7.1 最近最热门算法

T 表示当前时间时刻, 根据物品 Y_i 的流行度确定其中推荐结果中的重要程度。^[18]

$$\text{weight}(Y_k) = \sum_{m=1}^n \frac{1}{1 + \lambda (T - t(km))}$$

其中, λ 为时间衰减参数, m 表示第物品 Y_k 第 m 次被购买, $t(km)$ 表示物品 Y_k 第 m 次的购买时间。

在用户 X_h 的推荐列表中加入时间因素, $f(X_h,i)$ 表示 $f(X_h)$ 中物品 Y_i 分配的资源值, 则最

终的物品 i 的资源值为

$$g'(Y_k) = \text{weight}(Y_k) * g(Y_k)$$

3.7.2 用户购买周期算法

T 表示当前时间时刻, 根据物品的购买周期, 来进行推荐, 例如对物品 Y_i 平均一个月购买一次, 对物品 Y_j 平均一周购买一次, 当前时间距用户上一次购买 Y_i, Y_j 的时间刚好为一周, 那么用户购买 Y_j 的可能性要远远大于用户购买 Y_i 的可能性。当然此算法的前提条件是用户购买物品一次以上, 需要长时间的用户数据积累, 适用于重复购物推荐。^[18]

设 T_k 表示用户购买物品 Y_i 的平均时间间隔, 即为产品 Y_i 的购买周期, Δt 表示用户上一次购买 Y_i 距当前的时间, 则

$$\text{weight}(Y_k) = \frac{1}{1 + \lambda |\Delta t - T_k|}, \lambda \text{ 为时间衰减参数}$$

在用户 X_i 的推荐物品 Y_j 时加入时间因素, 则最终物品 Y_k 的资源值为

$$g'(Y_k) = \text{weight}(Y_k) * g(Y_k)$$

当前对于时间因素对推荐系统的影响, 部分学者进行了探讨, 但进行的实验也较少, 未来还值得进一步的实证研究。

3.8 组合法

组合推荐的基本思想是通过组合来弥补各类算法的缺点, 发挥各类算法的优点, 达到扬长避短的效果, 目前基于二分网络推荐算法有多种组合法, 大致可以划分为如下3类:

1) 前融合: 直接融合各种推荐方法, 例如后面将要介绍的将热传导算法和物质扩散算法融合到同一个框架下面。

2) 中融合: 以一种推荐算法为框架, 融合另外一种推荐算法, 如上文的基于内容的推荐算法, 考虑时间因素的推荐算法等。

3) 后融合: 融合两种或两种以上的推荐方法

各自产生的推荐结果，如将协同推荐的结果和基于二分网络的推荐结果进行线性融合，融合的列表决定最终的推荐结果。^[18] 由于前面已经有中融合组合推荐，下面将重点介绍前融合组合推荐和后融合组合推荐。

3.8.1 前融合推荐算法

物质扩散和热传导结合的混合推荐算法，物质扩散算法和热传导算法各具特点，物质扩散算法推荐的准确性更高，适合用户找到自己感兴趣的产品，而热传导算法具有较强的冷门商品推荐能力，推荐结果的多样性较好，更容易为用户提供个性化的物品。因此可以考虑将两者进行融合以提高推荐效果和推荐质量。

因此，文献^[9]物质扩散算法和热传导算法的资源转移方程的相似性提出了将二者融合的算法，资源转移方程如下：

$$W(kj) = \frac{1}{d(Yk)^\lambda d(Yj)^{1-\lambda}} \sum_{l=1}^m \frac{b(k,l)b(j,l)^*}{d(Xl)} \quad (0 \leq \lambda \leq 1)$$

该算法和热传导和物质扩散算法只是资源转移方程不同，其他的相同。

文献^[18]通过实验证明在最优情况下，物质扩散和热传导的混合推荐算法的精度、召回率和多样性要高于偏热传导算法。

3.8.2 后融合推荐算法

基于二分网络的多维度推荐技术，其基本思想是，采用多种推荐算法分别计算出用户 X_i 对物品 Y_k 的资源分配量，然后采用线性组合方法求出最大评分作为最佳组合，^[19] 例如，假设有算法 1 计算出用户 X_i 对物品 Y_j 的资源分配量 $g_1(Y_k)$ ，算法 2 计算出用户 X_i 对物品 Y_k 的资源分配量 $g_2(Y_k)$ ，…… 算法 n 计算出用户 X_i 对物品 Y_k 的资源分配量 $g_n(Y_k)$ 。 $g(Y_k) = a_1 * g_1(Y_k) + a_2 * g_2(Y_k) + \dots + a_n * g_n(Y_k)$ ， a_i 表示各个算法的权重，通过对 a_i 进行调节，求出最优的 r ，作为最

优推荐结果，使各种算法进行扬长避短，发挥最大作用。^[20]

文献^[15]采用基于用户的资源扩散，基于项目的资源扩散和基于属性投影的三种方法，进行了组合，采用 `movielens` 数据集进行了实验，结果显示，这混合算法较协同推荐方法提高了 10%，较基于二分网络的资源分配方法提高了 1.5%，较非均匀的初始资源分配算法提高了 2.8%，多样性和流行性也有较大的提高。

以上各种方法，根据推荐系统所面临的实际情况而提出，除基于时间因素的优化算法未能找到相关实验数据外，其他的优化算法均有学者进行了实验，证明了优化算法相对于标准的热传导算法和物质扩散算法有着更好的推荐效果，基于二分网络的推荐算法没有传统的协同过滤推荐系统所面临的数据稀疏性问题，同时基于内容的二分网络推荐算法有利于解决推荐系统所面临的冷启动问题。

4 基于二分网络推荐算法有价值的研究方向

在第三节中介绍了多种基于二分网络推荐的优化和改进算法，这些算法大多通过实验检验，可以提高推荐效果，但是当前研究还是存在一些不足，这些不足主要体现在下面的三个方面，未来值得重点研究。

4.1 推荐系统的脆弱性问题

受推荐系统在电子商务领域重大的经济利益的驱动，一些心怀不轨的用户通过提供一些作弊的行为，故意增加或者压制某些商品被推荐的可能性^[22]。对推荐系统进行攻击两种常见的方法是虚假购买和恶意评价。基于二分网络算法的推荐系统中，某个物品的连边越多其被推荐的可能性

越大。因此在很多商家就可以利用这一点来展开作弊行为,例如在淘宝推荐系统中,有不少卖家进行信用炒作,进行虚假的购买行为,诱导推荐系统进行推荐来提高销量,谋求更多的商业利益。也有商家为打击竞争对手,对对手是产品进行较低评分和差评,减少商品推荐的可能性。攻击者还通过将攻击对象和热销商品或特定用户群喜欢的商品绑定而提高攻击效果,甚至通过持续探测推荐系统的推荐算法,从而有针对性地开展攻击。^[23]加强推荐算法在对恶意攻击的健壮性,成为需要认真考虑的一个因素。以基于关联规则的推荐算法为例,Apriori 算法的健壮性就远胜于k近邻算法^[24]。有一些技术已经被设计出来提高推荐系统面对恶意攻击的健壮性,例如通过分析对比真实用户和疑似恶意用户之间行为模式的差异,提前对恶意行为进行判断,从而降低疑似恶意用户的权重进行攻击预防^{[25][26][27]}。总体来说,这方面的研究相对较少,系统性的分析还很缺乏。^[23]长期下去,必将引发消费者对推荐系统的信任问题。因此,如何识别出这种攻击行为和作弊行为,这也是未来基于二分网络的推荐算法值得进一步研究。

4.2 如何对假冒伪劣等违法物品的过滤

在现实的推荐系统中,我们可以看到一些假冒伪劣商品,用于犯罪的物品和色情、恐怖和反社会等低级趣和不良味影像等由于有需求存在而得到了的购买,而在二分网络中形成了大量的此类购买关系,而按照上述的推荐算法。它们会被像正常物品一样推荐给用户,这种推荐行为有悖于法律法规和伦理道德,并且易降低用户信任度和满意度,因此推荐系统应该识别出此类商品,并降低权重甚至不予以推荐。当前的算法在这方面研究不足,未来需研究将商品质量的同二分网络的推荐算法相结合的研究。

4.3 用户行为模式对推荐算法影响研究

有大量研究表明,用户的行为模式对推荐结果有着明显的影响。例如文献^[28]和^[29]研究表明新用户和老用户具有很不一样的选择模式:新用户倾向于选择热门的商品,而老用户对于小众商品关注更多,新用户所选择的商品相似度更高,老用户所选择的商品多样性较高。还比如上文我们提到的有些用户倾向于选度大的产品(大众产品),有些用户喜欢选择孤僻产品(小众产品),用户行为的时空统计特性也可以用于提高推荐或者设计针对特定场景的应用^[30]。例如,淘宝系统发现在双十一左右,用户喜欢购买衣服,在炎热的夏天用户喜欢购买饮料,因此,推荐系统应能够识别出用户的这种兴趣随着时空变化的特征。从时间数据中还可以分析出影响用户选择的长期和短期的兴趣,通过将这两种效应分离出来,可以明显提高推荐的精确度^[31]。事实上,简单假设用户兴趣随时间按照指数递减,也能够得到改进的推荐效果^[32]。在移动互联时代,推荐系统也有大量的地理位置数据可用,基于位置信息的推荐可能会成为个性化推荐的一个研究热点和重要的应用场景,而这个问题的解决需要能够对用户的移动模式有深入理解,包括预测用户的移动轨迹和判断用户在当前位置的状态以及是否适合进行推荐等^[33],同时还要有定量的办法去定义用户之间以及地点和场景之间的相似性^[34]。文献^[35]和文献^[36]研究了基于地理位置的推荐,发现考虑地理位置可以明显提高广告推荐和朋友推荐的精确度。通过收集用户的行为数据,运用消费者行为分析和社会学及心理学等知识,预测出不同用户在不同环境、条件、阶段、时间、空间等等行为模式所反映出来的用户需求,推荐系统根据需求预测结果来进行推荐,并且针对不同的需求匹配或组合不同的算法结果,将其推荐给用户。这需要更高级的数据挖掘和算法设计能力,还需要有丰富经

验了解业务逻辑和用户行为的研究者配合完成。^[23]然而如何对复杂的用户行为模式建模以及同推荐二分网络推荐算法相结合研究较少但是又有重要意义,未来值得进一步研究。^[23]

以上三个方面在推荐系统的实际应用中有较强的需求和重要的应用价值值得未来重点研究。

5 大数据环境下基于二分网络的推荐算法研究展望

在大数据时代,用户面对的信息过载问题将更加严重,推荐系统也必将发挥出更大的作用,因此,传统的同一些推荐算法可能不适应大数据时代的要求,需要进行改进和拓展,生成更加精准、效率更高、用户更加满意的推荐算法。大数据给推荐算法带来的影响主要体现在时空两个方面,空间上表现为数据的维度更高、来源多、异质性大、更大的冗余和噪音,更多的隐式数据和新增的地理数据;在时间上表现为数据的更新速度快,推荐结果是实时性要求高。这要求推荐算法处理速度快,精确性高,实时性好。^[37]基于二分网络推荐算法未来还有很大的发展空间,未来可能会朝哪些方向发展,笔者对大数据环境下基于二分网络的推荐算法研究进行展望。

5.1 由基于二分网络的推荐算法向基于多主体异构网络的推荐算法方向发展

大数据的一个重要特征就是数据的多样性,多样性一方面是指数据的来源多,另一方面是指涉及的主体多。不同于单一网站的用户购买数据,在大数据中,这些数据包含众多网站上的数据(例如交易数据,医疗数据,社交数据,求职数据,旅行数据,学习工作数据等),这些数据涉及不同主体以及相同用户的不同角色,这些数据中各个主体之间相互作用的关系,关系将这些关系投影

到网络上,就形成了多主体的异构数据。^[38]从推荐算法的数据输入上来看,

将有多种联系输入,不仅是用户—物品的二维数据输入,将会有用户—用户、用户—物品、物品—物品多维度数据输入,同时连接的链条可能更长,向三部图,四部图,五部图以及更高部数发展,对应的数据输出也更加多元,不仅是推荐商品,可能还推荐好友,推荐医生,推荐新闻,推荐合作伙伴,推荐工作等。这需要推荐算法对多源的数据进行用户身份的标识整合到统一的网络中,以及在实现异质性的网络中进行资源分配,不论是在理论上还是在实践上都要重要意义,在未来多主体异构网络上进行推荐是未来的重要发展方向。在多源异质性网络上进行推荐的将有可能解决困扰推荐算法的研究人员的一个重要难题:推荐系统的冷启动问题。冷启动问题是指推荐系统在新增加的用户和物品中,由于推荐系统没有其交互数据和使用记录,因此无法进行推荐,长期以来该问题众多学者进行了研究,但是没有彻底解决此问题。在大数据环境下,当推荐系统从其他来源获取到大量的用户偏好数据时,利用大数据所拥有多异质性源点数据,进行数据挖掘来获取新用户的偏好和新产品的配置信息,用户可能在打开网页的时候,推荐系统就可能知道用户的购买偏好,从而进行推荐。^[39]

5.2 推荐结果时空变化同大数据处理速度之间的矛盾

大数据环境下,一方面数据生成速度快,数据维度高,价值密度低,多源数据融合带来了更多的噪声和冗余,使得传统的基于二分网络算法的推荐算法的时间和空间复杂度较高;另一方面,不同于传统的 top-N 推荐,在大数据环境下,推荐系统拥有用户的地理数据,不论是商家还是用户都期望的是推荐结果随时空变化,对算法的实

时性、新颖性和多样性的要求比较高。这二者之间存在矛盾并将更加尖锐。对于传统的协同推荐也面对这一问题,但是已有学者进行了大量研究,并提出了并行计算、近似计算、矩阵分解和增量计算等有效的解决方案。例如成熟的并行化计算框架 Hadoop^[40] 使集群计算成为可能,通过对传统算法进行并行化编程^[41],则可以有效地提高推荐系统的计算效率; Mayer-SchÖnberger 等^[42] 提出的近似计算方式将逐步成为大数据环境下的主要计算方法,矩阵分解算法有奇异值分解(SVD, singular value decomposition)^[43]、非负矩阵分解(NMF, non-negative matrix factorization)^[44]、概率矩阵分解(PMF, probabilistic matrix factorization)^[45]等这些算法的共同特点是将高维矩阵分解成为2个或多个低维矩阵的乘积形式便于在一个低维空间研究高维数据的性质^[46]。矩阵分解算法能降低高维数据稀疏性对噪声和冗余不敏感,可扩展性好受到大量关注。增量更新逐渐成为一种主要的数据更新方式,当出现新的用户或项目数据时,只对新加入的项目以及产生关联的边进行更新计算,并对原有推荐结果进行微调,而不是对全部数据进行更新计算。每隔一段时间,就用自适应方法消除局部计算引入的误差,令推荐结果的偏差不上升^[42]。对于大数据流式数据,研究者改进算法消除相关限制,令数据不需要一次性加载进入内存,只使用个性化的小缓冲,这种方法每次只采样新数据进行计算,降低了计算复杂度^[47]。可以吸收这些针对大数据的处理的算法和技术,来提高基于二分网络推荐算法处理大数据的能力和效率。

5.3 用户面临的隐私和安全风险提高

保护用户隐私是法律和道德对于推荐系统的要求,也是用户的严重关切,也是推荐系统发展的关键问题。大数据环境下,大规模用户数据

包含更多用户隐私和安全信息,推荐系统可以分析用户的人口统计学信息、行为数据、上下文信息等^[48],这些单一碎片化的信息可能不涉及用户隐私,但是汇总起来可能就涉及到用户隐私,还有的数据现在没有用户隐私和安全问题,但并不代表未来随着数据量的增加没有隐私风险,这增加了推荐系统识别用户隐私的难度和侵犯用户隐私的风险。如何通过分布式的方式取代传统的集中式数据获取,在不改变用户行为习惯,不危害用户隐私安全的前提下,^[49]充分利用用户大数据生成精准推荐项目,成为大数据环境下推荐系统用户隐私保护和安全性研究的热点和难点。^[50]

参考文献

- [1] Tao Z, Jie R, Mat ús M, et al. Bipartite network projection and personal recommendation.[J]. Physical Review E, 2007, 76(4):70-80.
- [2] Yi C Z, Marcel B, Yi K Y. Heat conduction process on community networks as a recommendation model.[J]. Physical Review Letters, 2008, 99(15):12505-12508.
- [3] Liu J G, Guo Q, Zhang Y C. Information filtering via weighted heat conduction algorithm[J]. Physica A Statistical Mechanics & Its Applications, 2011, 390(12):2414-2420.
- [4] Zhou Y, Liu W, Zhang J. The power of ground user in recommender systems.[J]. Plos One, 2013, 8(8):57-57.
- [5] 关远. 推荐网络分析及个性化推荐算法研究[D]. 成都:电子科技大学,2014:12-40.
- [6] Zeng A, Chi H Y, Shang M, et al. The reinforcing influence of recommendations on global diversification[J]. Epl, 2012, 97(1):18005.
- [7] 贾春晓. 基于复杂网络的推荐算法和合作行为研究[D]. 合肥中国科学技术大学,2011.
- [8] Zhou T, Jiang L L, Su R Q, et al. Effect of initial configuration on network-based

- recommendation[J]. Epl, 2007, 81(5):15-18.
- [9] 韩腾跃. 基于二分网络的个性化推荐系统研究 [D]. 南昌: 南昌航空大学, 2013.
- [10] 杜晗. 基于网络结构的推荐算法的研究 [D]. 北京: 北京邮电大学, 2013.
- [11] Liu J G, Zhou T, Che H A, et al. Effects of high-order correlations on personalized recommendations for bipartite networks[J]. Physica A Statistical Mechanics & Its Applications, 2010, 389(4):881-886.
- [12] 江山. 基于复杂网络理论的推荐算法研究 [D]. 成都: 西南财经大学, 2012.
- [13] Zhou T, Su R Q, Liu R R, et al. Accurate and diverse recommendations via eliminating redundant correlations[J]. New Journal of Physics, 2009, 11(12):4652-4657.
- [14] 张新猛, 蒋盛益, 张倩生, 等. 基于用户偏好加权的混合网络推荐算法 [J]. 山东大学学报: 理学版, 2015(9):29-35.
- [15] 熊湘云. 基于二分网络的多维度推荐技术研究 [D]. 苏州: 苏州大学, 2013.
- [16] 全佳妮. 基于二分网络的协同推荐研究 [D]. 苏州: 苏州大学, 2012.
- [17] 项亮. 推荐系统实战 [M]. 北京: 人民邮电出版社, 2012.
- [18] Koren Y. Collaborative filtering with temporal dynamics[J]. Communications of the Acm, 2010, 53(4):89-97.
- [19] 许海玲, 吴潇, 李晓东, 等. 互联网推荐系统比较研究 [J]. 软件学报, 2009, 20(2):350-362.
- [20] 段双艳. 基于网络结构的个性化推荐算法研究 [D]. 重庆: 重庆大学, 2013.
- [21] Zhao G, Lee M L, Hsu W, et al. Increasing temporal diversity with purchase intervals[C]// Proceedings of the 35th international ACM SIGIR conference on Research and development in information retrieval. New York: ACM, 2012:165-174.
- [22] Mobasher B, Burke R, Bhaumik R, et al. Toward trustworthy recommender systems: An analysis of attack models and algorithm robustness.[J]. Acm Transactions on Internet Technology, 2007, 7(4):23.
- [23] 周涛. 个性化推荐十大挑战 [EB/OL] [2015-12-4]. <http://blog.sciencenet.cn/blog-3075-588779.html?COLLCC=2227031322>.
- [24] Sandvig J J, Mobasher B, Burke R. Robustness of collaborative recommendation based on association rule mining[C]// ACM Conference on Recommender Systems. New York: ACM, 2007:105-112.
- [25] Lam S K, Dan F, Riedl J. Do You Trust Your Recommendations? An Exploration of Security and Privacy Issues in Recommender Systems. [C]// Proceedings of International Conference on Emerging Trends in Information and Communication Security, 2006:14-29.
- [26] Resnick P, Sami R. The influence limiter: provably manipulation-resistant recommender systems[C]// Acm Conference on Recommender Systems, 2007:25-32.
- [27] Yu H, Shi C, Kaminsky M, et al. DSybil: Optimal Sybil-Resistance for Recommendation Systems[C]// 2009 30th IEEE Symposium on Security and Privacy. IEEE Computer Society, 2009:283-298.
- [28] Zhang C J, Zeng A. Behavior patterns of online users and the effect on information filtering[J]. Physica A Statistical Mechanics & Its Applications, 2012, 391(4):1822-1830.
- [29] Shang M, Lu L, Zhang Y C, et al. Empirical analysis of web-based user-object bipartite networks[J]. Epl, 2009, 90(4):1303-1324.
- [30] 刘怡君, 周涛. 社会动力学 [M]. 北京: 科学出版社, 2012.
- [31] Min S H, Han I. Detection of the customer time-variant pattern for improving recommender systems[J]. Expert Systems with Applications, 2005, 28(2):189-199.
- [32] Koren Y, Collaborative filtering with temporal dynamics [C] // Proceedings of the 15th ACM SIGKDD International Conference on Knowledge

Discovery and Data Mining. New York: ACM Press, 2009:447-456.

[33]Cho E, Myers S A, Leskovec J. Friendship and mobility: user movement in location-based social networks[C]// Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining. New York: ACM, 2011:1082-1090.

[34]Zheng V W, Zheng Y, Xie X, et al. Collaborative location and activity recommendations with GPS history data[C]// International Conference on World Wide Web, 2010:1029-1038.

[35]Dao T H, Jeong S R, Ahn H. A novel recommendation model of location-based advertising: Context-Aware Collaborative Filtering using GA approach[J]. Expert Systems with Applications, 2012, 39(3):3731-3739.

[36]Scellato S, Noulas A, Mascolo C. Exploiting place features in link prediction on location-based social networks[C]// ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2011:1046-1054.

[37]孟祥武, 纪威宇, 张玉洁. 大数据环境下的推荐系统[J]. 北京邮电大学学报, 2015(2):1-15.

[38]孙远帅. 基于大数据的推荐算法研究[D]. 厦门: 厦门大学, 2014.

[39]Bauer J, Nanopoulos A. A framework for matrix factorization based on general distributions[C]// ACM Conference on Recommender Systems. New York: ACM, 2014:249-256.

[40]Ferreira Cordeiro R L, Traina Junior C, Machado Traina A J, et al. Clustering very large multi-dimensional datasets with MapReduce[C]// Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining. New York: ACM, 2011:690-698.

[41]Katkar V D, Kulkarni S V. A novel parallel implementation of Naive Bayesian classifier for Big Data[C]// International Conference on Green

Computing, Communication and Conservation of Energy, 2013:847-852.

[42]Naimi A I, Westreich D J. Big Data: A Revolution That Will Transform How We Live, Work, and Think.[J]. Information Communication & Society, 2013, 17(1):181-183.

[43]Golub G, Kahan W. Calculating the Singular Values and Pseudo-Inverse of a Matrix[J]. Siam Journal on Numerical Analysis, 1965, 2(2):205-224.

[44]Lee D D. Algorithms for Non-negative Matrix Factorization[J]. Advances in Neural Information Processing Systems, 2001, 13(6):556--562.

[45]Hao Ma, Haixuan Yang, Michael R. Lyu, et al. SoRec: Social Recommendation Using Probabilistic Matrix Factorization[C]// Proceedings of the 17th ACM Conference on Information and Knowledge Management, 2008. 2010:931-940.

[46]何清, 李宁, 罗文娟, 等. 大数据环境下的机器学习算法综述[J]. 模式识别与人工智能, 2014, 27(4):327-336.

[47]Diaz-Aviles E, Drumond L, Schmidt-Thieme L, et al. Real-time top-n recommendation in social streams[C]// ACM Conference on Recommender Systems. New York: ACM, 2012:59-66.

[48]Bhagat S, Weinsberg U, Ioannidis S, et al. Recommending with an Agenda: Active Learning of Private Attributes using Matrix Factorization[C]// Proceedings of the 8th ACM Conference on Recommender systems. New York: ACM, 2013:65-72.

[49]Becchetti L, Bergamini L, Colesanti U M, et al. A lightweight privacy preserving SMS-based recommendation system for mobile users[J]. Knowledge & Information Systems, 2010, 40(1):191-198.

[50]Zheng H. A Survey of Trajectory Privacy-Preserving Techniques[J]. Chinese Journal of Computers, 2011, 34:1820-1830.