

基于信任网络随机游走模型的协同过滤推荐

何 明 刘伟世 魏 铮

(北京工业大学计算机学院 北京 100124)

摘 要 协同过滤是目前应用最广泛和最成功的推荐技术之一。然而,目前该技术的发展面临着严重的冷启动和稀疏性问题,降低了其推荐质量,因此提出了一种基于信任网络随机游走模型的协同过滤推荐方法。该方法融合了基于信任和项目的协同过滤推荐方法,并引入了信任因子作为引导推荐的重要因素。随机游走模型不仅考虑了信任用户对目标项目的评分,也考虑了他们对与目标项目相似的项目的评分。随着随机游走深度的增加,以相似项目的评分信息来代替目标项目的评分信息的概率也逐渐增大。在 Epinions 真实数据集上的验证结果表明,该方法在推荐评价指标上比其他算法具有更好的推荐结果。

关键词 协同过滤,推荐系统,随机游走,信任网络

中图法分类号 TP391 **文献标识码** A **DOI** 10.11896/j.issn.1002-137X.2016.6.051

Collaborative Filtering Recommendation Based on Random Walk Model in Trust Network

HE Ming LIU Wei-shi WEI Zheng

(College of Computer Science, Beijing University of Technology, Beijing 100124, China)

Abstract Collaborative filtering is one of the most widely used techniques for recommendation system which has been successfully applied in many applications. However, it suffers from serious problems of cold start and data sparsity. In addition, these methods can not indicate their confidence in recommendation. In this paper, we improved the random walk model combining trust-based and item-based collaborative filtering method for recommendation. The trust factor is introduced as an important factor of guiding recommendations. The random walk model considers not only the ratings of target item, but also those of the similar items. The probability of using the rating of a the similar item instead of a rating for the target item increases with increasing length of walk. Our framework contains both trust-based and item-based collaborative filtering recommendations as special cases. The empirical analysis on the Epinions dataset demonstrates that our method can provide better recommendation result in terms of evaluation metrics than other algorithms.

Keywords Collaborative filtering, Recommender system, Random walk, Trust network

1 引言

近年来,随着 Facebook、Twitter、新浪微博等社会化网络的日益发展,利用用户间的社会关系进行推荐的方法由于能更好地模拟现实社会中的推荐过程,更能体现出人在推荐过程中的作用,逐渐成为推荐领域的研究热点^[1]。在社会化网络中,人们更倾向于依赖他们信任的朋友的推荐,信任关系是目前社会化推荐系统中应用最广的一种社会网络关系^[2]。

协同过滤推荐^[3]是目前应用最广泛和最成功的推荐技术,其核心思想是根据用户对项目的评分或者偏好程度发现用户之间或项目之间的相关性,然后根据相关性进行推荐。然而,随着 Web2.0 的迅速发展和广泛应用,传统的基于用户或者基于项目的协同过滤推荐算法中存在数据稀疏性、冷启动等问题,这直接导致了其推荐质量的大大降低。为了解决这些问题,一些研究人员将社会化信息引入到推荐中,利用了

社会化网络用户聚类^[4,5]和社会化矩阵分解技术(Matrix Factorization, MF)^[6]。基于聚类的推荐方法的训练过程相对比较耗时,而矩阵分解技术对推荐结果的解释性较差。目前,在已有的社会化网络的典型信任模型和推理算法中,Advagato 算法^[7]和 Appleseed 算法^[8]把用户分成入门、进阶和高阶 3 种不同类型,每种类型间信任关系的强弱都不相同,前者采用网络流模型来计算用户的信任值,后者采用激活扩散机制来计算信任值。TidalTrust 算法^[9]和 MoleTrust 算法^[10]都是基于宽度优先搜索顺序迭代计算源用户和目标用户之间的信任值。Massa 和 Avesani 的研究显示通过在推荐系统中引入用户信任网络能够有效缓解数据稀疏性问题和冷启动问题^[11]。Jamali 和 Ester^[12]提出了一种随机游走模型,将基于信任的和基于项目的协同过滤方法相结合,在一定程度上缓解了推荐系统中的数据稀疏性和冷启动问题,同时也具有较高的推荐质量。

到稿日期:2015-05-26 返修日期:2015-10-08 本文受国家自然科学基金项目(60803086),国家科技支撑计划子课题(2013BAH21B02-01),北京市自然科学基金项目(4153058,4113076)资助。

何 明(1975—),男,博士,副教授,主要研究方向为推荐系统、数据挖掘、机器学习, E-mail: heming@bjut.edu.cn; 刘伟世(1989—),男,硕士生,主要研究方向为推荐系统、数据挖掘; 魏 铮(1990—),男,硕士生,主要研究方向为推荐系统、信息检索。

以上研究表明,在推荐系统的用户中建立社会化信任网络将有助于提高推荐性能。我们注意到,文献[12]中的随机游走模型假设从一个节点到另一个直接信任节点的跳转是等概率的,而我们认为在随机游走过程中对下一个节点的选择应该与当前节点的相似度和信任值有关。对此,基于文献[12],本文提出了一种基于信任因子随机游走模型的协同过滤推荐方法。我们认为那些与当前节点相比具有更高相似度、信任值高的节点在随机游走过程中应该以较大的跳转概率被选择,而不是与其他节点一样具有相同的跳转概率。

本文将协同过滤算法与社会化信任网络相结合,旨在建立一种基于信任关系的推荐系统,通过引入信任因素来缓解和解决传统协同过滤推荐系统中的数据稀疏性和冷启动问题所带来的负面影响。与此同时,在基于信任关系的推荐方法中由于用户对未知项目的预测评分来自于直接信任或间接信任的用户,因此可以提高推荐质量。本文的贡献主要包括以下4个方面:1)信任网络随机游走过程中选择下一跳用户节点的选择方式采用相似性与信任因子结合的方式。将Page-Rank思想与Sigmoid函数结合,对每个下游用户产生选择权重进行非等概率选择。一个用户被信任次数越多并且与初始目标用户相似性越大,则在游走过程中被选中的可能性也就越大。2)对游走轮数给出上下限的限定,以防止游走轮数过多而耗时和游走轮数过少对预测结果产生偏差。3)游走过程中记录的每次游走结果产生的总概率为 Pr_n 。当预测一个用户对某个项目的评分时,会进行一轮游走,每轮游走中包含了多次游走,而每一次游走对应一个概率值。概率值表示了每个节点的停止(继续)概率、下一跳用户选择概率以及最终的项目选择概率。4)最终每轮游走的结果会将该轮中多次结果进行加权计算,即根据每次游走的概率计算权重。通过实验观察,每轮的结果的总概率 Pr_n 的大小差异很大,如果仅仅使用 Pr_n 作为每轮结果权重会导致最终的预测评分产生严重的偏差,最终预测分值的大小将会与权重最大的几个评分极其相似。实验结果显示,本文采用 $-1/(\log(Pr_n))$ 作为每次游走结果权重,既考虑了概率大小对结果的影响,又避免了概率大小差异过大导致预测结果的严重偏差。

2 相关方法和模型

2.1 相似性度量方法

协同过滤推荐算法的核心是发掘用户之间和项目之间的相关性。所以需要采用一种合理的方法对相似性进行度量,相似性度量方法的合理与否也是整个协同过滤推荐是否成功的关键。通常做法是将用户评分数据用一个 $m \times n$ 阶矩阵 $R(m, n)$ 表示。 m 行表示 m 个用户, n 列表示 n 个项目, $r_{i,j}$ 代表用户 i 对项目 j 的评分。用户评分数据矩阵如表1所列。

表1 用户评分数据矩阵 $R(m, n)$

	Item ₁	...	Item _j	...	Item _n
User ₁	$r_{1,1}$...	$r_{1,j}$...	$r_{1,n}$
...
User _i	$r_{i,1}$...	$r_{i,j}$...	$r_{i,n}$
...
User _m	$r_{m,1}$...	$r_{m,j}$...	$r_{m,n}$

目前应用最为广泛的几种相似度计算公式有 Cosine、Pearson Correlation^[13]和 Constrained Pearson Correlation。本文选取 Pearson Correlation 方法作为基础的相似度计算方

法。项目 i 和项目 j 之间的皮尔逊相似度为 $corr_{i,j}$,则有:

$$corr_{i,j} = \frac{\sum_{u \in UC_{i,j}} (r_{u,i} - \bar{r}_i)(r_{u,j} - \bar{r}_j)}{\sqrt{\sum_{u \in UC_{i,j}} (r_{u,i} - \bar{r}_i)^2} \sqrt{\sum_{u \in UC_{i,j}} (r_{u,j} - \bar{r}_j)^2}} \quad (1)$$

其中, $UC_{i,j}$ 代表对项目 i 和项目 j 的共同评价用户集合, \bar{r}_i 与 \bar{r}_j 为项目 i 和项目 j 在 $UC_{i,j}$ 范围内评分的平均值。

皮尔逊相似度只考虑了项目之间评分数值的差异性,没有考虑项目之间共同评分用户数量对相似度的影响。例如:有3个项目 i, j 和 k ,项目 i 和项目 j 的皮尔逊相似度为0.70,它们之间有100个共同评分用户;项目 i 和项目 k 的皮尔逊相似度也同样等于0.70,但它们之间仅有10个共同评分用户。如果仅从皮尔逊相似度的数值上来看, $corr_{i,j} = corr_{i,k}$,但 i 和 j 之间却有着比 i 和 k 之间更多的共同评分用户。那么此时可以认为项目 i 和 j 之间的相似度是强于项目 i 和 k 之间的,从量化的角度来说应该有 $corr_{i,j} > corr_{i,k}$ 。而传统的皮尔逊相似度恰恰忽略了共同评分数对相似性的影响。所以引入共同评分数 $|UC_{i,j}|$,将其作为影响相似度的一个要素。项目 i 和 j 之间最终相似度为 $sim_{i,j}$:

$$sim_{i,j} = 1 / (1 + e^{-|UC_{i,j}|/2}) \times corr_{i,j} \quad (2)$$

其中,将Sigmoid函数^[14]与皮尔逊相似度相结合作为最终的相似度度量方式。其中Sigmoid因子 $1/(1 + e^{-|UC_{i,j}|/2})$ 取值范围为 $[0.5, 1)$,当共同评分数 $|UC_{i,j}|$ 足够大时Sigmoid因子将趋近于1。上文中讨论的均是以项目间相似度为例,用户相似度计算方式与项目类似,此处不再赘述。将式(2)中的相似度计算方式称为增强的皮尔逊相似性,在下文中如果没有特殊声明,所有的相似度公式均指此公式。

2.2 信任模型

在不同学科中对信任的概念都有着不同的定义,通常大多数是用来形容一个实体对另一个实体的能力、意图及可靠性的认可。Farag Azzedin曾对网络中的信任做出如下定义^[15]:“信任是一个实体对另一个实体将要做的事所需能力的坚定信念,并且这种信念是随实体行为以及时间的变化而变化的”。随着近年来社会化网络的兴起,关于社会化网络中信任关系的研究也越发地受到有关学者们的重视。本文正是根据社会化网络中信任的特性构建了信任模型。

2.2.1 信任特性

社会化网络中的信任关系有如下几个特性:

1)可度量性

信任是可度量的,可以用一个值来表示信任关系的紧密程度。信任的度量方式是信任模型建立的基础。通常可以用离散的布尔型变量或一段连续的取值区间作为信任的度量方式,本文采用布尔型变量来度量信任关系。

2)非对称性

在社会化网络中由于信任关系具有一定的主观性,因此两个对象之间的信任关系并不是对称的。例如:已知对象A信任对象B,并不能推断出对象B也信任对象A。根据此特性,本文用非对称有向图来表示对象之间的信任关系。

3)弱传递性

信任在对象之间存在一定的传递性,但信任关系的传递并不是绝对的,而是在一定条件下才能成立的。例如:对象A信任对象B,对象B又信任对象C,但对象A只有在一定条件

下才会信任对象 C。根据信任的弱传递性,本文将信任分为直接信任和间接信任两大类。直接信任是根据已有数据直接给出的,间接信任是由一定规则通过信任的传递性推算出来的。

4)多样性

信任关系的多样性指信任关系存在形式具有多样性。信任在对象之间可以是一对一的,也可以是多对一或一对多的。

2.2.2 信任网络

本文所述的推荐模型中用户之间的信任关系可以由一个用户信任评分矩阵表示。矩阵中的值 $b_{i,j}$ 为布尔型变量,当 $b_{i,j}$ 为 1 时表示 $User_i$ 信任 $User_j$,当 $b_{i,j}$ 为 0 时表示不信任。如表 2 所列。

表 2 用户信任评分矩阵 $B(m,m)$

	User ₁	...	User _j	...	User _m
User ₁	$b_{1,1}$...	$b_{1,j}$...	$b_{1,m}$
...
User _i	$b_{i,1}$...	$b_{i,j}$...	$b_{i,m}$
...
User _m	$b_{m,1}$...	$b_{m,j}$...	$b_{m,m}$

用户信任评分矩阵可以形象地表示成用户信任网络,如图 1 所示。

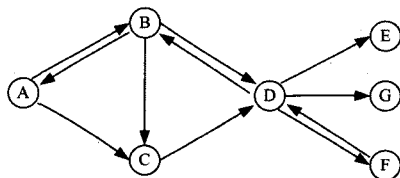


图 1 用户信任网络

在图 1 中,每个节点分别代表一个用户,每条边代表用户之间存在信任关系。由于信任关系是非对称的,因此信任网络中的边既有单向的也有双向的。每个用户节点可以有多个信任用户也可以被多个用户信任,体现了信任关系的多样性。

3 基于信任因子随机游走的随机模型

传统的协同过滤算法只考虑了用户对项目的历史评分数据,所以当出现稀疏数据或冷启动问题时推荐精度就会明显下降。本文引入信任网络中的随机游走机制来解决数据稀疏和冷启动问题。

3.1 单轮信任网络随机游走

假设模型中已存在用户集 $U = \{u_1, u_2, \dots, u_m\}$ 和项目集 $I = \{i_1, i_2, \dots, i_n\}$ 。用户 u 对项目 i 的评分为 $r_{u,i}$,本文中所有评分均为 1 到 5 之间的整数。如果用户 u 信任用户 v ,那么用 $t_{u,v}$ 表示用户 u 对 v 的信任值。由于本文中采用布尔型的信任度量方式,因此 $t_{u,v}$ 的取值用 0 或 1 表示。其中 0 表示不信任,1 表示信任。

另外定义 $TU_u = \{v \in U | t_{u,v} = 1\}$ 表示被用户 u 直接信任的用户集合,用图 $G = \langle U, TU \rangle$ 表示整个信任网络,图中有向边的集合定义为 $TU = \{(u,v) | u \in U, v \in TU_u\}$ 。将用户对项目的评分映射到信任网络后,如图 2 所示。图 2 中每个人物图标代表一个用户,每个用户头像上的图形代表该用户已评分过的不同项目,图形下边的数字代表具体评分。

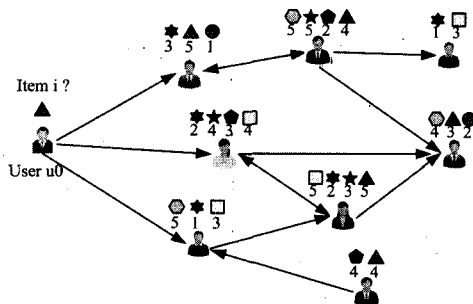


图 2 带有评分信息的用户信任网络

3.1.1 游走方式

如果预测一个用户 u_0 对项目 i 的评分为 $\bar{r}_{u_0,i}$,则会以用户 u_0 为根节点在信任网络中进行游走。当游走到信任网络中的某个用户节点 u 时,如果用户 u 对项目 i 已有评分,则停止随机游走并且返回项目评分 $r_{u,i}$ 作为随机游走结果;如果用户 u 对项目 i 没有评分,则有两种选择:

1) 以 $\alpha_{u,i,k}$ 的概率停止继续游走,并从当前用户 u 已评分的项目中选择一个与待评分项目 i 相似的项目 j ,将该用户的评分 $r_{u,j}$ 作为一轮的游走结果进行返回。从用户 u 已评分项目选择项目 j 的概率为 $Pr(\sigma_{u,i} = j) = sim_{i,j} / \sum_{h \in TU_u} sim_{i,h}$,其中 $\sigma_{u,i}$ 表示用户 u_0 预测项目 i 的评分停止游走时选择的项目。 $\alpha_{u,i,k}$ 为停止游走概率,将在后面章节专门讨论。

2) 另外以 $1 - \alpha_{u,i,k}$ 的概率从当前用户 u 的直接信任用户集合 TU_u 中选择一个用户作为下一跳用户 $v (v \in TU_u)$ 继续在信任网络中游走。选择下一跳用户节点之前,首先要计算 TU_u 中所有用户在这一次选择中被选中概率的权重 $\chi_{u,v}$ 。

$$\chi_{u,v} = (0.5 + 0.5 * sim(u_0, v)) * (1 / (1 + e^{-IN_v})) \quad (3)$$

$$IN_v = \sum_{w \in U \wedge w \neq v} t_{w,v} \quad (4)$$

根据 Google PageRank^[16] 的思想,在式(3)中引入了用户 v 的信任因子 IN_v , IN_v 表示用户集 U 中信任用户 v 的用户的总数。另外本文假设被选择用户的信任值越大并且与目标用户的相似性越高则在随机游走时被选中的概率就应该越高。所以将信任因子 IN_v 、用户相似度 $sim(u_0, v)$ 与 Sigmoid 函数相结合,使被选择用户 v 的概率权重 $\chi_{u,v}$ 与 IN_v 和 $sim(u_0, v)$ 均成正比。在计算完 TU_u 中每个用户的权重后,将所得权重归一化处理便可得出每个用户被选为下一跳用户的概率 $Pr(\chi_{u,v})$,具体方法如下:

$$Pr(\chi_{u,v}) = \chi_{u,v} / \sum_{z \in TU_u} \chi_{u,z} \quad (5)$$

3.1.2 一轮游走结束条件

当游走到每一个用户节点 u 时都会有一个对应的停止游走概率 $\alpha_{u,i,k}$,其中 u, i, k 分别表示当前游走节点为用户 u 、目标项目为 i 和游走深度为 k 。记用户节点 u 已评论的项目集合为 IU_u 。本文假设目标项目 i 与 IU_u 中项目的最大相似度越大并且游走步数 k 越大则在当前节点停止游走的概率 $\alpha_{u,i,k}$ 越高。 $\alpha_{u,i,k}$ 具体计算方式如下:

$$\alpha_{u,i,k} = \max_{j \in IU_u} sim_{i,j} \times 1 / (1 + e^{-k/2}) \quad (6)$$

至此,可以总结每轮游走停止的条件如下:

- 1) 游走过程中的某个节点对待预测目标项目 i 已做出直接评分并将该评分直接返回。
- 2) 当游走到节点 u 时, u 对目标项目 i 没有评分,但以概率 $\alpha_{u,i,k}$ 停止游走并返回一个类似项目的评分。

3)如果对游走深度不做任何限制,那么游走过程可能会存在一直游走的可能性。基于“六度分割”理论^[17],为了避免这种情况,本文规定游走深度最大为6步。当游走的深度为6时,无论该节点是否对项目*i*评分,都会停止游走并返回评分。

3.2 全局游走结束条件

为了得到一个可靠的预测评分,将会进行多轮随机游走,并根据多轮游走结果预测具体评分 $\bar{r}_{u_0,i}$ 。据此,需要对全局游走规定一个结束条件。当满足条件时全局游走结束,计算最终预测评分。

假设预测用户 u_0 对项目*i*的评分进行第*k*次随机游走返回的结果为 r_k ,共进行了*T*趟随机游走,*T*趟随机游走返回结果的平均值为 \bar{r} 。设 ϵ^2 为每轮游走结果的方差,则有

$$\epsilon^2 = \sum_{k=1}^T (r_k - \bar{r})^2 / T \quad (7)$$

定义 ϵ_i^2 为第*i*趟游走计算得到的方差。项目评分区间为[1,5],可证明经过若干轮游走后方差 ϵ^2 趋于收敛^[12],即 $|\epsilon_{i+1}^2 - \epsilon_i^2| \leq \zeta$ 。 ζ 为预先设定的一个极小的正数阈值,例如0.001。当方差收敛时即可结束游走,根据各轮结果计算预测分值 $\bar{r}_{u_0,i}$ 。

另外,对全局游走的总轮数*T*要规定上限和下限。当游走轮数小于下限时,无论方差是否收敛都会继续下轮游走。比如第一轮和第二轮游走返回结果都为5,此时方差为0小于阈值,但此结果并不可靠,需要得到更多轮游走结果来进行最终评分计算。当游走轮数等于上限时,无论方差是否收敛都会停止游走。因为当随机游走进行了一定轮数后,虽然方差可能并未收敛,但其结果集已足够用来计算预测分值,过多的游走轮数并没有必要。本文规定在信任网络中随机游走总轮数*T*的取值区间为[50,10000]。

3.3 推荐结果的生成

每轮游走结束后都会选择一个用户*v*对项目*j*的评分 $r_{v,j}$ 作为单轮随机游走的返回结果,其中项目*j*可能与目标项目*i*是相同的,也可能不同。当全局游走结束后,可以得到一个由每轮返回评分组成的结果集 $R_{u_0,i} = \{r_{v,j} | v \in U, j \in I\}$ 。假设总共进行*N*轮游走,那么最终的预测评分计算方法为:

$$\bar{r}_{u_0,i} = \sum_{n=1}^N W_n \cdot r_n \quad (8)$$

其中, W_n 为第*n*轮游走结果的权重, r_n 为第*n*轮游走的返回结果。本文假设权重 W_n 与各轮游走轨迹的总概率 Pr_n 相关,概率越高则权重越大,反之则越小。根据第*n*轮返回结果 r_n 的游走轨迹,可以计算从信任网络中初始节点 u_0 到最终节点 u_v 返回项目*j*的总概率 Pr_n :

$$Pr_n = (1 - \alpha_{u_0,u_1}) Pr(\chi_{u_0,u_1}) (1 - \alpha_{u_1,u_2}) Pr(\chi_{u_1,u_2}) \cdots \alpha_{u_v,i,k} Pr(\sigma_{v,i} = j) \quad (9)$$

但是在真实测试中发现,每轮的结果的总概率 Pr_n 的大小差异很大,如果仅仅使用 Pr_n 作为每轮结果权重,会导致最终的预测评分产生严重的偏移, $\bar{r}_{u_0,i}$ 的大小将会与权重最大的几个评分极其相似。因此此处使用 $-1/\log(Pr_n)$ 替代最初的 Pr_n 。将其进行归一化后得到每轮游走的最终权重:

$$W_n = (1/\log(Pr_n)) / (\sum_{i=1}^N \log(Pr_i)) \quad (10)$$

4 实验

为了验证提出的方法确实对推荐结果做出了优化,引入传统的基于协同过滤的推荐算法作为对比方法。另外实验中将测试随机游走推荐中不同的参数(游走步数的长短、方差收敛阈值)对推荐结果的不同影响。

本文中提到的所有方法均用Java语言实现。实验所用PC机配置为Intel core 2 duo,2.66GHz,4GB内存,操作系统为Windows 7。

4.1 数据集

实验中使用的数据集来自Epinions(<http://alchemy.cs.washington.edu/data/epinions>),其中包含了本文所需要的用户项目评分数据以及用户间的直接信任数据,其相对于其他一些未包含信任信息的数据集来说更适合本实验。

本文所采用实验数据中的评分数据极其稀疏。数据共包含了664824条评分数据,来自于40136名用户对139738个项目的评分,评分矩阵稀疏率接近于0.01%;此外,数据中还包括了442979条信任数据。我们定义对项目评分个数少于5的用户为冷启动用户,那么在40136名用户中,冷启动用户的比例占到了47%。从原始数据中提取出两个测试集分别为DS1和DS2。DS1中为从原数据中随机抽取的数据,DS2中为从原数据中抽取的冷启动用户数据。测试集数据量约为总数据的20%。原始数据摘要如表3所列。

表3 原始数据集

用户	项目	评分数据	信任数据
40136	139738	664824	442979

4.2 推荐结果评价标准

采用了较为常用的RMSE、MAE和F1-measure 3种指标去评价推荐结果的好坏。RMSE和MAE是评价预测评分的准确性标准,它们反映的是算法的预测评分与用户实际评分的近似度。其中,MAE的定义为 $MAE = \sum_{i,j} |r_{i,j} - \hat{r}_{i,j}| / T$, $r_{i,j}$ 表示用户 u_i 对项目*j*的实际评分, $\hat{r}_{i,j}$ 表示相应的预测评分,*T*表示测试样例的数量。与MAE的定义类似,RMSE的公式可以表示为 $RMSE = \sqrt{\sum_{i,j} (r_{i,j} - \hat{r}_{i,j})^2} / T$,并且MAE和RMSE的值越小,系统的推荐效果越好。

F1-measure是准确率precision和召回率recall的调和平均值。该值越高则表明推荐算法的综合性能越好。F1-measure定义为 $F1\text{-measure} = (2 \times \text{precision} \times \text{recall}) / (\text{precision} + \text{recall})$ 。

4.3 实验设计

为了证明基于信任因子随机游走的优越性,将不同参数下的随机游走与传统推荐算法在测试集DS1和DS2上进行实验对比。实验细节设计如下:

1)对比在信任网络中不同的游走深度对精确度的影响。首先将游走深度固定为1步,其次将游走深度固定为6步,最后按本文所介绍的方法进行随机概率游走,比较这3种方法的MAE值,分析不同游走深度对结果的影响。

2)找出最优的方差收敛阈值。将 ζ 分别设定为0.01,0.005,0.001,0.0005,并且以 ζ 为收敛阈值进行随机游走。观察不同的 ζ 对推荐精度的影响,找出最恰当的值。

3)将本文中提出的推荐方法与其他推荐算法进行比较。在两个测试集 DS1 和 DS2 上分别进行实验测试,观察各种算法在普通测试集和冷启动测试集中评价指标的表现。

4)判断相似性计算公式对推荐结果的影响。在 DS1 和 DS2 两个测试集中测试传统的皮尔逊相似度和改进的皮尔逊相似度对结果的影响。

5)测试游走过程中选择下一跳用户时概率权重 $\chi_{u,v}$ 对推荐结果的影响。除了本文中提到的根据权重概率 $Pr(\chi_{u,v})$ 选择的方法,还可以按照 $1/|TU_u|$ 等概率的方式选择。对比采用这两种方法所得的推荐结果的精度。

4.4 实验比较

实验比较中所参照的算法描述如下。

- User based:基于用户的传统协同过滤推荐算法。
- Item based:基于项目的传统协同过滤推荐算法。
- TrustWalker1:在信任网络中随机游走深度为 1 的基于信任因子随机游走模型的协同过滤推荐。
- TrustWalker6:在信任网络中随机游走深度为 6 的基于信任因子随机游走模型的协同过滤推荐。
- TrustWalker:文献[12]中基于随机游走模型的推荐算法。
- RTW:本文中提出的完整的基于信任因子随机游走模型的协同过滤推荐。

首先,比较 TrustWalker1, TrustWalker6 和 RTW 在 DS1 和 DS2 上的 MAE 值。从图 3 中可以看出,RTW 推荐结果的 MAE 要低于另外两种随机游走方法,这种优势在全部由冷启动用户组成的数据集 DS2 上表现得更加明显。在 DS1 上 TrustWalker1, TrustWalker6 和 RTW 3 种方法得出的 MAE 值分别为 0.879, 0.921 和 0.865, 在 DS2 上的 MAE 值为 0.998, 1.013 和 0.971。由此可以得出,改进后的游走深度控制方法对提升推荐的预测精度有明显效果。

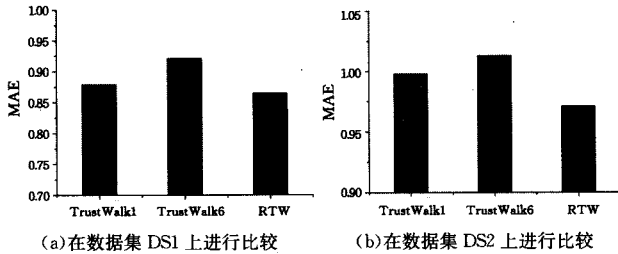


图 3 在两个数据集上比较不同的游走停止条件对 MAE 的影响 (方差收敛阈值 ζ 为 0.001)

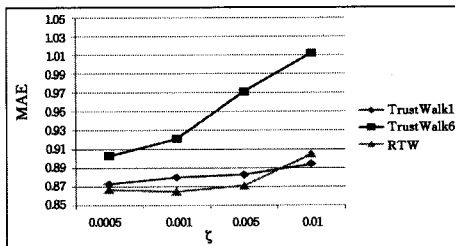


图 4 在数据集 DS1 上比较不同的方差收敛阈值 ζ 对推荐结果的影响

从图 4 中可以很清楚地看到收敛阈值 ζ 对推荐结果的影响。当 $\zeta > 0.001$ 时, MAE 随着 ζ 的减小而明显减小。但当 ζ 的值继续减小时, RTW 和 TrustWalker1 的 MAE 值并没有明显地下降, 其中 RTW 的 MAE 值反而有所上升。另外随着 ζ 下降, 为了得到最终的收敛结果, 在执行期间需要更多轮的

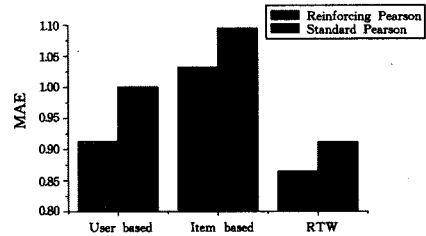
游走,这也导致了时间开销的增大。所以将 ζ 的最终值定为 0.001。

接着,将方差收敛阈值置为 0.001,在 DS1 和 DS2 两个数据集中对比所有方法的 MAE 值,结果如表 4 所列。在 DS1 中 User based 和 TrustWalker1 的 MAE 值与 RTW 比较接近,但在 DS2 上可以看出 RTW 的推荐精度要远远高于其他算法。

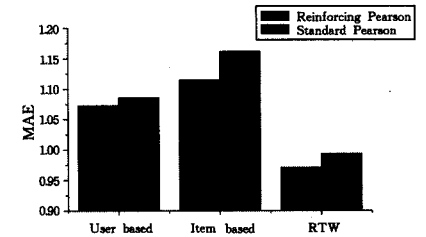
表 4 推荐在 DS1 和 DS2 上的 MAE 值比较(方差收敛阈值为 0.001)

Methods	DS1	DS2
TrustWalk1	0.879	0.998
TrustWalk6	0.921	1.013
RTW	0.865	0.971
User based	0.913	1.073
Item based	1.032	1.115

图 5 中对比的是传统的皮尔逊相似度与改进的皮尔逊相似度两种相似度计算方式对推荐的影响。从图中可以看出,改进的皮尔逊相似度确实使得推荐精度有所提高。但是由于 DS2 中的所有用户均由冷启动用户组成,用户之间的共同评分相对于 DS1 要少很多,因此改进的皮尔逊相似度对推荐的提升不如 DS1 中那么明显。



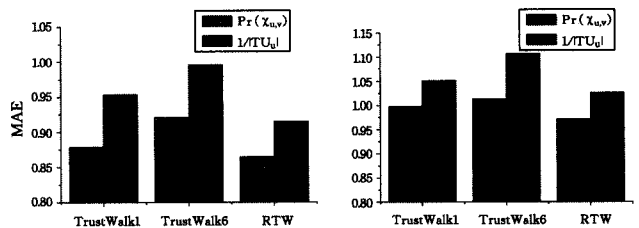
(a)在数据集 DS1 上进行比较



(b)在数据集 DS2 上进行比较

图 5 传统的皮尔逊相似度与改进后的相似度对推荐的影响

图 6 示出了采用可信权重概率 $Pr(\chi_{u,v})$ 选择和等概率选择两种选择方式对推荐结果的影响。可以看出在两个数据集中基于 $Pr(\chi_{u,v})$ 的方法的 MAE 值均小于等概率选择的方法。实验的结果证明了我们的假设:被选择用户的信任值越大并且与目标用户的相似性越高,那么在随机游走时被选中的概率就应该越高。



(a)在数据集 DS1 上进行比较

(b)在数据集 DS2 上进行比较

图 6 下一跳节点选择方式对推荐结果的影响

最后,在实验中比较 RTW 与社会化网络中基于信任推理的 TidalTrust 算法^[9]和 MoleTrust 算法^[10]、基于随机游走模型的 TrustWalker^[12]、基于用户的协同过滤推荐 User based 和基于项目的协同过滤推荐 Item based 分别在 RMSE 和 F_1 -measure 两个评价指数上的表现,实验结果如表 5 所列。

表 5 推荐在 DS1 和 DS2 上的 RMSE 和 F_1 -measure 值比较(方差收敛阈值 ζ 为 0.001)

Methods	DS1		DS2	
	RMSE	F_1 -measure	RMSE	F_1 -measure
TidalTrust	1.041	0.822	1.155	0.683
MoleTrust	1.038	0.801	1.353	0.645
TrustWalker	1.023	0.868	1.132	0.758
User based	1.261	0.721	1.420	0.293
Item based	1.182	0.733	1.443	0.355
RTW	1.012	0.887	1.123	0.826

从上述实验结果中可以看到,RTW 无论在推荐综合性能还是推荐质量上,都有一定的优势。这也验证了我们最初的设想是合理的,即在信任网络中进行随机游走时,引入节点信任因子,给予那些被更多其他用户所信任节点更高的权重对于提高推荐质量是有帮助的。

结束语 近年来,互联网上的信息呈几何级数增长。找出一种高效的信息筛选方法来帮助人们克服信息过载问题变得越来越重要。很多领域都在使用推荐系统来提供个性化推荐信息,其中最常用的是基于协同过滤的推荐系统。然而现有的推荐系统很难在“冷启动”情况下提供令人满意的推荐结果。另外,现有的大部分方法也没有利用用户之间的信任信息。

本文提出了一种基于信任因子的随机游走方法。此方法不仅考虑了用户对项目的评分信息同时还考虑了用户之间的信任信息。通过实验结果比较可以看出,基于信任因子的随机游走方法在推荐效果上优于其他推荐方法,尤其当评分数据极度稀疏的情况下,该方法的各评价指标将明显地优于其他推荐方法。

在未来的工作中可以从以下几个方面来继续完善推荐系统:首先,我们发现在传统单机环境下的计算能力并不足以支持基于信任因子的随机游走模型进行即时推荐,所以将方法扩展到分布式平台尤为必要。其次,本文中用的是布尔型的二值信任,在今后的工作中可以用实数型的变量来度量用户间的信任关系,使得信任的度量更加准确。最后,由于用户的兴趣随时间变化,因此我们希望在方法中引入时间因素,从而进一步提高推荐效果。

参 考 文 献

[1] Guo,Lei, Ma Jun, Chen Zhu-min, et al. Incorporating Item Relations for Social Recommendation[J]. Chinese Journal of Computers, 2014, 37(1): 219-228(in Chinese)
郭磊, 马军, 陈竹敏, 等. 一种结合推荐对象间关联的社会化推荐算法[J]. 计算机学报, 2014, 37(1): 219-228

[2] Meng Xiang-wu, Liu Shu-dong, Zhang Yu-jie, et al. Research on

social recommender systems[J]. Journal of Software, 2015, 26(6): 1356-1372(in Chinese)
孟祥武, 刘树栋, 张玉洁, 等. 社会化推荐系统研究[J]. 软件学报, 2015, 26(5): 1356-1372

[3] Goldberg D, Nichols D, Oki B M, et al. Using collaborative filtering to weave an information tapestry[J]. Communications of the ACM, 1992, 35(12): 61-70

[4] Huang Y X, Cui B, Zhang W, et al. TencentRec: Real-time stream recommendation in practice[C]// Proc. of SIGMOD'15. New York: ACM Press, 2015: 227-238

[5] Chen Ke-han, Han Pan-pan, Wu Jian. User Clustering Based Social Network Recommendation[J]. Chinese Journal of Computers, 2013, 36(2): 349-359(in Chinese)
陈克寒, 韩盼盼, 吴健. 基于用户聚类的异构社交网络推荐算法[J]. 计算机学报, 2013, 36(2): 349-359

[6] Koren Y. Factorization meets the neighborhood: a multifaceted collaborative filtering recommender systems [C] // Proc. of KDD'08. New York: ACM Press, 2008: 426-434

[7] Levien, Aiken. Advogato's trust metric[OL]. <http://advogato.org/trust-metric.html>

[8] Ziegler C N, Lausen G. Spreading activation models for trust propagation[C]// Proc. of the IEEE Int'l Conf. on e-Technology, e-Commerce, and e-Service. Washington: IEEE Computer Society, 2004: 83-97

[9] Golbeck J. Computing and applying trust in Web-based social network [D]. University of Maryland College Park, 2005

[10] Massa P, Avesani P. Trust-Aware recommender systems[C]// Proc. of RecSys 2007. New York: ACM Press, 2011: 257-260

[11] Massa P, Avesani P. Trust metrics in recommender systems [M]// Computing with Social Trust Human, Springer London, 2009: 259-285

[12] Jamali M, Ester M. TrustWalker: a random walk model for combining trust-based and item-based recommendation[C]// Proc of KDD'09. New York: ACM Press, 2009: 397-405

[13] Sarwar B, Karypis G, Konstan J, et al. Item-based collaborative filtering recommendation algorithms[C]// Proc. of WWW'01. New York: ACM Press, 2001: 285-295

[14] Fang Ran-ning, Guo Yun-fei, Hu Hong-chao, et al. Improved collaborative filtering recommender algorithm based on sigmoid function[J]. Application Research of Computers, 2013, 30(6): 1688-1691(in Chinese)
方耀宁, 郭云飞, 扈红超, 等. 一种基于 sigmoid 函数的改进协同过滤推荐算法[J]. 计算机应用研究, 2013, 30(6): 1688-1691

[15] Azzedin F, Maheswaran M. Evolving and Managing Trust in Grid Computing Systems[C]// Proc. of the IEEE Canadian Conference on Electrical & Computer Engineering. 2002: 1424-1429

[16] Page L, Brin S, Motwani R, et al. The pagerank citation ranking: Bringing order to the Web[R]. Stanford Digital Library Technologies Project, 1998

[17] Milgram S. The small world problem [J]. Psychology Today, 1967, 1(1): 61-67