

基于用户信任和兴趣的概率矩阵分解推荐方法^①

彭 鹏, 米传民, 肖 琳

(南京航空航天大学 经济与管理学院, 南京 211100)

摘 要: 传统协同过滤推荐算法存在数据稀疏性、冷启动、新用户等问题. 随着社交网络和电子商务的迅猛发展, 利用用户间的信任关系和用户兴趣提供个性化推荐成为研究的热点. 本文提出一种结合用户信任和兴趣的概率矩阵分解(STUIPMF)推荐方法. 该方法首先从用户评分角度挖掘用户间的隐性信任关系和潜在兴趣标签, 然后利用概率矩阵分解模型对用户评分信息、用户信任关系、用户兴趣标签信息进行矩阵分解, 进一步挖掘用户潜在特征, 缓解数据稀疏性. 在 Epinions 数据集上进行实验验证, 结果表明, 该方法能够在一定程度上提高推荐精度, 缓解冷启动和新用户问题, 同时具有较好的可扩展性.

关键词: 推荐系统; 协同过滤; 社交信任; 兴趣标签; 概率矩阵分解

引用格式: 彭鹏,米传民,肖琳.基于用户信任和兴趣的概率矩阵分解推荐方法.计算机系统应用,2017,26(9):1-9. <http://www.c-s-a.org.cn/1003-3254/5933.html>

Recommended Algorithm Based on User Trust and Interest with Probability Matrix Factorization

PENG Peng, MI Chuan-Min, XIAO Lin

(College of Economics and Management, Nanjing University of Aeronautics and Astronautics, Nanjing 211100, China)

Abstract: The traditional collaborative filtering recommendation algorithm has such problems as data sparseness, cold-start and new users. With the rapid development of social network and e-commerce, how to provide personalized recommendations based on the trust between users and user interest tag is becoming a hot research topic. In this study, we propose a probability matrix factorization model (STUIPMF) by integrating social trust and user interest. First, we excavate implicit trust relationship between users and potential interest label from the perspective of user rating. Then we use the probability matrix factorization model to conduct matrix decomposition of user ratings information, users trust relationship, user interest label information, and further excavate the user characteristics to ease data sparseness. Finally, we make experiments based on the Epinions dataset to verify the proposed method. The results show that the proposed method can to some extent improve the recommendation accuracy, ease cold-start and new user problems. Meanwhile, the proposed STUIPMF approach also has good scalability.

Key words: recommender system; collaborative filtering; social trust; interest tag; probability matrix factorization

随着网络和信息技术的不断普及, 人类产生的数据量正在呈指数级增长, “信息过载”(Information Overload)问题日益严重^[1]. 作为帮助用户发现其感兴趣的物品、解决信息过载问题的重要工具, 推荐系统应运而生. 越来越多的电子商务服务商包括 Amazon, Half.com, CDNOW, Netflix, 和 Yahoo!等都致力于使

用推荐系统为自身客户提供“量身定做”的购买建议^[2]. 然而传统的协同过滤(Collaborative Filtering Recommendation, CF)推荐算法本身存在的数据稀疏性^[3]、冷启动、新用户等问题, 影响了推荐的精度与质量.

近年来, 互联网用户数量呈现爆炸式增长, 社交网络异军突起. 2016年8月3日, 中国互联网络信息中心

^① 基金项目: 江苏省高校哲学社会科学基金(2015SJD039); 中央高校基本科研业务费专项资金(NS2016078)

收稿时间: 2016-12-23; 采用时间: 2017-01-12

(CNNIC)在京发布的第38次《中国互联网络发展状况统计报告》显示,截至2016年6月,中国网民规模达7.10亿,互联网普及率达到51.7%^[4]。美国著名的尼尔森调查机构调查了影响用户相信某个推荐的因素,调查显示,近百分之九十的用户相信朋友对他们的推荐^[5]。基于社交信任的推荐算法得到了广泛的研究,并且可以提高推荐的质量^[6]。然而,由于用户数目庞大,也存在用户间直接信任数据比较稀疏的问题。

在推荐系统中引入用户信任关系能够缓解冷启动问题,加入用户兴趣能够缓解数据的稀疏性,基于以上事实,本文从用户评分角度挖掘用户间的隐性信任关系和潜在兴趣,在概率矩阵分解模型的基础上,融入用户信任关系信息、用户兴趣相似信息,提出一种新的结合用户信任和兴趣的协同过滤 STUIPMF 推荐方法(Recommended algorithm combined with social trust and user interest based on probability matrix factorization model)。实验结果表明,本文所提出的方法能够在一定程度上提高推荐精度,同时缓解传统协同过滤推荐算法的冷启动和新用户等问题,具有较强的可扩展性。

1 相关工作

推荐算法(或叫推荐策略)是整个推荐系统中最为核心和关键的部分,在很大程度上决定了推荐系统性能的优劣^[7]。协同过滤推荐算法因其操作简单、解释性强、技术易于实现等优点成为应用最为广泛的推荐算法之一^[8],其主要根据用户对项目的评分计算相似度的高低来进行推荐,但有研究表明,在大型电子商务系统上,用户评分项目一般不会超过项目总数的1%^[9],因此其不可避免的会出现数据稀疏性、冷启动、新用户等问题。

为了解决这些问题,不少学者将信任机制引入基于模型的协同过滤推荐算法,大致分为两类:一类是基于邻域模型研究信任关系的推荐方法:Massa提出Mole Trust模型,利用深度优先搜索评分用户,通过信任在用户A的社会网络边上的传递,预测其对目标用户B的信任值^[10];与之类似,Golbeck提出Tidal Trust模型,改进宽度优先策略预测用户间的信任值^[11];Jamali提出TrustWalker模型将基于项目的推荐系统与基于信任的推荐系统相结合^[12]。杨雪梅综合考虑用户评分相似性和用户之间信任关系,利用层次分析法构建用户信任模型,提出一种融合用户信任模型的协

同过滤推荐算法^[13]。但是这些方法只考虑了近邻用户间的信任关系,忽视了对用户间的隐性信任关系的挖掘以及用户评分对推荐结果的影响。

另一类是基于矩阵分解模型考虑信任关系的推荐方法,Jamali提出融入用户信任信息的SocialMF(a matrix factorization based model for recommendation in social rating network)方法,引入信任传播的概念,考虑直连的信任用户和两步连接的用户的用户的信息来产生推荐,获得了较好的推荐效果,但计算复杂度较高,而且未采用不同的信任度量标准^[14]。Ma等提出SoRec(Social Recommendation)方法,通过共享用户隐特征向量空间把用户评分信息和用户社交信息联系起来进行研究^[15]。但這些方法更多是考虑直连的信任网络,而忽视了用户间隐性信任关系的挖掘。

由于基于矩阵分解的推荐算法利用的是隐因子,较难对推荐结果给出准确合理的解释。Salakhutdinov等从概率的角度描述了矩阵分解问题,提出概率矩阵分解模型(Probabilistic Matrix Factorization, PMF),通过给用户-推荐项目的特征矩阵加上先验分布,并最大化预测评价的后验概率来进行推荐,并在Netflix数据集上得到了十分优秀的预测结果^[16]。Noam Koenigstein等人在概率矩阵分解过程中集成了部分物品特征信息,并且在Xbox电影推荐系统上进行实验,验证了所提模型的有效性^[17]。

有研究表明,考虑用户兴趣进行建模有利于做出更精准的个性化推荐。姚平平提出一种基于用户偏好和项目属性的协同过滤推荐算法,但忽视了用户间信任关系对于推荐结果的影响^[18]。Lee提出将用户偏好信息和社交网络中的信任传播相结合,提高推荐质量^[19]。陶俊等提出了一种适应用户兴趣多样性的基于用户兴趣分类的协同过滤算法并利用改进的模糊聚类算法,搜索最近邻来改善推荐算法的准确性^[20]。嵇晓声等提出了一种基于用户兴趣度的相似性度量方法,利用用户对不同项目类别的兴趣程度与用户评分相结合进行用户之间的相似性计算^[21]。但这些方法大多数都只关注用户对项目评分值,没有考虑用户偏好以及用户评分与项目属性之间的关系对推荐精度的影响,也忽视了用户间信任关系的挖掘。

因此,本文综合考虑用户对项目的评分和用户间的隐性信任关系,在概率矩阵分解模型基础上加入用户间的信任关系和用户兴趣信息,进一步挖掘出隐藏在信任关系和用户评分背后的用户特征,提出一种新

的 STUIPMF 推荐算法, 实验结果表明: 该方法综合利用多方面信息, 能够提升推荐精度和模型的可扩展性.

2 结合用户信任和兴趣的概率矩阵分解推荐算法

2.1 概率矩阵分解模型

概率矩阵分解模型的原理是从概率的角度来预测用户对项目的评分. 为了便于形式化描述, 本文将用到的参数符号, 如表 1 所示.

表 1 参数列表

符号	意义
M, N, S	分别表示用户数目、商品数目、兴趣标签数目
K	隐特征向量维数
$U_{M \times K}$	用户隐特征向量
$V_{N \times K}$	商品隐特征向量
$F_{M \times K}$	信任隐特征向量
$L_{S \times K}$	兴趣标签隐特征向量
R_{ij}	用户 U_i 对商品 V_j 的真实评分
P_{ik}	用户 U_i 在兴趣标签 L_k 上的标注次数
T_{il}	用户 U_i 与朋友 F_l 之间的信任程度
\tilde{R}_{ij}	用户 U_i 对商品 V_j 的预测评分

PMF 的计算过程如下:

假设用户和物品的隐式特征向量都服从高斯先验分布:

$$P(U|\sigma_U^2) = \prod_{i=1}^M N(U_i|0, \sigma_U^2 I) \quad (1)$$

$$P(V|\sigma_V^2) = \prod_{j=1}^N N(V_j|0, \sigma_V^2 I) \quad (2)$$

再假设已获取的用户评分数据的条件概率也服从高斯先验分布:

$$P(R|U, V, \sigma_R^2) = \prod_{i=1}^M \prod_{j=1}^N \left[N(R_{ij} | g(U_i^T V_j), \sigma_R^2) \right]^{I_{ij}^R} \quad (3)$$

其中, I_{ij}^R 为指示函数, 如果用户 U_i 对物品 V_j 有过评分, 那么它的值等于 1, 否则为 0. $g(x)$ 将 $U_i^T V_j$ 的值映射到 $[0, 1]$ 区间内, 在本文中 $g(x) = 1/(1 + e^{-x})$.

通过贝叶斯推理, 可得用户和物品的隐式特征的后验概率:

$$IP(U, V | R, \sigma_R^2, \sigma_U^2, \sigma_V^2) \propto P(R|U, V, \sigma_R^2) \times P(U|\sigma_U^2) \times P(V|\sigma_V^2) \quad (4)$$

这样, 通过用户物品评分矩阵, 就可以学习到用户和物品的隐特征向量, 进而通过求内积的方式得到最近似的用户评分, 用公式表示即:

$$\tilde{R} \approx U_i^T V_j \quad (5)$$

相应的概率图模型如图 1 所示.

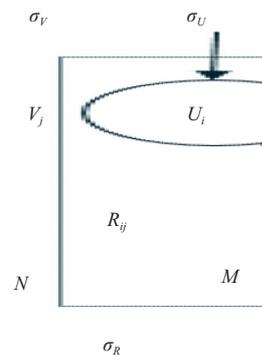


图 1 概率矩阵分解图模型

2.2 挖掘用户隐性信任关系

基于信任的推荐算法可以有效解决推荐系统的冷启动和数据稀疏性问题, 提高推荐覆盖率. 本文通过研究发现, 当前大多数算法只考虑直连的信任网络, 即用户间的显性信任关系, 而较少关注用户隐性信任关系的挖掘, 所以引入用户行为系数和用户信任度函数对用户信任关系衡量进行改进.

用户行为系数是基于用户评分信息进行信任推理, 通过计算评分准确度得到, 在用户评分相似度的基础上得到用户隐性信任关系; 本文中, 用户评分准确度利用用户与所有用户对项目平均评分的差值进行表示; 一般来讲, 用户评分的准确与否将直接影响其他用户对其信任的程度. φ_u 表示用户 u 的用户行为系数, 由用户评分准确度决定.

$$\varphi_u = \frac{1}{1 + \sum_{i=1}^N (R_{ui} - \bar{R}_i) \cdot I_{ui}} \quad (6)$$

R_{ui} 表示用户 u 对项目 i 的评分, \bar{R}_i 表示所有用户对于项目 i 的平均评分, 如果用户 u 对项目 i 有评分, $I_{ui}=1$, 没有评分, 则 $I_{ui}=0$.

评分相似度采用较为流行的 Pearson 相关系数进行度量. 它首先需要找到用户 i 和用户 j 共同评分的项目集合, 这个集合用表示, 如此两个用户的相似性 $sim_{i,j}$ 的计算公式为:

$$sim_{i,j} = \frac{\sum_{c \in U} (r_{i,c} - \bar{r}_i) \times (r_{j,c} - \bar{r}_j)}{\sqrt{\sum_{c \in U} (r_{i,c} - \bar{r}_i)^2 \times \sum_{c \in U} (r_{j,c} - \bar{r}_j)^2}} \quad (7)$$

其中 $r_{i,c}$ 和 $r_{j,c}$ 分别表示用户 i 和用户 j 对项目 c 的评分, \bar{r}_i 和 \bar{r}_j 分别表示用户 i 和用户 j 对项目 c 的平均评分.

本文中用户隐性信任关系用 TI 表示, 用户 i 和用户 j 的隐性信任关系:

$$TI_{ij} = \varphi_i \cdot sim_{ij} \quad (8)$$

本文用 t_{ij} 表示用户 i 和用户 j 的显式信任关系, 用户 i 信任用户 j 时, $t_{ij}=1$, 反之为 0. 由于信任关系具有非对称性, t_{ij} 并不能准确地反映用户之间的显性信任关系, 还应跟信任和被信任的用户个数有关. 比如, 当用户 t_i 信任很多用户时, 用户 t_i 和用户 t_j 之间的信任值 t_{ij} 应降低, 反之, 当用户 t_i 被很多用户信任时, 用户 t_i 和用户 t_j 之间的信任值应得到增强. 因此, 结合用户影响力对用户之间的显性信任值进行改进: TE_{ij} 表示改进后的显性信任值.

$$TE_{ij} = \sqrt{\frac{d^-(u_i)}{d^+(u_j) + d^-(u_i)}} \cdot t_{ij} \quad (9)$$

$d^-(u_i)$ 表示用户 u_i 被信任的用户数量, $d^+(u_j)$ 表示用户 u_j 信任的用户数量.

用户信任度函数用 T_{ij} 表示, 通过结合信任网络中声明的显性信任关系, 确定显性信任和隐性信任各自的权重系数之后计算获得.

$$T_{ij} = \alpha \cdot TE_{ij} + (1 - \alpha) \cdot TI_{ij} \quad (10)$$

其中 α 代表权重系数.

用户信任关系矩阵用 T 表示, 元素 T_{il} 表示用户 U_i 与朋友 F_l 之间的信任度, 已知用户信任条件概率分布函数, 如下所示:

$$P(T|U, F, \sigma_T^2) = \prod_{i=1}^M \prod_{l=1}^M \left[N(T_{il} | g(U_i^T F_l), \sigma_T^2) \right]^{I_{il}^T} \quad (11)$$

其中, I_{il}^T 表示指示函数, 如果用户 U_i 与用户 F_l 是朋友关系, 那么他的值等于 1, 否则为 0.

U 和 F 概率分布, 如下所示:

$$P(U|\sigma_U^2) = \prod_{i=1}^M N(U_i | 0, \sigma_U^2 I) \quad (12)$$

$$P(F|\sigma_F^2) = \prod_{l=1}^M N(F_l | 0, \sigma_F^2 I) \quad (13)$$

根据贝叶斯推理, 如下所示:

$$P(U, F | T, \sigma_T^2, \sigma_U^2, \sigma_F^2) \propto P(T | U, F, \sigma_T^2) \times P(U | \sigma_U^2) \times P(F | \sigma_F^2) \quad (14)$$

基于用户信任关系的模型, 相应的概率图模型如图 2 所示.

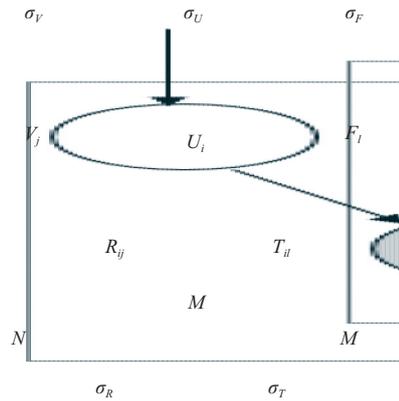


图 2 基于用户信任关系的概率图模型

2.3 挖掘用户兴趣相似关系

本文研究发现当前基于用户兴趣分类的算法较少关注用户偏好以及用户评分与项目属性之间的关系对推荐结果精度的影响, 所以本文基于用户—项目评分矩阵结合项目类型信息和用户评分阈值, 挖掘用户隐性标签, 得到用户—兴趣标签矩阵, 进而补充用户信息, 缓解数据稀疏性.

用户对所有项目的评分集合对应的中位数评分阈值集合为 $A = \{A_1, A_2, \dots, A_m\}$, 项目的属性标签集合 $L = \{L_1, L_2, \dots, L_k\}$, 当 $R_{ui} \geq A_i$ 时, 认为用户 u 喜欢项目 i , 项目 i 对应的项目属性标签项目 L_c 被标记为用户 u 的兴趣标签, 根据项目的属性和用户的评分阈值提取用户的兴趣标签. 某个用户可能多次被同一个兴趣标签标记, 对次数进行累加, 形成用户的兴趣标签矩阵 $L_{me} = \{L_{uy}\}$, L_{uy} 表示用户 u 对于项目属性 L 标记的兴趣标签次数. 然后将用户评分中低于评分阈值的评分记为 0, 得到用户-项目中位数评分矩阵; 结合项目-属性矩阵, 项目属于某一属性, 则记为 1, 反之为 0, 这样用户与项目属性之间就建立了联系就得到了用户-兴趣标签矩阵 P .

用户兴趣标签矩阵用 P 表示, 元素 P_{ik} 表示用户 U_i 在兴趣标签 L_k 上的标记次数, 已知用户兴趣标签概

率分布函数,如下所示:

$$P(P|U, L, \sigma_P^2) = \prod_{i=1}^M \prod_{k=1}^Q \left[N(P_{ik} | g(U_i^T L_k), \sigma_P^2) \right]^{I_{ik}^P} \quad (15)$$

其中, I_{ik}^P 表示指示函数, 如果用户 U_i 在兴趣标签 L_k 上至少标记过一次, 那么它的值等于 1, 否则为 0.

U_i 和 L 的概率分布, 如下所示:

$$P(U|\sigma_U^2) = \prod_{i=1}^M N(U_i | 0, \sigma_U^2) \quad (16)$$

$$P(L|\sigma_L^2) = \prod_{k=1}^Q N(L_k | 0, \sigma_L^2) \quad (17)$$

根据贝叶斯推理, 如下所示:

$$P(U, L | P, \sigma_P^2, \sigma_U^2, \sigma_L^2) \propto P(P | U, L, \sigma_P^2) \times P(U | \sigma_U^2) \times P(L | \sigma_L^2) \quad (18)$$

基于用户兴趣标签的模型, 相应的概率图模型如图 3 所示.

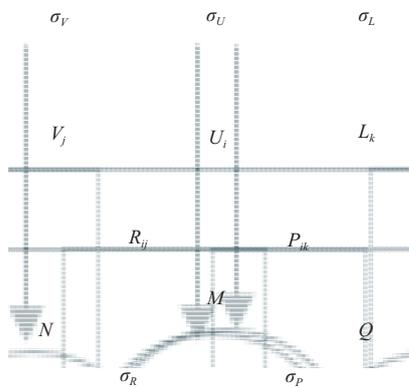


图 3 基于用户兴趣标签的概率图模型

3 融合用户信任和兴趣的 STUIPMF 模型具体实现

概率矩阵分解算法仅依据用户-项目评分矩阵, 学习相应的特征向量, 并没有考虑用户之间信任关系以及用户兴趣对推荐结果的影响, 为了体现这一影响, 本文对上述模型进行改进, 将用户信任关系矩阵、用户兴趣标签矩阵和用户-项目评分矩阵的分解整合起来, 通过用户潜在特征矩阵进行连接, 提出 STUIPMF 模型, 如图 4 所示.

联合之后的后验概率的对数值满足下式:

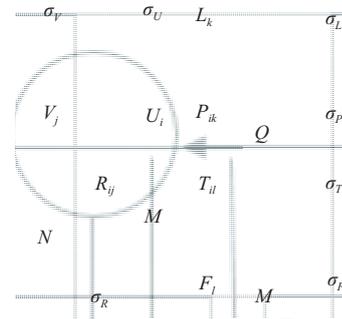


图 4 STUIPMF 概率图模型

$$\begin{aligned} \ln P(U, V, L, F | R, P, T, \sigma_R^2, \sigma_P^2, \sigma_T^2, \sigma_U^2, \sigma_V^2, \sigma_L^2, \sigma_F^2) \\ = -\frac{1}{2\sigma_R^2} \sum_{i=1}^M \sum_{j=1}^N I_{ij}^R (R_{ij} - g(U_i^T V_j))^2 \\ -\frac{1}{2\sigma_P^2} \sum_{i=1}^M \sum_{k=1}^Q I_{ik}^P (P_{ik} - g(U_i^T L_k))^2 \\ -\frac{1}{2\sigma_T^2} \sum_{i=1}^M \sum_{l=1}^M I_{il}^T (T_{il} - g(U_i^T F_l))^2 \\ -\frac{1}{2\sigma_U^2} \sum_{i=1}^M U_i^T U_i - \frac{1}{2\sigma_V^2} \sum_{j=1}^N V_j^T V_j - \frac{1}{2\sigma_L^2} \sum_{k=1}^Q L_k^T L_k \quad (19) \\ -\frac{1}{2\sigma_F^2} \sum_{l=1}^M F_l^T F_l - \frac{1}{2} \left(\sum_{i=1}^M \sum_{j=1}^N I_{ij}^R \right) \ln \sigma_R^2 \\ -\frac{1}{2} \left(\left(\sum_{i=1}^M \sum_{k=1}^Q I_{ik}^P \right) \ln \sigma_P^2 + \left(\sum_{i=1}^M \sum_{l=1}^M I_{il}^T \right) \ln \sigma_T^2 \right) \\ -\frac{1}{2} \left((M \times K) \ln \sigma_U^2 + (N \times K) \ln \sigma_V^2 \right) \\ -\frac{1}{2} \left((M \times K) \ln \sigma_F^2 + (Q \times K) \ln \sigma_L^2 \right) + C \end{aligned}$$

最大化其对数后验等价于最小化如下目标函数:

$$\begin{aligned} S(U, V, L, F, R, P, T) \\ = \frac{1}{2} \sum_{i=1}^M \sum_{j=1}^N I_{ij}^R (R_{ij} - g(U_i^T V_j))^2 \\ + \frac{\lambda_P}{2} \sum_{i=1}^M \sum_{k=1}^Q I_{ik}^P (P_{ik} - g(U_i^T L_k))^2 \\ + \frac{\lambda_T}{2} \sum_{i=1}^M \sum_{l=1}^M I_{il}^T (T_{il} - g(U_i^T F_l))^2 \quad (20) \\ + \frac{\lambda_U}{2} \sum_{i=1}^M \|U_i\|_F^2 + \frac{\lambda_V}{2} \sum_{j=1}^N \|V_j\|_F^2 + \frac{\lambda_L}{2} \sum_{k=1}^Q \|L_k\|_F^2 \\ + \frac{\lambda_F}{2} \sum_{l=1}^M \|F_l\|_F^2 \end{aligned}$$

$\lambda_P = \sigma_R^2 / \sigma_P^2$, $\lambda_T = \sigma_R^2 / \sigma_T^2$, $\lambda_U = \sigma_R^2 / \sigma_U^2$, $\lambda_V = \sigma_R^2 / \sigma_V^2$, $\lambda_L = \sigma_R^2 / \sigma_L^2$, $\lambda_F = \sigma_R^2 / \sigma_F^2$ 都是正则化固定参数, $\|\cdot\|_F$ 表示矩阵的 Frobenius 范数.

本文采用随机梯度下降法学习得到相应的潜在特征矩阵:

$$\frac{\partial S}{\partial U_i} = \sum_{j=1}^N I_{ij}^P g'(U_i^T V_j) (g(U_i^T V_j) - R_{ij}) V_j + \lambda_P \sum_{k=1}^Q I_{ik}^P g'(U_i^T L_k) (g(U_i^T L_k) - P_{ik}) L_k + \lambda_T \sum_{l=1}^M I_{il}^T g'(U_i^T F_l) (g(U_i^T F_l) - T_{il}) F_l + \lambda_U U_i \quad (21)$$

$$\frac{\partial S}{\partial V_j} = \sum_{i=1}^M I_{ij}^R g'(U_i^T V_j) (g(U_i^T V_j) - R_{ij}) U_i + \lambda_V V_j \quad (22)$$

$$\frac{\partial S}{\partial L_k} = \lambda_P \sum_{i=1}^Q I_{ik}^P g'(U_i^T L_k) (g(U_i^T L_k) - P_{ik}) U_i + \lambda_L L_k \quad (23)$$

$$\frac{\partial S}{\partial F_l} = \lambda_T \sum_{i=1}^M I_{il}^T g'(U_i^T F_l) (g(U_i^T F_l) - T_{il}) U_i + \lambda_F F_l \quad (24)$$

为降低计算复杂度, 令 $\lambda_U = \lambda_V = \lambda_T = \lambda_L = \lambda$, λ_P 和 λ_F 的取值后面会讨论到.

每次迭代时, U_i, V_j, L_k, F_l 调整如下:

$$U_i \leftarrow U_i - \gamma \cdot \frac{\partial S}{\partial U_i}, V_j \leftarrow V_j - \gamma \cdot \frac{\partial S}{\partial V_j} \quad (25)$$

$$L_k \leftarrow L_k - \gamma \cdot \frac{\partial S}{\partial L_k}, F_l \leftarrow F_l - \gamma \cdot \frac{\partial S}{\partial F_l}$$

其中, γ 为预先定义的步长.

重复上述训练过程, 每次迭代后, 计算并验证平均绝对误差, 当目标函数 S 值的变化小于预先定义的很小的常数后终止迭代过程. 得到迭代终止后的 U_i, V_j, L_k, F_l 之后, 就可以预测用户 U_i 对商品 V_j 的未知评分, 对于每一个用户, 根据计算得到的预测评分值由高到低对候选商品进行排序, 产生 Top- N 推荐列表, 然后推荐给用户.

4 实验验证及结果分析

4.1 数据集介绍

Epinions 作为 1999 年在美国成立的一个点评性质的网站, 用户在上面能够浏览到其他用户对于琳琅满目的物品的评论, 进而为其购物、选择公司、观看节目或者电影时提供参考意见. 网站定位是“社会化商

务”, 它建立了一种“信任机制”, 即用户能够对他人对物品点评质量的好坏做出自己的判定, 假如相信某个用户, 就可以把他加入自己的信任列表, 反之也有权利加入不信任列表. 这样 Epinions 网站的数据中就同时包含用户评分信息和用户信任关系, 因此在当前推荐系统研究领域, 都会把 Epinions 数据集作为一个基准数据库, 来进行基于信任的推荐方法的研究.

本文的数据集就是由 Paolo Massa 和 Paolo Avesani 提供的, 爬取自“Epinions.com”的网站. 关于这个数据集的统计情况如表 2 所示.

表 2 Epinions 数据集的统计信息

数据集	Epinions
用户数	49290
商品数	139738
评分数	664813
信任关系数	487181
兴趣标签数	154

4.2 评价指标

为了评价推荐系统预测的准确性, 本文采用常用的评价指标: 平均绝对误差(MAE, 全称 Mean Absolute Error)和 RMSE 均方根误差(Root Mean Squared Error, RMSE), 分别定义如下:

$$MAE = \frac{1}{n} \sum_{i=1}^n |p_i - r_i| \quad (26)$$

$$RMSE = \sqrt{\left(\sum_{i=1}^n (p_i - r_i)^2 \right) / n} \quad (27)$$

其中 n 为参与预测的项目数, p_i 是系统预测的目标用户的评分, r_i 是目标用户实际评分.

将本文所提算法与文献[16]提到的 PMF 模型、文献[14]提到的 SocialMF 模型、文献[15]提到的 SoReg 这三种典型推荐算法的效果进行对比.

4.3 参数 λ_P, λ_F 的影响

在本文所提的方法中, 参数 λ_P, λ_F 的设置显得至关重要, 它们起到平衡的作用. 当 λ_P 设置为 0 时, 系统推荐时就不考虑用户之间的信任关系, 而只考虑用户的评分矩阵和用户的隐性兴趣标签. 当 λ_P 设置为无穷大时, 系统推荐时就只考虑用户之间的信任关系, 而不考虑其他因素; 同理, 当 λ_F 设置为 0 时, 系统推荐时就不考虑用户的隐性兴趣标签, 而只考虑用户的评分矩阵和用户之间的信任关系. 当 λ_F 设置为无穷大时, 系统推

荐时就只考虑用户隐性兴趣标签,而不考虑其他因素.

图 5 显示的是在其他参数设置不变的情况下,参数 λ_p 在潜在特征数分别为 5 和 10 的情况下对 MAE 和 RMSE 的影响. 由图 5 的(a)和(b)可知,随着 λ_p 的增加, MAE 和 RMSE 都在降低,即预测的准确性在提高;当 λ_p 达到一定阈值时,随着 λ_p 的增加, MAE 和 RMSE 都在提高,即预测的准确性在降低. 由此可得出, $\lambda_p \in [0.01, 0.1]$ 时,推荐的准确度比较高,因此,在后面的实验中我们均采用这个区间的平均值作为近似最优值,即 $\lambda_p = \lambda_f = 0.005$ 进行实验. 图 6 展示的 λ_f 取值的影响,同理,在此不再赘述.

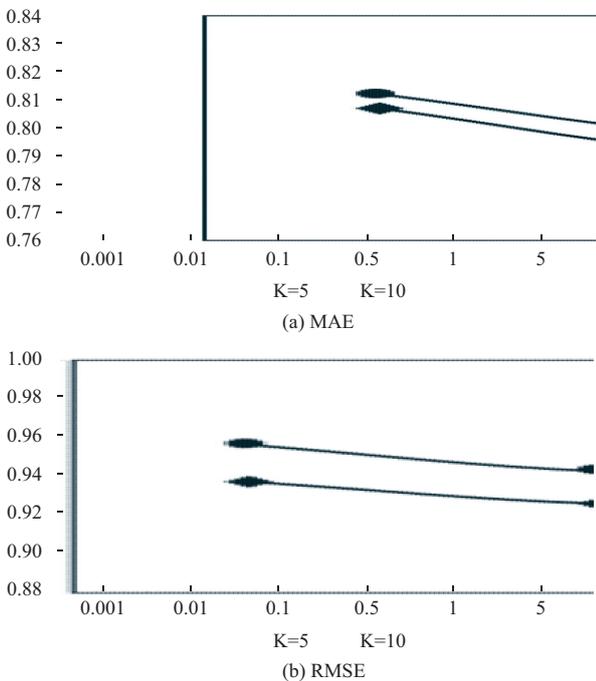


图 5 参数 λ_p 对 MAE 和 RMSE 的影响

为了验证实验效果,分别选择其中 80% 的用户作为训练集, 20% 的用户作为测试集和 90% 的用户作为训练集, 10% 的用户作为测试集进行实验.

在实验过程中,相关参数主要是根据实验效果进行最优选择的. 本文的 STUIPMF 推荐方法的参数设置如下: $\lambda_U = \lambda_V = \lambda_T = \lambda_L = \lambda = 0.001$, $\lambda_p = \lambda_f = 0.005$, 特征向量 K 的取值分别为 5 和 10; 其他方法的参数设置: 在 PMF 方法中, $\lambda_U = \lambda_V = 0.001$; 在 SocialMF 方法中, $\lambda_U = \lambda_V = 0.001$, $\lambda_T = 0.5$; 在 SoReg 方法中, $\lambda_U = \lambda_V = 0.001$, $\alpha = 0.1$. STUIPMF 方法与其他方法的试验效果对比如图 7、图 8 所示.

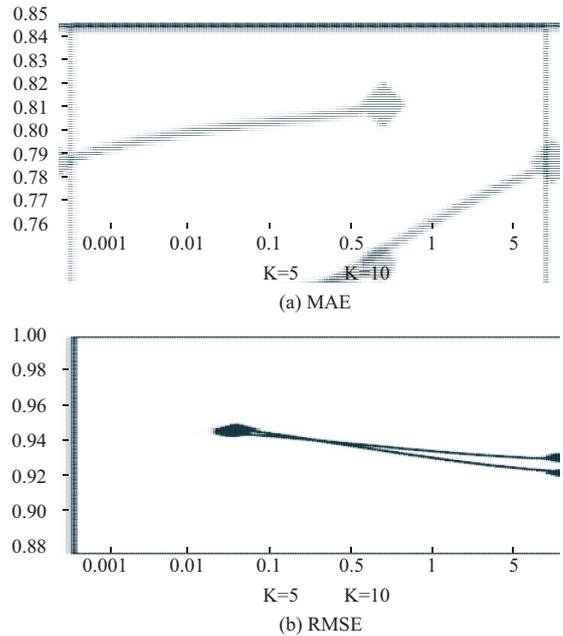


图 6 参数 λ_f 对 MAE 和 RMSE 的影响

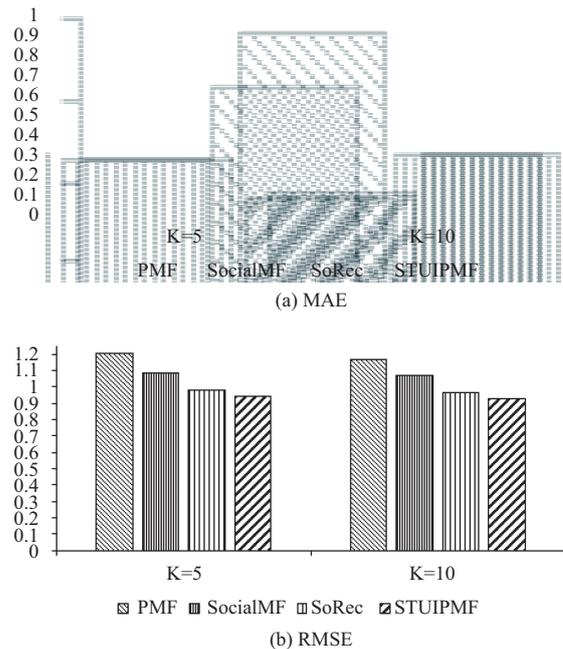


图 7 STUIPMF 方法与其他方法实验效果对比(80% 训练集)

根据图 7、8 所示,我们可以得出如下结论:

1) 本文所提出的 STUIPMF 模型综合考虑了用户评分信息、用户信任关系和兴趣信息,在所有实验参数都选择最优的情况下,80% 作为训练集,20% 作为测试集时,与 PMF、SocialMF、SoReg 相比, MAE 值分别下降 17%、5.8%、5.3%, RMSE 值分别下降 21%、

13%、4%；90% 训练集，10% 测试集时，与 PMF、SocialMF、SoReg 相比，MAE 值分别下降 16.2%、4.1%、3.7%，RMSE 值分别下降 20.8%、13.5%、4.1%；说明本文所提方法在推荐准确性上有所提高。

2) 随着隐特征向量维数的增加，推荐的精度有所提高，但另一方面可能出现过拟合和计算复杂度增加的问题。

3) 对用户信任关系矩阵和用户兴趣标签矩阵进行概率矩阵分解后，能增加用户特征的先验信息，从而缓解推荐系统中的冷启动和新用户问题。

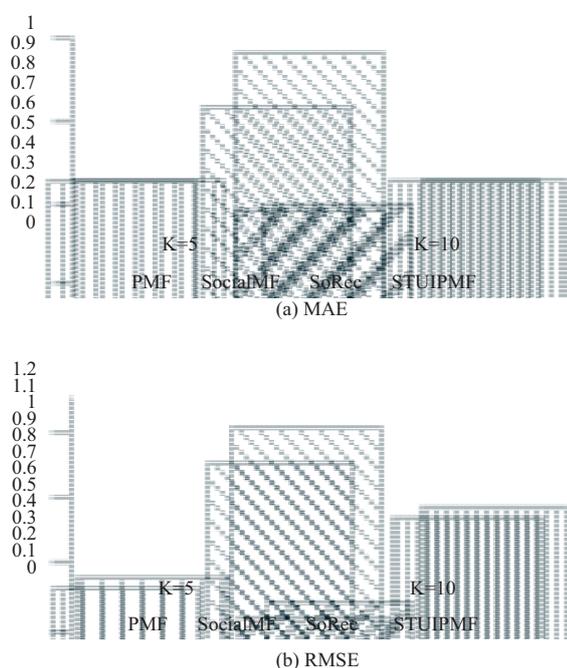


图8 STUIPMF方法与其他方法实验效果对比(90%训练集)

5 结论及展望

随着个性化服务在当今经济和社会生活中的地位和重要性日益突出，如何通过用户行为准确把握用户真实兴趣与需求，提供高质量的个性化推荐成为当前的研究热点。针对传统协同过滤方法存在的冷启动和数据稀疏性等问题，本文提出一种结合用户信任和兴趣的概率矩阵分解推荐方法(STUIPMF方法)。首先从用户评分角度挖掘用户间的隐性信任关系和潜在兴趣标签，然后利用概率矩阵分解模型对用户评分信息、用户信任关系、用户兴趣标签信息进行矩阵分解，进一步挖掘用户潜在特征，缓解数据稀疏性，产生更精准的推荐。在 Epinions 数据集上进行实验验证，结果表明，本文所提出的 STUIPMF 方法综合利用用户评分、用

户信任关系、用户兴趣等多方面信息，能够在一定程度上提高推荐精度，缓解冷启动和新用户问题，同时具有较强的可扩展性。

但其中也存在一些问题，比如模型中 λ 的取值我们是使用的是近似最优值，接下来将进一步确定 λ 的最优值以及动态 λ 值的变化，提高推荐准确度；另一方面会考虑更多的信息融入所提模型中，如文本信息、位置信息、时间因素等，关注用户信任关系和兴趣的更新，以及加入对用户之间不信任关系的考量。

参考文献

- 1 Borchers A, Herlocker J, Konstan J, *et al.* Ganging up on information overload. *Computer*, 1998, 31(4): 106–108. [doi: 10.1109/2.666847]
- 2 Linden G, Smith B, York J. Amazon.com recommendations: Item-to-item collaborative filtering. *IEEE Internet Computing*, 2003, 7(1): 76–80. [doi: 10.1109/MIC.2003.1167344]
- 3 Vozalis E, Margaritis KG. Analysis of recommender systems' algorithms. Proc. of the the 6th Hellenic-European Conference on Computer Mathematics and its Applications. Athens, Greece. 2003. 732–745.
- 4 中国互联网络信息中心. 第38次《中国互联网络发展状况统计报告》. 北京: 中国互联网络信息中心, 2016.
- 5 Bobadilla J, Ortega F, Hernando A, *et al.* Recommender systems survey. *Knowledge-Based Systems*, 2013, 46: 109–132. [doi: 10.1016/j.knosys.2013.03.012]
- 6 Guo GB. Integrating trust and similarity to ameliorate the data sparsity and cold start for recommender systems. Proc. of the 7th ACM Conference on Recommender Systems. Hong Kong, China. 2013. 451–454.
- 7 王国霞, 刘贺平. 个性化推荐系统综述. *计算机工程与应用*, 2012, 48(7): 66–76.
- 8 张学钱, 林世平, 郭昆. 协同过滤推荐算法对比分析与优化应用. *计算机系统应用*, 2015, 24(5): 100–105.
- 9 Sun XH, Kong FS, Ye S. A comparison of several algorithms for collaborative filtering in startup stage. Proc. of 2005 IEEE Networking, Sensing and Control. Tucson, AZ, USA. 2005. 25–28.
- 10 Massa P, Avesani P. Trust-aware recommender systems. Proc. of the 2007 ACM Conference on Recommender Systems. Minneapolis, MN, USA. 2007. 17–24.
- 11 Golbeck J. Personalizing applications through integration of inferred trust values in semantic web-based social networks. Proc. of Semantic Network Analysis Workshop. Galway, Ireland. 2005.

- 12 Jamali M, Ester M. TrustWalker: A random walk model for combining trust-based and item-based recommendation. Proc. of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. Paris, France. 2009. 397–406.
- 13 杨秀梅, 孙咏, 王丹妮, 等. 融合用户信任模型的协同过滤推荐算法. 计算机系统应用, 2016, 25(7): 165–170. [doi: [10.15888/j.cnki.csa.005229](https://doi.org/10.15888/j.cnki.csa.005229)]
- 14 Jamali M, Ester M. A matrix factorization technique with trust propagation for recommendation in social networks. Proc. of the 4th ACM Conference on Recommender Systems. Barcelona, Spain. 2010. 135–142.
- 15 Ma H, Yang HX, Lyu MR, *et al.* SoRec: Social recommendation using probabilistic matrix factorization. Proc. of the 17th ACM Conference on Information and Knowledge Management. Napa Valley, California, USA. 2008. 931–940.
- 16 Salakhutdinov BR, Mnih A. Probabilistic matrix factorization. Proc. of the 20th International Conference on Neural Information Processing Systems. 2015. 1257–1264.
- 17 Koenigstein N, Paquet U. Xbox movies recommendations: Variational bayes matrix factorization with embedded feature selection. Proc. of the 7th ACM Conference on Recommender Systems. Hong Kong, China. 2013. 129–136.
- 18 姚平平, 邹东升, 牛宝君. 基于用户偏好和项目属性的协同过滤推荐算法. 计算机系统应用, 2015, 24(7): 15–21.
- 19 Lee WP, Ma CY. Enhancing collaborative recommendation performance by combining user preference and trust-distrust propagation in social networks. Knowledge-Based Systems, 2016, (106): 125–134. [doi: [10.1016/j.knosys.2016.05.037](https://doi.org/10.1016/j.knosys.2016.05.037)]
- 20 陶俊, 张宁. 基于用户兴趣分类的协同过滤推荐算法. 计算机系统应用, 2011, 20(5): 55–59.
- 21 嵇晓声, 刘宴兵, 罗来明. 协同过滤中基于用户兴趣度的相似性度量方法. 计算机应用, 2010, 30(10): 2618–2620.