

# A Collaborative Filtering Recommender Algorithm Based on Privacy Preserving

Chenchen Li, Lefeng Zhang, Hui Hui, Ping Xiong

School of Information and Security Engineering, Zhongnan University of Economics and Law, Wuhan Hubei  
Email: lichenchen@zuel.com

Received: Jul. 7<sup>th</sup>, 2016; accepted: Jul. 26<sup>th</sup>, 2016; published: Jul. 29<sup>th</sup>, 2016

Copyright © 2016 by authors and Hans Publishers Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

---

## Abstract

Privacy preserving in recommender system is a hot research area currently. With the premise that recommender system server is untrusted, we propose a privacy-preserving collaborative filtering algorithm based on substitution encryption. Users encrypt their rating information at the client side and submit it to the recommender server. With the encrypted ratings collected from users, the recommender server predicts the ratings for users on unrated items with collaborative filtering algorithm. We represent a method for computing the similarity of users without knowing the meaning of the ratings, which is used for identifying the nearest neighbors of each user in collaborative filtering and predicting. The experimental results demonstrate the superiority of the proposed method comparing to the traditional collaborative filtering recommender algorithms.

## Keywords

Privacy Preserving, Collaborative Filtering, Recommender System, Substitution Encryption

---

# 一种基于隐私保护的协同过滤推荐算法

李晨晨, 张乐峰, 惠 慧, 熊 平

中南财经政法大学信息与安全工程学院, 湖北 武汉  
Email: lichenchen@zuel.com

收稿日期: 2016年7月7日; 录用日期: 2016年7月26日; 发布日期: 2016年7月29日

## 摘要

推荐系统中的用户隐私保护问题是当前的一个研究热点。以推荐系统服务器不可信为前提,提出了一种基于代换加密的隐私保护协同过滤算法。用户在客户端对评分信息进行代换加密并提交给推荐服务器,服务器则根据收集的评分密文信息进行协同过滤推荐。提出了一种无语义条件下的用户模式相似度计算方法,用以在隐私保护协同过滤中确定每个用户的近邻,进而对用户的评分密文进行预测。实验结果验证了该方法相对于传统协同过滤推荐算法的优越性。

## 关键词

隐私保护, 协同过滤, 推荐系统, 代换加密

## 1. 引言

随着计算机网络和电子商务的迅猛发展,推荐系统成为企业提高市场竞争力的重要工具。作为一种信息过滤技术[1],推荐系统能够利用客户对产品项目的历史评价来分析客户的消费模式,从而实现对客户的个性化信息推送,在帮助客户摆脱信息过载困境的同时,为企业创造更多的商业利益。例如,美国亚马逊公司至少有 20% 的销售利润来自于推荐算法[2]。协同过滤(Collaborative Filtering, CF)是目前应用最为广泛的一种推荐算法,它利用用户消费行为上的相似性来实现个性化的推荐。通常,用户的历史信息越详尽,推荐结果则越精准。然而,对用户数据的深度分析与挖掘会对用户的隐私造成严重的威胁。Calandrino 等[3]的研究证明了利用项目评分和用户相似度矩阵来进行推演攻击的可行性。这导致了用户对个人隐私信息的普遍担忧。Kobsa 的研究[4]表明,70%~89.5%的互联网用户认为个人隐私信息面临泄露风险。因此,如何在保护用户隐私的前提下实现准确的推荐,成为目前推荐系统领域的一个研究热点。

针对这一问题,一系列隐私保护方法近年来被提出,主要包括随机干扰方法[5][6]和分组匿名方法[7]。这些方法的核心思想是把推荐服务提供者(service provider, SP)作为不可信的实体,用户的信息必须经过相应的处理后才能提交给 SP。然而,虽然随机干扰方法能够有效保护用户的历史信息,但却无法防止 SP 根据产生的推荐结果来推测用户的行为特征;分组匿名方法虽然能够将个人的行为模式泛化,但却需要用户之间建立信任关系并充分交换信息。因此,这些方法在实际应用中都具有一定的局限性。

本文提出一种基于隐私保护的协同过滤推荐算法,用代换密码来保证用户评分信息的语义安全,同时设计无语义前提下的相似度计算方法,使得推荐算法在无法获取用户评分语义的情况下,仍然能够为用户提供较为精准的推荐。实验结果验证了该方法的可行性。

## 2. 基于隐私保护的推荐系统

### 2.1. 基本框架

本文提出的隐私保护推荐系统采用客户-服务器结构。用户和推荐算法分别处于客户端和服务端,并各自完成相应的数据处理过程。系统的基本流程框架如图 1 所示,主要包括以下几个步骤:

- (1) 首先,每个用户在客户端自主地定义评分符号集合(密文空间),符号所代表的评分语义不对外公开;
- (2) 每个用户用自定义的评分符号对项目进行评分,并将评分信息上传到服务器;

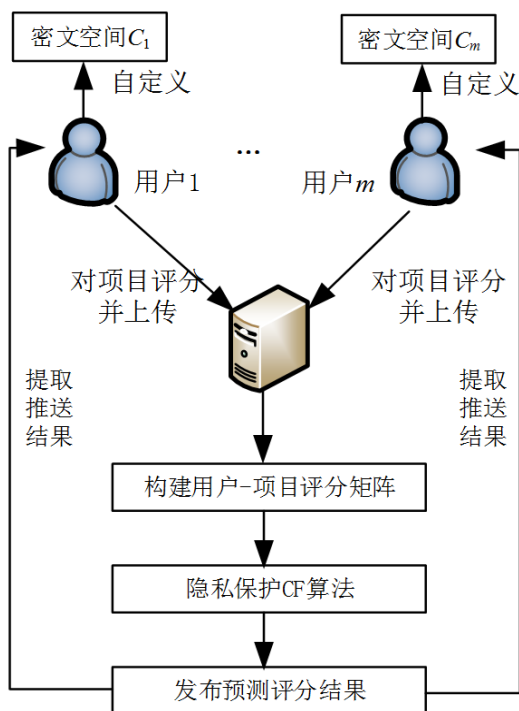


Figure 1. The framework of algorithm CF  
图 1. 隐私保护 CF 算法框架

- (3) 服务器收集到所有用户加密后的评分信息，建立用户 - 项目评分矩阵；
  - (4) 服务器上的推荐算法根据用户 - 项目评分矩阵计算用户相似度矩阵，并利用本文所提算法产生各用户对未评分项目的预测评分，预测结果对外发布；
  - (5) 每个用户根据服务器发布的预测结果，根据自定义的符号语义，提取出推荐信息。
- 本文后续内容对以上关键步骤进行详细阐述。

## 2.2. 评分信息的代换加密

协同过滤的基本原理是根据所有用户的历史评分信息来推荐用户可能感兴趣的项目。设用户集合为  $U = \{u_1, u_2, \dots, u_m\}$ ，项目集合  $I = \{i_1, i_2, \dots, i_n\}$ 。用户  $u$  对项目  $i$  的评分记为  $r_{ui}$ 。设项目评分符号明文空间为  $P$ ，任意元素  $p \in P$  具有指定的语义。以 Netflix 网站的电影评分为例，明文空间为  $P = \{1, 2, \dots, 5\}$ ，其中“1”表示非常不喜欢，“5”表示非常喜欢。

在隐私保护推荐框架下，用户首先在客户端自行定义自己的密文空间  $C$ ，其元素数量与明文空间相同，即  $|C|=|P|$ 。对于每个用户， $C$  中的元素完全由用户自主确定。例如，以电影评分为例，一种可能的明文-密文对代换映射为： $\{(1, \bullet), (2, \diamond), (3, \square), (4, \blacktriangle), (5, \blacksquare)\}$ 。由于每个用户的密文空间是自由设定的，因此确定这种映射关系的密钥是完全随机的。代换加/解密的另一个优点是计算代价极小，不会对客户端的性能产生明显的影响。

采用以上自定义的评分符号，用户对各项目进行评分，并上传至服务器。服务器在收集到所有用户的评分信息之后，即可构建用户 - 项目评分矩阵，如表 1 所示。

需要指出的是，推荐算法的预测结果也是用各用户自定义的评分符号来表示的，最终由用户自己在客户端提取推荐结果的语义。不可信的服务器由于不掌握密钥(映射关系)，因此无法推测用户的评分倾向。

**Table 1.** User-Item rating matrix  
**表 1.** 用户 - 项目评分矩阵

	$i_1$	$i_2$	.....	$i_n$
$u_1$	X	X		Y
$u_2$	#	&		&
$\vdots$				
$u_m$	▲	▲		◇

### 2.3. 用户相似度计算

为了预测用户  $u$  对项目  $i$  的评分  $r_{ui}$ ，协同过滤算法首先要确定与用户  $u$  最相似的  $K$  个邻居，然后综合这些邻居的评分，给出对  $r_{ui}$  的预测。这就涉及到用户相似度计算的问题。由于评分矩阵中所有评分的语义都被加密屏蔽，因此传统的相似度计算方法(例如余弦相似度、Tanimoto系数和皮尔逊相关系数等)都不再适用。

针对这种无语义条件下的相似度，本文提出模式相似度的概念及其计算方法：给定两个用户  $u_1$  和  $u_2$ ， $I_{u_1 \cap u_2}$  是两个用户共同评价的项目集合。 $R_{u_1}(I_{u_1 \cap u_2})$  和  $R_{u_2}(I_{u_1 \cap u_2})$  分别为  $u_1$  和  $u_2$  对  $I_{u_1 \cap u_2}$  的评分序列。用户  $u_1$  和  $u_2$  的模式相似度分别为：

$$s_{u_1, u_2} = \frac{1 + CR_{u_1, u_2}}{2} \times \left( 1 + \log_2 \left( |I_{u_1 \cap u_2}| \right) \right) \quad (1)$$

其中， $|I_{u_1 \cap u_2}|$  是两个用户共同评价的项目数量。 $CR_{u_1, u_2}$  是两个序列  $R_{u_1}(I_{u_1 \cap u_2})$  和  $R_{u_2}(I_{u_1 \cap u_2})$  之间的 CR 系数(Correct Rand Index) [8]。CR 系数是一种用来考察聚类算法性能的指标，其实质是度量两个聚类结果序列的匹配程度。根据给定的数据集，设  $C_1 = \{\lambda_1, \lambda_2, \dots, \lambda_n\}$ ， $C_2 = \{\gamma_1, \gamma_2, \dots, \gamma_l\}$ ，则两者的 CR 系数为：

$$CR_{u_1, u_2} = \frac{\sum_{h=1}^{|C_1|} \sum_{l=1}^{|C_2|} \binom{t_{hl}}{2} - \binom{|I_{u_1 \cap u_2}|}{2}}{\frac{1}{2} \left[ \sum_{h=1}^{|C_1|} \binom{t_h}{2} + \sum_{l=1}^{|C_2|} \binom{t_l}{2} \right] - \binom{|I_{u_1 \cap u_2}|}{2}} \quad (2)$$

其中  $t_{hl}$  表示被评分为  $\lambda_h$  和  $\gamma_l$  的项目数量， $t_h$  表示被  $u_1$  评分为  $\lambda_h$  的项目数量， $t_l$  表示被  $u_2$  评分为  $\gamma_l$  的项目数量。CR 系数的取值范围为  $[-1, 1]$ ，越接近 1，表明两个评分序列的一致程度越高。例如，设  $u_2 = \{X, X, Y, X, Z, Z\}$ ， $u_6 = \{\blacktriangle, \blacktriangle, \blacksquare, \blacktriangle, \diamond, \diamond\}$ ，那么这两个评分序列的 CR 系数为 1。

另外，我们认为  $|I_{u_1 \cap u_2}|$  越大，评分序列携带的信息量也越大。因此，在 CR 系数相同的情况下，序列长度越大，则用户相似度越高。从公式(1)可以看出，两个序列的一致性越高、序列长度越大，则它们的模式相似度越高。

计算任意两个用户的模式相似度，可以得到用户模式相似度矩阵：

$$S = \begin{pmatrix} 0 & s_{u_1, u_2} & \cdots & s_{u_1, u_m} \\ s_{u_2, u_1} & 0 & \cdots & s_{u_2, u_m} \\ \vdots & \vdots & \ddots & \vdots \\ s_{u_m, u_1} & s_{u_m, u_2} & \cdots & 0 \end{pmatrix}$$

主对角线上的元素被置为 0。显然， $S$  是一个对称阵。

## 2.4. 评分预测

为了预测用户  $u$  对项目  $i$  的评分  $r_{ui}$ ，首先根据用户模式相似度矩阵选择出  $u$  的  $K$  个邻居。每个邻居都对  $r_{ui}$  给出自己的预测，最后将  $K$  个结果聚合为最终结果。

给定用户  $u$  及其邻居  $u'$  的评分序列，首先要找出二者密文符号之间的对应关系，才能根据  $u'$  对  $i$  的评分来预测  $u$  的评分。这里，我们利用“熵”来确定这些对应关系。设  $A$  为一随机变量，熵的计算公式为：

$$H(A) = -\sum_{i=1}^n p(a_i) \log p(a_i) \quad (3)$$

其中  $p_{xi}$  表示  $A$  的值等于  $a_i$  的概率。熵越小，代表随机变量的确定性越高。

例如，设  $u_1$  和  $u_9$  的评分序列如表 2 所示。

对于  $u_9$  的每个评分符号，统计其对应  $u_1$  的各评分符号的频次，则可建立频次矩阵如表 3 所示。

对频次矩阵的每一行，计算相应评分符号的熵，可得  $H(\bullet)=0$ ， $H(\blacktriangle)=1$ ， $H(\blacksquare)=0.918$ 。

选择熵最小的行，并将行符号与频次最高的列符号建立对应关系。在本例中第一行熵最小，“ $\bullet$ ”对应符号“ $Z$ ”。在建立一对符号的对应关系后，对频次矩阵进行更新，并从频次矩阵中删除该行。重复以上步骤，直到频次矩阵的所有行均被删除，两个用户的评分符号间的一一对应关系即可建立。本例中，根据熵最小首先可确定“ $\bullet$ ”对应“ $Z$ ”，更新对应关系表后可确定“ $\blacktriangle$ ”对应“ $X$ ”、“ $\blacksquare$ ”对应“ $Y$ ”。最后，由于  $u_9$  对  $i_{32}$  的评分为“ $\bullet$ ”，而“ $\bullet$ ”对应“ $Z$ ”，因此  $u_9$  对  $i_{32}$  的预测为“ $Z$ ”。

根据以上预测方法， $K$  个邻居均给出对目标项目的评分预测。设  $K$  个邻居为  $N_1, N_2, \dots, N_K$ ，他们与用户  $u$  的模式相似度分别为  $s_{N_1,u}, s_{N_2,u}, \dots, s_{N_K,u}$ 。  $K$  个预测值为  $\hat{r}_1, \hat{r}_2, \dots, \hat{r}_K$ ，均属于用户  $u$  的密文空间  $C$ 。对  $C$  中的每个元素  $c$ ，统计在  $c$  在  $K$  个预测值中出现的加权次数：

$$\text{count}(c) = \sum_{k=1}^K (s_{N_k,u} \times l_k), \quad \text{其中 } l_k = \begin{cases} 1 & \text{if } \hat{r}_k = c \\ 0 & \text{其它} \end{cases}$$

最后，加权次数最多的元素即为最终的预测结果： $\hat{r}_{u,i} = \arg \max_{c \in C} (\text{count}(c))$ 。

具体实现算法流程如下：

输入：带有缺失值的用户项目评分矩阵  $R$ ，邻居数  $K$

输出：更新的评分矩阵  $R$

Table 2. The example of rating sequence

表 2. 评分序列示例

	$i_1$	$i_3$	$i_7$	$i_{15}$	$i_{18}$	$i_{20}$	$i_{23}$	$i_{39}$	$i_{90}$	$i_{32}$
$u_9$	X	X	Y	X	Z	X	X	X	X	?
$u_1$	$\bullet$	$\bullet$	$\blacktriangle$	$\bullet$	$\blacksquare$	$\blacksquare$	$\bullet$	$\blacksquare$	$\bullet$	$\bullet$

Table 3. Frequency matrix

表 3. 频次矩阵

	X	Y	Z
$\bullet$	0	0	4
$\blacktriangle$	1	0	1
$\blacksquare$	0	1	2

- 1) 对任意两个用户, 根据式(1)计算用户模式相似度, 得到用户模式相似度矩阵  $S$ ;
- 2) for 每个待预测评分  $r_{ui}$ {
- 3) 在  $S$  中选出与  $u$  具有最大相似度的  $K$  个邻居;
- 4) 根据式(2)确定  $u$  和  $N_k$  的评分符号对应关系, 并给出预测值  $\hat{r}_k$ ;
- 5)  $\hat{r} = \hat{r} \cup \hat{r}_k$ ;
- 6) 对  $u$  的密文空间  $C$  中任意元素  $c$ ;
- 7) 根据式(3)统计  $c$  在  $\hat{r}$  中的加权次数  $\text{count}(c)$ ;
- 8) 根据式(4)输出最终预测  $\hat{r}_{u,i}$  并填充入  $R$ ;
- 9) 输出  $R$ ;

根据算法输出的用户项目评分矩阵  $R$  对外发布。用户自行读取  $R$  中的相应行, 然后对预测评分进行逆置换, 即可读取推荐信息。

### 3. 实验及结果

本文采用 MovieLens 数据集和 Netflix 数据集对所提出的方法进行性能评估。MovieLens 数据集包含 943 名用户对 1682 部电影的 100,000 个评分, 每个用户至少参与 20 部电影的评分。Netflix 数据集从 Netflix Prize 数据集中抽样而得, 包括 1000 个用户对 712 部电影的评分。在每组实验中, 我们取 80% 数据为训练数据, 20% 为测试数据, 进行 5 折交叉验证。我们以传统的非隐私保护 CF 算法为对比基准, 对本文提出的隐私保护 CF 算法在预测准确率(Precision)和平均绝对误差(MAE)两个方面进行了测试和比较。预测准确率为正确预测的项目数与总预测数的比值。平均绝对误差定义是所有单个观测值与算术平均值的偏差的绝对值的平均, 在本实验中, 表示预测值与实际值的平均差。另外, 在传统 CF 算法中, 我们分别采用了余弦相似度、Tanimoto 系数和皮尔逊相关系数作为用户相似度的度量,  $K$  的取值为 5, 10,  $\dots$ , 50。

对 MovieLens 数据集的实验结果如图 2 和图 3 所示。由图 2 可以看出, 预测准确率随着  $K$  值的增大逐渐提高, 当  $K$  大于 30 后则趋于稳定, 说明相似度高的邻居越多, 预测准确率越高。另外, 隐私保护 CF 算法在预测准确率上(大于 0.4)明显高于传统 CF 算法(小于 0.38), 这说明本文提出的基于密文序列匹配的用户相似度计算方法以及评分预测算法在预测准确率上比传统 CF 算法具有更好的性能。但是, 从平均绝对误差来看, 如图 3 所示, 隐私保护 CF 算法的 MAE 比传统 CF 算法略高一些, 这是因为在隐私

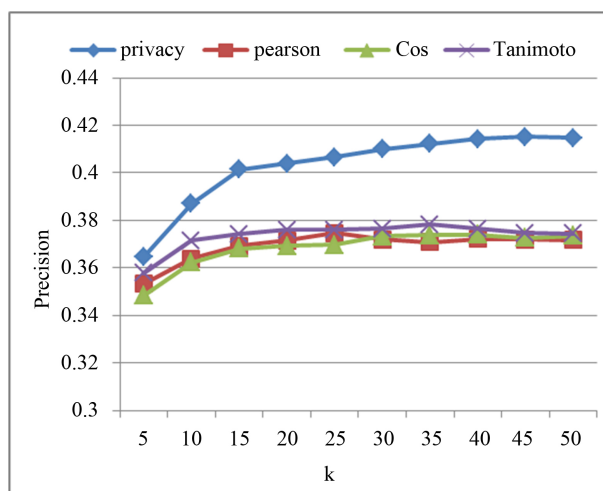


Figure 2. Comparison of predicting accuracy (MovieLens)

图 2. 预测准确率对比(MovieLens)

保护 CF 算法中, 评分信息的含义被屏蔽, 使得用户之间的评分失去的可比较性, 从而导致总体预测误差增大。从实验结果来看, 隐私保护 CF 算法的 MAE 平均比传统 CF 算法高 5%~7%, 这是为实现隐私保护而付出的必要的和可接受的代价。

实验在 Netflix 数据集上得到了相似的结果(如图 4 和图 5 所示)。隐私保护 CF 算法在预测准确率上优于传统 CF 算法, 但平均绝对误差则高于传统 CF 算法。例如当  $K = 30$  时, 隐私保护 CF 算法获得最好的性能, 其准确率比传统 CF 算法高约 8%~10%, 但平均绝对误差也比传统 CF 算法大 8%~10%。从总的趋势来看, 随着  $K$  值的增长, 预测准确率逐渐增长, 平均绝对误差则逐渐降低, 说明推荐算法的性能越来越好, 当  $K$  值增长到 20~25 时趋于平稳。

在两个数据集上的一致结果证明了隐私保护 CF 算法的有效性和鲁棒性。

#### 4. 结束语

本文基于推荐系统服务器不可信的前提, 提出了一种实现隐私保护的协同过滤推荐框架及其实现算

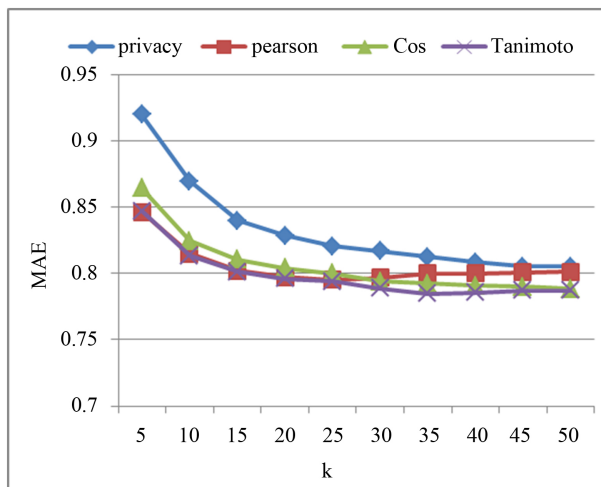


Figure 3. Comparison of MAE (MovieLens)

图 3. MAE 对比(MovieLens)

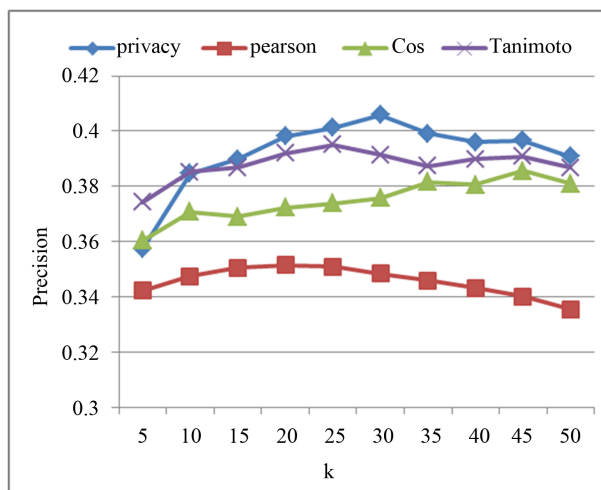


Figure 4. Comparison of predicting accuracy (Netflix)

图 4. 预测准确率对比(Netflix)



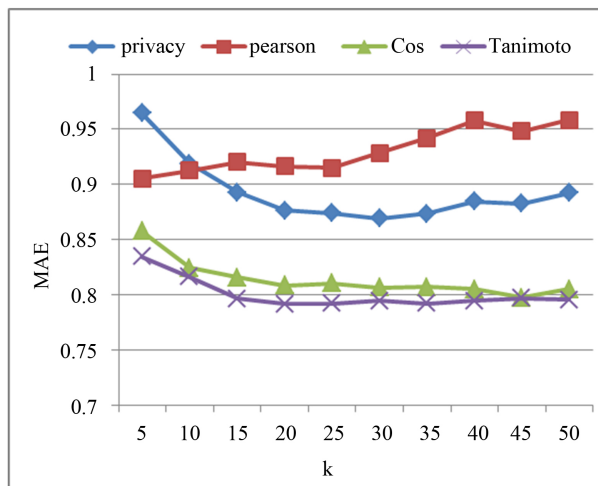


Figure 5. Comparison of MAE (Netflix)

图 5. MAE 对比(Netflix)

法。用户通过定义自己专用的代换密码机制，实现对评分信息的加密。推荐服务器则通过本文提出的隐私保护推荐算法来预测用户的项目的评分。最后用户再根据推荐服务器的预测结果来提取推荐信息。本文提出的隐私保护协同过滤推荐方法的优点在于，用户评分的真实含义仅能由用户自己定义和读取，推荐服务器或其他用户无法根据用户的评分来推测用户的偏好特征。同时，该方法能够以较小的计算代价和通信代价来实现隐私保护和精确推荐。本文提出的基于隐私保护的协同过滤推荐算法能够在保护用户隐私的前提下广泛应用于多种评分系统和推荐系统，具有广阔的应用前景。

## 基金项目

湖北省自然科学基金项目(2014CFB354)。

## 参考文献 (References)

- [1] Melville, P. and Sindhvani, V. (2011) Recommender Systems. In: *Encyclopedia of Machine Learning*, Springer, 829-838.
- [2] 项亮. 推荐系统实践[M]. 北京: 人民邮电出版社, 2012.
- [3] Calandrino, J.A., Kilzer, A., Narayanan, A., Felten, E.W. and Shmatikov, V. (2011) You Might Also Like: Privacy Risks of Collaborative Filtering. *IEEE Symposium on Security & Privacy*, Oakland, 11 December 2011, 231-246.
- [4] Kobsa, A. (2007) Privacy-Enhanced Web Personalization. *Lecture Notes in Computer Science*, **4321**, 628-670. [http://dx.doi.org/10.1007/978-3-540-72079-9\\_21](http://dx.doi.org/10.1007/978-3-540-72079-9_21)
- [5] Bilge, A. and Polat, H. (2013) A Scalable Privacy-Preserving Recommendation Scheme via Bisecting k-Means Clustering. *Information Processing & Management*, **49**, 912-927. <http://dx.doi.org/10.1016/j.ipm.2013.02.004>
- [6] Zhu, T., Li, G., Zhou, W., Xiong, P. and Yuan, C. (2015) Privacy-Preserving Topic Model for Tagging Recommender Systems. *Knowledge & Information Systems*, **46**, 1-26.
- [7] Shokri, R., Pedarsani, P., Theodorakopoulos, G. and Hubaux, J.-P. (2009) Preserving Privacy in Collaborative Filtering through Distributed Aggregation of Offline Profiles. *Proceedings of the Third ACM Conference on Recommender systems*, **2**, 157-164.
- [8] Hubert, L. and Arabie, P. (1985) Comparing Partitions. *Journal of Classification*, **2**, 193-218. <http://dx.doi.org/10.1007/BF01908075>



**期刊投稿者将享受如下服务：**

1. 投稿前咨询服务 (QQ、微信、邮箱皆可)
2. 为您匹配最合适的期刊
3. 24 小时以内解答您的所有疑问
4. 友好的在线投稿界面
5. 专业的同行评审
6. 知网检索
7. 全网络覆盖式推广您的研究

投稿请点击：<http://www.hanspub.org/Submission.aspx>