

协同过滤推荐算法综述

Survey of Collaborative Filtering Recommendation Algorithm

黄正 HUANG Zheng

(华南理工大学计算机科学与工程学院, 广州 510640)

(School of Computer Science and Engineering, South China University of Technology, Guangzhou 510640, China)

摘要: 推荐系统是电子商务系统中最重要技术之一, 协同过滤推荐技术是目前应用最广泛和最成功的推荐技术。本文首先介绍了协同过滤的基本概念和原理, 然后总结了协同过滤推荐算法中的关键问题和相关解决方案, 最后介绍了协同过滤推荐算法需要进一步解决的问题和可能的发展方向。

Abstract: Recommendation system is one of the most important technologies in E-commerce. Collaborative filtering is the most widely used and the most successful recommendation technology. This paper first introduces the basic concept and principle of collaborative filtering. And then, this paper summarizes the key problems and related solutions of the collaborative filtering recommendation algorithm. Finally, this paper introduces the collaborative filtering recommendation algorithm need to be further solved problems and possible development direction.

关键词: 协同过滤; 推荐算法; 稀疏性; 扩展性

Key words: collaborative filtering; recommendation algorithm; sparsity; scalability

中图分类号: TP39

文献标识码: A

文章编号: 1006-4311(2012)21-0226-03

0 引言

随着电子商务和互联网的迅速发展, 人们能够获得的资讯越来越多, 但是想快捷地获得自己感兴趣的资讯越来越难。近年来兴起的个性化推荐系统成为解决这些问题的一个重要途径。个性化推荐系统很多, 协同过滤推荐是当前最成功的推荐技术^[1]。

协同过滤的概念是由 Goldberg、Nicolis、Oki 以及 Terry 在 1992 年首次提出的^[2], 应用于 Tapestry 系统, 其基本思想是: 通过对用户的显式输入或隐式输入的历史数据收集并统计计算, 预测与此用户兴趣相似的用户, 并将其相似用户感兴趣的项目推荐给此用户。

目前, 主要的协同过滤推荐算法有两类: 基于用户的协同过滤推荐算法和基于项目的协同过滤推荐算法。本文对协同过滤推荐算法中的关键技术进行了详细的介绍, 并对协同过滤推荐算法中存在的问题进行了分析, 同时也介绍了一些解决问题的改进方法和算法评估方法, 最后对协同过滤推荐算法的发展进行了展望。

1 传统的协同过滤推荐算法

传统的协同过滤推荐算法可以分为三个步骤:

1.1 构建用户-项目评价矩阵 用户评分数据可以用一个 $m \times n$ 的用户-项目评价矩阵 R_{mn} 来表示 (如图 1 所示)。其中, m 行表示用户的数量, n 列表示项目的数量, 第 i 行第 j 列元素 R_{ij} 表示用户 i 对项目 j 的评分。

User/Item	I_1	I_2	...	I_j	...	I_n
U_1	R_{11}	R_{12}	...	R_{1j}	...	R_{1n}
U_2	R_{21}	R_{22}	...	R_{2j}	...	R_{2n}
...
U_i	R_{i1}	R_{i2}	...	R_{ij}	...	R_{in}
...
U_m	R_{m1}	R_{m2}	...	R_{mj}	...	R_{mn}

图 1 用户-项目评价矩阵

1.2 最近邻居集的形成 这一步骤是基于用户的协同

过滤推荐算法的核心步骤。对目标用户 u , 算法找到与其相似度最高的 K 个用户, 这 K 个用户就是该目标用户的最近邻居集合 $N(u) = \{u_1, u_2, \dots, u_k\}$, $u \notin N(u)$ 且 $N(u)$ 中的用户 u_k 按其与用户 u 之间的相似度 $\text{sim}(u_k, u)$ ($1 \leq k \leq K$) 由大到小排序。

传统的计算用户相似度的方法有三种^[3]:

①余弦相似性 (Cosine Similarity): 把用户对项目的评价看做 n 维项目空间的向量, 两个用户间的相似性通过二者向量夹角的余弦度量。设用户 i 和用户 j 在 n 维项目空间的评分分别用向量 \vec{i} 和向量 \vec{j} 表示, 则两者间的相似性为:

$$\text{sim}(i, j) = \cos(i, j) = \frac{\vec{i} \cdot \vec{j}}{\|\vec{i}\| \times \|\vec{j}\|}$$

②相关相似性 (Correlation Similarity): 相关相似性也称皮尔森 (Pearson) 系数相关, 通过皮尔森 (Pearson) 相关系数来度量用户间的相似性。设用户 i 和用户 j 共同评分的项目集合用 I_{ij} 表示, 则两者间的相似性为:

$$\text{sim}(i, j) = \frac{\sum_{c \in I_{ij}} (R_{i,c} - \bar{R}_i)(R_{j,c} - \bar{R}_j)}{\sqrt{\sum_{c \in I_{ij}} (R_{i,c} - \bar{R}_i)^2} \sqrt{\sum_{c \in I_{ij}} (R_{j,c} - \bar{R}_j)^2}}$$

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 与 \bar{R}_j 分别表示用户 i 和用户 j 对项目的平均评分。

③修正的余弦相似性 (Adjusted Cosine Similarity): 由于余弦相似性度量方法中没有考虑到用户不同的评价尺度问题, 在修正的余弦相似性度量方法中通过减去用户对项目的平均评分来改善上述缺陷。设用户 i 和用户 j 共同评分的项目集合用 I_{ij} 表示, I_i 和 I_j 分别表示用户 i 和用户 j 评分的项目集合, 则两者间的相似性为:

$$\text{sim}(i, j) = \frac{\sum_{c \in I_{ij}} (R_{i,c} - \bar{R}_i)(R_{j,c} - \bar{R}_j)}{\sqrt{\sum_{c \in I_i} (R_{i,c} - \bar{R}_i)^2} \sqrt{\sum_{c \in I_j} (R_{j,c} - \bar{R}_j)^2}}$$

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 与 \bar{R}_j 分别表示用户 i 和用户 j 对项目的平均评分。

1.3 产生推荐 目标用户的“最近邻居”集产生后, 可

作者简介: 黄正 (1983-), 男, 湖北武汉人, 研究方向为计算机应用。

以计算两类结果:用户对任意项的兴趣度的预测值和 Top-N 形式的推荐集。

①目标用户对任意项的兴趣度的预测值。目前大部分协同过滤推荐系统都采用平均加权策略来产生推荐,目标用户 i 对项目 c 的预测评分为:

$$P_{i,c} = \bar{R}_i + \frac{\sum \text{sim}(i,j) \times (R_{j,c} - \bar{R}_j)}{\sum |\text{sim}(i,j)|}$$

其中, $\text{sim}(i,j)$ 为用户 i 和用户 j 的相似度, $R_{j,c}$ 为最近邻居中的用户 j 对项目 c 的评分, \bar{R}_i 与 \bar{R}_j 分别表示用户 i 和用户 j 对项目的平均评分。

②Top-N 形式的推荐集。分别统计“最近邻居”集中的用户 i 对不同项的兴趣度的加权平均值,取其中 N 个排在最前面且不属于 I_i (I_i 表示用户 i 评分的项目集合)的项作为 Top-N 推荐集。

2 协同过滤推荐算法存在的问题及解决方法

协同过滤推荐技术在个性化推荐系统中取得了巨大的成功,并得到了广泛的应用。然而,随着互联网的发展和普及,电子商务系统规模的扩大,站点结构、内容的复杂程度和用户人数的不断增加,以及协同过滤本身的特点,协同过滤推荐面临着一些难以解决的问题,比较典型的有数据稀疏问题、冷启动问题、算法的可扩展性等问题。

2.1 数据稀疏性问题 这是推荐系统面临的主要问题,也是导致系统推荐质量下降的主要原因。在众多推荐系统中,用户评分过的项目数量往往很有限,相关研究指出这个比例在 3% 左右^[4],使得用户-项目评分矩阵极端稀疏性,导致用户最近邻居和项目最近邻居的计算准确性降低,使得推荐系统的推荐质量急剧下降。

目前,解决数据稀疏性问题的一种方法是矩阵填充技术。最简单的填充办法就是将用户未评分项目的评分设定为一个固定的缺省值,或者设置为其他用户对该项目的平均评分^[5]。然而,用户对为评分项目的评分不可能完全相同,这种方法也就不能从根本上解决稀疏性问题。比较理想的方法是采用预测评分对用户-项目评分矩阵进行填充,主要有以下几种:BP 神经网络、Naive Bayesian 分类方法、基于内容的预测。

另一种解决数据稀疏性问题的方法是矩阵降维技术。通过降低用户-项目评分矩阵的维数,删除那些不重要的或噪音用户和项目,将两个用户投影到一个低维空间上,计算两者间的相似度,提高了算法的效率。比较典型的降维技术是奇异值分解 (Singular Value Decomposition, SVD) 技术,将大小为 $m \times n$ (假设 $m \geq n$) 的用户-项目评分矩阵 R 分解为大小分别为 $m \times m$, $m \times n$, $n \times n$ 的三个矩阵。

2.2 冷启动问题 冷启动问题包括用户冷启动问题 (New user problem) 和项目冷启动问题 (New item problem),可以说是数据稀疏性问题的极端情况。①在基于用户的协同过滤推荐系统中,对于一个新加入的用户来说,系统中没有该用户的任何浏览或购买信息,甚至连该用户的浏览信息都没有,因此无法找到该用户的最近邻居,从而无法对其进行推荐。②在基于项目的协同过滤推荐系统中,当加入一个新的项目,由于系统中没有对该项目的评分,从而无法找出该项目的最近邻居,也就无法将该项目推荐出去。而且,在新项目出现早期,用户评价较少,推荐的准确性也比较差^[6]。

为了解决冷启动问题,普遍采用基于内容的最近邻居查找技术^[7],其基本思想是:①利用聚类技术将用户按照属性相似性聚类,从项目属性的角度找到新项目的最近邻居;②用新项目 k 的所有最近邻居的平均评分来代替已有评分的平均值 \bar{R}_k 。

2.3 可扩展性问题 传统的协同过滤推荐算法的运算时间是随着用户和项目数目的增加而急剧增长的,面对越来越多的用户和项目数据,传统的推荐系统将面临严重的可扩展性问题。通常采用聚类技术来解决这一问题,典型的方法有以下几种: k -means 聚类算法、EM (Expectation-Maximization) 算法、Gibbs Sampling 方法和模糊聚类。

上面讨论的都是目前协同过滤推荐系统面临的几个关键问题,事实上,推荐系统还面临着许多其他问题,例如统一性问题、安全问题、隐私问题等。虽然很多研究者对传统的协同过滤推荐系统作了改进,但都很难从根本上解决问题,这也是未来研究协同过滤推荐系统的重点。

3 算法性能评估标准

个性化推荐系统的性能是由用户对其推荐结果的满意度决定的。评价推荐系统推荐质量的度量主要包括统计精度度量方法和决策支持精度度量方法两类^[8]。统计精度度量方法中常用的是平均绝对偏差 MAE (Mean Absolute Error); 决策支持精度度量方法中主要有召回率 (Recall)、准确率 (Precision) 等两种方法。

3.1 平均绝对偏差 MAE 通过计算用户或项目的预测评分与实际评分之间的偏差量来衡量推荐结果的准确度, MAE 越小,推荐的质量就越高。假设预测的用户评分集合表示为 $\{p_1, p_2, \dots, p_N\}$, 对应的实际用户评分集合为 $\{q_1, q_2, \dots, q_N\}$, 则平均绝对偏差 MAE 定义^[9]为: $MAE = \frac{\sum_{i=1}^N |p_i - q_i|}{N}$ 。

3.2 召回率 (Recall) 反映的是待推荐项目被推荐的比率: $\text{recall} = \frac{|\text{test} \cap \text{top-N}|}{|\text{test}|}$ 。

其中, test 表示测试数据集中的项目数量, top-N 表示系统推荐给用户的 N 个项目。

3.3 准确率 (Precision) 表示算法推荐成功的比率,即预测准确的项目总数与预测的所有项目总数之比:

$$\text{precision} = \frac{|\text{test} \cap \text{top-N}|}{N}$$

其中, test 表示测试数据集中的项目数量, top-N 表示系统推荐给用户的 N 个项目。

4 总结及展望

本文围绕协同过滤推荐这一主题,对协同过滤推荐算法进行了详细的介绍,对存在的数据稀疏性问题、冷启动问题和可扩展性问题进行了深入的研究,并介绍了一些经典的解决方法。最后对算法的评估方法作了详细的介绍。

今后,如何解决协同过滤推荐算法存在的问题将是这一课题研究的重点,研究者也提出了很多改进方法。但是,这些方法只能在一定程度上或某些特定场合选才能有效解决问题。随着电子商务的进一步发展,这些问题必定会进一步凸显。而且,协同过滤推荐还会面临一些新的挑战,例如安全问题、隐私问题等。另外,随着推荐算法的层出不穷,对算法的评估标准的研究也会越来越重要。同时,将其

基于以太网的雷达发射机控制与保护系统设计

The Design of System to Control and Protect Radar Transmitter Based on Ethernet

洪远刚 HONG Yuan-gang; 肖剑 XIAO Jian; 邓凤军 DENG Feng-jun

(中国卫星海上测控部, 江阴 214431)

(China Department of Satellite Marine Track & Control, Jiangyin 214413, China)

摘要: 随着大功率发射机技术的发展, 发射机的控制和保护要求也越来越高, 各模块的参数采集与传输的实时性、高数据率和可靠性要求越来越高, 传统的基于串行通讯的控保系统面临挑战。为了解决上述问题, 本文设计了一种基于以太网的控保系统, 在工程实践中, 取得了很好的应用效果。

Abstract: With the development of high power radar transmitter technology, the requirement to control and protect radar transmitter has become higher and higher, and the real-time and high-data and high-credibility requirement of parameter collection begin higher and higher. The routine design method of Serial/Parallel communication faces challenge. In order to solve the above problem, the composition of Ethernet and its application to the control and protection system of Radar transmitter are discussed in this paper, and the good effect have gained in engineering practice.

关键词: 雷达发射机; 以太网; 发射机控制和保护

Key words: radar transmitter; Ethernet; transmitter control and protection

中图分类号: TN959.51

文献标识码: A

文章编号: 1006-4311(2012)21-0228-02

0 引言

在发射机的工作过程中, 需要对各种参数进行监测, 以便系统操作员及时了解发射机的工作状况。一旦发现故障, 需自动或人工干预, 并及时进行保护。参数监测一般是通过发射机监控台上的各种电压表或电流表、指示灯或示波器等仪表显示。现代发射机大多遥控操作, 为便于主控制操作员及时了解发射机的工作状况, 传统的控保系统通常采用串口通信接口将被监测参数传送至雷达主控制台^[1], 并通过串口接收主控制台发送来的各个指令。由于采用的是基于串口通信的信息传输模式, 系统的数据率及可靠性不高。

随着发射机技术的发展, 发射机数字化、集成化、模块化和自动化程度越来越高, 需要监测和控制参数的种类与个数有时可能达到几千个以上; 传统的串口通信是一种适合于点对点、数据率较低的通讯方式, 特别是多套发射机并行工作时, 更凸现了其与现代发射机所要求的高数据率已不相适应。高发射功率也是发射机的趋势, 随之带来的

高压、高温和高电磁辐射严重干扰了基于串口通信控保系统的可靠性和有效性。近来, 高速发展的以太网以其成熟的技术、高效的组网效率和高可靠性被成功地应用于各种处于严酷环境的工业自动化系统的组网设计中, 本文设计的基于以太网的发射机控保系统, 充分发挥了以太网的优势, 在工程上取得了良好的效果。

1 系统组成

发射机控保系统按其功能可分为监控单元、数据传输单元、主控制台三部分组成, 如图 1 所示。

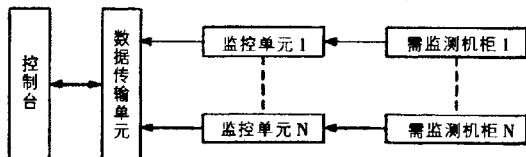


图 1 基于以太网的控保系统组成框图

1.1 监控单元 目前, 发射机正往高度集成化、数字化和模块化方向发展, 需要监测的参数随之大幅度增多, 使用单片机已不能满足要求, 必须采用分散检测和数据采集、集中显示的多机分层检测方案。即将所有检测点分为若干个监控单元, 每个单元有自己的 CPU, 完成本单元的

作者简介: 洪远刚(1986-), 男, 四川绵阳人, 助理工程师, 学士, 研究方向为雷达测试与信号处理。

他领域技术与推荐技术相结合也将会成为未来的一个重点研究方向。

参考文献:

- [1] HERLOCKER J, KONSTAN J, TERVEEN L, et al. Evaluating collaborative filtering recommendation systems [J]. ACM Transactions on Information Systems, 2004, 22(1): 5-53.
- [2] Goldberg D, Nichols D, Oki B M, et al. Using collaborative filtering to weave an information tapestry [J]. Communications of the ACM, December, 1992, 35(12): 61-70.
- [3] Lu Pei. An Improved Project Clustering-Based Collaborative Filtering Recommendation Algorithm [J]. Public Communication of Science & Technology, 2011, No. 1: 205-206.
- [4] 孙小华. 协同过滤推荐系统的稀疏性与冷启动问题研究 [D]. 浙江大学, 2006.

- [5] Deng Ai-lin, Zhu Yang-yong, Shi Bai-le. A collaborative filtering recommendation algorithm based on item rating prediction [J]. Journal of Software, 2003, 14(9): 1621-1628.
- [6] 许敏, 邱玉辉. 电子商务中推荐系统存在的问题及其对策研究 [J]. 计算机科学, 2001, 28(4): 122-124.
- [7] Ma Hong-wei, Zhang Guang-wei, Li Peng. Survey of Collaborative Filtering Algorithms [J]. Journal of Chinese Computer Systems, 2009, Vol. 30, No. 7: 1282-1288.
- [8] 吴发青, 贺樑, 夏薇薇, 等. 一种基于用户兴趣局部相似性的推荐算法 [J]. 计算机应用, 2008, 28(8): 1981-1985.
- [9] Sarwar B, Karypis G, Konstan J, et al. Item-based collaborative filtering recommendation algorithms [C]// Proceedings of the 10th International World Wide Web Conference, 2001: 285-295.

协同过滤推荐算法综述

作者: [黄正, HUANG Zheng](#)
作者单位: [华南理工大学计算机科学与工程学院, 广州, 510640](#)
刊名: [价值工程](#) **ISTIC**
英文刊名: [Value Engineering](#)
年, 卷(期): 2012, 31 (21)

参考文献(9条)

1. [HERLOCKER J;KONSTAN J;TERVEEN L](#) [Evaluating collaborative filtering recommend systems](#) 2004(I)
2. [Goldberg D;Nichols D;Oki B M](#) [Using collaborative filtering to weave an information tapestry](#)[外文期刊] 1992(12)
3. [Lu Pei](#) [An Improved Project Clustering-Based Collaborative Filtering Recommendation Algorithm](#) 2011(01)
4. [孙小华](#) [协同过滤推荐系统的稀疏性与冷启动问题研究](#) 2006
5. [Deng Ai-lin;Zhu Yang-yong;Shi Bal-le](#) [A collaborative filtering recommendation algorithm based on item rating prediction](#)[期刊论文]-[Journal of Software](#) 2003(09)
6. [许敏;邱玉辉](#) [电子商务中推荐系统存在的问题及其对策研究](#)[期刊论文]-[计算机科学](#) 2001(04)
7. [Ma Hong-wei;Zhang Guang-wei;Li Peng](#) [Survey of Collaborative Filtering Algorithms](#)[期刊论文]-[Journal of Chinese Computer Systems](#) 2009(07)
8. [吴发青;贺操;夏薇薇](#) [一种基于用户兴趣局部相似性的推荐算法](#)[期刊论文]-[计算机应用](#) 2008(08)
9. [Sarwar B;Karypis G;Konstan J](#) [Item -baaed collaborative fil-tering recommendation algorithms](#) 2001

本文链接: http://d.wanfangdata.com.cn/Periodical_jzgc201221110.aspx