

融合知识图谱表示学习的栈式自编码器推荐算法

王卫红 冯倩 吕红燕 曹玉辉

(河北经贸大学信息技术学院 河北 石家庄 050061)

摘要 针对目前协同过滤推荐算法中数据稀疏和语义信息欠缺问题,提出一种融合知识图谱表示学习的栈式自编码器推荐算法(SAEKG-CF)。将评分矩阵作为栈式自编码器的输入,训练得到项目的隐性特征向量,并据此计算特征相似性矩阵;利用知识图谱表示学习算法将项目中的实体映射到低维向量空间,并计算出低维向量空间中实体间的语义相似性矩阵;将特征相似性矩阵与语义相似性矩阵相融合,得到融合相似性矩阵,进而依据最优融合相似性矩阵产生 top-k 推荐列表。实验结果表明,该算法能有效地同时解决数据稀疏与语义信息欠缺问题,提高推荐的准确率。

关键词 协同过滤 栈式自编码器 知识图谱 推荐系统

中图分类号 TP3 文献标志码 A DOI:10.3969/j.issn.1000-386x.2021.02.043

RECOMMENDATION ALGORITHM BASED ON REPRESENTATION LEARNING OF KNOWLEDGE GRAPH AND STACK AUTOENCODER

Wang Weihong Feng Qian Lü Hongyan Cao Yuhui

(School of Information Technology, Hebei University of Economics and Business, Shijiazhuang 050061, Hebei, China)

Abstract Aiming at the problem of data sparseness and lack of semantic information in the current collaborative filtering recommendation algorithm, a stacked autoencoder recommendation algorithm based on knowledge graph representation learning (SAEKG-CF) is proposed. It used the rating matrix as the input of stack self-encoder, trained the implicit feature vector of the project, and then calculated the feature similarity matrix; the knowledge graph representation learning algorithm was used to map the entities in the project to the low-dimensional vector space, and calculated semantic similarity matrix between entities; the feature similarity matrix was merged with the semantic similarity matrix to obtain a fusion similarity matrix; according to the optimal fusion similarity matrix, a top-k recommendation list was generated. The experimental results show that the proposed algorithm can effectively solve the problem of sparse data and lack of semantic information, and improve the accuracy of recommendation.

Keywords Collaborative filtering Stack autoencoder Knowledge graph Recommendation system

0 引言

爆炸性增长的信息数据,为人们的工作、生活、学习提供巨大便利的同时也带来了信息迷航^[1]与信息过载^[2]问题,为此推荐系统应运而生。推荐系统主要是通过对用户行为数据进行分析与处理、构建用户兴趣模型,从而主动向用户推荐最合适的商品,帮助用户快速地做出决策^[3]。目前,因其可以挖掘用户隐性兴趣

并提供个性化服务,已成为一个重要的研究领域。协同过滤算法因其通用性、易理解性等优势已成为推荐系统中使用最为广泛的算法之一,但实际应用中因其存在数据稀疏、语义信息欠缺等问题,导致该算法推荐性能较低^[4]。

因此,本文提出一种融合知识图谱表示学习的栈式自编码器推荐算法。该算法利用栈式自编码器获取项目深层特征,并与知识图谱表示学习算法得到的语义信息进行融合,能够同时解决数据稀疏性和语义信

息欠缺问题,得到更为准确的推荐结果。

1 相关工作

目前,国内外研究学者针对协同过滤推荐算法存在的数据稀疏性问题进行了大量研究。文献[5]综合考虑用户评分行为、评分偏好和属性特征,提出一种混合相似度计算模型;文献[6]采用EMD方法对传统的相似度计算进行改进,并与用户信任关系相融合;文献[7]将评分相似度、兴趣相似度、属性相似度与置信度相融合,更为准确地为用户找到所需项目,有效地解决了数据稀疏问题;文献[8]提出一种融合多源异构数据的矩阵分解模型,缓解数据稀疏性问题;文献[9]提出一种新的相似性强化机制来增强基于内存的协同过滤方法,从而解决数据稀疏性问题;文献[10]将属性进行聚类,利用相似特征选择过滤冗余属性,降低了数据稀疏性;文献[11]使用关联检索框架中的传播激活算法探索用户之间的关联关系,帮助构建用户兴趣模型有效处理稀疏性问题;文献[12]提出基于邻域的CF的相似性度量方法,提高处理稀疏数据的处理能力。上述研究从评分矩阵填充、相似性方法改进等多个角度探讨了解决数据稀疏性的方法,但未考虑语义信息问题。

语义信息对于推荐系统而言至关重要。2006年Loizou A博士在推荐系统研讨会(ECAI2006)上指出:传统的推荐算法由于没有考虑语义信息,使得在实际应用中,这些算法在实时性、鲁棒性和推荐质量等方面存在严重的不足^[13]。因此,将语义信息融入协同过滤推荐算法中,利用语义知识对用户兴趣和项目内容进行描述以提高算法的推荐质量便成为近年来协同过滤算法研究的一个重要方向。目前针对语义推荐算法的研究较少,文献[14]将基于用户显式数据的协同过滤算法与知识图谱表示学习方法相融合,来增强项目的语义信息;文献[15]将语义网中的本体技术运用到社交网络的服务推荐系统中,得到更加准确的推荐结果;文献[16]引入WordNet词汇结构,分析用户标签与项目之间的相似度,在语义方面更好地理解用户偏好;文献[17]提出混合多准则的语义增强协同过滤算法,具有较好的推荐效果;文献[18]构建基于语义网的推荐模型,应用于农业学习资源推荐。由此可知,这些探讨语义问题的研究工作中并未考虑数据稀疏性问题。

本文为了同时解决数据稀疏与语义信息欠缺问题,提出一种融合知识图谱表示学习的栈式自编码器推荐算法。

2 SAEKG-CF

2.1 基本思想

本文提出的算法旨在同时解决数据稀疏与语义信息欠缺问题,针对数据稀疏问题采用栈式自编码器进行降维和特征提取,针对语义信息欠缺问题采用知识图谱表示学习算法获得项目的知识语义信息,并将两者相结合,得到优化的推荐结果。

栈式自编码器^[19]是一种深度模型结构,通过多层编码器和解码器可以自动地抽取深层隐式特征,相当于输入数据的压缩表示,这种表示能更好地代替原始数据。栈式自编码器强大的特征提取能力能够有效地缓解数据稀疏性问题,但缺少对结果的语义解释。

知识图谱融合多源异构数据信息将用户与用户、用户与项目、项目与项目之间相互连接起来,且可以结合推理得到数据的语义信息^[20]。现有研究表明,知识图谱表示学习方法能将知识图谱中的实体和关系嵌入到一个低维语义空间,利用连续数值向量高效地计算实体间的语义联系^[21]。

因此,SAEKG-CF将栈式自编码器和知识图谱表示学习方法进行融合,既可以对栈式编码器产生的结果增加语义解释,又可以缓解协同过滤推荐算法中数据稀疏性和语义信息欠缺的问题。该算法的工作原理如图1所示。

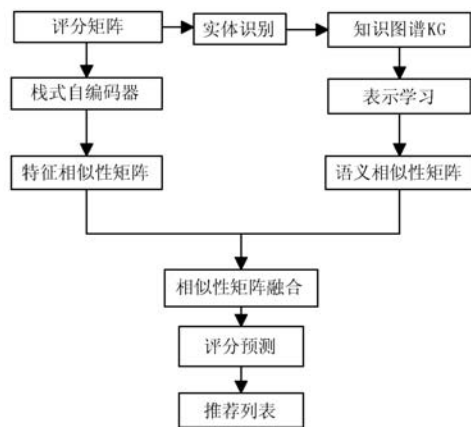


图1 SAEKG-CF 工作原理

2.2 算法设计

2.2.1 基于栈式自编码器的特征相似性度量

设用户个数为 M ,项目个数为 N ,用户对项目的评分矩阵记为 \mathbf{R}_{MN} ,用户集合为 $U = \{u_1, u_2, \dots, u_i, \dots, u_M\}$,项目集合为 $V = \{v_1, v_2, \dots, v_i, \dots, v_N\}$,对于任意一个项目,根据用户对项目的评分可产生一个项目向量 $\mathbf{R}_j = (r_{1j}, r_{2j}, \dots, r_{Mj})$ 。本文算法将原有的三层自编码器网络结构设置为七层的栈式自编码器,其中包括一个

输入层,五个隐含层和一个输出层,其结构如图 2 所示。

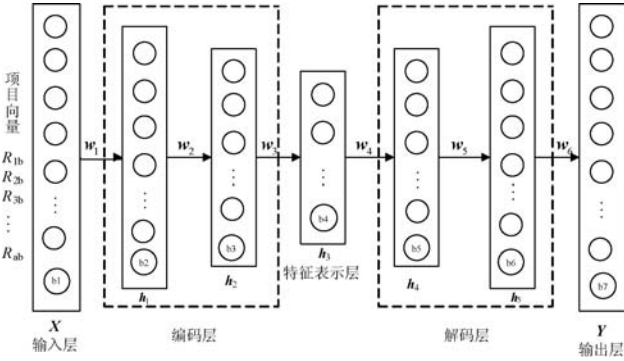


图 2 栈式自编码器结构图

可以看出,前三层为编码阶段,后三层为解码阶段,中间为特征表示层,各层之间采用的是全连接,各层之间的输出为:

$$h_i = f(w_i h_{i-1} + b_i) \quad (1)$$

式中: w_i 为权重矩阵; b_i 为偏置向量; $f(\cdot)$ 为激活函数。为防止过拟合,加入了正则项的损失函数:

$$L = \sum_{i=1}^N (h_{1i} - h_{7i})^2 + \frac{\lambda}{2} \sum (\|w_k\|_2^2 + \|b_k\|_2^2) \quad (2)$$

式中: h_{1i} 为输入向量的第 i 个分量; h_{7i} 表示重构后向量的第 i 个分量; λ 为正则化系数。

本文算法通过最小化损失函数来训练整个网络,输出训练好的隐含层数据 h_4 记为 R_{v_i} ,表示项目的特征向量。根据得到的项目特征向量进行相似性度量,采用余弦相似度计算:

$$sim_1(v_a, v_b) = \frac{R_{v_a} \cdot R_{v_b}}{|R_{v_a}| \times |R_{v_b}|} \quad (3)$$

2.2.2 基于知识图谱表示学习的语义相似性度量

知识图谱的表示学习算法能够将知识图谱三元组映射到低维向量空间,从而实现数值表示。本文使用在语义表达方面具有较好性能的 TransE^[22] 算法,来增强用户-项目评分矩阵中项目的语义信息。

TransE 通过不断调整头实体向量 h 、关系向量 r 和尾实体向量 t ,使得 $h + r$ 尽可能地与 t 相等,即 $\|h + r\| \approx \|t\|$, $\|\cdot\|$ 可以选择 L1 或者 L2 范数,其模型如图 3 所示。

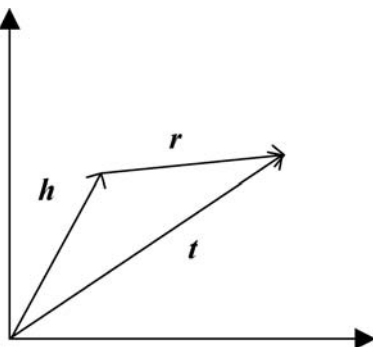


图 3 TransE 模型

TransE 的优化目标函数定义为:

$$L = \sum_{(h,r,t) \in S} \sum_{(h',r',t') \in S'(h,r,t)} [\gamma + \|h + r - t\| - \|h' + r' - t'\|] \quad (4)$$

式中: γ 是间距参数; S 为知识图谱中三元组集合; S' 为负例样本三元组组合,该类负例样本由 S 中的头实体或者尾实体随机替换成其他实体所得。

TransE 算法在处理单一关系上具有高效性,但在多关系上存在局限性^[23],因此,本文算法采用随机梯度下降进行不断迭代,使目标函数达到最优,针对知识图谱中存在的多关系实体进行单独训练,得到每个关系对应的该实体向量,取均值作为该实体的语义向量,从而改善 TransE 算法在多元关系上的性能。将项目语义向量表示为:

$$V = (e_1, e_2, \dots, e_d)^T \quad (5)$$

式中: e_n 表示项目 V 语义向量在第 n 维上的值。

由于在对实体进行映射时,采用的是欧氏距离,因此为了准确描述项目间的语义相似性,本文算法采用同等范数的欧氏距离来计算项目之间相似性,计算公式如下:

$$sim_2(v_a, v_b) = \frac{1}{1 + \sqrt{\sum_{k=1}^d (e_k^{v_a} - e_k^{v_b})^2}} \quad (6)$$

可以看出计算结果越接近于 1,两个项目之间的语义相似性越高;反之计算值越接近于 0,说明两个项目之间关系疏远。

2.2.3 融合项目相似性计算与评分预测

(1) 融合相似性矩阵的计算。基于评分矩阵,通过式(3)和式(6)可以得到项目的特征相似性矩阵和语义相似性矩阵,然后对两个矩阵进行融合,依据推荐结果的准确率得到最优融合相似性矩阵,原理示意如图 4 所示。

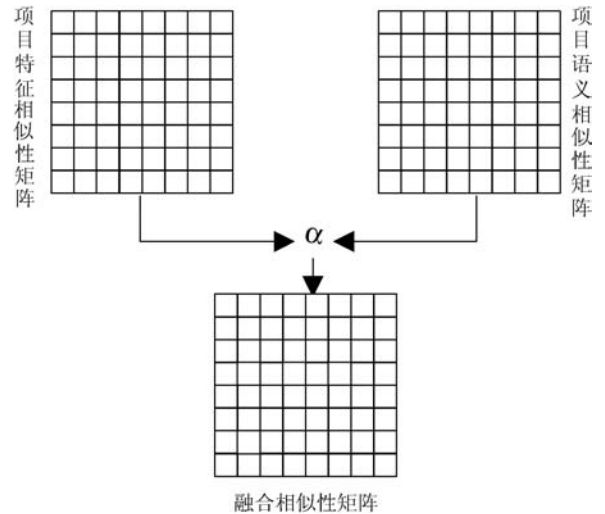


图 4 融合原理图

融合公式如下:

$$sim(v_a, v_b) = \alpha sim_2(v_a, v_b) + (1 - \alpha) sim_1(v_a, v_b) \quad (7)$$

式中: α 为融合因子, 实验中通过不断调整 α 的值, 找到最优融合相似性矩阵。

(2) 评分预测。根据得到的最优融合相似性矩阵, 预测用户对未评价物品的评分, 计算公式如下:

$$P_{u,v_a} = \sum_{j \in T_u} R_{u,v_b} sim(v_a, v_b) \quad (8)$$

式中: R_{u,v_b} 为用户 u 对项目 v_b 的评分; $sim(v_a, v_b)$ 表示项目 v_a 和 v_b 的相似度; T_u 为用户交互的项目列表。

3 实验

3.1 数据集与评价指标

实验数据集采用 MovieLens-1M 和 Book-Crossing Dataset, MovieLens-1M 数据集主要包括 6 040 个用户对 3 952 部电影的 100 多万条评分记录^[24]。Book-Crossing Dataset^[25] 主要包括 17 860 个用户对 14 967 本书籍的 100 多万条评分记录。本文知识图谱采用文献^[26]开源的电影和书籍知识图谱, 数据集 80% 作为训练集, 20% 作为测试集, 并用 5 折交叉验证的方法进行实验, 将平均值作为最终的实验结果。本文采用准确率 (Precision)、召回率 (Recall) 和 F1 值作为评价标准。计算公式如下:

$$Precision = \frac{M(u) \cap N(u)}{M(u)} \quad (9)$$

$$Recall = \frac{M(u) \cap N(u)}{N(u)} \quad (10)$$

$$F1 = \frac{2 \times Precision \times Recall}{Precision + Recall} \quad (11)$$

式中: $M(u)$ 为推荐算法产生的推荐列表; $N(u)$ 为真实的行为列表。准确率与召回率两者之间是相互影响的, 理想情况是两者取值都高, 但实际情况通常是准确率高时, 召回率的值较低, 因此采用 F1 值来权衡这两个指标, 作为推荐效果好坏的依据。

3.2 实验结果与分析

在 MovieLens-1M 数据集上栈式自编码器采用的网络规模为 6 040 - 4 000 - 2 000 - 300 - 2 000 - 4 000 - 6 040, 正则化系数 λ 取 0.1, 并在该数据集上对参数 α 和知识图谱表示学习嵌入维度 dim 进行调优。Book-Crossing 数据集中栈式自编码器采用的网络规模为 17 860 - 8 000 - 3 000 - 300 - 3 000 - 8 000 - 17 860, 其他参数采用 MovieLens-1M 数据集调优得到的参数值。

3.2.1 融合因子的优化及 K 值确定

融合因子 α 取值不同, 所得到的推荐效果也会有

所差异, 本文在区间 $[0, 1]$ 内, 以 0.2 为间隔选取 α 值, 当 α 为 0 时表示基于栈式自编码器的推荐算法, 当 α 为 1 时表示基于知识图谱表示学习的推荐算法, 选取表示学习嵌入维度为 200, 图 5 和图 6 分别为准确率和召回率曲线。

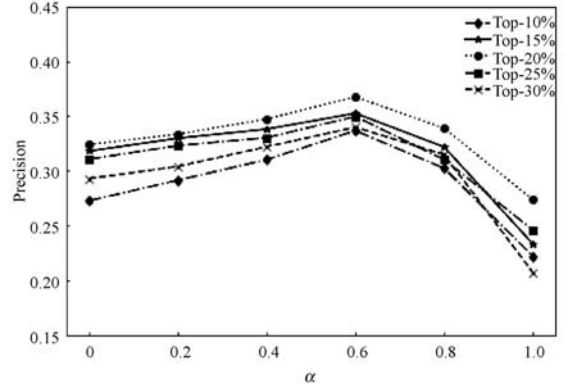


图 5 Precision 比较

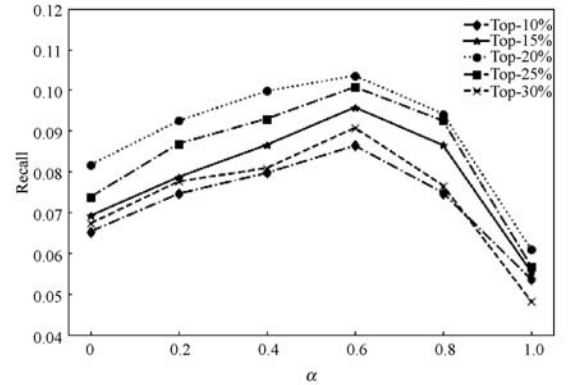


图 6 Recall 比较

可以看出, 随着 α 的增加, 准确率和召回率均呈现先上升后下降的趋势, 当 α 值为 0.6, 准确率和召回率的值达到最高。当推荐列表长度为 20 时, 整体推荐效果最佳。因此本文算法将 α 值设置为 0.6, k 设置为 20。

3.2.2 表示学习嵌入维度确定

知识图谱表示学习将实体嵌入低维一个空间, 维度的取值同样会对推荐效果产生一定的影响, 本文分别选取了 50 ~ 300 维进行实验, 图 7 - 图 9 分别为准确率、召回率和 F1 值随维度的变化情况。

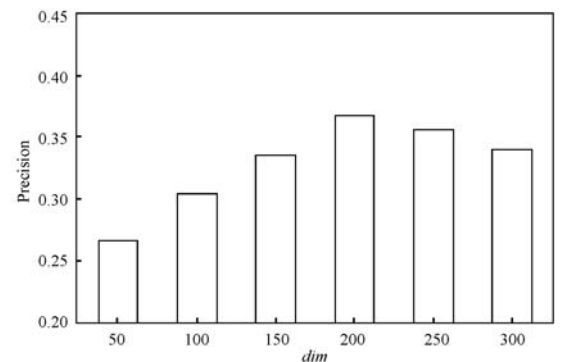


图 7 Precision 随维度的变化情况

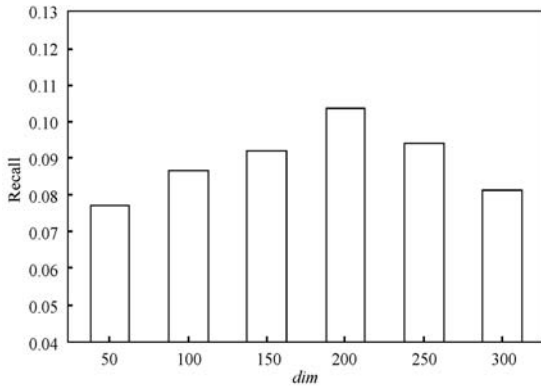


图8 Recall 随维度的变化情况

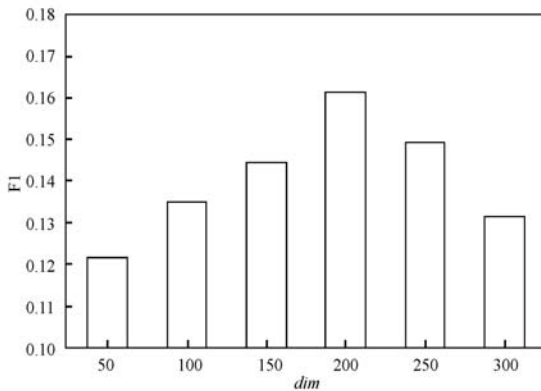


图9 F1 随维度的变化情况

可以看出,随着嵌入维度的增加,准确率、召回率与 F1 值先逐步增加,当维度为 200 时取值最高,然后逐步减小,因此选取知识图谱嵌入维度为 200。

3.2.3 结果分析

为了验证该算法的有效性,将其与协同过滤推荐算法(Item-CF)、奇异值分解(SVD)、栈式自编码器算法(SAE)、知识图谱表示学习算法(KG-CB)在不同稀疏程度的数据集 MovieLens-1M 和 Book-Crossing 上进行对比实验,表 1 为数据集稀疏程度,表 2 和表 3 为对比实验结果。

表 1 数据集稀疏程度

数据集	稀疏程度/%
MovieLens-1M	95.81
Book-Crossing	99.57

表 2 MovieLens-1M 对比实验结果

算法	Precision	Recall	F1
Item-CF	0.240 1	0.051 4	0.084 6
SVD	0.254 1	0.058 2	0.094 7
SAE	0.323 9	0.081 6	0.130 3
KG-CB	0.273 4	0.060 8	0.099 5
SAEKG-CF	0.367 4	0.103 5	0.161 5

表 3 Book-Crossing 对比实验结果

算法	Precision	Recall	F1
Item-CF	0.175 4	0.051 0	0.079 0
SVD	0.182 6	0.059 1	0.089 3
SAE	0.219 6	0.071 4	0.107 8
KG-CB	0.207 9	0.065 4	0.099 5
SAEKG-CF	0.230 8	0.086 3	0.125 6

从表 2、表 3 可以看出,SAEKG-CF 性能在不同稀疏程度的数据集上均优于其他四种算法,准确率、召回率与 F1 值有所提升,因此算法在一定程度上有效地缓解了协同过滤推荐算法中数据稀疏性和语义信息欠缺的问题。

4 结 语

针对协同过滤推荐算法中数据稀疏和语义信息欠缺问题,本文提出一种融合知识图谱表示学习的栈式自编码器推荐算法,利用栈式自编码器获取项目深层特征,并与知识图谱表示学习算法得到的语义信息进行融合。实验结果表明,该算法能够同时解决数据稀疏性和语义信息欠缺问题,得到更为准确的推荐结果。未来考虑在不同领域的数据集上进行实验,并寻找更优的网络层数,其次考虑引入其他深度学习模型来提高推荐效果。

参 考 文 献

- [1] Li L, Zheng L, Yang F, et al. Modeling and broadening temporal user interest in personalized news recommendation[J]. Expert Systems with Applications, 2014, 41 (7): 3168 - 3177.
- [2] Xiong H T, Liu Z B. A situation information integrated personalized travel package recommendation approach based on TD-LDA model[C]//2015 International Conference on Behavioral, Economic and Socio-cultural Computing (BESCI). IEEE, 2015: 32 - 37.
- [3] 路春霞. 个性化推荐中协同过滤算法研究[D]. 北京: 北京交通大学, 2016.
- [4] 何明, 刘伟世, 魏铮. 基于信任网络随机游走模型的协同过滤推荐[J]. 计算机科学, 2016, 43(6): 257 - 262.
- [5] 廖志芳, 符本才, 孔令远, 等. 一种新颖的混合相似度计算模型[J]. 计算机应用与软件, 2018, 35(1): 175 - 182.
- [6] 胡勋, 孟祥武, 张玉洁, 等. 一种融合项目特征和移动用户信任关系的推荐算法[J]. 软件学报, 2014, 25(8): 1817 - 1830.
- [7] 王三虎, 王丰锦. 融合用户评分和属性相似度的协同过滤

- 推荐算法[J]. 计算机应用与软件, 2017, 34(4): 305-308.
- [8] 吴宾, 娄铮铮, 叶阳东. 一种面向多源异构数据的协同过滤推荐算法[J]. 计算机研究与发展, 2019, 6(5): 1034-1047.
- [9] Hu Y, Shi W S, Li H, et al. Mitigating data sparsity using similarity reinforcement-enhanced collaborative filtering[J]. ACM Transactions on Internet Technology, 2017, 17(3): 31.
- [10] Zhang L M, Ma J F, Lu D, et al. Attribute clustering based collaborative filtering[C]//2014 International Conference on Advances in Materials Science and Information Technologies in Industry (AMSITI), 2014: 965-968.
- [11] Huang Z, Chen H, Zeng D, et al. Applying associative retrieval techniques to alleviate the sparsity problem in collaborative filtering[J]. ACM Transactions on Information Systems, 2004, 22(1): 116-142.
- [12] Patra B K, Launonen R, Ollikainen V, et al. A new similarity measure using Bhattacharyya coefficient for collaborative filtering in sparse data[J]. Knowledge Based Systems, 2015, 82(C): 163-177.
- [13] 黄震华, 张佳雯, 张波, 等. 语义推荐算法研究综述[J]. 电子学报, 2016, 44(9): 2262-2275.
- [14] Gradgyenge L, Kiss A, Filzmoser P, et al. Graph embedding based recommendation techniques on the knowledge graph[C]//Adjunct Publication of the 25th Conference on User Modeling, Adaptation and Personalization. ACM, 2017: 354-359.
- [15] 陈鹤. 基于语义本体的社交网络服务推荐系统[D]. 长春: 吉林大学, 2014.
- [16] Chen H, Zhang M F. Improve tagging recommender system based on tags semantic similarity[C]//2011 IEEE 3rd International Conference on Communication Software and Networks. IEEE, 2011: 94-98.
- [17] Shambour Q, Lu J. A Hybrid multi-criteria semantic-enhanced collaborative filtering approach for personalized recommendations[C]//2011 IEEE/WIC/ACM International Conferences on Web Intelligence and Intelligent Agent Technology. IEEE, 2011: 71-78.
- [18] 王俊红. 基于语义网的农业学习资源推荐系统研究[J]. 计算机应用与软件, 2013, 30(8): 233-235.
- [19] Wang H, Shi X J, Yeung D Y. Relational stacked denoising autoencoder for tag recommendation[C]//Proceedings of the 29th AAAI Conference on Artificial Intelligence. ACM, 2015: 3052-3058.
- [20] 常亮, 张伟涛, 古天龙, 等. 知识图谱的推荐系统综述[J]. 智能系统学报, 2019, 14(2): 207-216.
- [21] 吴玺煜, 陈启买, 刘海, 等. 基于知识图谱表示学习的协同过滤推荐算法[J]. 计算机工程, 2018, 44(2): 226-232.
- [22] Bordes A, Usunier N, Garcia-Duran A, et al. Translating embeddings for modeling multi-relational data[C]//Proceedings of the 26th International Conference on Neural Information Processing Systems. ACM, 2013: 2787-2795.
- [23] 方阳, 赵翔, 谭真, 等. 一种改进的基于翻译的知识图谱表示方法[J]. 计算机研究与发展, 2018, 55(1): 139-150.
- [24] Harper F M, Konstan J A. The movieLens datasets: History and context[J]. ACM Transactions on Interactive Intelligent Systems, 2015, 5(4): 19.
- [25] Ziegler C, McNeel S M, Konstan J A, et al. Improving recommendation lists through topic diversification[C]//Proceedings of the 14th International Conference on World Wide Web. ACM, 2005: 22-32.
- [26] Wang H W, Zhang F Z, Wang J L, et al. RippleNet: propagating user preferences on the knowledge graph for recommender systems[C]//Proceedings of the 27th ACM International Conference on Information and Knowledge Management. ACM, 2018: 417-426.
- ~~~~~
- (上接第204页)
- [22] Minhas R, Mohammed A A, Wu Q M J. Incremental learning in human action recognition based on snippets[J]. IEEE Transactions on Circuits & Systems for Video Technology, 2012, 22(11): 1529-1541.
- [23] Wu X X, Xu D, Duan L X, et al. Action recognition using context and appearance distribution features[C]//2011 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). IEEE, 2011: 489-496.
- [24] 范晓杰, 宣士斌, 唐凤. 基于Dropout卷积神经网络的行为识别[J]. 广西民族大学学报(自然科学版), 2017, 23(1): 76-82.
- [25] Ji S W, Xu W, Yang M, et al. 3D convolutional neural networks for human action recognition[J]. IEEE Transactions on Pattern Analysis & Machine Intelligence, 2013, 35(1): 221-231.
- [26] Wang Y, Mori G. Max-margin hidden conditional random fields for human action recognition[C]//2009 IEEE Conference on Computer Vision and Pattern Recognition. IEEE, 2009: 872-879.
- [27] Yeffet L, Wolf L. Local trinary patterns for human action recognition[C]//2009 IEEE 12th International Conference on Computer Vision. IEEE, 2009: 492-497.
- [28] Zhang Y, Liu X, Chang M C, et al. Spatio-temporal phrases for activity recognition[C]//The 2012 European Conference on Computer Vision. Springer, 2012: 707-721.
- [29] Liu J, Yang Y, Shah M. Learning semantic visual vocabularies using diffusion distance[C]//2009 IEEE Conference on Computer Vision and Pattern Recognition. IEEE, 2009: 461-468.
- [30] Wang H, Klaser A, Schmid C, et al. Dense trajectories and motion boundary descriptors for action recognition[J]. International Journal of Computer Vision, 2013, 103: 60-79.