

基于用户分类的隐含因子模型研究*

黎新志, 高茂庭

(上海海事大学 信息工程学院, 上海 201306)

摘要: 针对现有隐含因子模型存在的新用户冷启动问题, 提出基于用户分类的隐含因子模型, 将用户分类信息融入到隐含因子的矩阵分解当中。先在原评分矩阵和用户分类信息的基础上使用指示函数和数据归一化等方法构建一个分类评分矩阵; 再将分类评分矩阵融入到隐含因子模型的评分预测中。通过与传统隐含因子模型等方法在多个不同隐含因子个数上的实验比较分析, 实验结果表明, 改进模型不仅能够解决新用户和项目的冷启动问题, 还能有效降低预测评分的均方根误差, 并提高预测推荐的准确度。

关键词: 推荐系统; 隐含因子模型; 冷启动; 用户分类; 随机梯度下降法

中图分类号: TP181

文献标志码: A

文章编号: 1001-3695(2018)08-2289-04

doi: 10.3969/j.issn.1001-3695.2018.08.012

Research on latent factor model based on user classification

Li Xinzhi, Gao Maoting

(College of Information Engineering, Shanghai Maritime University, Shanghai 201306, China)

Abstract: To solve the problem of user cold start, the paper proposed a latent factor recommendation algorithm based on the user classification, which integrated user classification information into matrix factorization. Firstly, it created the classification rating matrix by using the method of index function and unifying based on the user-item rating matrix and user classification information, then it combined the data of classification rating matrix into the rating prediction stage in the latent factor model, and finally it used a number of latent factors to compare the new method with the traditional latent factor model. The experimental results show that the proposed model can not only solve the cold start problem, but also reduce the root mean squared error between predicted ratings and real ones, and improves the accuracy of prediction effectively.

Key words: recommendation system; latent factor model; cold start; user classification; stochastic gradient descent

0 引言

个性化推荐系统是一种通过挖掘用户的兴趣模式来指导用户选择自己喜欢的产品或服务的技术, 它解决了大数据背景下信息过载所导致的用户在海量数据中选择困难的问题。目前, 推荐系统主要采用基于协同过滤(collaborative filtering)^[1,2]的推荐技术, 协同过滤一般分为基于内存的协同过滤^[3]和基于模型的协同过滤^[4]。基于模型的协同过滤是通过训练原数据中的用户偏好信息得出一个推荐模型, 根据该模型进行预测和推荐。

在基于模型的推荐系统中, 推荐系统通常要完成三个环节的任务: 构建评分矩阵、构建推荐模型、预测和推荐。通过原数据中的用户信息、项目信息以及用户对项目的评分信息来构建一个评分矩阵去表示用户对项目的喜爱程度, 但是原始评分矩阵中的元素通常存在大量的空缺项。因此, 推荐问题可以被简化为预测评分矩阵中的空缺项——即预测用户对某一项目的评分, 再依据该预测评分值的大小向用户进行推荐。

隐含因子模型是一种基于矩阵分解^[5]的机器学习算法, 通过将原始高维矩阵分解成两个低维隐含特征矩阵, 再将分解后的矩阵相乘得到用户对项目的预测评分。该模型具有良好的可扩展性, 能够融合用户和项目的偏见信息, 比如, 某些用户习惯于评高分等。Koren 等人^[5]研究表明, 融合了评分偏见信息的隐含因子模型在对项目进行评分预测时, 其预测的准确度要高于传统协同过滤推荐。然而, 隐含因子模型仍然存在一些不足: a) 忽略了用户的分类属性信息, 没有充分考虑分类信息对预测评分的影响; b) 由于新用户和新项目无历史评分记录, 无法进行建模分析, 导致其无法得到推荐系统的推荐, 存在冷启动问题。

用户的分类属性信息是用户固有的一种属性信息, 这些属性信息是不会轻易改变的, 并且在推荐系统中的作用日益明显。通过将原数据中的用户信息进行分类处理后融入推荐模型中, 一方面, 可以对新加入的无任何评分记录的新用户和新项目提供推荐服务; 另一方面, 可以提高个性化推荐的准确度。

为此, 在隐含因子模型中结合用户分类信息, 如用户职业分类, 提出基于用户分类的隐含因子模型, 以解决新用户和项目的冷启动问题并提高个性化推荐的准确度。

1 隐含因子模型相关研究

1.1 用户评分矩阵分解

基于评分预测的推荐系统的核心问题是对评分矩阵 R 的空缺项进行合理的填充, 而矩阵分解技术可以实现这一目标。对该类问题的研究, Liu 等人^[6]通过融合用户的社会关系和项目内容提出了一种改进的基于贝叶斯概率矩阵分解(Bayesian probabilistic matrix factorization)模型的推荐算法; 王建芳等人^[7]利用奇异值分解(singular value decomposition, SVD)初始化用户和项目的隐含因子向量, 避免随机赋值带来的局限性, 提出融合偏置信息的概率矩阵分解方法(recommendation with SVD initialization PMF and bias information, RSPB); 刘慧婷等人^[8]利用用户间相似性等信息, 提出基于用户偏好的矩阵分解方法(users' preference PMF, USPMF)。

奇异值分解^[9]是一种普遍用于样本特征提取和降维的矩阵分解方法: $R = A \Sigma B^T = A' B'^T$ 。其中, A 和 B 是正交矩阵, Σ 是对角矩阵并且对角元素是从大到小排列的奇异值, 这些奇异值代表着原始矩阵的重要特征, 也可将对角矩阵 Σ 吸收到矩阵 A 和 B 中实现两个矩阵相乘的形式。

Funk^[10]从最小化评分误差的角度描述了矩阵分解, 提出

收稿日期: 2017-04-11; 修回日期: 2017-05-18 基金项目: 国家自然科学基金资助项目(61202022)

作者简介: 黎新志(1991-), 男, 硕士研究生, 主要研究方向为数据挖掘、数据分析(95107541@qq.com); 高茂庭(1963-), 男, 系统分析员, 教授, 博士, 主要研究方向为智能信息处理、数据库与信息系统。

了 Funk-SVD 方法,该方法的核心思想是将传统的评分矩阵分解转换为求解最优问题。

一方面,评分矩阵 R 是一个 $U \times I$ 的矩阵,表示 U 个用户对 I 个项目的评分,每个评分值用 r_{ui} 表示。从矩阵分解的角度将评分矩阵 R 分解为两个低维的隐含特征矩阵相乘的形式: $R_{U \times I} = P_{U \times K} Q_{K \times I}$,其中,参数 K 为隐含的特征因子个数,并且 K 小于用户个数 U 和项目个数 I ,实现矩阵分解降维的作用,因此,用户 $u \in U$ 对项目 $i \in I$ 评分可用式(1)预测。

$$\bar{r}_{ui} = p_u^T q_i \quad (1)$$

其中: p_u 是分解后的矩阵 P 中的第 u 个行向量,表示用户 u 的隐含特征因子; q_i 是分解后的矩阵 Q 中的第 i 个列向量,表示项目 i 的隐含特征因子。两个特征因子向量的点积即为预测评分,并且所有的预测评分构成了 $U \times I$ 维的矩阵 \bar{R} ,即 U 个用户对 I 个项目的预测评分。

另一方面,实际评分值和预测评分值之间存在一定的误差 $e_{ui} = r_{ui} - \bar{r}_{ui}$,如果预测误差 e_{ui} 足够小的话,就能得到一个较为准确的预测评分矩阵 $\bar{R} = PQ$,于是,预测评分矩阵的计算就转换为一个目标函数为式(2)的优化问题。

$$E = \sum_{ui} e_{ui}^2 = \sum_{ui} (r_{ui} - \bar{r}_{ui})^2 \quad (2)$$

由于评价推荐系统预测好坏的一个评价指标是均方根误差(root mean square error, RMSE),所以,可以采用随机梯度下降法,通过对训练集进行不断的迭代学习,最小化式(2)的目标函数来得到 P 和 Q 矩阵。进而获得用户对项目的预测评分矩阵 \bar{R} ,再结合预测评分矩阵 \bar{R} 和测试集评分矩阵就可计算求得均方根误差值,并且所求得的均方根误差值为最小值,均方根误差小说明预测推荐的误差越小,表明推荐的准确度越高。

1.2 隐含因子模型

Koren 等人研究发现在式(1)的预测公式中融入原数据中的用户评分偏见和项目评分偏见能提高评分预测的准确度,改进后的模型被称为隐含因子模型(latent factor model, LFM)^[11,12]。Funk 将传统矩阵分解方法转换为求解最优问题,但是该方法在计算预测评分时却如同传统的矩阵分解方法一样局限在了用户对项目的评分上,未考虑到其他信息对评分预测的影响。然而,现实世界中,样本数据在特定的应用环境下有其自身的属性,比如对同一个项目有些用户倾向于对项目评高分而有些用户却倾向于对项目评低分,称这种现象为用户固有的偏见^[11]。同理,也存在项目偏见。可以初始化用户 u 和项目 i 的偏见分别为用户 u 对所有项目评分的平均值和项目 i 的所有评分的平均值。因此,考虑用户评分偏见 b_u 和项目评分偏见 d_i 后,用户对项目评分预测计算如式(3)所示。

$$\bar{r}'_{ui} = p_u^T q_i + \mu + b_u + d_i \quad (3)$$

其中: μ 为常数表示所有评分数据的平均值。

采用随机梯度下降法对可微的目标函数进行迭代训练时一般按如下步骤进行:

- 对目标函数中的变量赋初值,初始值可以为 0 或随机值。
- 迭代更新目标函数中变量的值,使得目标函数向梯度下降的方向逐步减小。

在对训练集进行迭代训练求解局部最优解时,可能会出现数据的过度拟合(over-fitting)导致模型的测试效果不好,进而影响到预测推荐的准确度。常用的避免过拟合方法是正则化(regularizer),因为文中的目标函数为平方损失,所以正则项可以使用变量的 L_2 范数,二范数的正则化项一般具有如下形式:

$$\lambda \sum_{j=1}^n \theta_j^2$$

参数 θ 是待惩罚的变量,参数 λ 是正则项系数用来权衡正则项与原目标函数项的比重。 λ 的值需要根据具体的应用场景反复实验得到。相应地,目标函数修改为式(4)。

$$E' = \sum_{(u,i) \in K} (r_{ui} - \bar{r}'_{ui})^2 + \lambda (\|p_u\|^2 + \|q_i\|^2 + b_u^2 + d_i^2) \quad (4)$$

考虑用户评分偏见和项目评分偏见的隐含因子模型可以较好地消除用户和项目的评分偏见。但是,仍然存在冷启动问题:当一个新用户加入系统之后,由于推荐系统对该用户除了注册的用户信息外一无所知,也就无法对其进行建模分析。这

时,隐含因子模型也就无法对该用户进行个性化的推荐。

2 基于用户分类的隐含因子模型

为了给用户提供更好的服务,网站往往需要用户注册成为会员,如淘宝网、豆瓣网等,让用户在账号注册过程中填写一些基本信息,如性别、年龄、职业、兴趣等。推荐系统正好可以利用这些信息对用户进行分类,当有新用户加入系统之后,根据新用户所属的分类就可对该用户进行推荐服务。基于此,隐含因子模型可以使用指示函数和归一化处理等方法对用户—项目评分矩阵 R 中的用户维度进行用户分类处理,分为 M 类($M \ll U$),比如按用户职业进行分类;同样地,再对用户—项目评分矩阵 R 中的项目维度进行分类处理,分为 N 类($N \ll I$)。从而得到一个 M 行 N 列的用户类别—项目类别评分矩阵 S ,矩阵元素值 s_{mn} 表示用户类别 m 对项目类别 n 的总体评分。因此,基于用户分类信息对传统隐含因子模型的改进主要体现在以下几个方面:

a) 使用指示函数以及数据归一化方法对原数据中的用户信息进行分类处理,得出分类评分矩阵 S 。

b) 将分类评分矩阵 S 中的分类评分 s_{mn} 进一步融入到考虑用户评分偏见和项目评分偏见的隐含因子模型的预测公式中,设用户 u 属于用户类别 m ,项目 i 属于项目类别 n ,得出改进后的评分预测公式:

$$\bar{r}''_{ui} = p_u^T q_i + \mu + b_u + d_i + s_{mn} \quad (5)$$

c) 当有新的用户加入推荐系统之后,利用该用户的用户信息计算出该用户所属的用户类别,根据该用户类别对项目的预测评分进行预测和推荐。

基于用户分类的隐含因子模型的主要处理流程如图 1 所示。

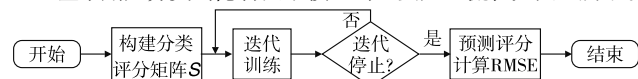


图1 基于用户分类的隐含因子模型工作流程

2.1 构建分类评分矩阵 S

对原数据进行预处理和提取得到用户对项目评分矩阵 R 并作为推荐模型输入数据的一部分。矩阵 R 是非常稀疏的矩阵,非空值表示用户对项目的实际评分,空缺值是推荐系统需要预测的值。因此,补全矩阵中的空缺项是推荐系统需要做的主要任务。

对于分类评分矩阵 S 的构建是基于用户评分矩阵 R 中的评分值和原数据中的用户和项目信息,一般来说,用户对某个项目有过评分行为(即评分非 0 的情况),往往表明该用户或该用户所属类别对该项目及其所属的项目类别是感兴趣的。比如,电商平台上一个用户对某个具体品牌的奶粉有过评分行为——不管是最低分 1 分还是最高分 5 分,都表明该用户对奶粉类别有一定的兴趣度。那么,有如下所示的指示函数 I_{ui} ,表明用户 u 对项目 i 评分与否:

$$I_{ui} = \begin{cases} 0 & \text{用户 } u \text{ 对项目 } i \text{ 无评分} \\ 1 & \text{用户 } u \text{ 对项目 } i \text{ 有评分} \end{cases} \quad (6)$$

通过式(6)可以将原始评分矩阵 R 映射到一个元素只为 0 或 1 的指示矩阵 I 上。在此基础上,再结合原数据中用户所属的类别进行统计分析获得用户类别—项目类别评分矩阵 S ,用户类别 m 对项目类别 n 的评分 s_{mn} 按式(7)计算。

$$s_{mn} = \frac{\sum_{u \in G(u), n \in H(i)} I_{ui}}{m} \quad (7)$$

其中: $G(u)$ 表示用户 u 所属的用户类别; $H(i)$ 表示项目 i 所属的项目类别。

将评分矩阵 R 映射到分类评分矩阵 S 之后还需要作进一步处理,这是由于基于用户分类的各个类别下的用户基数往往不尽相同,并且有时会出现各个类别下用户基数的两极分化情况。比如,某些用户类别内的用户基数很大,常常会导致该用户类别的评分很高,而有一些用户类别则恰恰相反,这种情况会影响数据分析的结果,所以需要奇异样本数据进行归一化处理。

常用的归一化处理的方法有 min-max 标准化方法以及 z-

score 标准化方法两种,而 z-score 标准化方法处理后的结果常常会出现负数的情况,这与本文使用的数据集评分制为非负整数不相符(文中实验部分使用的 MovieLens 数据集采用的是 1~5 评分制)。因此,本文采用 min-max 标准化方法,并对原始的 min-max 标准化方法进行改进,使该方法的转换函数适用于采用 1~5 评分制的数据集。修改后的 min-max 标准化方法的转换函数如式(8)。

$$s_{mn} = \frac{(s_{mn} - \min S_{MN}) \times 5}{\max S_{MN} - \min S_{MN}} \quad (8)$$

其中: $\min S_{MN}$ 表示所有用户类别对项目类别评分的最低分; $\max S_{MN}$ 表示所有用户类别对项目类别评分的最高分。采用式(8)即可将总评分矩阵的评分值规范化到[0-5]上,从而构建出分类评分矩阵 S 。

2.2 对矩阵进行迭代分解

利用机器学习的方法,通过不断地迭代学习来拟合原始评分矩阵 R 可以有效提高矩阵分解的效率并获得较好的训练结果。在基于用户分类的隐含因子模型中,采用随机梯度下降法^[13]进行训练学习,最小化预测评分与实际评分之间的误差 $e_{ui} = r_{ui} - \bar{r}_{ui}$ 求解最近似预测评分,从而得到矩阵 R 中的空缺项。

隐含因子模型不同于传统 SVD 最重要的一点是,隐含因子模型具有很好的扩展性,Koren 等人^[5]消除了用户和项目评分偏见对预测的干扰。在此基础上,进一步充分考虑分类信息对预测评分的影响,将式(3)改进为式(5),并得到相应的目标函数,如式(9)所示。

$$E'' = \sum_{(u,i) \in K} (r_{ui} - \mu - b_u - d_i - s_{mn})^2 + \lambda (\|p_u\|^2 + \|q_i\|^2 + b_u^2 + d_i^2) \quad (9)$$

采用随机梯度下降法最小化式(9)的目标函数:首先分别对 p_u, q_i, b_u, d_i 四个变量求偏导得到目标函数在各个变量处的梯度;然后,将变量分别向各自梯度的负方向下降。变量的迭代更新一般具有如下形式:

$$\theta \leftarrow \theta + \gamma \cdot \nabla^{(-)}$$

其中: $\nabla^{(-)}$ 表示梯度负方向,在变量下降的过程中需要通过参数 γ 来控制下降的速率,即学习的步长。步长 γ 的取值一般在 0~1,步长太大会导致目标函数值来回振荡得不到局部最优解;而步长太小也会导致学习的速率太慢,在迭代一定次数后同样得不到局部最优解。因此,四个变量的迭代更新公式如下所示:

$$\begin{cases} p_u \leftarrow p_u + \gamma(e_{ui}q_i - \lambda p_u) \\ q_i \leftarrow q_i + \gamma(e_{ui}p_u - \lambda q_i) \\ b_u \leftarrow b_u + \gamma(e_{ui} - \lambda b_u) \\ d_i \leftarrow d_i + \gamma(e_{ui} - \lambda d_i) \end{cases} \quad (10)$$

其中:参数 γ 为步长,需要根据具体的应用场景反复实验获得;参数 λ 的意义如前所述,为用来控制正则化程度的常数。通过多次的迭代训练以更新评分预测矩阵 R ,迭代停止的条件可以是指定迭代的次数或迭代到均方根误差小于一个预先设定的值时停止,然后,输出预测评分矩阵 R 。最后,根据预测评分进行个性化推荐服务。

3 实验及结果分析

为了全面地验证本文所研究的基于用户分类的隐含因子模型在推荐系统中的有效性,通过在规模不同的两个数据集上进行多次实验来验证改进模型的推荐质量,并讨论随着隐含因子个数的不同对推荐结果的影响,同时与其他推荐模型进行对比分析。

3.1 数据集

实验数据选用 MovieLens 数据集中 MovieLens-100K 和 MovieLens-1M 两个数据集,该数据集由美国 GroupLens 项目组建,采用的评分度量是 1~5 分的评分制,评分越高表明用户的喜爱程度越高,0 分表示用户未评分,即用户对该项目没有相应的历史行为记录,用户从事的职业共分为 21 种,电影共分

为 19 个类别。

MovieLens-100K 数据集包含 943 位用户对 1 682 部电影总计 100 000 条评分数据。其中,u. data 文件记录用户对电影的评分信息;u. user 文件记录了用户信息;u. item 记录了电影信息。

MovieLens-1M 数据集包含 6 040 个用户对 3 900 部电影的评分记录。其中 ratings. dat 文件记录了用户对电影的评分情况;users. dat 和 movies. dat 文件分别记录了用户和电影的信息。

实验中,将数据集按 80% 和 20% 的比例进行划分,其中 80% 的数据作为训练样本进行模型训练,20% 的数据作为测试样本进行测试并统计相应的评价指标。

3.2 推荐质量的度量标准

在推荐系统研究领域中,推荐系统预测的准确性是最重要的一个评估指标^[14],比如,文献[15]通过比较推荐算法的预测准确性来比较 Funk-SVD 和基于项目的协同过滤的推荐质量。

准确性既可以用来评估预测评分与实际评分之间的误差值,也可以用来评估预测排名与实际排名之间的误差值。在推荐系统领域通常采用均方根误差 RMSE 来度量预测评分的准确性。RMSE 越大,则表明推荐系统的误差就越大,准确性越低;相反地, RMSE 越小,则表明推荐系统的误差越小,准确性越高。

$$RMSE = \sqrt{\frac{\sum_{(u,i) \in T} (r_{ui} - \bar{r}_{ui})^2}{|T|}} \quad (11)$$

其中: r_{ui} 表示用户 u 对项目的实际评分; \bar{r}_{ui} 表示用户 u 对项目 i 的预测评分; T 表示测试集。

3.3 实验及结果分析

实验中,先通过将评分矩阵进行迭代训练实现矩阵的分解,再将训练后输出的预测评分矩阵 R 与测试集进行评价计算,得出该模型的推荐精度。为了较为全面地掌握基于用户分类的隐含因子模型的性能,从以下两个方面进行对比实验:

a) 讨论 $R_{U \times I} = P_{U \times K} Q_{K \times I}$ 矩阵分解模型中的 K 值(即隐含因子的个数)对结果的影响,取多个不同的 K 值,检测在取 5、10、20、40 等多个不同值时 RMSE 值的变化。

b) 将本文提出的基于用户分类的隐含因子模型与传统的隐含因子模型 LFM 和基于用户偏好的矩阵分解模型 USPMF^[8] 对比,观察 RMSE 值的变化。

随机梯度下降法迭代的学习速率 γ 和目标函数中的正则化参数 λ 通过实验确定,针对 MovieLens-100K 的数据集,其参数值分别初始化为: $U = 943, I = 1\ 682, M = 21, N = 19, \gamma = 0.003, \lambda = 0.1$,迭代停止的条件是在训练集上迭代 100 次后停止。实验结果如图 2 所示。

在 MovieLens-1M 数据集上初始化输入参数为: $U = 6\ 040, I = 3\ 900, M = 21, N = 19, \gamma = 0.01, \lambda = 0.1$,同样,最大迭代次数 $\max Iteration = 100$ 。实验结果如图 3 所示。

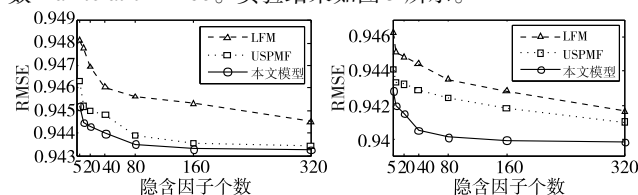


图2 MovieLens-100K集实验结果

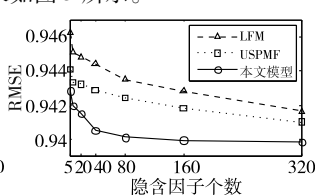


图3 MovieLens-1M集实验结果

通过对图 2 和 3 结果的分析对比,可以得到以下几个方面的结论。

3.3.1 模型的稳定性

在 MovieLens-100K 和 MovieLens-1M 两个数据集上的实验中,当隐含因子个数增加到大约 160 之后,本文提出的模型的 RMSE 值均不再有较大变化,趋于稳定,即针对不同规模的数据集, RMSE 变化趋势具有相似的规律,表明该模型良好的稳定性。

- [5] Yildirimoglu M, Geroliminis N. Experienced travel time prediction for congested freeways[J]. *Transportation Research, Part B: Methodological*, 2013, 53(4) :45-63.
- [6] Salamanis A, Meladianos P, Kehagias D, *et al.* Evaluating the effect of time series segmentation on STARIMA-based traffic prediction model[C]//Proc of the 18th IEEE International Conference on Intelligent Transportation Systems. Piscataway, NJ: IEEE Press, 2015: 2225-2230.
- [7] Nath R P D, Lee H J, Chowdhury N K, *et al.* Modified K-means clustering for travel time prediction based on historical traffic data [C]//Proc of the 14th IEEE International Conference on Knowledge-Based and Intelligent Information and Engineering Systems. Berlin: Springer, 2010: 511-521.
- [8] Zhang Ting, Xia Yingjie, Zhu Qianqian, *et al.* Mining related information of traffic flows on lanes by K-medoids[C]//Proc of the 11th IEEE International Conference on Fuzzy Systems and Knowledge Discovery. Piscataway, NJ: IEEE Press, 2014: 390-396.
- [9] Lin Lei, Wang Qian, Sadek A W. A novel variable selection method based on frequent pattern tree for real-time traffic accident risk prediction[J]. *Transportation Research, Part C: Emerging Technologies*, 2015, 55(6) :444-459.
- [10] Chen Ying, Kim J, Mahmassani H S. Pattern recognition using clustering algorithm for scenario definition in traffic simulation-based decision support systems[C]//Proc of the 17th IEEE International Conference on Intelligent Transportation Systems. Piscataway, NJ: IEEE