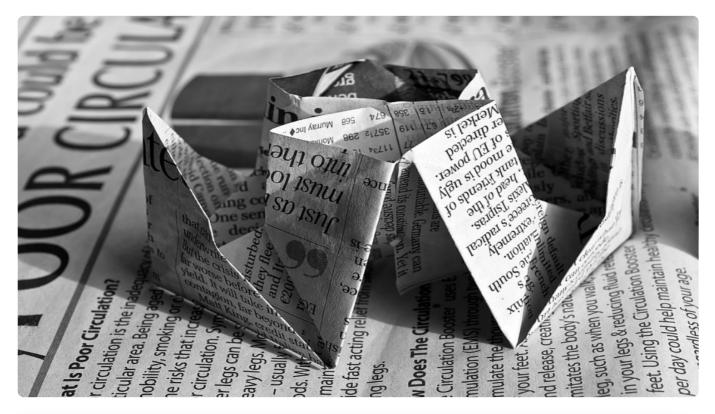
(2)

# 102 | The Web 2018论文精读:如何改进经典的推荐算法BPR?

2018-05-28 洪亮劼

AI技术内参 进入课程>



**讲述:初明明** 时长 05:46 大小 2.65M D

今天,我们来看万维网大会上的一篇优秀短论文。在万维网大会上,主要发表两类论文。一类是 10 页的长论文,一类是 2 页的短论文或称作展板论文。短论文主要是发表短小的成果或者是还在研究过程中的重要成果。每一届的万维网大会,都会评选出一篇最佳短论文奖。

今天我和你分享的论文,题目是《利用查看数据,贝叶斯个性化排序的一种改进的取样器》(An Improved Sampler for Bayesian Personalized Ranking by Leveraging View Data)。这篇论文也有六位作者,和我们介绍的上一篇论文一样,都来自清华大学和新加坡国立大学。

## 贝叶斯个性化排序

要想理解这篇论文的内容,我们必须要讲一下什么是"**贝叶斯个性化排序**"(Bayesian Personalized Ranking),或者简称是**BPR**。有关 BPR 的详细介绍,可以阅读参考文献 [1]。我们在这里仅对 BPR 进行一个高维度的总结。

简单来说,BPR 是推荐系统中的一个配对排序(Pairwise)学习算法。在我们前面介绍搜索算法的时候,曾经提到了各种配对排序学习算法。配对排序学习不是针对每一个数据实例来学习其标签或者响应变量,而是学习一个相对的顺序,希望能够把所有的正例都排列到负例之前。也就是说,对于配对排序来说,每一个数据实例的预测值本身并不重要,排序算法在意的是对于一正一负的一个配对来说,是否能够把正例给准确地排列到负例之上。这其实就要求 BPR 在数值上对正例的预测值能够比负例的预测值高。

BPR 主要是解决了在推荐系统中长期以来只对单个数据点进行预测,比如需要对用户物品的喜好矩阵建模的时候,之前的大多数算法都无法有效地对没有观测到的数据进行建模。而BPR 是配对算法,因此我们只需要关注观测的数据以及他们之间的关系,从而能够对用户的喜好,特别是有"**隐反馈**"(Implicit Feedback)数据的时候,取得更加明显的效果。这里的隐反馈指的并不是用户告诉系统其对每一个物品的喜好程度,而是用户在和系统的交互过程中通过一些行为表达出的喜好。这些用户的行为往往并不全面,因此需要算法和模型能够对这些行为进行有效建模。

#### 论文的主要贡献和核心方法

了解了 BPR 大概是怎么回事以后,我们来看一看这篇论文的主要贡献和核心方法。

首先我们刚才讲到 BPR 的核心是学习一个配对的排序问题。那么在训练的时候,我们需要对一个正例和一个负例的配对进行学习,更新参数。然而在一个自然的用户隐反馈数据集里,正例相对来说往往是少数,负例则是绝大多数。因此,一个传统的方法就是在组成一个配对的时候,相对于一个正例来说,我们都"均匀地" (Uniformly) 选取负样本来组成配对,这个过程有时候也叫"采样" (Sampling)。

这篇论文有两个主要贡献。第一个贡献是,作者们发现,如果在全局均匀地采样负样本,第一没有必要,第二可能反而会影响最后学习的效果。第二个贡献是,针对电子商务的应用,作者们发明了一种负样本采样的方法,使得学习算法可以利用到更多的用户"浏览"(View)信息,从而能够对算法的整体训练效果有大幅度的提升。

### 方法的实验效果

这篇论文的数据集分别使用了母婴产品"贝贝网"和天猫的数据。其中,贝贝网有约 16 万用户、12 万商品、260 万次购买和 4600 万次浏览;天猫的数据则有 3 万用户、3 万多商品、46 万次购买和 150 多万次浏览。两个数据集都呈现了大于 99% 的"稀疏度"(Sparsity)。

首先,作者们实验了不从全局中选取负样本而仅仅采样一部分,而且是相比于原来的空间非常小的样本,比如仅仅几百个负样本而不是几万个的情况。实验效果在贝贝网上不仅没有影响算法的精确度,算法的精确度反而还有提升。而在天猫的数据集上,算法效果没有提升,而有一些小幅度的下降,但是作者们认为这样的代价还是值得的,因为数据集的减少,算法的训练时间会大幅度降低。从这个实验中,作者们得出了不需要从全局进行采样的结论。

紧接着,作者们提出了一个新的概念,那就是,对用户的数据集合进行划分,把用户的行为分为"购买集"(C1)、"浏览但没有购买集"(C2)、"剩下的数据"(C3)这三个集合。作者们提出,BPR 要想能够达到最好的效果,需要对这三种数据集进行采样。也就是说,我们需要组成 C1 和 C2、C1 和 C3 以及 C2 和 C3 的配对来学习。

具体来说,用户在贝贝网和天猫的数据中尝试了不同的比例来对这三种集合进行采样。总体的经验都是 C3 中采样的数据要大于 C2 中的,然后要大于 C1 中的。这其实就是说训练算法要更好地学习到用户不喜欢某件东西的偏好。采用这样的采样方式,作者们展示了模型的效果比传统的 BPR 或仅仅使用"最流行的物品"作为推荐结果要好 60% 左右。

## 小结

今天我为你讲了今年万维网大会的一篇优秀短论文。文章介绍了如何对一个经典的推荐算法 BPR 进行改进,从而提高效率并且大幅度提升算法有效度。

一起来回顾下要点:第一,我们从高维度介绍了BPR的含义;第二,我们简要介绍了论文的主要贡献和思路;第三,我们简单分享了论文的实验成果。

最后,给你留一个思考题,除了这篇论文提出的组成正例和负例的配对思路以外,你能不能想到在用户浏览网站的时候,还有哪些信息可以帮助我们组成更多的配对呢?

欢迎你给我留言,和我一起讨论。

## 参考文献

1. Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. BPR: Bayesian personalized ranking from implicit feedback. Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence (UAI '09). AUAI Press, Arlington, Virginia, United States, 452-461, 2009.



© 版权归极客邦科技所有,未经许可不得传播售卖。 页面已增加防盗追踪,如有侵权极客邦将依法追究其法律责任。

上一篇 101 | The Web 2018论文精读:如何对商品的图片美感进行建模?

下一篇 103 | The Web 2018论文精读:如何从文本中提取高元关系?

## 精选留言(1)



<sub>L</sub>



吴文敏

2018-07-20

对这个思路稍作拓展,我们只要定义隐反馈行为间的偏序关系,就可以基于多种隐反馈 (浏览、点击、购买、加入购物车)进行配对采样