

# 社交媒体内容可信性分析与评价

刘 波<sup>1,3</sup> 李 洋<sup>2,3</sup> 孟 青<sup>1</sup> 汤小虎<sup>1</sup> 曹玖新<sup>2</sup>

<sup>1</sup>(东南大学计算机科学与工程学院 南京 211189)  
<sup>2</sup>(东南大学网络空间安全学院 南京 211189)  
<sup>3</sup>(计算机网络与信息集成教育部重点实验室(东南大学) 南京 211189)  
(bliu@seu.edu.cn)

## Evaluation of Content Credibility in Social Media

Liu Bo<sup>1,3</sup>, Li Yang<sup>2,3</sup>, Meng Qing<sup>1</sup>, Tang Xiaohu<sup>1</sup>, and Cao Jiuxin<sup>2</sup>

<sup>1</sup>(School of Computer Science and Engineering, Southeast University, Nanjing 211189)  
<sup>2</sup>(School of Cyber Science and Engineering, Southeast University, Nanjing 211189)  
<sup>3</sup>(Key Laboratory of Computer Network and Information Integration (Southeast University), Ministry of Education, Nanjing 211189)

**Abstract** With the rapid development of social media in recent years, the access to information has been broadened, but the spreading of incredible information has been facilitated at the same time, which brings a series of negative impacts to cyber security. Compared with the traditional online media, the information in social media is more open and complicated, giving rise to great challenges to judge online information credibility for individuals. How to filter the incredible information becomes an urgent problem. In the existing research on the assessment of information credibility in social media, lots of effort has been involved in extracting the useful factors for credibility assessment, but the processing of noisy data is neglected, and a large number of useless tweets can be included in the evaluation process, resulting in the deviation of the information credibility assessment. So it is particularly important to select the significant tweets for information credibility assessment. This paper takes the topic factor and conformity of users into consideration to relieve the impact of noisy data, such as conformity retweeting, on information credibility assessment, and uses Bayesian network to establish an evaluation model for information credibility in social media. Then we verify the effectiveness of our model using a real dataset.

**Key words** social media; content credibility; topic factor; conformity; probabilistic graphical model

**摘 要** 近年来社交媒体在拓宽人们获取信息渠道的同时,也方便了虚假信息的传播,并造成了严重的负面影响.与传统互联网媒体相比,社交媒体包含的信息更加复杂多样,为内容可信性的判断带来了新的

收稿日期:2018-09-12;修回日期:2019-04-09  
基金项目:国家重点研发计划项目(2017YFB1003000);国家自然科学基金项目(61370208,61472081,61320106007,61272531);国家“八六三”高技术研究发展计划基金项目(2013AA013503);江苏省网络与信息安全重点实验室基金项目(BM2003201);江苏省计算机网络技术重点实验室基金项目(BE2018706)  
This work was supported by the National Key Research and Development Program of China (2017YFB1003000), the National Natural Science Foundation of China (61370208, 61472081, 61320106007, 61272531), the National High Technology Research and Development Program of China (863 Program) (2013AA013503), the Jiangsu Provincial Key Laboratory of Network and Information Security Foundation (BM2003201), and the Jiangsu Provincial Key Laboratory of Computer Network Technology Foundation (BE2018706).

挑战.已有研究在分析社交媒体内容可信性时,对挖掘可信性影响因素进行了很多工作,但缺乏对噪音数据的处理,大量的无用推文会对推文可信性判断造成干扰,进而会影响事件层面的可信性判断,从大量噪音数据中筛选出真正有用的推文数据就显得尤为重要.在推文层面同时考虑用户的主题因素和从众行为,减少了从众转发等噪音数据在可信性判断过程中的作用,对社交媒体内容的可信性进行研究,采用贝叶斯网络建立了社交媒体内容可信性评价模型,并通过新浪微博公开数据集验证了模型的有效性.

**关键词** 社交媒体;内容可信性;主题因素;从众行为;概率图模型

**中图法分类号** TP393

随着移动互联网以及智能移动终端的普及,新浪微博<sup>①</sup>、Facebook<sup>②</sup>、Instagram<sup>③</sup>等社交媒体平台将人们的生活和互联网越来越紧密地联系在一起.由全球领先的社交媒体数字营销机构 DataReportal 发表的 2019 全球数字报告显示,全球 77 亿人口中活跃的社交媒体用户已经达到 45%,Facebook 的月活跃用户达到了 22.71 亿,中国的新浪微博月活跃用户也已经达到了 4.46 亿<sup>[1]</sup>.此外,Kantar Media CIC 在 2017 年中国社会化媒体格局概览中指出中国的社会网络如新浪微博等社交媒体平台,已经覆盖了人们生活的方方面面<sup>[2]</sup>.这是因为社交媒体具有快捷、方便、双向、开放等特点,给人们消费信息带来了巨大的便利.然而,社交媒体的这些特点也使它成为了孕育不可信信息的温床.一方面由于在社交媒体中内容的发布几乎是零门槛,用户自身认识局限性导致的错误观点或者是用户出于某种目的而设计的片面新闻、虚假新闻都能轻易地发布在社交媒体平台上.另外一方面,由于社交媒体中信息的交换十分频繁,不可信的内容能够很快传播开来,覆盖大量的用户,给社会和个人带来严重的负面影响.

传统的内容可信性判断是通过人工来实现的,对于社交媒体中海量的内容,这种方法已经不可行.如今随着数据挖掘、机器学习等技术的发展,采用计算机评估内容可信性成为了主流.该方法的最大优势在于能够从全局角度去评价内容可信性,避免了人工评判中信息不对称的问题.现在大部分社交媒体平台都有自动化的信息过滤机制,如点评类网站 Yelp<sup>④</sup>对垃圾评论进行过滤,问答互动型网站 Quora<sup>⑤</sup>会隐藏劣质答案而向用户推送最佳答案.本文将从社交媒体中用户的特点出发,考虑用户的主题因素和从众因素,提出一种基于概率图模型的方法来对社交媒体中的内容可信性进行判断.

## 1 相关研究综述

从 20 世纪 90 年代中期开始,互联网内容的可信性研究就成为了一个重要的研究领域<sup>[3]</sup>.随着社交媒体的兴起,研究社交媒体中内容的可信性变得尤为重要.对于计算机领域中的可信性,Fogg 等人<sup>[4]</sup>给出了被大部分研究者所认同的解释.他们认为可信性包含 2 个基本维度:可信赖度(trustworthiness)和专业度(expertise).可信赖度包含无恶意(well-intentioned)、真实(truthful)、公正(unbiased)3 个方面,侧重于描述信息本身;专业度包含经验丰富(experienced)、知识渊博(knowledgeable)、能力突出(competent)3 个方面,侧重于描述信息源.

根据上述 2 个维度,可以将社交媒体中内容的可信性研究分为面向信息源的可信性研究和面向信息可信性的研究.考虑到传统网络媒体中内容可信性的研究方法也适用于社交媒体,本研究把传统网络媒体当作特殊的社交媒体也纳入到社交媒体内容可信性研究的讨论中,那么信息源就体现为传统网络媒体中的网站和社交媒体中的用户.信息则是网站或者用户发布于传统网络媒体或社交媒体上的多媒体内容.

面向信息源的可信性研究可以分为 2 类:基于网络拓扑结构的信息源可信性研究和基于信息源特征的信息源可信性研究.基于网络拓扑结构的研究以信息源为节点、信息源之间的关系为边构造网络模型,根据信息源在网络中所处的位置,对信息源的可信性进行计算.PageRank<sup>[5]</sup>算法是其中最为经典的算法,之后出现了许多基于 PageRank 的改进算法,比如 Appleseed<sup>[6]</sup>,TrustRank<sup>[7]</sup>,CredibleRank<sup>[8]</sup>,

① <http://weibo.com>

② <http://www.facebook.com>

③ <http://www.instagram.com>

④ <http://www.yelp.com>

⑤ <http://www.quora.com>

VoteTrust<sup>[9]</sup>等算法.基于信息源特征的研究是寻找影响信息源可信性的因素,比如信息源的活跃度、权威度、与其他信息源的关系、历史行为和信息源发布内容语义信息、传播范围、时效性等,研究这些因素如何影响信息源的可信性,采用合适的模型对这些因素进行组合,从而得到信息源的可信性<sup>[10-11]</sup>.虽然在信息源的可信性计算中加入了很多因素,但这些研究大部分都忽略了信息源的主题因素,默认信息源在所有主题下是一样的,不符合常理,比如说人们更倾向于相信一个医生发布的关于药品的内容,而不相信他发布的关于天文的内容.

面向信息的可信性研究方面通常考虑多种因素,采用迭代模型、优化模型和概率图模型3种模型来进行研究.使用迭代模型度量信息可信性的研究利用影响信息可信性的因素和信息可信性之间的相互影响关系,通过影响信息可信性因素计算信息可信性,然后通过信息可信性量化影响信息可信性的因素,不断重复这个过程直至收敛.采用迭代模型的最简单情形是利用信息源可信性和信息可信性之间的相互影响来计算信息可信性<sup>[12-14]</sup>.也有一些迭代模型考虑多种因素,如信息源之间的关系、信息的语义等<sup>[15-16]</sup>.采用优化模型的信息可信性计算方法主要目的是寻找一个映射把影响信息可信性因素和信息可信性联系起来,有2种实现方式:一种是回归<sup>[17-20]</sup>,另一种是分类<sup>[21-22]</sup>,区别在于前者得到的是连续值,后者得到的是离散值.采用回归方法时通常会使用逻辑回归、最大似然估计等算法,或者是根据具体应用场景设计相应的回归算法;分类方法中会采用支持向量机、决策树等算法.基于概率图模型的信息可信性研究认为:信息源做出的判断、信息源的特征、信息本身的特征等可观测变量的分布依赖于信息可信性、信息源可信性等随机变量,通过建立随机变量和可观测变量之间的关系得到概率图模型.大部分研究采用了贝叶斯网络<sup>[23-25]</sup>,也有研究使用的是条件随机场模型<sup>[26]</sup>.

随着社交媒体的普及,近几年研究重心逐渐从传统网络转移到Twitter<sup>①</sup>、新浪微博等社交媒体平台.目前国外研究涉及到的媒体平台主要包括Twitter<sup>[14,21]</sup>、新浪微博<sup>[20,26]</sup>以及Yelp<sup>[27]</sup>,研究对象可划分为事件层面的信息可信性研究<sup>[14,20]</sup>以及推文层面的可信性研究<sup>[21,26-27]</sup>.文献[14]通过构建推文与信息源、信息源之间的关系图,将推文的可信

性作为隐含变量通过最大期望(expectation maximum, EM)算法进行求解,进而通过投票思想获得事件的可信性.文献[20]通过构建推文间的关系图来将事件的可信性计算转化为图优化问题.社交媒体中往往存在很多噪音数据,比如大量的从众转发等现象,对推文可信性判断带来干扰进而使事件的可信性判断出现不可忽视的偏差.由此可见,从大量噪音数据中筛选出真正有用的推文数据就显得十分重要.推文层面的研究大多依赖于数据集的标注标签以使用传统机器学习方法,考虑到训练集规模较大、人工标注耗费成本较高,我们更倾向于使用无需人工标注的方法,如使用带标签数据集或无监督学习方法.Fontanarava等人<sup>[27]</sup>使用Yelp带标签数据集采用了集成学习的方法,混合了多个模型对Yelp上特定领域的评论可信性进行了研究.他们从评论的语言学特征入手,采用判别模型支持向量机和生成模型循环神经网络对评论内容的可信性进行分析,另一方面采用随机森林,根据用户和评论元数据的特征,对评论的可信性进行了分类.最后将3个模型得到的结果采用线性内插法结合到一起,得到最终的结果.Yelp等评价类网站与新浪微博等内容导向的社交平台的元数据特征存在明显差别,如评价类平台特有的星级等,所以针对内容导向的社交平台仍需挖掘有用特征.

目前国内对社交媒体信息可信性评价的相关工作较少.谢柏林等人<sup>[28]</sup>在2016年的研究中,侧重于及早发现微博中的虚假信息,将转发以及评论内容的观点倾向,结合用户对信息的识别度作为观测值,使用状态持续时间概率为Gamma分布的隐半马尔可夫模型计算原创微博的可信性.除此以外,任亚峰等人<sup>[29]</sup>在2015年针对虚假评论检测进行了研究,该研究考虑到人工标注数据集后采用监督学习的不合理性,基于少量已知正例样本采用PU(positive and unlabeled)学习算法标注未知标签数据,最后在标注数据集中构建多核分类器来检测虚假评论.虽然这些研究开始重视社交网络信息的可信性,但近几年国内在该方面的研究还很少见.

面向信息可信性研究中大多忽略了主题因素,然而用户在不同主题下具有不同的可信性<sup>[30]</sup>.默认信息源在所有主题下具有相同可信性,一方面削弱了信息源在其擅长主题下的可信性,另一方面也增强了信息源在其不擅长主题下的可信性,从而

① <http://twitter.com>



影响最终可信性计算结果的准确性.文献[21]在研究 Twitter 平台中的信息可信性时考虑了用户主题对信息可信性的影响,认为用户的主题与其参与的推文的主题偏差越大,推文和用户的可信性就越低.该研究针对 Twitter 平台的 10 个话题爬取了 2 000 条相关推文,采用人工进行标注分析,一方面模型对标注信息依赖性较高,另一方面数据规模较小,无法充分挖掘潜在特征.

在解决冲突数据相关问题中,需要保证数据源之间的独立性.数据源之间的依赖关系如频繁的拷贝行为,会对最终数据准确性的分析产生影响<sup>[31]</sup>.Dong 等人<sup>[13]</sup>在该问题的研究中,通过贝叶斯建模数据源之间的依赖关系,并据此调整数据源可信性在数据准确性分析中的权重.社交媒体内容可信性分析也应当考虑同样的问题.此外,频繁拷贝信息的信息源不仅经常出现在不可信信息的发布者中,也经常出现在可信的发布者中,并不能从他们发布信息的行为中得到所发布信息可信性的倾向.

综上所述,本文将同时考虑用户的主题因素和拷贝因素对社交媒体中内容可信性进行进一步研究.由于社交媒体中缺乏内容和用户可信性的标记,人工标记难度很大,成本很高,比较适合采用无监督的方法进行研究,而概率图模型比较适合无监督的学习<sup>[32]</sup>,同时具有直观易于理解的特点,所以本文在考虑用户主题和拷贝因素的基础上,使用了贝叶斯网络对社交媒体中内容可信性进行分析和评价.本文的主要贡献在于同时考虑了用户的主题特性和从众行为特性,一方面将可信性评价与用户的擅长领域联系起来,另一方面也降低了社交平台中拷贝内容等噪音数据给可信性评价带来的干扰,最终在新浪微博真实数据集的实验结果表明本文提出的社交媒体内容可信性评价模型相比其他模型更具有适用性.

## 2 内容可信性评价模型

为评价社交媒体信息可信性,本文提出社交媒体内容可信性评价模型 LCEM(latent credibility evaluation model).首先描述模型背后的思想,简要介绍用户的主题因素、从众因素以及各因素与内容可信性之间的关系,然后给出模型的构建过程.

### 2.1 模型思想

在社交媒体中,用户发表或者转发一条内容的行为可以看作是一次投票行为.对于转发微博,其投

票对象是转发微博对应的原始微博.对于原创微博,其投票对象是发表微博所承载的内容信息,可以认为原创微博是一种特殊的转发微博,是将抽象内容转发为具体文本,而不是文本对文本的转发.为了将原创微博与转发微博统一起来,近似认为原创微博的投票对象也是原始微博.如果一个用户发表了原创内容或者单纯转发了他人发表的内容,可以看作是該用户对其发表或者转发的内容投了一次赞成票,表示其认为原始内容是可信的.如果在转发的同时加上了自己对内容的观点,当观点的情感极性是正向的,那么可以当作用户相信原始内容,投出了赞成票;反之,如果评论的情感极性是负向的,那么可以当作用户不认可原始内容,投出了反对票.很显然,不同用户投票对于人们判断内容可信性的参考价值是不一样的.

首先,如果一个投票是在用户从众的情况下产生的,那么意味着这个投票的产生未经过用户的判断,投票中没有赞成和反对的倾向,其产生独立于内容的可信性,所以不具备参考价值.如果用户在非从众的情况下做出了一次投票,表明用户是通过自己的思考,利用相关的知识经验进行了判断.由于知识经验和内容是相关的,所以投票也就与内容的可信性联系在一起,具有参考价值.CNNIC2016 年中国互联网新闻市场研究报告<sup>[33]</sup>中显示,超过 60% 的用户在转发新闻内容的时候并未对内容的可信性进行判断,这些大量的从众投票会严重干扰人们对内容可信性进行判断,所以依据用户的从众行为过滤没有价值的投票显得十分必要.

此外,用户在非从众情况下投票的参考价值也有着很大的差异.如果一个用户在一个主题下比较活跃,那么用户对该主题相关的知识掌握的也就相对较多,也就越容易做出正确的判断,用户在该主题下的投票参考价值也就越大;相反,在用户不熟悉的主题下,用户缺乏判断该主题下内容可信性的知识,不容易做出准确判断,所以这时候用户做出投票的参考价值很小.总的来说,用户非从众情况下投票的参考价值很大程度上取决于用户在投票对象主题下的专业程度.本文将用户的活跃程度视为用户的专业程度.

综合考虑用户的从众行为和主题分布可以很大程度上过滤掉没有价值的投票,提升具有参考价值投票的作用,从而提高对内容可信性判断的准确度.下面从这 2 方面出发,以新浪微博平台为例,从用户视角阐述新浪微博中投票的产生.首先用户打开其

微博主页会看到最新发表或转发的微博,如果用户倾向于从众的话,他很有可能直接转发看到的热门微博.如果该用户独立思考能力比较强,那么他会选择自己感兴趣的微博进行转发,并且会考虑微博的可信性,以一定的概率转发微博,做出投票.这个过程中涉及到 2 种投票的概率:一种是从众情况下的概率;另一种是非从众情况下的概率.对于前者,用户呈现出的态度可能是支持也可能是反对.可以认为用户是从已有的转发微博中随机挑选了一条进行转发,所以他的态度取决于他转发微博的态度.那么从众用户投出赞成票的概率就是用户所处环境下赞成票数占所有票数的比例,即表示支持的转发数占总转发数的比例,投出反对票的概率则是反对票数占总票数的比例.对于后者,用户也会投出赞成票或者反对票,这取决于用户自身的属性.用户在非从众情况下可能赞成了可信的内容(真阳性),也有可能支持了不可信的内容(假阳性),同样也会出现反对可信内容(假阴性)和反对不可信内容(真阴性)的情况.所以用户投赞成票和反对票的概率就是它们在内容可信性下的边缘概率,也就是在内容可信与否 2 种情况下的投票概率之和,其中边缘概率的

计算公式为

$$P(X=x_i)=\sum_{j=1}^{\infty}P(X=x_i|Y=y_j)P(Y=y_j). \tag{1}$$

2.2 模型建立

考虑到缺乏带标记的社交媒体内容可信性数据,本文基于生成模型的思想,采用贝叶斯网络建立了社交媒体内容可信性评价图模型 LCEM,利用盘式记法简化表示为图 1,模型中的各个符号含义如表 1 所示.

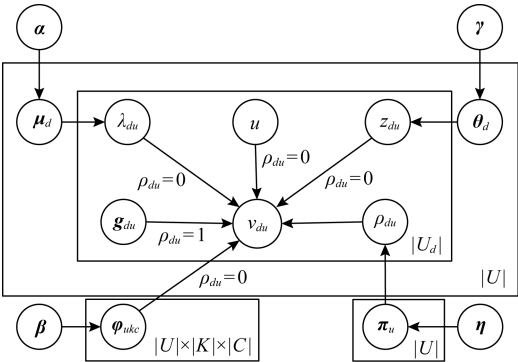


Fig. 1 Latent credibility evaluation model  
图 1 社交媒体内容可信性评价模型

Table 1 The Description of Symbols  
表 1 模型符号说明

| Symbol    | Description   |
|-----------|---|
| $K$       | Number collection of content topic $\langle \{1, 2, \dots,  K \} \rangle$   |
| $U$       | User collection   |
| $D$       | Content collection  |
| $U_d$     | User collection that vote on content $d$  |
| $C$       | The credibility label collection $\{0, 1\}$ . 0 means incredibility, 1 means credibility.   |
| $X$       | The conformity status collection $\{0, 1\}$ . 0 means non-conformity, 1 means conformity.   |
| $Y$       | Voting status collection $\{0, 1\}$ . 0 means negative voting, 1 means positive voting.   |
| $\alpha$  | Parameter of content credibility prior distribution.  |
| $\mu$     | Parameter of content credibility distribution.  |
| $\lambda$ | Content credibility labels. $\lambda = (\lambda_{du}   d \in D, u \in U), \lambda_{du} \in C$ .   |
| $\gamma$  | Parameter of content topic prior distribution.  |
| $\theta$  | Parameters of content topic distribution. $\theta = (\theta_d   d \in D)$ .   |
| $z$       | Content Topic numbers. $z = (z_{du}   d \in D, u \in U), z_{du} \in K$ .  |
| $\beta$   | In non-conformity situation ( $\rho_{du}=0$ ), given the credibility and topic of content, the parameter of the prior distribution of user voting status.   |
| $\phi$    | In non-conformity situation ( $\rho_{du}=0$ ), given the credibility and topic of content, the parameters of user's voting status distribution. $\phi = (\phi_{ukc}   u \in U, k \in K, c \in C)$ . |
| $v$       | Users voting statuses. $v = (v_{du}   d \in D, u \in U), v_{du} \in Y$ .  |
| $u$       | User number. $u \in U$ .  |
| $\eta$    | Parameter for the prior distribution of user's conformity status distribution.  |

Continued (Table 1)

| Symbol   | Description   |
|----------|---|
| $\pi$    | Parameter of the user's conformity status distribution, $\pi = (\pi_d   d \in D)$ .   |
| $\rho$   | User conformity statuses, $\rho = (\rho_{du}   d \in D, u \in U), \rho_{du} \in X$ .  |
| $g_{du}$ | Context for voting record.  |
| $w$      | Voting record, a triple consists of observable variables, $w = (g_{du}, u, v_{du})$ . |
| $W$      | Voting record collection  |
| $W_d$    | Voting record collection of content $d$   |

下面详细描述模型的建立过程,从变量之间的关系建模到整个贝叶斯网络的构建.

1)  $\rho_{du}$  表示用户  $u$  在内容  $d$  中的从众行为.从众行为分为不从众和从众,分别用 0 和 1 表示, $\rho_{du}$  服从伯努利分布,即单次二项分布,假设分布参数  $\pi_u = (\pi_u^0, \pi_u^1)$ ,其中  $\pi_u^0$  和  $\pi_u^1$  分别表示  $\rho_{du}$  取 0 或 1 的概率,那么有:

$$\rho_{du} \sim \text{Bino}(1, \pi_u). \quad (2)$$

2)  $z_{du}$  表示用户  $u$  投票给内容  $d$  时内容  $d$  的主题.其中用户的主题分布可以通过其参与的所有内容的主题来体现.假设有  $|K|$  个主题,那么  $z_{du}$  服从单次的多项分布,假设分布服从的参数为  $\theta_d = (\theta_d^1, \theta_d^2, \dots, \theta_d^{|K|})$ ,其中各分量分别表示内容对应主题的概率,那么有:

$$z_{du} \sim \text{Mutl}(1, \theta_d). \quad (3)$$

3)  $\lambda_{du}$  表示用户  $u$  投票给内容  $d$  时内容  $d$  的可信性,内容可信性分为不可信和可信,分别用 0 和 1 表示, $\lambda_{du}$  服从单次二项分布,假设分布参数为  $\mu_d = (\mu_d^0, \mu_d^1)$ ,其中各分量分别表示内容不可信和可信的概率,那么有:

$$\lambda_{du} \sim \text{Bino}(1, \mu_d). \quad (4)$$

4)  $v_{du}$  表示用户  $u$  对内容  $d$  的投票,投票分为赞成票和反对票,分别用 1 和 0 表示.其服从的分布分为 2 种情况:一种是用户  $u$  在从众情况下产生;另一种是用户  $u$  在非从众情况下产生.

4.1) 若投票是用户  $u$  在非从众情况下产生(如图 1,  $\lambda_{du} \rightarrow v_{du}, z_{du} \rightarrow v_{du}, \phi_{ukc} \rightarrow v_{du}$ ):用户  $u$  做出的投票结果  $v_{du}$  受到内容  $d$  的可信性  $\lambda_{du}$  和主题  $z_{du}$  影响,此时  $v_{du}$  取决于用户  $u$  在内容可信性为  $\lambda_{du}$  以及主题为  $z_{du}$  时做出反对或者支持的概率,本文使用  $\phi_{ukc} = (\phi_{ukc}^0, \phi_{ukc}^1)$  来刻画这个概率,其中  $k$  表示  $z_{du}$  的取值结果, $c$  表示  $\lambda_{du}$  的取值结果,那么有:

$$v_{du} = \text{Bino}(1, \phi_{ukc}). \quad (5)$$

4.2) 若投票是用户  $u$  在从众情况下产生(如图 1,  $g_{du} \rightarrow v_{du}$ ):投票结果  $v_{du}$  受上下文环境  $g_{du} =$

$(g_{du}^0, g_{du}^1)$ ,也就是用户  $u$  投票时内容  $d$  已有赞成票和反对票比例的影响,那么有:

$$v_{du} \sim \text{Bino}(1, g_{du}), \quad (6)$$

上下文环境变量  $g_{du}$  的建模方法为

$$g_{du} = \left( \frac{n_y^0}{n_y^0 + n_y^1}, \frac{n_y^1}{n_y^0 + n_y^1} \right), \quad (7)$$

其中,  $n_y^0$  是投票时已有的反对票数,  $n_y^1$  是投票时已有的赞成票数.但是这样处理会有一个问题,就是当一个转发序列中尚未出现反对票时,用户投赞成票的概率都为 1,这显然是不合理的,可以优化为

$$g_{du} = \left( \frac{n_y^0 + \tau}{n_y^0 + n_y^1 + 2 \times \tau}, \frac{n_y^1 + \tau}{n_y^0 + n_y^1 + 2 \times \tau} \right), \quad (8)$$

其中,参数  $\tau$  是新引入的一个超参数,用于平衡用户投赞成票和反对票的概率.

上述 1)~4) 所有变量中,如图 1,用户从众行为  $\rho_{du}$ 、内容主题  $z_{du}$  和内容可信性  $\lambda_{du}$  将作为模型的隐含变量;投票结果  $v_{du}$ 、用户  $u$  以及上下文环境  $g_{du}$  则作为可观测变量;剩下的用户从众概率分布参数  $\pi_u$ 、内容主题分布参数  $\theta_d$ 、内容可信性分布参数  $\mu_d$  以及用户投票行为分布参数  $\phi_{ukc}$  是待估计参数,也就是需要求解的变量.

为了提高模型的灵活性和进行平滑处理,为每个待估计参数引入相应的先验分布,先验分布的参数就是超参数.首先用户  $u$  在内容  $d$  的从众行为  $\rho_{du}$  服从单次二项分布,那么用户的所有投票的从众行为  $\rho_u = (\rho_{d_1u}, \rho_{d_2u}, \dots)$  服从二项分布,那么有:

$$\rho_u \sim \text{Bino}(n_u, \pi_u), \quad (9)$$

其中,  $n_u$  表示用户  $u$  的投票次数.为便于计算,  $\pi_u$  的分布满足二项分布的共轭先验贝塔分布,也就是:

$$\pi_u \sim \text{Beta}(\eta), \quad (10)$$

其中,超参数  $\eta = (\eta^0, \eta^1)$ ,每个分量表示 0 和 1 的个数.同理有:

$$\theta_d \sim \text{Dir}(\gamma), \quad (11)$$

其中,超参数  $\gamma = (\gamma^0, \gamma^1, \dots, \gamma^{|K|})$ ,

$$\mu_d \sim \text{Beta}(\alpha), \quad (12)$$

其中,超参数  $\alpha = (\alpha^0, \alpha^1)$ ,

$$\varphi_{ukc} \sim \text{Beta}(\beta), \quad (13)$$

其中,超参数  $\beta = (\beta^0, \beta^1)$ .

模型中,  $\{u, v_{du}, g_{du}\}$  是可观测变量,  $\{\mu_d, \varphi_{ukc}, \pi_u, \theta_d\}$  是待估计参数,  $\{\rho_{du}, z_{du}, \lambda_{du}\}$  是隐含变量. 模型的输入是所有投票记录对应的可观测变量和超参数的值, 输出是所有隐含变量以及待估计参数的值.

图 1 中各变量的联合概率分布的抽象表达为

$$P(W, \lambda, z, \rho, \mu, \theta, \varphi, \pi; \alpha, \gamma, \beta, \eta). \quad (14)$$

根据上面提出的社交媒体内容可信性评价模型, 投票产生的具体过程为

1) 对于每一个用户  $u$ 、每一个内容主题  $k$  和每一种内容可信性  $c$ , 从贝塔分布  $\text{Beta}(\varphi_{ukc} | \beta)$  中取样生成非从众情况下用户  $u$  在内容主题为  $k$  和可信性为  $c$  情况下投票行为的分布参数  $\varphi_{ukc}$ ;

2) 对于每个用户  $u$ , 从贝塔分布  $\text{Beta}(\pi_u | \eta)$  中取样生成用户  $u$  的从众行为分布参数  $\pi_u$ ;

3) 对于每条内容  $d$ :

3.1) 从狄利克雷分布  $\text{Dir}(\theta_d | \gamma)$  取样生成内容  $d$  的主题分布  $\theta_d$ ;

3.2) 从贝塔分布  $\text{Beta}(\mu_d | \alpha)$  中取样生成内容的可信性分布  $\mu_d$ ;

3.3) 对于每个投票给内容  $d$  的用户  $u$ :

3.3.1) 从二项分布  $\rho_{du} \sim \text{Bino}(1, \pi_u)$  中取样生成用户  $u$  的从众行为  $\rho_{du}$ ;

3.3.2) 从二项分布  $\lambda_{du} \sim \text{Bino}(1, \mu_d)$  中取样生成内容的可信性标签  $\lambda_{du}$ ;

3.3.3) 从多项分布  $z_{du} \sim \text{Multi}(1, \theta_d)$  中取样生成内容的一个主题  $z_{du}$ ;

3.3.4) 若  $\rho_{du} = 0$ , 则从二项分布  $v_{du} \sim \text{Bino}(1, \varphi_{ukc})$  中取样生成投票  $v_{du}$ , 其中  $k$  表示  $z_{du}$  的取值结果,  $c$  表示  $\lambda_{du}$  的取值结果; 若  $\rho_{du} = 1$ , 则从二项分布  $v_{du} \sim \text{Bino}(1, g_{du})$  中取样生成投票  $v_{du}$ .

### 3 模型求解

在完成概率图模型建立后, 需要针对其中的待估计参数进行求解. 本文在参数估计的过程中采用了吉布斯采样算法. 吉布斯采样作为马尔可夫蒙特卡洛方法的一种特殊情况, 适用于高维数据的采样, 普遍应用于概率图模型中. 采用吉布斯采样求解模型, 最主要的工作是推导隐含变量的采样规则. 根据采样结果可以很容易地计算待估计参数.

#### 3.1 隐含变量联合概率推导

首先给出隐含变量在已知数据, 即数据集和超参下的联合概率分布形式, 表示为

$$P(\lambda, z, \rho | W; \alpha, \gamma, \beta, \eta). \quad (15)$$

引入隐含变量分布参数, 即待估计参数后, 式(15)可表示为

$$\int_{\mu} \int_{\theta} \int_{\varphi} \int_{\pi} P(\lambda, z, \rho, \mu, \theta, \varphi, \pi | W; \alpha, \gamma, \beta, \eta) d\mu d\theta d\varphi d\pi. \quad (16)$$

那么要计算隐含变量的联合概率分布, 需要先计算  $P(\lambda, z, \rho, \mu, \theta, \varphi, \pi | W; \alpha, \gamma, \beta, \eta)$ , 即隐含变量和待估计参数在已知信息下的联合概率分布. 根据贝叶斯公式以及 D-分离规则有:

$$\begin{aligned} & P(\lambda, z, \rho, \mu, \theta, \varphi, \pi | W; \alpha, \gamma, \beta, \eta) \propto \\ & P(W, \lambda, z, \rho | \mu, \varphi, \pi) P(\theta | \gamma) P(\mu | \alpha) \times \\ & P(\varphi | \beta) P(\pi | \eta). \end{aligned} \quad (17)$$

根据概率图模型中各条生成路线, 式(17)可以整理得到:

$$P(\lambda, z, \rho, \mu, \theta, \varphi, \pi | W; \alpha, \gamma, \beta, \eta) \propto \prod_{u \in U} P(\pi_u | \eta) \prod_{d \in D} \prod_{w \in W_d} P(\rho_w | \pi_{u_w}) \times \quad (18-1)$$

$$\prod_{u \in U} \prod_{k \in K} \prod_{c \in C} P(\varphi_{ukc} | \beta) \times \prod_{d \in D} \prod_{w \in W_d} P(v_w | \varphi_{u_w z_w \lambda_w})^{1-\rho_w} \times \quad (18-2)$$

$$\prod_{d \in D} P(\theta_d | \gamma) \prod_{d \in D} \prod_{w \in W_d} P(z_w | \theta_d) \times \quad (18-3)$$

$$\prod_{d \in D} P(\mu_d | \alpha) \prod_{d \in D} \prod_{w \in W_d} P(\lambda_w | \mu_d) \times \quad (18-4)$$

$$\prod_{d \in D} \prod_{w \in W_d} P(v_w | g_w)^{\rho_w}. \quad (18-5)$$

(18)

表达式(18-1)对应图 1 中  $\eta \rightarrow \pi_u \rightarrow \rho_{du}$  生成路线, 表示用户从众行为的先验分布中采样生成用户从众行为的分布, 然后从该分布中采样出用户是否从众. 同理表达式(18-2)对应路线  $\beta \rightarrow \varphi_{ukc} \rightarrow v_{du}$ ; 表达式(18-3)对应路线  $\gamma \rightarrow \theta_d \rightarrow z_{du}$ ; 表达式(18-4)对应路线  $\alpha \rightarrow \mu_d \rightarrow \lambda_{du}$ ; 表达式(18-5)对应路线  $g_{du} \rightarrow v_{du}$ . 对于投票结果  $v_w$ , 其生成路径分别对应表达式(18-2)和表达式(18-5)这 2 种不同的情况, 由公式的指数上标也就是投票记录对应的用户从众行为  $\rho_w$  决定. 当该投票是在用户非从众 ( $\rho_w = 0$ ) 情况下产生时, 其生成路径对应表达式(18-2), 当该投票在用户从众 ( $\rho_w = 1$ ) 情况下产生时, 其生成路径对应表达式(18-5).

将式(18)带入式(16)计算隐含变量的概率分布, 并且将多重积分根据积分变量进行转化来简化



计算复杂度,整理为

$$P(\lambda, z, \rho | W; \alpha, \gamma, \beta, \eta) \propto$$

$$\int_{\pi} \prod_{u \in U} P(\pi_u | \eta) \prod_{d \in D} \prod_{w \in W_d} P(\rho_w | \pi_{u_w}) d\pi \times \quad (19-1)$$

$$\int_{\varphi} \prod_{u \in U} \prod_{k \in K} \prod_{c \in C} P(\varphi_{ukc} | \beta) \times \prod_{d \in D} \prod_{w \in W_d} P(v_w | \varphi_{u_w z_w \lambda_w})^{1-\rho_w} d\varphi \times \quad (19-2)$$

$$\int_{\theta} \prod_{d \in D} P(\theta_d | \gamma) \prod_{d \in D} \prod_{w \in W_d} P(z_w | \theta_d) d\theta \times \quad (19-3)$$

$$\int_{\mu} \prod_{d \in D} P(\mu_d | \alpha) \prod_{d \in D} \prod_{w \in W_d} P(\lambda_w | \mu_d) d\mu \times \quad (19-4)$$

$$\prod_{d \in D} \prod_{w \in W_d} P(v_w | g_w)^{\rho_w}. \quad (19-5)$$

(19)

对于表达式(19-1),假设不同用户的从众行为相互独立.其中  $n_u^x$  表示观察到的用户  $u$  做出从众行为  $x$  的次数,  $n_u = (n_u^0, n_u^1)$ .  $B(\cdot)$  表示贝塔函数:

$$\int_{\pi} \prod_{u \in U} P(\pi_u | \eta) \prod_{d \in D} \prod_{w \in W_d} P(\rho_w | \pi_{u_w}) d\pi = \left( \frac{1}{B(\eta)} \right)^{|U|} \prod_{u \in U} \int_{\pi_u} \prod_{x \in X} \pi_{u,x}^{\eta^x + n_u^x - 1} d\pi_u \propto \prod_{u \in U} B(\eta + n_u). \quad (20)$$

对于表达式(19-2),假设不同投票的产生之间相互独立.其中,  $n_{ukc}^y$  表示观察到的用户  $u$  在主题  $k$  和可信性为  $c$  下做出投票  $y$  的次数,  $n_{ukc} = (n_{ukc}^0, n_{ukc}^1)$ .

$$\int_{\varphi} \prod_{u \in U} \prod_{k \in K} \prod_{c \in C} P(\varphi_{ukc} | \beta) \times \prod_{d \in D} \prod_{w \in W_d} P(v_w | \varphi_{u_w z_w \lambda_w})^{1-\rho_w} d\varphi = \left( \frac{1}{B(\beta)} \right)^{|U| \times |K| \times |C|} \prod_{u \in U} \prod_{k \in K} \prod_{c \in C} \int_{\varphi_{ukc}} \prod_{y \in Y} \varphi_{ukc,y}^{\beta^y + n_{ukc}^y - 1} d\varphi_{ukc} \propto \prod_{u \in U} \prod_{k \in K} \prod_{c \in C} B(\beta + n_{ukc}). \quad (21)$$

对于表达式(19-3),假设不同内容的主题相互独立.其中,  $n_d^k$  表示观察到的内容  $d$  中所有用户出现在主题  $k$  下的次数,  $n_d = (n_d^1, n_d^2, \dots, n_d^{|K|})$ .

$$\int_{\theta} \prod_{d \in D} P(\theta_d | \gamma) \prod_{d \in D} \prod_{w \in W_d} P(z_w | \theta_d) d\theta = \left( \frac{1}{B(\gamma)} \right)^{|D|} \prod_{d \in D} \int_{\theta_d} \prod_{k \in K} \theta_{d,k}^{\gamma^k + n_d^k - 1} d\theta_d \propto \prod_{d \in D} B(\gamma + n_d). \quad (22)$$

对于表达式(19-4),假设不同内容的可信性相互独立.其中  $n_d^c$  表示观察到内容  $d$  中所有用户出现在可信性  $c$  下的次数,  $n_d' = (n_d'^0, n_d'^1)$ .

$$\int_{\mu} \prod_{d \in D} P(\mu_d | \alpha) \prod_{d \in D} \prod_{w \in W_d} P(\lambda_w | \mu_d) d\mu = \left( \frac{1}{B(\alpha)} \right)^{|D|} \prod_{d \in D} \int_{\mu_d} \prod_{c \in C} \mu_{d,c}^{\alpha^c + n_d'^c - 1} d\mu_d \propto \prod_{d \in D} B(\alpha + n_d'). \quad (23)$$

并且表达式(19-5)可以转化为

$$\prod_{d \in D} \prod_{w \in W_d} P(v_w | g_w)^{\rho_w} = \prod_{w \in W} g_w^{\rho_w}. \quad (24)$$

至此,结合式(20)~(24),可以得到隐含变量的联合概率分布:

$$P(\lambda, z, \rho | W; \alpha, \gamma, \beta, \eta) \propto F_1 \times F_2 \times F_3 \times F_4 \times F_5.$$

$$\begin{aligned} F_1 &= \prod_{u \in U} B(\eta + n_u); \\ F_2 &= \prod_{u \in U} \prod_{k \in K} \prod_{c \in C} B(\beta + n_{ukc}); \\ F_3 &= \prod_{d \in D} B(\gamma + n_d); \\ F_4 &= \prod_{d \in D} B(\alpha + n_d'); \\ F_5 &= \prod_{w \in W} g_w^{\rho_w}. \end{aligned} \quad (25)$$

### 3.2 采样算法

在3.1节隐含变量的联合概率分布的推导基础上,继续阐述隐含变量的状态转移分布推导过程,并给出 LCEM 的吉布斯采样算法.

根据吉布斯采样算法,LCEM 的转移概率为

$$P(\lambda_o, z_o, \rho_o | \lambda_{\neg o}, z_{\neg o}, \rho_{\neg o}, W; \alpha, \gamma, \beta, \eta) \propto \frac{P(\lambda, z, \rho | W; \alpha, \gamma, \beta, \eta)}{P(\lambda_{\neg o}, z_{\neg o}, \rho_{\neg o} | W_{\neg o}; \alpha, \gamma, \beta, \eta)} = \frac{F_1}{F_{1,\neg o}} \times \frac{F_2}{F_{2,\neg o}} \times \frac{F_3}{F_{3,\neg o}} \times \frac{F_4}{F_{4,\neg o}}, \quad (26)$$

其中  $(\lambda_o, z_o, \rho_o)$  表示与一个投票  $v_o$  对应的隐含变量,  $\{\lambda_{\neg o}, z_{\neg o}, \rho_{\neg o}\}$  表示剔除该投票  $v_o$  对应的隐含变量后剩余投票对应的隐含变量.可以看出需要采样的隐含变量的转移概率同所有隐含变量的联合概率与剔除该组变量的隐含变量的联合概率比值成正比.并且可以使用式(15)的形式来表示联合概率,整个转移概率公式推导可拆分成对每一部分的推导.下面具体推导式(26)中  $F_1$  和  $F_{1,\neg o}$  的关系,其中  $\Gamma(\cdot)$  表示伽玛函数:

$$\begin{aligned} F_1 &= \prod_{u \in U} B(\eta + n_u) = \\ &= B(\eta + n_{u_o}) \prod_{u \in U/u_o} B(\eta + n_u) = \\ &= \frac{\eta^{\rho_o} + n_{u_o}^{\rho_o, \neg o}}{\sum_{x \in X} (\eta^x + n_{u_o}^x, \neg o)} \frac{\Gamma(\eta^{\rho_o} + n_{u_o}^{\rho_o, \neg o})}{\Gamma(\sum_{x \in X} (\eta^x + n_{u_o}^x, \neg o))} \times \\ &= \prod_{x \in X/\rho_o} \Gamma(\eta^x + n_{u_o}^x) \prod_{u \in U/u_o} B(\eta + n_u) = \end{aligned}$$



$$\frac{\eta^{\rho_o} + n_{u_o, \neg o}^{\rho_o}}{\sum_{x \in X} (\eta^x + n_{u_o, \neg o}^x)} \prod_{u \in U} B(\boldsymbol{\eta} + \mathbf{n}_{u, \neg o}) = \frac{\eta^{\rho_o} + n_{u_o, \neg o}^{\rho_o}}{\sum_{x \in X} (\eta^x + n_{u_o, \neg o}^x)} F_{1, \neg o}. \quad (27)$$

同式(27)推导过程,式(26)中  $F_2$  和  $F_{2, \neg o}$  的关系为

$$F_2 = \frac{\beta^{v_o} + n_{u_o z_o \lambda_o, \neg o}^{v_o}}{\sum_{y \in Y} (\beta^y + n_{u_o z_o \lambda_o, \neg o}^y)} F_{2, \neg o}. \quad (28)$$

式(26)中  $F_3$  和  $F_{3, \neg o}$  的关系为

$$F_3 = \frac{\gamma^{z_o} + n_{d_o, \neg o}^{z_o}}{\sum_{k \in K} (\gamma^k + n_{d_o, \neg o}^k)} F_{3, \neg o}. \quad (29)$$

式(26)中  $F_4$  和  $F_{4, \neg o}$  的关系为

$$F_4 = \frac{\alpha^{\lambda_o} + n_{d_o, \neg o}^{\lambda_o}}{\sum_{c \in C} (\alpha^c + n_{d_o, \neg o}^c)} F_{4, \neg o}. \quad (30)$$

式(26)中  $F_5$  和  $F_{5, \neg o}$  的关系为

$$F_5 = (g^{v_o})^{\rho_o} F_{5, \neg o}. \quad (31)$$

综合式(27)~(31),一组隐含变量转移概率的具体表达形式为

$$P(\lambda_o, z_o, \rho_o \mid \boldsymbol{\lambda}_{\neg o}, \mathbf{z}_{\neg o}, \boldsymbol{\rho}_{\neg o}, W; \boldsymbol{\alpha}, \boldsymbol{\gamma}, \boldsymbol{\beta}, \boldsymbol{\eta}) \propto \frac{\eta^{\rho_o} + n_{u_o, \neg o}^{\rho_o}}{\sum_{x \in X} (\eta^x + n_{u_o, \neg o}^x)} \times \left( \frac{\beta^{v_o} + n_{u_o z_o \lambda_o, \neg o}^{v_o}}{\sum_{y \in Y} (\beta^y + n_{u_o z_o \lambda_o, \neg o}^y)} \right)^{1-\rho_o} \times \frac{\gamma^{z_o} + n_{d_o, \neg o}^{z_o}}{\sum_{k \in K} (\gamma^k + n_{d_o, \neg o}^k)} \times \frac{\alpha^{\lambda_o} + n_{d_o, \neg o}^{\lambda_o}}{\sum_{c \in C} (\alpha^c + n_{d_o, \neg o}^c)} \times (g^{v_o})^{\rho_o}, \quad (32)$$

其中  $\rho_o \in \{0, 1\}$ , 若当前隐含变量对应的投票记录在从众情况下产生, 即  $\rho_o = 0$ , 最终的概率与上下文  $g_{v_o}$  无关, 同理  $\rho_o = 1$ , 最终概率与第 2 项无关。

对其中某个隐含变量进行采样时, 另外 2 个变量作为隐含变量的固定值. 所以该隐含变量的采样概率只和式(32)中的相关项有关, 其他项在当前采样过程中作为常量. 所以各隐含变量的采样规则为

$$P(\lambda_o \mid \boldsymbol{\lambda}_{\neg o}, \mathbf{z}_{\neg o}, \boldsymbol{\rho}_{\neg o}, W; \boldsymbol{\alpha}, \boldsymbol{\gamma}, \boldsymbol{\beta}, \boldsymbol{\eta}) \propto \left( \frac{\beta^{v_o} + n_{u_o z_o \lambda_o, \neg o}^{v_o}}{\sum_{y \in Y} (\beta^y + n_{u_o z_o \lambda_o, \neg o}^y)} \right)^{1-\rho_o} \times \frac{\alpha^{\lambda_o} + n_{d_o, \neg o}^{\lambda_o}}{\sum_{c \in C} (\alpha^c + n_{d_o, \neg o}^c)}, \quad (33)$$

$$P(z_o \mid \boldsymbol{\lambda}_{\neg o}, \mathbf{z}_{\neg o}, \boldsymbol{\rho}_{\neg o}, W; \boldsymbol{\alpha}, \boldsymbol{\gamma}, \boldsymbol{\beta}, \boldsymbol{\eta}) \propto \left( \frac{\beta^{v_o} + n_{u_o z_o \lambda_o, \neg o}^{v_o}}{\sum_{y \in Y} (\beta^y + n_{u_o z_o \lambda_o, \neg o}^y)} \right)^{1-\rho_o} \times \frac{\gamma^{z_o} + n_{d_o, \neg o}^{z_o}}{\sum_{k \in K} (\gamma^k + n_{d_o, \neg o}^k)}, \quad (34)$$

$$P(\rho_o \mid \boldsymbol{\lambda}_{\neg o}, \mathbf{z}_{\neg o}, \boldsymbol{\rho}_{\neg o}, W; \boldsymbol{\alpha}, \boldsymbol{\gamma}, \boldsymbol{\beta}, \boldsymbol{\eta}) \propto \left( \frac{\beta^{v_o} + n_{u_o z_o \lambda_o, \neg o}^{v_o}}{\sum_{y \in Y} (\beta^y + n_{u_o z_o \lambda_o, \neg o}^y)} \right)^{1-\rho_o} \times \frac{\eta^{\rho_o} + n_{u_o, \neg o}^{\rho_o}}{\sum_{x \in X} (\eta^x + n_{u_o, \neg o}^x)} \times (g^{v_o})^{\rho_o}. \quad (35)$$

根据隐含变量的采样公式对隐含变量进行采样, 将采样得到的结果作为后验知识, 结合事先设定的先验知识, 利用先验分布和后验分布的共轭关系, 可以得到各个待估计参数的计算规则:

$$\hat{\boldsymbol{\mu}}_d = E(\boldsymbol{\mu}_d) = \left( \frac{\boldsymbol{\alpha}^{\lambda'} + n_d^{\lambda'}}{\sum_{c \in C} (\alpha^c + n_d^c)} \right)_{\lambda' \in C}, \quad (36)$$

$$\hat{\boldsymbol{\theta}}_d = E(\boldsymbol{\theta}_d) = \left( \frac{\boldsymbol{\gamma}^{z'} + n_d^{z'}}{\sum_{k \in K} (\gamma^k + n_d^k)} \right)_{z' \in K}, \quad (37)$$

$$\hat{\boldsymbol{\pi}}_u = E(\boldsymbol{\pi}_u) = \left( \frac{\boldsymbol{\eta}^{\rho'} + n_u^{\rho'}}{\sum_{x \in X} (\eta^x + n_u^x)} \right)_{\rho' \in X}, \quad (38)$$

$$\hat{\boldsymbol{\phi}}_{ukc} = E(\boldsymbol{\phi}_{ukc}) = \left( \frac{\beta^{v'} + n_{ukc}^{v'}}{\sum_{y \in Y} (\beta^y + n_{ukc}^y)} \right)_{v' \in Y}. \quad (39)$$

根据这些规则就可以得到本文提出的可信性评价模型 LCEM 的吉布斯采样算法. 算法输入是所有内容对应的投票记录集合  $W$ 、内容所有主题类别  $K$ 、内容可信性先验分布参数  $\boldsymbol{\alpha}$ 、内容主题先验分布参数  $\boldsymbol{\gamma}$ 、用户从众行为先验分布参数  $\boldsymbol{\eta}$ 、用户在不同主题和可信性下投票行为的先验分布参数  $\boldsymbol{\beta}$ 、上下文环境变量平衡参数  $\tau$ , 以及采样迭代次数  $I$ . 算法输出包括内容可信性分布  $\boldsymbol{\mu}$ 、内容主题分布  $\boldsymbol{\theta}$ 、用户从众行为分布  $\boldsymbol{\pi}$ 、用户在不同主题和内容可信性下投票行为分布  $\boldsymbol{\phi}$ , 以及所有隐含变量  $\{\boldsymbol{\lambda}, \mathbf{z}, \boldsymbol{\rho}\}$ . 详细过程如算法 1 所示。

**算法 1.** LCEM 吉布斯采样算法.

输入:  $\{W, K, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\eta}, \tau, I\}$ ;

输出:  $\{\boldsymbol{\mu}, \boldsymbol{\theta}, \boldsymbol{\pi}, \boldsymbol{\phi}, \boldsymbol{\lambda}, \mathbf{z}, \boldsymbol{\rho}\}$ .

① /\* 隐含变量初始化 \*/

② for all  $d \in D$  do

③ for all  $w \in W_d$  do

④ /\* 从均匀分布中对隐含变量随机赋值 \*/

⑤  $\lambda_w \sim U(0, 1)$ ;

⑥  $z_w \sim U(1, |K|)$ ;

⑦  $\rho_w \sim U(0, 1)$ ;

⑧ end for

⑨ end for

⑩ /\* 吉布斯采样过程 \*/

```

⑪ for  $i=1$  to  $I$  do
⑫   for all  $d \in D$  do
⑬     for all  $w \in W_d$  do
⑭        $\lambda'_w \sim \text{式}(33)$ ;
⑮        $z'_w \sim \text{式}(34)$ ;
⑯        $\rho'_w \sim \text{式}(35)$ ;
⑰       assign  $\{\lambda'_w, z'_w, \rho'_w\}$  to  $\{\lambda_w, z_w, \rho_w\}$ ;
⑱     end for
⑲   end for
⑳ end for
㉑ /* 更新待估计参数 */
㉒ 利用式(36)~(39)更新参数  $\{\mu_d, \theta_d, \pi_u,$   

    $\varphi_{ukc}\}$ , 得到  $\{\hat{\mu}_d, \hat{\theta}_d, \hat{\pi}_u, \hat{\varphi}_{ukc}\}$ .
```

4 模型评价

本节采用真实社交媒体平台的数据来验证本文提出的模型.采用的数据来自于文献[34]的新浪微博公开数据集,数据中有 3 万条原创微博、3 700 万条转发微博、140 万个用户.从数据中可以提取出可观测变量的值,从而得到模型的输入.

4.1 参数设定

需要设定的参数包括迭代次数  $I$ 、微博主题类别  $K$ ,以及各先验分布的超参.  
首先对迭代次数  $I$  的设定.由于吉布斯采样是

一个随机化求解方法,无法保证迭代确定次数后收敛.本文将迭代次数设定为一个较大值 1 000,通过观察困惑度(perplexity)来判断是否收敛.困惑度是一个用于衡量概率模型拟合程度的量,值越小表示拟合效果越好.随着采样进行,困惑度会不断减少,当困惑度变化范围小于一定阈值,则认为其收敛,实验设定阈值为 0.001.困惑度计算方法为

$$Perplexity(V)=\exp\left\{-\frac{\mathcal{L}(V)}{|V|}\right\}, \tag{40}$$

其中, $V$  表示所有记录集合,本文表示所有投票集合, $\mathcal{L}(V)$ 是似然函数,计算方法为  
$$\mathcal{L}(V)=\ln P(V|\Phi,\Psi)=\sum_{i \in W} \ln P(v_i|\Phi,\Psi), \tag{41}$$
其中  $\Phi$  表示待估计参数集合, $\Psi$  是超参数集合.

本文中模型的似然函数可以表示为

$$\mathcal{L}(V)=\sum_{i \in W} \ln(P(0|\pi_{u_i})(\sum_{k \in K} \sum_{c \in C} P(v_i|\varphi_{ukc}) P(k|\theta)P(c|\mu)))+P(1|\pi_{u_i})P(v_i|g_i)). \tag{42}$$

对于主题类别  $K$  的设定,也就是主题类别个数的设定,本文采用 HDP(hierarchical Dirichlet processes)模型<sup>[35]</sup>. HDP 模型可以看做是 LDA(latent Dirichlet allocation)模型<sup>[36]</sup>的扩展,是非参数化的 LDA 模型,可以自动调整主题个数,达到不用人工确定主题个数的目的.本文将微博文本输入 HDP 模型,经过 5 天 32 182 次迭代,得到如图 2 所示的主题与困惑度关系:

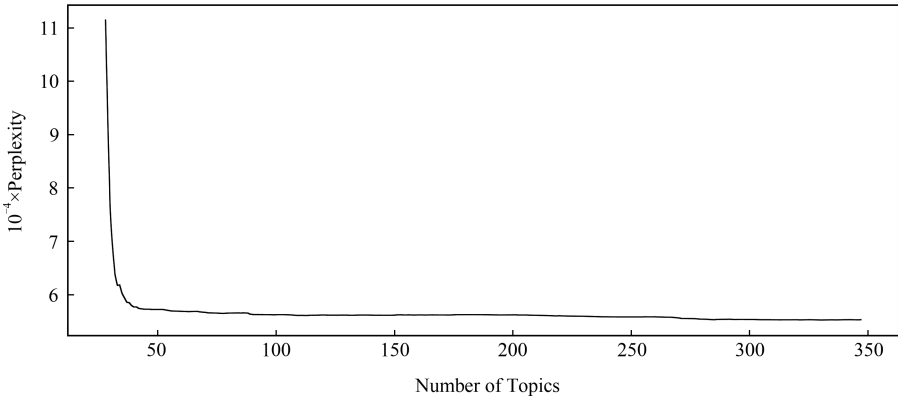


Fig. 2 The relation between number of topics and perplexity  
图 2 主题数和困惑度的关系

从图 2 中可以发现主题数 100 之后,困惑度趋于平稳,所以本文将主题数确定为 100.  
对于  $\{\alpha, \gamma, \eta, \beta\}$  这些先验分布的超参数,假设它们的各个分量都相等,即  $\alpha=\alpha' \times (1,1), \gamma=\gamma' \times$

$(1,1,\cdots,1), \eta=\eta' \times (1,1), \beta=\beta' \times (1,1)$  那么对向量的设定就可以转化为对标量的设定,即对系数  $\{\alpha', \gamma', \eta', \beta'\}$  的设定.设定  $\{\alpha', \gamma', \eta', \beta', \tau\}$  这些参数时,本文利用了贝叶斯优化工具<sup>①</sup> 搜寻合适的

① <https://github.com/fmfn/BayesianOptimization>

参数,设定的搜索区间为 $\alpha' \in [0.01, 2]$ , $\gamma' \in [0.01, 2]$ , $\eta' \in [0.01, 2]$ , $\beta' \in [0.01, 2]$ , $\tau \in [0.01, 100]$ ,搜寻结果为 $\alpha' = 0.01$ , $\gamma' = 0.01$ , $\eta' = 0.01$ , $\beta' = 0.01$ , $\tau = 94.47$ .

4.2 实验结果

为了验证 LCEM 模型的有效性,本文选取了 6 个模型进行对比.

1) LCEM/H.本文提出的社交媒体内容可信性模型除去用户从众因素,只利用用户的主题因素来辅助内容可信性的判断.

2) LCEM/T.本文提出的社交媒体内容可信性模型除去用户主题因素,只利用用户的从众因素来辅助内容可信性的判断.

3) TruthFinder<sup>[37]</sup>.该方法是一种迭代模型,通过信息源(source)建立事实(fact)之间的联系,采用类似于 PageRank 的方法计算 fact 的可信性.

4) LTM<sup>[24]</sup>.该模型也是概率图模型,其思想是各个 fact 中每个 source 做出的声明(claim)受到 fact 可信与否的影响,利用这个影响关系来判断 fact 的可信性.

5) KDEm<sup>[17]</sup>.该模型是一种回归模型,采用了核密度估计的思想,将同一 fact 的所有 claim 映射到函数空间,将用户的可信性作为权重,对 fact 的可信性进行拟合.

6) CATD<sup>[38]</sup>.该模型也是回归模型,通过 fact 可信性与 source 之间的关联,建立优化目标,在计算 source 权重时考虑了 source 发表的 claim 数服从幂律分布,每个 source 权重的置信度会有很大差

别,根据置信度来修正权重.

其中 TruthFinder, KDEm, CATD 有开源代码<sup>①</sup>,由文献[17]提供.这 6 个模型中的 source, claim, fact 分别对应着本文研究场景中的用户、投票、微博.

由于本文采用的公开数据集中并不携带内容可信或者不可信的标签,常规的 F1 值评价方法并不适用.本文将采用的评价方法为:取实验结果中可信性最高的 100 条微博和可信性最低的 100 条微博,采用人工的方式判断前 100 条中可信微博的数量和后 100 条微博中不可信微博的数量,将前 100 条中可信微博的比例和后 100 条中不可信微博的比例作为评价指标.

对各个模型中输出的内容可信性评分排序,提取出可信性最高的 100 条微博和可信性最低的 100 条微博,得到的对比结果如图 3 所示.本文提出模型的准确程度都要高于其他模型,即使除去用户主题因素的考虑,相比其他模型也具有一定的优势.不考虑从众因素的情况下,效果也和其他模型中最好的相差无几.其中的原因是对比模型是建立在用户行为差异比较大的基础上,即所有用户投出的赞成票数和反对票数差别较小.但是在社交媒体中反对票数本来就远小于赞成票数,加上从众用户的存在,它们悬殊更加巨大.在本文使用的数据集中,根据情感分析得到的赞成票数和反对票数的比值达到了 900.而本文从用户的从众因素和主题因素 2 个角度弱化了这种负面影响,得到了相对于其他模型较好的结果.虽然 TruthFinder 也考虑到用户之间存在着拷贝,

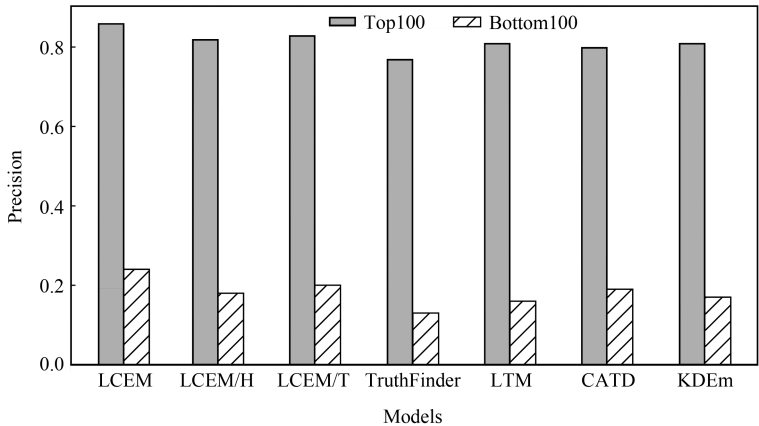


Fig. 3 The precision of credibility top100 and bottom100 microblogs

图 3 可信性 Top100 以及 Bottom100 微博的精确率

① <https://github.com/MengtingWan/KDEm>

但是只是单纯地为所有内容的可信性加上了一个相同衰减系数,并不影响最终可信性的排名.图3也体现了用户从众因素对内容可信性评价的影响大于用户的主题因素,原因在于用户参与其不熟悉的主题往往反映了一定的从众倾向,即从众因素中包含了部分主题因素.

同时,从图3中可以明显看到,在Top100中可信微博的比例都比较高,而在Bottom100中不可信微博的比例都很低.究其原因,一方面数据集中可信内容数要远大于不可信内容数;另一方面,用户在参与负面新闻时往往持批判的态度,即根据情感极性分析得到的是反对票,但实际上是赞成票.这样就导致了在Bottom100中负面新闻占据了很大一部分,使得真正不可信的内容减少.

图4是LCEM,LCEM/H,LCEM/T这3个模型困惑度随着迭代次数变化的情形.在进行100次迭代之后,困惑度的变化幅度已经不明显.可以发现LCEM的困惑度是最大的,其次是LCEM/H,最后是LCEM/T.虽然之前提到了困惑度可以用于评价概率图模型对数据拟合的好坏程度,但是跨模型对比困惑度是没有意义的,首先是不同模型超参数设置会有差别,另外不同模型中的变量含义也不同.所以困惑度大的模型不一定效果差,反之困惑度小的模型效果不一定就好.图4中,LCEM/H和LCEM/T的困惑度小于LCEM的是因为减少了投票产生的约束条件,使得投票出现概率增大,从而减小了困惑度.另外可以看到LCEM/T困惑度的降幅要大于LCEM/H的,这进一步说明了用户从众行为对内容可信性判断的影响要大于用户主题分布对内容可信性判断的影响.

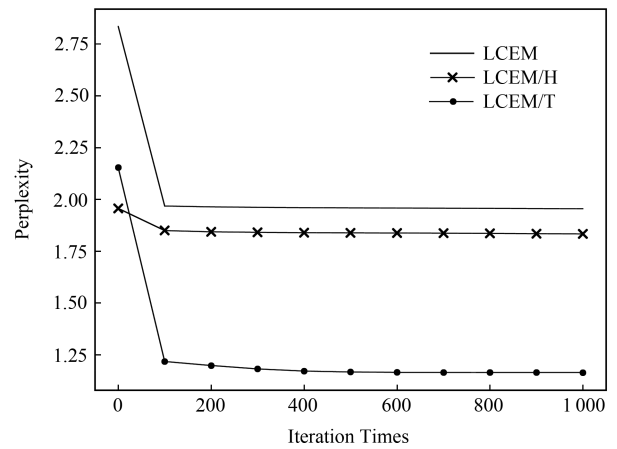


Fig. 4 The perplexity of LCEM, LCEM/H and LCEM/T  
图4 模型LCEM,LCEM/H,LCEM/T的困惑度

5 总结与展望

本文致力于解决的问题是社交媒体中内容可信性判断的问题.针对该问题,考虑到在社交媒体中用户在消费信息时有跟风的倾向和选择自己感兴趣信息的倾向,本文从用户的从众因素和主题因素以及内容的可信性因素出发,对用户发表或传播内容时持有的支持或反对态度进行分析建模,从而实现对内容可信性的评价.实验结果表明,本文提出的模型更加适合社交媒体中内容可行性的评价.

虽然相比现有的内容可信性评价模型,本文提出的模型具有较好的效果,但是本文模型在以下方面仍有改进的空间:提高评论支持或反对的计算准确程度;更加准确地衡量用户转发内容时的上下文环境;除了考虑用户的从众行为,加入用户对特定用户的依赖能够提高可信性判断的准确程度.

参 考 文 献

[1] DataReportal. Digital 2019: Global digital overview [EB/OL]. (2019-01-31) [2019-02-28]. <https://datareportal.com/reports/digital-2019-global-digital-overview>

[2] Xu Lingbei. 2017 China social media landscape [EB/OL]. (2017-07-12) [2019-02-28]. <https://cn.kantar.com/媒体动态/社交/2017/2017年中国社会化媒体格局概览/> (in Chinese) (徐凌蓓. 2017年中国社会化媒体格局概览[EB/OL]. (2017-07-12) [2019-02-28]. <https://cn.kantar.com/媒体动态/社交/2017/2017年中国社会化媒体格局概览/>)

[3] Shah A A, Ravana S D, Hamid S, et al. Web credibility assessment: Affecting factors and assessment techniques [J]. Information Research, 2015, 20(1): 365-391

[4] Fogg B J, Tseng H. The elements of computer credibility [C] //Proc of the SIGCHI Conf on Human Factors in Computing Systems. New York: ACM, 1999: 80-87

[5] Page L, Brin S, Motwani R, et al. The PageRank citation ranking: Bringing order to the Web [R]. Palo Alto, CA: Stanford InfoLab, 1999

[6] Ziegler C N, Lausen G. Propagation models for trust and distrust in social networks [J]. Information Systems Frontiers, 2005, 7(4/5): 337-358

[7] Gyöngyi Z, Garcia-Molina H, Pedersen J. Combating Web spam with trustrank [C] //Proc of the 30th Int Conf on Very Large Data Bases. San Francisco, CA: Morgan Kaufmann, 2004: 576-587

[8] Caverlee J, Liu Ling. Countering Web spam with credibility-based link analysis [C] //Proc of the 26th Annual ACM Symp on Principles of Distributed Computing. New York: ACM, 2007: 157-166



- [9] Xue Jilong, Yang Zhi, Yang Xiaoyong, et al. VoteTrust: Leveraging friend invitation graph to defend against social network sybils [C] //Proc of IEEE INFOCOM'13. Piscataway, NJ; IEEE, 2013: 2400-2408
- [10] Momeni E, Cardie C, Diakopoulos N. A survey on assessment and ranking methodologies for user-generated content on the Web [J]. ACM Computing Surveys, 2016, 48(3): 41:1-41:49
- [11] Pei Qingqi, Yan Dingyu, Ma Lichuan, et al. A strong and weak ties feedback-based trust model in multimedia social networks [J]. The Computer Journal, 2015, 58(4): 627-643
- [12] Pasternack J, Roth D. Knowing what to believe (when you already know something) [C] //Proc of the 23rd Int Conf on Computational Linguistics. Stroudsburg, PA; ACL, 2010: 877-885
- [13] Dong X L, Berti-Equille L, Srivastava D. Integrating conflicting data: The role of source dependence [J]. Proceedings of the VLDB Endowment, 2009, 2(1): 550-561
- [14] Wang Dong, Amin M T, Li Shen, et al. Using humans as sensors: An estimation-theoretic perspective [C] // Proc of the 13th Int Symp on Information Processing in Sensor Networks. Piscataway, NJ; IEEE, 2014: 35-46
- [15] Vydiswaran V G, Zhai Chengxiang, Roth D. Content-driven trust propagation framework [C] //Proc of the 17th ACM SIGKDD Int Conf on Knowledge Discovery and Data Mining. New York; ACM, 2011: 974-982
- [16] Zhao Liang, Hua Ting, Lu C T, et al. A topic-focused trust model for Twitter [J]. Computer Communications, 2016, 76: 1-11
- [17] Wan Mengting, Chen Xiangyu, Kaplan L, et al. From truth discovery to trustworthy opinion discovery: An uncertainty-aware quantitative modeling approach [C] //Proc of the 22nd ACM SIGKDD Int Conf on Knowledge Discovery and Data Mining. New York; ACM, 2016: 1885-1894
- [18] Popat K, Mukherjee S, Strötgen J, et al. Credibility assessment of textual claims on the Web [C] //Proc of the 25th ACM Int on Conf on Information and Knowledge Management. New York; ACM, 2016: 2173-2178
- [19] Wu Shu, Liu Qiang, Liu Yong, et al. Information credibility evaluation on social media [C] //Proc of the 30th AAAI Conf on Artificial Intelligence. Menlo Park, CA; AAAI, 2016: 4403-4404
- [20] Jin Zhiwei, Cao Juan, Zhang Yongdong, et al. News verification by exploiting conflicting social viewpoints in microblogs [C] //Proc of the 30th AAAI Conf on Artificial Intelligence. Menlo Park, CA; AAAI, 2016: 2972-2978
- [21] Ito J, Song Jing, Toda H, et al. Assessment of tweet credibility with LDA features [C] //Proc of the 24th Int Conf on World Wide Web. New York; ACM, 2015: 953-958
- [22] Castillo C, Mendoza M, Poblete B. Information credibility on Twitter [C] //Proc of the 20th Int Conf on World Wide Web. New York; ACM, 2011: 675-684
- [23] Zhao Bo, Han Jiawei. A probabilistic model for estimating real-valued truth from conflicting sources [C/OL]. //Proc of the 10th Int Workshop on Quality in Databases. 2012 [2019-02-28]. [http://hanj.cs.illinois.edu/pdf/qdb12\\_bzhao.pdf](http://hanj.cs.illinois.edu/pdf/qdb12_bzhao.pdf)
- [24] Zhao Bo, Rubinstein B I P, Gemmell J, et al. A Bayesian approach to discovering truth from conflicting sources for data integration [J]. Proceedings of the VLDB Endowment, 2012, 5(6): 550-561
- [25] Pasternack J, Roth D. Latent credibility analysis [C] //Proc of the 22nd Int Conf on World Wide Web. New York; ACM, 2013: 1009-1020
- [26] Guo Qiaozhen, Huang Wei (Wayne), Huang Kai, et al. Information credibility: A probabilistic graphical model for identifying credible influenza posts on social media [C] //Proc of the Int Conf on Smart Health 2015. Berlin; Springer, 2015: 131-142
- [27] Fontanarava J, Pasi G, Viviani M. An ensemble method for the credibility assessment of user-generated content [C] // Proc of the Int Conf on WI'17. New York; ACM, 2017: 863-868
- [28] Xie Bailin, Jiang Shengyi, Zhou Yongmei, et al. Misinformation detection based on gatekeepers' behaviors in microblog [J]. Chinese Journal of Computers, 2016, 39(4): 730-744 (in Chinese)  
(谢柏林, 蒋盛益, 周咏梅, 等. 基于把关人行为的微博虚假信息及早检测方法[J]. 计算机学报, 2016, 39(4): 730-744)
- [29] Ren Yafeng, Ji Donghong, Zhang Hongbin, et al. Deceptive reviews detection based on positive and unlabeled learning [J]. Journal of Computer Research and Development, 2015, 52(3): 639-648 (in Chinese)  
(任亚峰, 姬东鸿, 张红斌, 等. 基于 PU 学习算法的虚假评论识别研究[J]. 计算机研究与发展, 2015, 52(3): 639-648)
- [30] Gupta M, Sun Yizhou, Han Jiawei. Trust analysis with clustering [C] //Proc of the 20th Int Conf Companion on World Wide Web. New York; ACM, 2011: 53-54
- [31] Li Yaliang, Gao Jing, Meng Chuishi, et al. A survey on truth discovery [J]. ACM SIGKDD Explorations Newsletter, 2016, 17(2): 1-16
- [32] Ng A Y, Jordan M I. On discriminative vs. generative classifiers: A comparison of logistic regression and naive Bayes [C] //Proc of the 14th Int Conf on Neural Information Processing Systems: Natural and Synthetic. Cambridge, MA; MIT Press, 2001: 841-848
- [33] 2016 China Internet news market research report [EB/OL]. [2019-02-28]. <http://www.cnnic.net.cn/hlwfzyj/hlwxbzg/mtbg/201701/P020170112309068736023.pdf> (in Chinese)  
(2016 年中国互联网新闻市场研究报告 [EB/OL]. [2019-02-28]. <http://www.cnnic.net.cn/hlwfzyj/hlwxbzg/mtbg/201701/P020170112309068736023.pdf>)

[34] Zhang Jing, Liu Biao, Tang Jie, et al. Social influence locality for modeling retweeting behaviors [C] //Proc of the 33rd Int Joint Conf on Artificial Intelligence. Menlo Park, CA: AAAI, 2013: 2761-2767

[35] Teh Y W, Jordan M I, Beal M J, et al. Hierarchical Dirichlet processes [J]. Journal of the American Statistical Association, 2006, 101(476): 1566-1581

[36] Blei D M, Ng A Y, Jordan M I. Latent Dirichlet allocation [J]. Journal of Machine Learning Research, 2003, 3(Jan): 993-1022

[37] Yin Xiaoxin, Han Jiawei, Philip S Y. Truth discovery with multiple conflicting information providers on the Web [J]. IEEE Transactions on Knowledge and Data Engineering, 2008, 20(6): 796-808

[38] Li Qi, Li Yaliang, Gao Jing, et al. A confidence-aware approach for truth discovery on long-tail data [J]. Proceedings of the VLDB Endowment, 2014, 8(4): 425-436



**Liu Bo**, born in 1975. PhD, associate professor in the Southeast University. Her main research interests include social network, social big data.



**Li Yang**, born in 1995. Master candidate in the Southeast University. His main research interests include information credibility analysis.



**Meng Qing**, born in 1990. PhD candidate in the Southeast University. His main research interests include social influence, user behavior modeling, and machine learning.



**Tang Xiaohu**, born in 1990. Master from the Southeast University. His main research interests include big data analysis, and information credibility analysis.



**Cao Jiuxin**, born in 1967. PhD, professor in the Southeast University. His main research interests include computer network, social network, big data, privacy protection, and cloud resource scheduling.