

在线社交网络中的多主题谣言溯源

戴树兴, 夏正友

(南京航空航天大学 计算机科学与技术学院, 江苏 南京 211106)

摘要:随着通信技术的快速发展,用户之间的信息可以很快地流通,同时也导致谣言在社交网络中传播,因此亟需对谣言来源进行检测以确保社交网络的公信力。目前关于谣言溯源的研究方向基本注重于单主题谣言传播,然而社交网络中存在大量不同主题的谣言,谣言源头以及谣言主题数量越多,产生的不良影响越大。针对多主题谣言同时存在的情况,信息的传播过程需要被重新定义。因此,该文提出了一种多主题独立级联模型,并在该模型的基础上定义了谣言溯源问题。从已感染的网络子图中,基于影响力最大化的原则找出前 k 个可疑节点,这组节点被认为是最可能的谣言来源。并证明了该问题是NP难的,以及目标函数是单调且子模的。在此基础上,提出了一种基于影响力最大化的近似比为 $(1 - 1/e)$ 的贪婪算法。在大型真实数据集上的实验表明,平均误差距离控制在1跳之内。而且与其他算法相比,该算法具有更高的准确性以及有效性。

关键词:多主题;社交网络;谣言溯源;谣言来源;独立级联

中图分类号:TP39

文献标识码:A

文章编号:1673-629X(2024)01-0030-07

doi:10.3969/j.issn.1673-629X.2024.01.005

Tracing to Source of Multi-topic Rumors in Online Social Networks

DAI Shu-xing, XIA Zheng-you

(School of Computer Science and Technology, Nanjing University of Aeronautics and Astronautics, Nanjing 211106, China)

Abstract: With the rapid development of communication technology, information between users can flow quickly, which also leads to the spread of rumors in social networks, so there is an urgent need to detect the source of rumors to ensure the credibility of social networks. At present, the research on rumor traceability basically focuses on the spread of single-topic rumors. However, there are a large number of rumors with different topics in social networks. The more the source of rumors and the number of rumor topics, the greater the adverse effects. In view of the fact that multi-topic rumors exist at the same time, the process of information dissemination needs to be redefined. Therefore, a multi-topic independent cascade model is proposed, and the rumor traceability problem is defined on the basis of this model. From the infected network subgraph, the first k suspicious nodes are identified based on the principle of maximizing influence, and this group of nodes is considered to be the most likely source of rumors. It is proved that the problem is NP-hard and the objective function is monotone and submodular. On this basis, a greedy algorithm based on influence maximization with approximate ratio $(1 - 1/e)$ is proposed. Experiments on large real data sets show that the average error distance is controlled within 1 hop. And compared with other algorithms, the proposed algorithm has higher accuracy and effectiveness.

Key words: multi-topic; social network; rumor tracing; rumor source; independent cascade

0 引言

现今人们的生活离不开网络,网络给用户提供了信息交流的平台。随着手机、电脑等通信设备的快速普及,社交平台也逐渐开始兴起,例如 Twitter、Facebook、微博。用户之间通过社交平台进行信息分享。同时,信息的快速传播也会产生不同的影响,

例如用户可以通过网络获取天气预报、股票市场情况变化等信息,但同样也会受到虚假信息的影响,由于虚假信息具有新颖性,抓住人们的猎奇心理,从而导致虚假信息比真实信息具有更强的传染性。此外,对各种各样信息的真实性进行检测是困难的,这导致了谣言在社交网络中传播的问题,恶意用户发散各种不实信

收稿日期:2023-03-02

修回日期:2023-07-05

基金项目:国家重点研发计划“宽带通信和新型网络”重点专项中项目(2018YFB1800600, 2018YFB1800602)

作者简介:戴树兴(1999-),男,研究生,研究方向为计算机网络和复杂网络;通信作者:夏正友(1973-),男,博士,副教授,CCF会员(54073M),研究方向为计算机网络和信息安全。

息,可能会影响社会的稳定,产生严重的后果^[1]。

为了确保社交平台的公信力,以及减少错误信息在社交网络中的影响,追踪识别散播谣言的源头是非常重要的步骤,通过了解这些谣言源头有助于平台设计有效的策略遏制谣言的散播。而且感染源头检测技术在其他领域有着很多成功的应用,例如找出污水网络中的病毒、流行传染病的最初感染者^[2]。

目前谣言溯源的工作集中于对单一主题的谣言进行源头检测,且大部分为单源检测,这些工作并不全面。考虑到一个更为现实的场景,在同一社交网络中,会有多个不同主题的谣言从多个源头同时传播,每个用户可以同时接收到这些不同主题的谣言。因此,该文将单主题谣言溯源的工作拓展到多主题谣言溯源。在这种情况下如何高效地对谣言源头进行追溯是一个很有挑战性的课题。经过一段时间后,在只知道底层的社会网络结构以及时间 t 的感染子图情况下,如何确定谣言的来源。要解决这个问题,需要解决几个关键挑战。首先,在多主题多源谣言传播的场景下,信息是在网络中如何扩散的;其次,如何确定每个感染节点为源头的可能性;最后,对谣言源头进行识别的方法是否有近似保证,是否会产生较大的偏差。

在单主题谣言传播中,独立级联模型^[3]被广泛应用于各项研究中。而多主题谣言传播已经有研究证明其影响是异质的^[4]。因此,需要重新证明激活节点目标函数的子模性质,子模能够保证解的良好逼近。该文提出了一种谣言传播的多主题独立级联模型,基于该模型定义了多主题谣言溯源问题,并证明了该问题是 NP 难的,以及目标函数具有单调、子模的性质。从目标函数的性质出发,提出了一种基于影响力最大化的贪婪算法来进行谣言源头识别,该算法能保证 $(1 - 1/e)$ 的比例逼近最优解。

1 相关工作

Shah 和 Zaman 首次提出了谣言溯源问题^[5],他们假设单一感染源在树状网络下满足 SI 模型传播,并提出谣言中心性的概念估计谣言来源。Cai 等人^[6]研究了在一般传播时间分布条件下,SI 模型对单一源多次独立网络快照的检测概率。

Xu 和 Chen 提出一种新的源检测方法^[7],设置一些监视器节点来获得谣言传播的速度,以此提出一种多项式算法计算节点的可达性并进行重要性排序。谣言源头检测概率取决于监视器数量。

文献[8]中作者在经典传染病模型的基础上进行改进,加入了新的“辟谣者”状态,并基于贪婪算法识别源头。文献[9]提出一种 SIOR 传染病模型,并研究了在该模型下的谣言溯源问题。文献[10]提出一种

SEIR 模型,研究了基于网络观测快照下的单一来源检测问题。

Choi 等人^[11]提出了基于查询的方法,首先向节点进行一次简单的查询并根据回答生成网络。然后使用交互式查询,询问节点从谁接收到谣言。该方法保证了在规则树网络中的检测概率。文献[12]中,作者在基于监视的观察下,通过监视器节点发送“反谣言”,并用最大后验估计器来检测谣言来源。

在多源检测工作中,Wang^[13]首次将谣言中心度拓展到了多源检测。Dong 等人^[14]提出了一种基于深度学习的谣言溯源模型,在缺乏底层网络信息传播模型的先验知识下,依然能检测到多个谣言来源。

Nguyen 提出了基于排名和优化的方法将感染节点按可疑性排序,找出前 k 个可疑节点^[15]。李城等人^[16]基于最长公共子序列改进了 LCS 算法。

叶增炜等^[17]提出了一种基于有责量和免责量的谣言溯源方法。廖艺等人^[18]基于谱优化社区划分算法将感染子图划分为两个社区后寻找谣言来源。

2 谣言来源检测和问题描述

在这一节内容中,首先将介绍谣言的基本传播模型,然后给出了改进后的多主题谣言传播模型,接着描述了相关问题的描述以及证明。该文将独立级联模型拓展到多主题独立级联模型。

2.1 独立级联模型

在独立级联模型中,社交网络被视作为一个有向加权图 $G = (V, E)$,每个节点 v 代表不同的用户,每条有向边 $e = (u, v) \in E$ 代表用户 u 和用户 v 之间的关系,每条边会被分配一个权重 $p(u, v) \in [0, 1]$,代表用户 u 对用户 v 的影响程度。

在谣言传播过程中,每个节点只有两种状态:活跃和非活跃。在时间步 t 时,当节点 u 被激活为活跃状态时,该节点会依次向其每个邻居 v 以概率 $p(u, v)$ 进行谣言传播,如果激活成功,邻居 $t + 1$ 在第 u 时刻转化为活跃状态。在后续传播过程中,节点 u 将不再尝试激活其邻居。当没有新的节点被激活时,谣言传播过程结束。

2.2 多主题独立级联模型

独立级联模型研究单个主题的信息传播过程。但是考虑到多个不同主题的谣言可以同一时间传播,同一用户可以在短时间内同时接收到不同主题的谣言。不同主题的谣言不仅内容不同,传染力也可能不同。娱乐八卦类的谣言比农业、军事等类型的谣言传播更广,影响人数更多。而且同一用户传播不同类型的谣言时,对邻居产生的影响也可能不同。

因此,在多主题谣言传播的情况下,需要重新对信

息扩散过程进行建模。而独立级联模型并不能处理这种情况,因为它很难体现不同主题之间的复杂相关性。已经有相关研究证明,当采用多主题信息级联时,计算激活节点的目标函数不再是子模的^[19-20]。

在多主题独立级联模型中,在线社交网络用 $G = (V, E, P)$ 表示,其中 V 为节点集, E 为边集,且 $|V| = n$, $|E| = m$ 。每个节点 v 代表不同的用户,每条有向边 $e = (u, v) \in E$ 代表用户 u 和用户 v 之间的关系,每条边会被分配一个初始影响权重 $p_0(u, v) \in [0, 1]$ 代表用户 u 对用户 v 的影响程度。用 $N_i(v)$ 表示节点 v 的传入邻居节点集, $N_o(v)$ 表示传出邻居节点集。

假设在社交网络中有 q 个主题的谣言,例如娱乐、体育、农业等。被不同主题谣言激活的节点集为 $S = \{S_1, S_2, \dots, S_q\}$ 。其中 S_i 表示被第 i 个主题的谣言所激活的节点集合,可以知道节点感染的谣言主题类型,而且可以被多个主题谣言同时感染。在 q 个主题上谣言传播的感染节点集 $S = \bigcup_{i=1}^q S_i$ 。

对每个节点 $v \in V$,由于同一节点可以被不同主题谣言多次感染,因此节点 v 在时间步 t 有 $q+1$ 个状态,即 $Q_v^t = \{\text{非活跃}, \text{主题 1 激活}, \text{主题 2 激活}, \dots, \text{主题 } q \text{ 激活}\}$ 。每个向量分量用 0 和 1 表示,当且仅当节点 v 没有被任何谣言感染时,非活跃分量的值为 1,其他分量的值为 0。当被第 i 个主题的谣言所激活时,主题 i 激活分量的值为 1。

在现实中,节点 u 对节点 v 的影响可能还会受到谣言主题的影响,因此在给每个节点 v 分配一个影响系数向量 $\delta_v = \{\delta_v^1, \delta_v^2, \dots, \delta_v^q\}$,其中 $\delta_v^i \in [0, 1]$,表示主题 i 对其邻居的影响。之前已给节点 u 对节点 v 分配初始感染概率 $p_0(u, v)$,故新的感染概率 $p(u, v) = \{p_v^1, p_v^2, \dots, p_v^q\}$,其中 $p_v^i = \delta_v^i \cdot p_0(u, v)$ 。

令 $I_0 \subseteq V$ 为初始恶意节点集合,其初始状态为 Q_v^0 。谣言在多主题独立级联模型下随着离散时间步 t 在网络中传播,当节点 u 为活跃状态时,该节点会依次向其每个邻居 v 以概率 $p(u, v)$ 进行谣言传播,无论邻居 v 是否被激活成功,每个主题的谣言只传播一次。如果节点 u 被其他新主题 i 的谣言激活,那么节点 u 会向每个邻居 v 传播主题 i 的谣言。当没有新的节点被任意主题的谣言感染时,谣言传播过程结束。

2.3 问题描述

该文的目标是在已知底层网络结构以及第 t 个时间步感染子图的情况下找出一组节点,而这组节点被认为是最可能的谣言来源。现实中谣言源头想通过影响尽可能多的用户以达到某种目的,因此在直觉中如果能找出影响最大的一组节点,那么这组节点中为谣言来源的可能性非常大。利用这一特点,基于影响力

最大化的原理将活跃集合 S 中每个节点根据可疑性程度进行排序,其中前 k 个节点被视作最可疑的节点。

假设 $\sigma_i^t(G, A)$ 为初始恶意节点 A 的情况下,在第 t 时间步所有被主题 i 感染的节点集合, $\sigma_i(G, A)$ 则表示为最终所有被主题 i 感染的节点集合。 $\sigma(G, A)$ 表示最终所有主题下的感染人数: $\sigma(G, A) = \sum_{i=1}^q \sigma_i(G, A)$ 接下来给出 k -可疑节点问题的定义。

定义 1 (k -可疑节点):在线社交网络用有向加权图 $G = (V, E, P)$ 表示以及给出带有 q 个主题的感染节点集 $S = \{S_1, S_2, \dots, S_q\}$,利用多主题独立级联模型来模拟信息扩散过程。目标是在已知感染子图 $G' = (V', E', P')$ 时找出一组不超过 k 个可疑节点集 $A \subseteq S$,使得激活节点数均值 $\varphi(G', A) = E(|\sigma(G', A) \cap S|)$ 最大。

一些关于影响最大化问题是 NP 难的证明可以在文献[3, 20-21]中找到。接下来给出在多主题独立级联模型下的 k -可疑节点问题是 NP 难的证明。

定理 1:基于多主题独立级联模型的 k -可疑节点问题是 NP 难的。

证明:为了证明 k -可疑节点问题是 NP 难的,用背包问题归约到 k -可疑节点问题。而背包问题是公认被证明的 NP 完全问题。

背包问题:假设有 n 个物品,物品 i 的重量为 w_i ,价值为 c_i 。背包的容量为 W 。其中重量和价值的值都是非负的。问题的目标是怎么样找出一组物品,使其价值最大。即找出向量 $X = \{x_1, x_2, \dots, x_n\}$,满足 $\sum_i x_i w_i \leq W$,使 $\sum_i x_i c_i \geq C$ 。

令 $\pi_1 = (X, W)$ 是背包问题的一个实例,令 $\pi_2 = (G', S, k)$ 是 k -可疑节点问题的一个实例。其中 S 是谣言的源头节点, k 是确定的前 k 个可疑节点。对于接下来构造从 π_1 到 π_2 的一个归约,如图 1 所示。

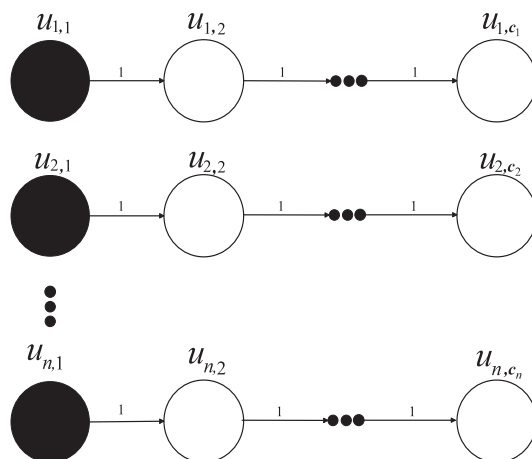


图 1 k -可疑节点问题到背包问题的归约
归约:为了构造归约,首先给出一个感染子图 $G' =$

(V, E', P') 使其满足以下条件: 其中存在谣言源头集 S , 对每个物品的价值 c_i 构造一条含有 $c_i + 1$ 个节点的简单路径: $S \rightarrow u_{i,1} \rightarrow u_{i,2} \rightarrow \dots \rightarrow u_{i,c_i}$, 并设路径的每条有向边权重为 1。对任意的 $u_{i,1}$ 都有 $w_i = 1$, 令 $k = W$ 和 $M = C$ 。 M 是一个正整数。接下来证明 π_1 有解 $X = \{x_1, x_2, \dots, x_n\}$ 当且仅当 π_2 也有对应的解 $A = \{u_{i,1} \mid i = 1\}$ 使 $\varphi(G, A) > M$, 反之亦然。

(\rightarrow) 假设 $X = \{x_1, x_2, \dots, x_n\}$ 是 π_1 的解, 对于 π_2 的解选择这样的 $A = \{u_{i,1} \mid x_i = 1\}$ 且 $|A| \leq k$, 有

$\sum_{i \mid x_i = 1} x_i w_i = \sum_i w_i \leq W = k$ 。根据多主体独立级联模型可知, 当选择节点 $u_{i,1}$ 时, 其所在路径后面的子节点 $u_{i,1}, u_{i,2}, \dots, u_{i,c_i}$ 都会被感染。因此有 $\varphi(G, A) = \sum_{i \mid x_i = 1} x_i c_i = \sum_i c_i \geq C = M$, 所以 A 是 π_2 的解。

(\leftarrow) 假设 A 是 π_2 的解, 那么 A 中一定不包括节点 $u_{i,j>1}$ 。因为若 A 中存在节点 $u_{i,j>1}$, 那么有集合 A' 包含节点 $u_{i,j'<j}$ (即节点 $u_{i,j}$ 的祖先), 使 $\varphi(G, A') > \varphi(G, A)$, 这与 A 是 π_2 的解矛盾。对于 π_1 的解 $X = \{x_1, x_2, \dots, x_n\}$, 如果 $u_{i,1} \in A$, 则令 $x_i = 1$, 否则令 $x_i = 0$ 。因此有 $\sum_i x_i w_i = \sum_{i \mid u_{i,1} \in A} w_i \leq W = k$ 且 $\varphi(G, A) = \sum_{i \mid u_{i,1} \in A} c_i = \sum_i x_i c_i \geq M = C$ 。

综上所述, 背包问题的解是 k -可疑节点问题的解, 而 k -可疑节点问题的解也是背包问题的解。因此 k -可疑节点问题是 NP 难的。

接下来需要证明目标函数是单调且子模的。

定理 2: 多主题独立级联模型下的目标函数 $\varphi(\cdot)$ 是单调递增的子模函数。

证明: 如果函数满足“边际效益递减”规律, 即对所有元素 u 和集合 $S \subseteq T$ 有 $f(S \cup \{u\}) - f(S) \geq f(T \cup \{u\}) - f(T)$, 则称函数 f 是子模函数。对于函数 $\varphi(G, A)$, 有集合 $A = \sum_{i=1}^q A_i$, $G = \sum_{i=1}^q G_i$, 其中 A_i 表示被主题 i 谣言感染的集合, G_i 表示只含有主题 i 谣言的感染子图。因此有 $\varphi(G, A) = \sum_{i=1}^q \varphi(G_i, A_i)$ 。在文献 [6] 中已经证明可以分解 $\varphi(G_i, A_i) = \sum_g \text{prob}(g) \varphi_g(G_i, A_i)$, 其中 g 和 prob 分别表示随机产生的活跃节点子图和产生该子图的概率, 且文献 [15] 证明了 $\varphi(G_i, A_i)$ 是单调递增的子模函数。由 $\varphi(G, A) = \sum_{i=1}^q \varphi(G_i, A_i)$ 可知, 子模函数的非负线性组合依然还是子模函数, 因此目标函数 $\varphi(G, A)$ 是单调递增的子模函数。

2.4 基于影响力最大化的贪婪算法

前面已经证明多主题独立级联模型下的 k -可疑节点问题是 NP 难的, 因此想要求最优解是需要花费指数增长的时间, 成本太高。因此, 该文提出了一种求解该问题的贪婪算法。由于子模函数的性质, 该算法能以 $(1 - 1/e)$ 的比例逼近最优解, 能够在一定程度内保证解的质量和精度。

在已知某一时刻的感染子图情况下, 对任意节点 u , 如果能感染尽可能多 S 中的节点, 那么该节点为谣言源头的可能性越大。以此计算并排序将所有节点的可疑程度, 最后输出前 k 个最可疑的节点。算法过程如下:

算法 1: 贪婪算法

输入: 感染子图 $G = (V, E', P')$, 以及感染节点集 S , 模拟次数 R , 谣言主题数 q

输出: k 个最可疑节点

```

 $I \leftarrow \varnothing$ 
for  $i = 1$  to  $q$ 
  for  $j = 1$  to  $k$ 
    for  $u \in S_i \setminus I_i$ 
       $\varphi_i^u \leftarrow 0$ 
      for  $m = 1$  to  $R$ 
         $\varphi_i^u = \varphi_i^u + |\sigma_i(G, (I_i \cup \{u\})) \cap S_i|$ 
      end
       $\varphi_i^u \leftarrow \varphi_i^u / R$ 
    end
     $I_i \leftarrow I_i \cup \arg\max_{u \in S_i \setminus I_i} \{\varphi_i^u\}$ 
  end
end
Return  $I$ 

```

为了尽可能提高算法的准确性, 使用蒙特卡罗方法对多主题谣言传播过程进行 R 次模拟, 对节点可疑程度取平均值并取前 k 个节点。在时间复杂度方面, 假设模拟一次多主题谣言传播过程的时间需要 $O(m)$, 谣言主题数量为 q , 那么总花费时间为 $O(Rmqk | S |)$ 。

3 实验

该文分别在三个不同的真实在线社交网络上进行了谣言溯源的实验。结果表明, 与其他的算法相比, 贪婪算法有着更高的准确性。

3.1 实验设置

实验分别在 Slashdot, Epinions 以及 gemsec 三个真实数据集上进行。数据集从 [http://snap.stanford.edu/data/] 获得。此外, 该文对数据集做了相关处理, 去除了所有自环边。关于它们的具体数据可以从该网站中找到, 表 1 给出了关于这些数据集的相关信息。

实验中, 将每条边 (u, v) 的初始感染概率设为

$p(u, v) = 1 / |N_o(v)|$ 。对每个主题,每个节点还会被分配一个影响系数向量 $\delta_v = \{\delta_v^1, \delta_v^2, \dots, \delta_v^q\}$, δ_v 的每个分量在范围 $(0, 1]$ 内随机取值。为了减少误差,用蒙特卡罗方法对多主题级联过程重复 1 000 次模拟,最后取平均检测概率。实验分别展示了谣言主题数量 $q = 2$ 和 $q = 3$ 情况下的结果。每个主题谣言的初始恶意节点数量为 $|S_i| = 10$ 。

表 1 数据集

数据集	节点	边	类型
Slashdot	77 360	828 161	有向图
Epinions	75 888	508 837	有向图
gemsec	47 528	222 887	有向图

实验将贪婪算法与最大度算法和随机算法进行了比较。结果表明,在三个不同的真实数据集上以及在不同谣言主题数量的情况下,贪婪算法的表现都要优于最大度算法以及随机算法。

3.2 实验结果

表 2 和表 3 分别显示谣言主题数量 $q = 2$ 和 $q = 3$

表 2 三种算法在不同数据集上的性能对比(谣言主题数量 $q = 2$)

算法	Epinions			gemsec			Slashdot		
	检测概率 /%	误差距离 (1 跳之内) /%	平均误差 距离(跳)	检测概率 /%	误差距离 (1 跳之内) /%	平均误差 距离(跳)	检测概 率/%	误差距离 (1 跳之内) /%	平均误差 距离(跳)
贪婪算法	70	80	0.60	55	85	0.50	70	80	0.55
最大度算法	30	45	1.40	45	75	0.80	35	30	1.70
随机算法	20	65	1.20	40	45	1.30	20	40	1.60

表 3 三种算法在不同数据集上的性能对比(谣言主题数量 $q = 3$)

算法	Epinions			gemsec			Slashdot		
	检测概率 /%	误差距离 (1 跳之内) /%	平均误差 距离(跳)	检测概率 /%	误差距离 (1 跳之内) /%	平均误差 距离(跳)	检测概 率/%	误差距离 (1 跳之内) /%	平均误差 距离(跳)
贪婪算法	43.3	66.7	0.90	53.3	86.7	0.47	50	86.7	0.70
最大度算法	23.3	36.7	1.43	26.7	86.7	0.93	20	36.7	1.47
随机算法	20	53.3	1.47	43.3	45	0.93	16.7	43.3	1.60

当 $q = 3$ 时,由表 3 可以看出,在 Epinions, gemsec 和 Slashdot 数据集上,贪婪算法的检测概率能分别达到 43.3%, 53.3% 和 50%。最大度算法的检测概率分别达到 23.3%, 26.7% 和 20%。而随机算法则分别达到 20%, 43.3% 以及 16.7%。随着谣言主题数量增加,贪婪算法依然明显优于最大度算法和随机算法。

3.2.2 误差距离

图 2 显示了不同算法在不同真实数据集下的误差距离频率。其中(a1, b1, c1)为谣言主题数量 $q = 2$ 的情况,在 Epinions 数据集上,贪婪算法有 70% 的把握能够正确找到谣言来源,且有 80% 的把握保证检测的

的情况下,通过检测概率,误差距离以及平均误差距离用来评估算法的性能。检测概率是指检测谣言来源节点数量与真实源节点数量之比。误差距离是指检测谣言来源节点与真实源节点的最小距离。从表 2 和表 3 可以看出,在不同谣言主题数量的情况下,贪婪算法在检测概率、误差距离(1 跳之内)以及平均误差距离上的表现都要优于其他两种算法。

3.2.1 检测概率

随着 k 的增大,检测概率提高。理论上,如果 k 增大到与感染子图的节点数量一致的话,那么一定包含所有真实谣言来源,但这种做法并不现实。所以实验中将 k 的最大值设置为 10。

由表 2 可以看出,在谣言主题数量 $q = 2$ 的情况下,在 gemsec 数据集上,贪婪算法能达到 55% 的成功检测概率,而最大度算法和随机算法分别只有 45% 和 40%。在 Slashdot 和 Epinions 数据集上,贪婪算法的检测概率能达到 70%,最大度算法有 35% 和 30% 的检测概率,而随机算法则都只有 20% 的检测概率。

谣言节点离真实谣言来源节点的距离在 1 跳之内。而随机算法和最大度算法分别只有 65% 和 45% 的把握。在 gemsec 和 Slashdot 数据集上,贪婪算法则分别有 85% 和 80% 的把握正确找到真实谣言来源或者只差 1 跳的误差距离,随机算法则只有 45% 和 40%,最大度算法有 75% 和 30%。

在谣言主题数量 $q = 3$ 的情况下,即图 2(a2, b2, c2)中可以看出,贪婪算法的表现依然出色且稳定,在三种数据集上分别有 66.7%, 86.7%, 86.7% 的概率能够在 1 跳的距离之内找到真实谣言来源。而随机算法则分别达到 53.3%, 45%, 43.3%。最大度算法则

分别达到 36.7%、86.7%、36.7%。随着谣言主题数

量增加,贪婪算法优于其他两种算法。

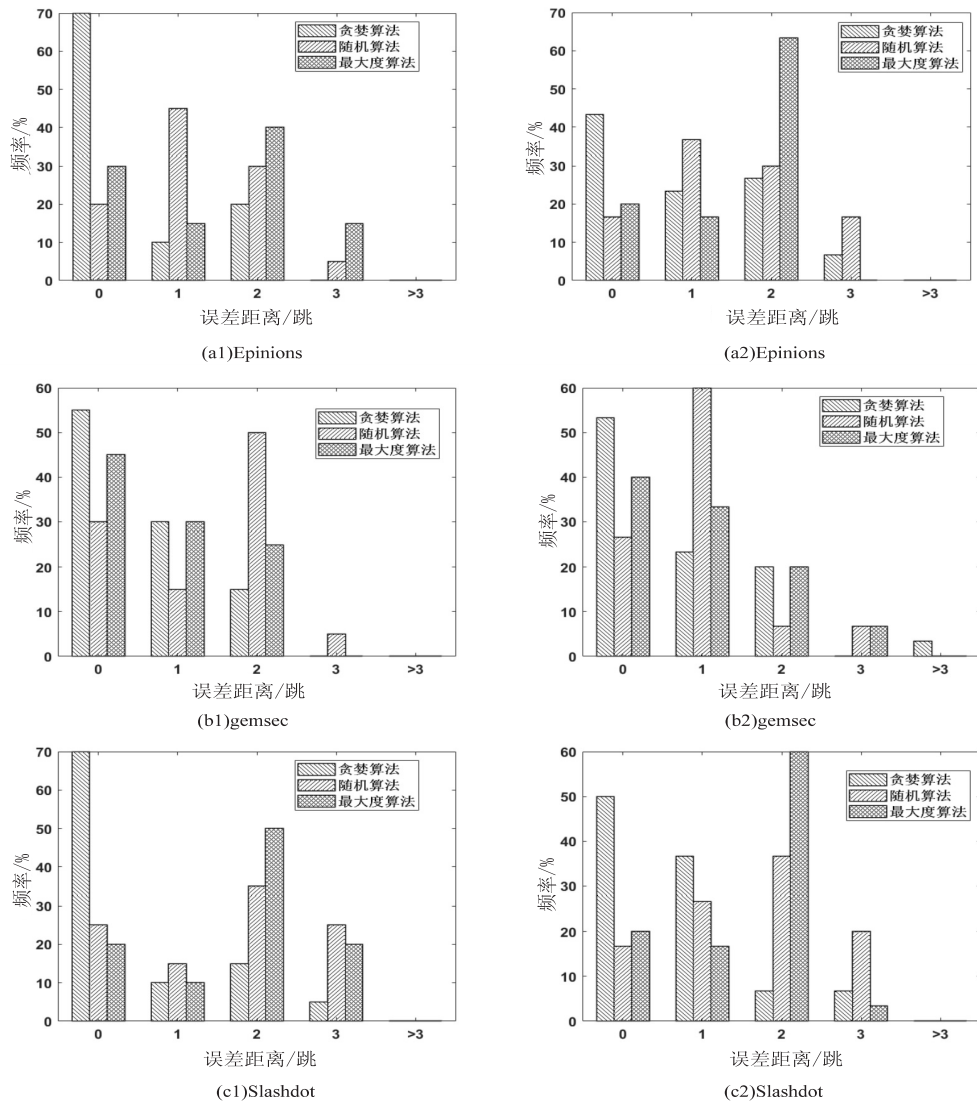


图2 不同数据集和算法下的误差距离

相比于其他两种算法,贪婪算法检测到的大部分谣言节点离真实谣言节点的距离在1跳以内,检测距离达到3跳或者超过3跳的节点不超过10%,表现稳

定且高效。随机算法与最大度算法所检测的谣言节点与真实谣言节点的距离随机分布在0~3跳之间,误差较大。

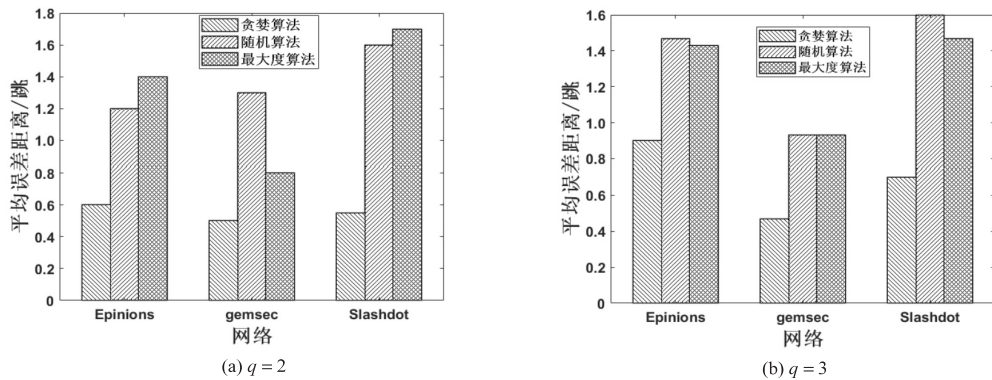


图3 不同算法在不同数据集下的平均误差距离

图3显示了不同算法在不同数据集下的平均误差距离,进行1000次蒙特卡罗模拟后取平均值。当 $q=2$ 时,在三种数据集下的平均误差距离分别为0.6跳、

0.5跳、0.55跳。而最大度算法分别为1.40跳、0.80跳、1.70跳。随机算法则分别为1.20跳、1.30跳、1.60跳。当 $q=3$ 时,贪婪算法平均误差距离分别为

0.90 跳、0.47 跳、0.70 跳。最大度算法分别为 1.43 跳、0.93 跳、1.47 跳。随机算法分别为 1.47 跳、0.93 跳、1.60 跳。可以看出贪婪算法明显优于其他算法。

4 结束语

考虑到每个节点可能在同一时间收到多种类型的谣言并传播给邻居的情况,在传统独立级模型的基础上进行了扩展,提出一种多主题谣言传播的独立级联模型。在该模型下,每个节点可以被不同主题的谣言多次感染,因此需要重新定义影响力最大化的溯源问题以及相应的目标函数。

在该模型的基础上,提出了考虑影响力的 k -可疑节点谣言溯源问题,即找出前 k 个影响力最大的节点,这些节点被认为是最有可能的谣言来源。并证明了多主题独立级联模型下的 k -可疑节点谣言溯源问题是 NP 难的,以及对应的目标函数 $\varphi(\cdot)$ 是单调递增的子模函数。

提出了一种近似比为 $(1 - 1/e)$ 的贪婪算法。实验结果表明,在不同大型网络上,贪婪算法可以有效地识别出不同谣言主题的来源。

参考文献:

- [1] 王超,倪静.融入心理因素的在线社交网络谣言传播模型[J].计算机技术与发展,2022,32(11):190-197.
- [2] 黄春林,刘兴武,邓明华,等.复杂网络上疾病传播溯源算法综述[J].计算机学报,2018,41(6):1376-1399.
- [3] KEMPE D, KLEINBERG J, TARDOS É. Maximizing the spread of influence through a social network[C]//Proceedings of the ninth ACM SIGKDD international conference on knowledge discovery and data mining. Washington: ACM, 2003:137-146.
- [4] FAN J, QIU J, LI Y, et al. OCTOPUS: an online topic-aware influence analysis system for social networks[C]//2018 IEEE 34th international conference on data engineering (ICDE). Paris: IEEE, 2018:1569-1572.
- [5] SHAH D, ZAMAN T. Rumors in a network: who's the culprit? [J]. IEEE Transactions on Information Theory, 2011, 57(8):5163-5181.
- [6] CAI K, XIE H, LUI J C S. Information spreading forensics via sequential dependent snapshots[J]. IEEE/ACM Transactions on Networking, 2018, 26(1):478-491.
- [7] XU W, CHEN H. Scalable rumor source detection under independent cascade model in online social networks[C]//2015 11th international conference on mobile ad-hoc and sensor networks (MSN). Shenzhen: IEEE, 2015:236-242.
- [8] 陈一新,陈馨悦,刘奕,等.基于SIDR模型的谣言传播与源头检测研究[J].数据分析与知识发现,2021,5(1):78-89.
- [9] 吴杨,吴国文,张红,等.基于扩展传染病模型的谣言溯源[J].计算机与现代化,2022(1):113-119.
- [10] ZHOU Y, WU C, ZHU Q, et al. Rumor source detection in networks based on the SEIR model[J]. IEEE Access, 2019, 7:45240-45258.
- [11] CHOI J, MOON S, WOO J, et al. Rumor source detection under querying with untruthful answers[C]//IEEE INFOCOM 2017-IEEE conference on computer communications. Atlanta: IEEE, 2017:1-9.
- [12] CHOI J, MOON S, SHIN J, et al. Estimating the rumor source with anti-rumor in social networks[C]//2016 IEEE 24th international conference on network protocols (ICNP). Singapore: IEEE, 2016:1-6.
- [13] WANG Z, DONG W, ZHANG W, et al. Rooting our rumor sources in online social networks: the value of diversity from multiple observations[J]. IEEE Journal of Selected Topics in Signal Processing, 2015, 9(4):663-677.
- [14] DONG M, ZHENG B, QUOC VIET HUNG N, et al. Multiple rumor source detection with graph convolutional networks[C]//Proceedings of the 28th ACM international conference on information and knowledge management. Beijing: ACM, 2019:569-578.
- [15] NGUYEN D T, NGUYEN N P, THAI M T. Sources of misinformation in online social networks: who to suspect? [C]//MILCOM 2012-2012 IEEE military communications conference. Orlando: IEEE, 2012:1-6.
- [16] 李城,沙俊淞,武文.基于最长公共子序列的微博谣言溯源研究[J].计算机与现代化,2018(1):107-112.
- [17] 叶增炜,王友国,柴允.基于有责量和免责量的谣言溯源算法[J].计算机技术与发展,2022,32(1):40-46.
- [18] 廖艺,王友国,朱亮.基于谱优化社区划分的双信源溯源算法[J].计算机技术与发展,2020,30(12):72-76.
- [19] TONG G, WU W, DU D Z. Distributed rumor blocking with multiple positive cascades[J]. IEEE Transactions on Computational Social Systems, 2017, 5(99):468-480.
- [20] PHAM D V, NGUYEN G L, NGUYEN T N, et al. Multi-topic misinformation blocking with budget constraint on online social networks[J]. IEEE Access, 2020, 8:78879-78889.
- [21] PHAM CANH V, THAI MY T, DUONG HIEU V, et al. Maximizing misinformation restriction within time and budget constraints[J]. Journal of Combinatorial Optimization, 2018, 35(4):1202-1240.