

基于标签传播的社区挖掘算法研究综述

王 庚 宋传超 盛玉晓 王童童 李盛恩

(山东建筑大学 计算机科学与技术学院 山东 济南 250101)

摘 要: 社会网络由于其流行程度已经成为众多学者的研究热点。通过社区挖掘算法可以发现存在于社会网络中的潜在社区,而重叠社区挖掘则可以挖掘出更具有现实意义的社区结构。但是在研究中社会网络所包含的庞大数据量又会为之带来种种不便,因此快速的社区挖掘算法就受到了越来越多的重视。基于标签传播的社区挖掘算法具有近乎线性的时间复杂度。文中将从多方面研究目前基于标签传播的社区挖掘算法的优劣,并且详细分析基于标签传播算法在以后研究中的改进思路。

关键词: 社会网络; 标签传播; 社区挖掘; 重叠社区

中图分类号: TP311.13

文献标识码: A

文章编号: 1673-629X(2013)12-0069-05

doi: 10.3969/j.issn.1673-629X.2013.12.017

Research Summary on Communities Mining Algorithm Based on Label Propagation

WANG Geng SONG Chuan-chao SHENG Yu-xiao WANG Tong-tong LI Sheng-en

(College of Computer Science and Technology, Shandong Jianzhu University, Jinan 250101, China)

Abstract: Social networks have been a hot area of research because of its popularity. Discover potential communities in social networks through community mining and find community structures that have more realistic significance through detecting overlapping communities. However there is lot of inconvenience because of the sheer amount of data in social networks. So fast algorithm for mining community are getting more and more attention. The algorithms based on the thoughts of label propagation have nearly linear time complexity. In this paper study the algorithms based on the thoughts of label propagation from various aspects and analyze those algorithms' improvement ideas in the future research.

Key words: social networks; label propagation; community mining; overlapping community

0 引言

近些年随着 Web2.0 的兴起,社会网络也以各种各样的形式出现在了人们的视野中,它的迅速流行也使得很多研究者开始了对社会网络的研究。在这些研究中,社会网络通常被抽象成一个图的形式,图中用顶点表示用户,边用来表示用户之间存在的关系。

在社会网络的用户之间还存在一些潜在的社区结构,社区结构是由一组相似的顶点互相连接而成的,同一社区内部之间连接稠密,不同社区之间连接较为稀疏。比如在社会网络中共同喜好足球这项运动的用户就可以划分为一个社区。通过一些社区挖掘算法就能得到这些社区结构。比较经典的社区挖掘算法有 GN 算法^[1]、谱分析思想的算法^[2]、层次距离算法^[3]、边集

聚系数法^[4]等。另外 Fortunato 在文献[5]中也对各种社区挖掘算法进行了更为详细的介绍研究。由于当前社会网络中海量的数据信息,因此要求社区挖掘算法应具有尽量低的时间复杂度。在众多的社区挖掘算法中,由 Raghava 等人^[6]于 2007 年率先提出的基于标签传播(LPA)思想的 RAK 算法是目前所知的效率最高的社区挖掘算法,它具有的线性的时间复杂度极大地提高了社区挖掘的效率。在此之后,基于标签传播的社区挖掘算法也一直被研究者所关注和改进,文中将对此做进一步的探讨。

在真实的社会网络中由于事物往往存在多样性的特点,所以一个节点往往不仅仅从属于一个社区,而是同时可以划分到多个社区中,比如一个用户喜欢足球

收稿日期: 2013-03-04

修回日期: 2013-06-15

网络出版时间: 2013-09-29

基金项目: 国家自然科学基金资助项目(61170052)

作者简介: 王 庚(1987-),男,硕士研究生,研究方向为社会网络、数据库;李盛恩,教授,研究方向为数据库、数据挖掘。

网络出版地址: <http://www.cnki.net/kcms/detail/61.1450.TP.20130929.1541.033.html>

这项运动的同时也可能有音乐、电影等其他方面的兴趣爱好,因此一个用户可能同时带有多个属性标签。由此可以看到社会网络中存在重叠社区结构。但是之前所说的社区挖掘算法所挖掘得到的社区结构中每个节点只可能从属于某一个社区,这不符合现实社会网络中的实际情况。进而又有很多研究者开始了对重叠社区挖掘算法的研究,目前的重叠社区挖掘算法主要有派系过滤算法(CPM)^[7]、EAGLE 算法^[8]、GCE 算法^[9]、LFK 算法^[10]等。由 Steve^[11]在 RAK 算法的基础上提出的 COPRA 算法是一种基于标签传播的重叠社区挖掘算法。根据大量的研究工作,认为基于标签传播的社区挖掘算法中,最核心的两个算法便是 RAK 算法和 COPRA 算法,因此将根据这两个算法来展开对基于标签传播社区挖掘算法的总结和研究。

1 非重叠 RAK 算法及其改进算法研究

1.1 标签传播思想及 RAK 算法

标签传播算法(LPA)最早是由 Zhu 等人^[12]于 2002 年提出的,它是基于图的一种半监督学习方法,其算法基本思想是用已标记了节点的标签信息去预测其他还未标记节点的标签信息。首先利用节点样本间的关系来建立关系完全图模型,在初始化的完全图中,节点分为已标注标签的节点和未标注标签的节点,通过节点之间的边来表示两个节点的相似度。然后已标注标签的节点按相似度来传递标签给其他节点。

由 LPA 算法基本理论,可以看出每个节点的标签是按照节点间的相似度来传播标签信息的,在标签传播的每一步,未标记节点根据相邻节点的标签信息来更新自己的标签,同该节点相似度越大,其邻接点对其标注的影响权值也就越大,当相似节点的标签越趋于一致时,其标签就越容易传播。

RAK 算法^[6]是基于标签传播算法的思想来快速地发现社区结构。在该算法中,首先需要为每个节点分配一个唯一的标签,由于每个节点不是单独存在的,都会有自己的相邻节点,在迭代过程中,节点就会根据它相邻节点标签的情况不断地迭代更新自己的标签信息,直到所有节点的标签不再发生任何变化,最后再根据每个节点的标签把它划分到相应的社区中。上述是算法的基本思想,接下来将该算法描述如下:

(1) 对于图中任意的节点 $x, x \in G(x)$, $G(x)$ 表示的是网络中所有的节点集合,都为其分配一个唯一的标签 c_x ,用于表示节点 x 所存在的社区。

(2) 节点 x 将会根据其邻接点集 $N(x)$ 中的标签情况来迭代更新自己的标签 c_x ,这个过程一直持续到 c_x 等于 $N(x)$ 中包含的最多的那个标签为止。在此期间如果存在不止一个可选的标签的情况时,那么随机

选择多个标签中的一个来作为最后的更新结果。在经过几次迭代之后,各个节点的标签变化情况也将趋于稳定。

(3) 对于任意的节点 $x, y \in G(x)$,如果存在 $c_x = c_y$,那么判定节点 x, y 属于同一个社区。

RAK 算法的时间复杂度为 $O(km)$,其中参数 k 用来表示迭代次数,参数 m 表示网络中的边数。由此可以看出其时间复杂度接近线性,在处理大规模的网络时也能保证很好的效率,从作者得到的实验结果中也可以看出通过该算法挖掘得到的社区结构具有比较理想的模块化系数。但是,由于 RAK 算法有很多随机性的因素,这会导致其结果很不稳定。特别是当节点的标签在进行更新时,如果存在多个可选标签的情况,算法便会进行一次随机选择,这会导致算法每次产生的社区结构都会存在一定的差异。在处理大型网络的数据集时,这种现象出现的更加明显和频繁。因为在社会网络的结构中,节点都不是单独存在的,一个节点的标签情况不仅仅影响到它本身,还会对它的相邻节点甚至是一个范围内的网络节点造成不小的影响。

在 RAK 算法中,Raghavan 等人区别于同步(synchronous)更新提出了一种标签的异步(asynchronous)更新策略。在这里异步指的是在第 t 次标签传播的迭代过程中,当节点 x 要根据其相邻节点进行更新时,对于还未更新过的邻接点根据其 $t-1$ 次的标签情况对节点 x 进行更新,而对于已经更新过的邻接点就根据它第 t 次的标签更新结果进行更新。相对异步更新策略同步更新指的是:在标签传播的第 t 次迭代过程中,全部根据其邻接点的第 $t-1$ 次迭代的结果进行标签更新。最终通过实验结果研究者发现异步更新策略相对同步来说可能需要更多的迭代次数,但是得到的社区结构也相对更加稳定。由于 RAK 算法是首次将 LPA 的思想应用于社区挖掘中,所以在之后针对社区发现的研究中,也有很多研究者经常把 RAK 算法直接叫做 LPA 算法。

1.2 基于 RAK 算法的改进算法

在之后的研究中,Leung 等人^[13]对 RAK 算法进行了一些改进。在研究中他们指出模块最大化方法不是一个无标度区间的测度方法,仅依靠它来检测社区有它本身的局限性,于是又提出了一种扩展的 LPA 用于实时检测社区,通过采用启发式方法提高了其平均检测性能和适应性。在实验部分,研究者也分别利用同步和异步更新策略来检测网络社区,并对这两种更新策略进行了进一步的测试。另外通过简单的修改参数可使算法具有一定的可扩展性,从而适用于不同规模的网络。

Barber 等人^[14]为了避免在标签传播的过程中所

有顶点都分配到同一个社区,改进 RAK 算法,最终提出了一种模块化标签传播算法(LPAm),即基于约束的 LPA 社区挖掘算法。通过引入跳跃衰减(hop attenuation)和节点倾向性选择(node preference)防止 RAK 算法在一些网络上会出现挖掘的社区结构很大的情况。根据跳跃衰减的情况,每个标签会有一个评分:

$$S_i(l) = (\max_{j \in N(i)} s_j(l_j)) - \delta \quad (1)$$

式中 $s_i(l)$ 表示的是节点 i 上标签 l 的评分; l_j 表示的是节点 j 的标签; $N(i)$ 表示的是节点 i 的邻接点中标签为 l 的节点集合; δ 为衰减因子。初始化时节点的评分为 1。随着标签传播的过程评分会逐渐地减小,当评分降低为 0 的时候,这个标签就无法再传递给其他节点。通过引入跳跃衰减的策略就能有效地避免大社区的形成。通过后续的研究,实验结果发现跳跃衰减结合节点倾向性选择将能取得更好的效果,在这两种策略结合的情况下,节点 i 的新标签 l_i' 为:

$$l_i' = \operatorname{argmax}_{j \in N(i)} s_j(l_j) f(j)^m w(i, j) \quad (2)$$

其中 $w(i, j)$ 表示节点 i 和节点 j 之间边的权重; $f(j)$ 表示的是节点倾向性选择函数; m 参数起到一个平衡因子的作用。Barber 等人在实验中初始 $f(j) = 1$, 即表示节点的度数为 1。这样的话,如果取 $m > 0$, 就表示算法倾向于选择邻接点中度数较大的节点所带有的标签。若一开始取 $m < 0$ 则表示相反的选择。在通常情况下, m 取 0.1。

Liu 等人^[15]经过继续研究发现上述 LPAm 算法可能在模块空间中陷入局部极大值,从而会导致社区挖掘不准确的问题,在分析了多种方案后,提出了将 LPAm 算法与多步贪婪凝聚算法(MSG)相结合,并最终设计出了一种模块化专业化的标签传播算法(LPAm+)。该算法利用 MSG 算法可以同时合并多个相似社区,比较有效地避免了 LPAm 算法中陷入局部最大值的问题,从而保证可以更加精确地对社区结构进行挖掘。

2011 年,金弟等人^[16]通过进一步的研究提出了基于遗传算法的复杂网络社区挖掘方法,该算法主要是在分析网络模块性函数 Q 局部单调性的基础上,结合遗传算法给出了一种快速、有效的局部搜索策略——局部搜索的遗传算法。在算法的开始,研究者利用 LPA 作为初始种群的生成方法,提供了高精度和多样性的初始种群,最后再将标签传播给未标注节点,从而发现高质量的种群社区。

还有研究者从比较独特的角度改善了 LPA 算法在社区挖掘中的执行速度。比如 Xie 等人^[17]通过实验发现 LPA 算法在经过大概五次迭代后,95% 的节点已经可以正确地聚集。而后面的迭代过程主要是对社

区内节点的更新,相对来说没有太大的意义。根据这个实验结果他们改进了 LPA 算法中的更新准则和迭代规则,大大地减少了原算法中不必要的更新和迭代过程,特别是在处理较为复杂的网络结构时,该算法的执行效率有了显著的提高。

在后续的研究中 Cordasco 等人^[18]也提出了一种半同步的 LPA 算法,该算法通过对网络节点并行着色,使任何两个相邻的节点都不拥有相同的颜色,并行同时进行传播标签。该算法中所使用的更新策略,能够比较好地克服网络中存在的节点振荡问题,而且每一步传播的迭代过程都并行作业,从而进一步提高了 LPA 算法的执行速度。

2 重叠 COPRA 算法及其改进算法研究

2.1 可挖掘重叠社区的 COPRA 算法描述

以上所提到的算法虽然从很多方面对 LPA 算法进行了改进,但是这些算法所挖掘得到的社区结构都是非重叠的,也就是说,每个节点只能从属于某一个社区,社区之间不会存在重叠结构。显然这与现实社会网络中的情况是不相符的,针对这种情况 Steve 基于 Raghava 等人的研究进一步提出了 COPRA 算法,通过该算法就可以挖掘得到重叠的社区结构。

在 COPRA 算法中,为了能够挖掘得到重叠社区,Steve 在研究中引入了新的标签结构 (c, b) , 对每个节点 $x, x \in G(x)$ 都拥有一组这样的标签集。其中,参数 c 表示社区的标识符;参数 b 用来表示节点 x 在社区 c 中的从属系数(belonging coefficient),并且 $0 \leq b \leq 1$ 。在此基础上通过 $b_t(c, x)$ 表示在第 t 次迭代过程中节点 x 与社区 c 之间的从属系数,具体的计算方法如公式(3)所示:

$$b_t(c, x) = \frac{\sum_{y \in N(x)} b_{t-1}(c, y)}{|N(x)|} \quad (3)$$

其中, $N(x)$ 表示节点 x 的邻接点集合。通过对上面式子的观察可以发现 COPRA 算法在标签传播过程中所使用的是同步更新策略。接下来借助图 1 来简单描述一下 COPRA 算法的标签更新过程并且给出其相应的算法描述。

在图 1 中节点 5 需要根据其相邻节点 1、2、3、4 的标签情况来更新自己的标签,在这个过程中首先要计算其邻接点中所存在的社区标签的和,接下来选择 v 个数值最大的标签(v 表示设定的重叠社区个数,图中 v 的值为 2)。若存在多个可选标签的话,这时就需要进行标签的随机选择,最后再将所选的标签进行标准化得到节点 5 的标签更新结果。

现在将 COPRA 算法过程描述如下:

(1) 首先要进行初始化,对于任意的节点 $x, x \in G(x)$,一开始为其分配一个唯一的标签,并且初始它的从属系数为1,如果用参数 c_x 表示 x 所存在的社区,那么根据前面的介绍,现在就可以用 $(c_x, 1)$ 表示这种标签结构。

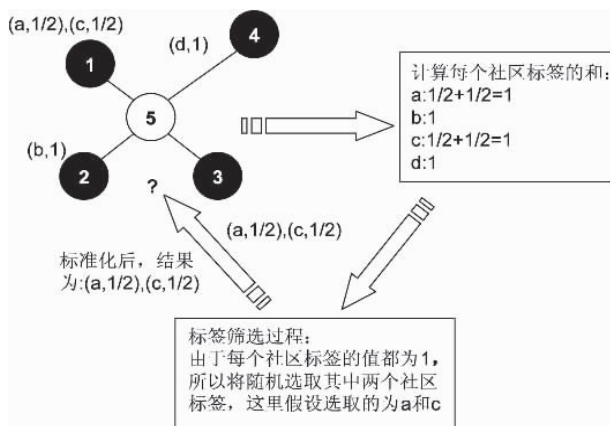


图1 标签传播过程示例

(2) 下面是算法的迭代过程,节点 x 会根据它的邻接点集的标签情况来更新自己的标签。在标签更新的过程中若存在多个可选标签的话,算法就会随机选择其中的 v 个标签来作为更新结果。作者设置了参数 v 来限制每个节点可以拥有的标签数目,是为了避免所有的标签都更新成相同的结果,其中 $0 < |C(x)| \leq v$ 。每次标签更新完成后,还需根据相应的策略对标签的从属系数进行一次标准化,从而保证在第 t 次迭代的过程中,节点的从属系数符合如下的表达式:

$$\sum_{c_x \in C(x)} b_i(c_x, x) = 1 \quad (4)$$

其中 $C(x)$ 表示节点 x 可能属于的社区集合。

(3) 上述的迭代过程会在满足设定的结束条件后所终止。

(4) 在最后将具有相同社区标签节点合并为同一社区。

另外从图1中节点5更新标签的例子中也可以看到,在遇到存在多个标签可选的情况下,该算法会随机进行选择,这个过程也会导致该算法的不稳定性。

2.2 基于 COPRA 算法的改进算法

在 COPRA 算法的基础上, Wu 等人^[19]于2012年对其进行改进并提出了基于多标签传播的平衡社区发现算法(BMLPA),该算法重新设计了 COPRA 中的标签更新策略,通过设定阈值 p 来控制每个节点可以拥有的标签数目,从而使得在新算法中不需要初始 v 的值,也就是说挖掘的重叠社区的数目不再受参数 v 的限制。

另外该算法也从某些角度增加了算法的稳定性,但随之带来的问题是该算法挖掘得到的社区结构比较固定,使得算法的适应性有所下降。

3 基于标签传播社区挖掘算法改进思路总结

经过对基于标签传播的社区挖掘算法及其改进算法的研究,发现这些算法还可以从以下三个方面进行改进:

(1) 标签初始阶段:由于标签传播的第一步是先为网络中所有节点分配一个唯一的标签,之后这些标签会不断地更新,当这个网络的规模越大时,更新标签所需要的资源消耗就越大,可以注意到这么一点,网络中有很多节点其实是可以绑定到一起的,如果这样的话就可以从整体上大大减少初始的标签数目,从而提高标签更新时算法执行的效率,但是这样就需要预先对所有节点进行简单的处理分类,考虑到在之后的标签更新过程中给算法带来的收益,这种方法也是可以值得一试的。另外也可以考虑简单的将基于标签传播的社区挖掘算法同其他经典的社区挖掘算法相结合,最后平衡地取舍多种算法的优劣,从而得到效率、性能都比较让人满意的算法。

(2) 标签传播阶段:在这个阶段,首先可以考虑从同步和异步更新策略来进行改进,根据之前的研究知道同步更新策略的优势是所需的迭代次数少,但是劣势是得到的社区结构没有异步更新所挖掘得到的稳定。综合两种更新策略的优劣,提出半异步的更新方法,也就是说在对标签进行更新时,一部分节点采用同步更新的办法,另一半则采用异步更新,而如何进行选择则可以根据节点的度数大小或者其他的属性来进行判断。另外也可以如同 BMLPA 算法一样,巧妙地利用阈值参数来重新设定更新策略。

(3) 标签选择阶段:基于标签传播的算法之所以存在不稳定性的原因,都是由于如果存在多个标签可选的情况时需要进行一次随机选择。可以通过节点的度数或者其他合理的策略去人为地来干预这个随机选择的过程,比如增大选到这多个标签中影响最大的节点的可能性等等。

另外值得一提的是,目前对基于标签传播的重叠社区挖掘算法的研究相对来说还比较少,而且发现 COPRA 算法和 BMLPA 算法都是在 RAK 算法的基础上提出的,因此这两种算法没有吸收 RAK 算法之后很多的改进算法的优势。在以后的研究中也可以将几种改进算法放到一块来得到新的算法。总体来说,在重叠社区挖掘这方面还有比较大的研究空间。

4 结束语

文中主要结合 RAK 算法和 COPRA 算法分析研究了基于标签传播的社区挖掘算法各自的特点及优劣,

并在此基础上提出了新算法在以后研究中的改进思路。由于真实社会网络中存在的重叠社区结构,所以基于 COPRA 算法的研究所能挖掘得到的社区结构是更具有现实意义的。在研究中还发现,标签传播算法在其他的领域也有很多应用,比如文本信息检索、多媒体信息检索等方面^[20]。也可以尝试从这些研究领域中借鉴一些新的研究思路。

参考文献:

- [1] Girvan M, Newman M E J. Community structure in social and biological networks [J]. Proc natl acad sci USA, 2002, 99 (12): 7821 – 7826.
- [2] Capocci A, Servidio V D P, Caldarelli G, et al. Detecting communities in large networks [J]. Physica A: Statistical and theoretical physics, 2005, 352(2 – 4): 669 – 676.
- [3] Boccaletti S, Latora V, Moreno Y, et al. Complex networks: Structure and dynamics [R]. [s. l.]: [s. n.], 2006.
- [4] Radicchi F, Castellano C, Cecconi F, et al. Defining and identifying communities in networks [J]. Proceedings of the national academy of sciences, 2004, 101(9): 2658 – 2663.
- [5] Fortunato S. Community detection in graphs [R]. [s. l.]: [s. n.], 2010.
- [6] Raghavan U N, Albert R, Kumara S. Near linear time algorithm to detect community structures in large – scale networks [J]. Phys rev, 2007, E76(3): 036106.
- [7] Palla G, Derényi I, Farkas I. Uncovering the overlapping community structure of complex networks in nature and society [J]. Nature, 2005, 435(6): 814 – 818.
- [8] Shen H W, Cheng X Q, Cai K. Detect overlapping and hierarchical community structure in networks [J]. Physica A, 2009, 388(8): 1706 – 1712.
- [9] Lee C, Reid F, McDaid A, et al. Detecting highly overlapping community structure by greedy clique expansion [J/OL]. 2010. arXiv: 1002.1827v2.
- [10] Lancichinetti A, Fortunato S, Kertész J. Detecting the overlapping and hierarchical community structure of complex networks [J]. New J phys, 2009, 11: 033015.
- [11] Gregory S. Finding overlapping communities in networks by label propagation [J]. New J physics, 2010, 12(10): 103 – 118.
- [12] Zhu Xiaojin, Ghanramani Z. Learning from labeled and unlabeled data with label propagation [R]. Pittsburghers: Carnegie Mellon University, 2002.
- [13] Leung I X Y, Hui P, Liò P, et al. Towards real – time community detection in large networks [J]. Phys rev, 2009, 79(6): 066107.
- [14] Barber M J, Clark J W. Detecting network communities by propagating labels under constraint [J]. Physical review E, 2009, 80(2): 129 – 139.
- [15] Liu Xin, Murata T. Advanced modularity – specialized label propagation algorithm for detecting communities in networks [EB/OL]. 2010 – 03 – 19 [2012 – 05 – 01]. <http://arxiv.org/pdf/0910>.
- [16] 金 弟, 刘 杰, 杨 博, 等. 局部搜索与遗传算法结合的大规模复杂网络社区探测 [J]. 自动化学报, 2011, 37(7): 873 – 882.
- [17] Xie Jierui, Szymanski B K. Community detection using a neighborhood strength driven label propagation algorithm [C]// Proc of IEEE network science workshop (NSW). [s. l.]: [s. n.], 2011: 188 – 195.
- [18] Cordasco G, Gargano I. Community detection via semi synchronous label propagation algorithms [EB/OL]. 2011 – 04 – 23 [2012 – 03 – 05]. <http://arxiv.org/abs/1103>.
- [19] Wu Zhihao, Lin Youfang, Gregory S. Balanced multi – label propagation for overlapping community detection in social networks [J]. JCST, 2012, 27(3): 468 – 479.
- [20] 张俊丽, 常杨丽, 师 文. 标签传播算法理论及其应用研究综述 [J]. 计算机应用研究, 2013, 30(1): 21 – 25.

(上接第 68 页)

- computational sciences. Hangzhou: Zhejiang University Press, 2006: 780 – 786.
- [3] 朱晓武. 商务智能的理论和应用研究综述 [J]. 计算机系统应用, 2007(1): 114 – 117.
- [4] 王 娟. 数据挖掘技术在旅游业中的应用——以黄山市为例 [J]. 国土与自然资源研究, 2012(4): 72 – 73.
- [5] 谢 炜, 徐晓飞, 刘 昊, 等. 商务智能: 新一代决策支持领域 [J]. 计算机科学, 2001, 28(4): 9 – 12.
- [6] 党安荣, 张丹明, 陈 杨. 智慧景区的内涵与总体框架研究 [J]. 中国园林, 2011(9): 15 – 21.
- [7] Inmon W H, Strauss D, Neushloss G. DW2.0: The architecture for the next generation of data warehousing [M]. [s. l.]: Morgan Kaufmann, 2003.
- [8] 宋丽丽, 王嵘冰. 商务智能系统的数据体系结构研究 [J]. 辽宁大学学报(自然科学版), 2009, 36(1): 55 – 58.
- [9] 江光中. 湖北烟草商务智能决策支持系统设计 [D]. 北京: 北京工业大学, 2009.
- [10] 汪祖丞, 刘 玲. 旅游客流量预测模型的比较及其实证研究——以黄山风景区为例 [J]. 安徽师范大学学报(自然科学版), 2010, 33(3): 286 – 290.
- [11] 周 瑾. 基于商务智能的税务征管决策支持研究 [J]. 中国管理信息化, 2009(9): 88 – 90.
- [12] 裴盈盈, 袁国宏. 智慧旅游浅析 [J]. 当代经济, 2012(5): 46 – 47.