

基于关系数据模型的犯罪网络挖掘研究*

李万彪^{1,2}, 余志¹, 龚峻峰³, 陈锐祥¹

(1. 中山大学工学院//智能交通研究中心, 广东 广州 510275;

2. 广东机电职业技术学院信息工程学院, 广东 广州 510515;

3. 华南理工大学土木与交通学院, 广东 广州 510641)

摘要: 提出了一种基于关系数据模型的犯罪网络挖掘方法, 基于侦查办案分析需求, 建立关系数据模型生成犯罪网络, 运用该模型可挖掘已知嫌疑人的其他团伙成员, 也可在未知嫌疑人的条件下生成所有的犯罪网络, 具有共同对象的不同模型还可组合进行跨模型网络挖掘。利用为期半年的通话与转账数据作为数据源进行实验, 实验结果表明, 文中模型不仅可以直观准确地实现已知嫌疑人和未知嫌疑人的犯罪网络挖掘, 还可通过不同模型的组合应用深度挖掘犯罪团伙关系。

关键词: 关系数据; 关系数据模型; 犯罪网络; 社会网络分析

中图分类号: TP391 **文献标志码:** A **文章编号:** 0529-6579(2014)05-0001-07

An Approach of Crime Network Analysis Based on Association Data Model

LI Wanbiao^{1,2}, YU Zhi¹, GONG Junfeng³, CHEN Ruixiang¹

(1. Research Center of ITS//School of Engineering, Sun Yat-sen University, Guangzhou 510275, China;

2. College of Information and Engineering of Guangdong Jidian Polytechnic, Guangzhou 510515, China;

3. School of Civil Engineering and Transportation, South China University of Technology, Guangzhou 510641, China)

Abstract: In order to uncover the criminal gang, an approach of Crime Network Analysis based on association data model is proposed. The proposed method can discover all the members of a gang when a certain member is known, and it can also generate all the gangs based on known data source even none of the members are known. Especially, different models can be combined together to excavate the crime network. Experiments are executed based on the data of “communication” and “transfer” during 6 months. The results show that the proposed algorithm can demonstrate criminal gang no matter its members are known or not, and the same occurs as combining different models.

Key words: association data; association data model; crime network; social network analysis

近年来, 随着社会经济的迅猛发展, 各类经济案件和刑事案件的发案率不断上升。据公安部门统计, 在这些案件中, 多人犯罪或团伙犯罪的比重越来越大, 而且, 犯罪团伙的隐蔽性越来越高, 犯罪手段的科技含量越来越高。因此, 有效地发掘并

遏制打击犯罪团伙成为公安侦查办案部门的重要工作。

有效发掘犯罪团伙, 归根结底为两个问题: 一是已知嫌疑人, 发掘出该特定嫌疑人的其他犯罪团伙; 二是在未知嫌疑人的前提下识别不同的犯罪团

* 收稿日期: 2013-12-13

基金项目: 广东省科技计划资助项目(2011B090400620); 中央高校基本科研业务费专项资助项目

作者简介: 李万彪(1982年生), 男; 研究方向: 智能交通系统、数据挖掘; 通讯作者: 龚峻峰; E-mail: ctgongj@scut.edu.cn

伙。随着信息化技术的不断发展,各个部门都积累了大量的数据资源,比如公安的户政数据、银行的账户转账数据、电信运营商的通讯数据,可通过分析这些数据资源来挖掘线索从而协助侦破案件,然而,面对海量的数据,仅依靠人工分析识别犯罪团伙几乎是不可能的,需要有高效智能化的分析手段来识别犯罪团伙。

“9.11”事件后,国内外学术界对犯罪团伙的识别作了大量的研究工作,最典型的是将社会网络分析(Social Network Analysis SNA)的方法引入犯罪侦查领域^[1],基于社会网络分析的手段,研究发掘犯罪团伙,称为犯罪网络分析(Crime Network Analysis CNA)^[2-4]。

一些学者,特别是 Sparrow^[5], Coles^[6], Klerks^[7] 和 Williams^[8], 在犯罪网络分析的研究中,提出了犯罪网络结构的一些特征(网络大小、密度、关联强度、中心性等),以及分析这些特征的方法,为犯罪网络的理论发展做出了重大贡献。Krebs^[4]通过搜集和整理有关“9.11”恐怖袭击事件成员之间的社会关系,构建恐怖袭击网络并进行分析。XU^[3]等人根据犯罪网络相关分析方法,设计出一个犯罪网络知识发现体系 CrimeNet Explorer,该体系能够建立犯罪网络,并进行犯罪网络结构分析和可视化。高建强等^[9]提出的以“某个嫌疑人为团伙犯罪中的一员”为前提假设进行犯罪网络的团伙分析,是犯罪网络分析中已知嫌疑人分析犯罪团伙的重要应用。四川大学的唐常杰及其团队等人基于最短路径算法,提出了网络核心挖掘算法和子网络分析算法等^[10-19]。马方^[20]介绍了社会网络分析方法在我国犯罪团伙挖掘中的应用及挑战。黄慧霞^[10]分析我国近年来的犯罪数据,发现了贩毒组织的小世界的组织特征。潘芳,张自力等将模糊处理技术与层次聚类算法相结合^[21-22],应用于犯罪网络分析,提出了一种新的基于模糊层次聚类算法(FHCM)的犯罪网络分析方法,基于9.11犯罪数据进行实验,划分犯罪网络并找到不同网络的中间人。李亮^[17]从公安业务的角度出发,基于已掌握犯罪团伙某一个或多个嫌疑人资料的情况,建立犯罪网络识别系统,采用 Radicchi 快速分裂算法对社会网络进行分解,并基于 SPLINE 的 KMM 算法对子网络进行分析,突出犯罪团伙的核心成员,协助侦查办案部门进行犯罪团伙筛选并应用于苏州公安部门,取得了不错的应用效果。

然而,上述研究大部分集中于社会网络结构的理论研究或小型的犯罪网络研究,基于公开的数据

进行实验(如9.11事件相关数据),但应用到海量数据挖掘的较少。本文针对实际可能获取的海量数据与分析需求,利用社会网络分析手段进行犯罪网络分析,分析数据特点,建立关系数据模型,开展已知嫌疑人和未知嫌疑人的犯罪团伙及成员挖掘,从各类数据资源中发掘犯罪团伙信息。

1 背景及基本概念

1.1 问题背景

通常情况下,团伙犯罪案件中的犯罪嫌疑人都有着较密切的联系,如通过分析部分案件发现几名涉案嫌疑人相互之间通信频繁、银行账户转账密切。

基于团伙犯罪的一般规律,希望可以快速、准确、直观地从海量的数据资源中发掘犯罪团伙信息。而发掘犯罪团伙信息存在两种情况,一种情况是已知某个嫌疑人是潜在的犯罪团伙的一员,根据已知的团伙成员发掘与其有关联关系的其他团伙成员,从而发掘整个犯罪团伙;另一种情况是从已有的数据资源中发掘未知的犯罪网络和犯罪网络中的成员,从而可以实现犯罪预警。

1.2 基本概念

1.2.1 对象 对象是有可区别性且独立存在的某种事物,是一具有相同属性描述的实体集合,对象可用来指人、动物、植物、地名、机构名、汽车等事物。在本文中,对象指人、电话、账户等。

1.2.2 属性数据 对象所具有的特征称为属性。如人有身份证号码、姓名、年龄、性别等属性;又如,电话有电话号码、归属地、机主身份证号码等属性。

属性数据指用来描述对象的属性的数据,仅与对象本身的特征有关,不涉及对象与其他对象之间的联系。属性数据可以通过二维矩阵来描述,即以行、列的形式来描述,行代表对象,列代表每一个对象所具有的不同属性。属性数据可以很清晰地描述对象自身的特征,但是无法表征对象与对象之间的关联关系。

1.2.3 关系数据 关系数据除了具备属性数据的基本特征外,还包括对象与对象之间的关联关系(也称联系:Association),如人与电话间存在所有者关系,两人如若通话则存在“通话”关系;人与银行账户间存在所有者关系,两人账户有经济往来则存在“转账”关系。关系数据不仅表征对象本身的属性,还能表征对象与其他对象之间的联系。

1.2.4 关系数据模型 基于对象、属性及对象间

的关联关系等要素，建立一个四元组网络模型 $M = (O, K, P, A)$ 。O 表示对象的集合；K 表示对象的主键，用以唯一识别对象；P 表示对象的属性集；A 表示对象与对象之间的关联关系。利用关系数据模型可以搜索特定对象并挖掘对象之间的关系，也可以组合多个模型搜索，发掘关系网络。

1.2.5 社会网络与犯罪网络 文献 [13] 指出，社会网络 (Social Network) 是社会行动者及其间的关系集合。社会网络是社会个体成员之间互动形成的相对稳定的关系体系，任何一种用于建立个体之间联系的自然现象、社会活动或技术体制都能形成一个网络。社会网络分析利用社会学、数学和图论等相关方法，分析对象节点、对象连结之间的社会关系，对象节点表示网络的行动者，结点间的边表示行动者之间的关系。

社会网络分析的手段应用于犯罪侦察领域，称为犯罪网络分析 (SNA)，相应的社会网络称为犯罪网络 (Crime Net)。在犯罪网络中，行动者通过通话、转账等发生关联关系，进而形成关系网络。

2 研究方法

2.1 数据预处理

本文利用关系数据，建立模型分析对象与对象间的关联关系。要对数据进行预处理，处理成对象与关系的格式，如“人”与“电话”间的“所有者”关系。

再者，本文研究采用的实验数据格式与质量参差不齐，需要进行必要的预处理，包括数据清洗、转换、隐私信息加密处理。首先，由于采集来源、采集手段和录入水平不一，导致部分数据存在错误的取值、空值或者重复的取值，这些数据记录在应用中属于“脏”数据，需要清洗。如，身份证号码为“12345678”类似的值，为明显错误取值，清洗时需将记录删除；又如，账户转账时间为空值，则可考虑平滑处理，以记录的录入时间来代替。其次，不同数据来源存在将同一对象定义为不同的数据类型的情况，比如有些数据源将日期定义为 DATETIME 类型，但另一些数据源却定义为 VARCHAR 类型，表征同一对象的数据字段类型要转换成统一的数据类型。此外，由于部分数据涉及隐私信息，因此在数据分析前需要对数据进行加密处理，对姓名、身份证号码、手机号、账户号等信息要作映射处理，比如将身份证号后四位及姓名的最后一位转换成“A”到“Z”的大写字母。

本文实验采用的是通话记录 (手机号码数据、

通话数据) 和转账记录 (账户数据、转账数据)，按照规则预处理后的部分数据格式如表 1-4 所示。

表 1 手机号码数据格式

Table 1 Sample data format of mobile number

ID	机主号码	机主身份号码	归属地	……
00001	13570	41012219	021	……
	46ABCD	871006ABCD		

表 2 账户数据格式

Table 2 Sample data format of account

ID	账户号码	户主身份号码	开户行	……
00001	44520219	41012219	001	……
	79ABCD	871006ABCD		

表 3 通话数据格式

Table 3 Sample data format of communication

ID	机主号码	被叫号码	通话时段	时长	…
00001	13570	13879	201010	05: 22	…
	46ABCD	73ABCD	2702		

表 4 转账数据格式

Table 4 Sample data format of transfer

ID	账户号码	收款方号码	转账时段	转账金额	…
00001	44520219	47520249	201010	1 000.00	…
	79ABCD	85 ABCD	2703		

2.2 模型构建

2.2.1 建模步骤 基于关系数据特点，针对不同的业务需求建立关系数据模型，运行模型生成犯罪网络，如“通话模型” (模型中成员具备通话关系)、“转账模型” (模型中的成员具备经济往来关系) 的关系数据模型。不同模型的建模算法相同，本文以“通话”模型为例，介绍建模的基本算法。

1) 选择数据源。本模型以通话记录 (手机数据、通话数据) 为数据源。

2) 确定模型对象。本模型包括两个对象：人与手机。

3) 确定对象主键字段。主键唯一识别对象，以身份证号作为对象“人”的主键，手机编号作为对象“手机”主键。

4) 确定对象的属性。对象“人”的属性包括姓名、性别、身份证号等；对象“手机”的属性包括手机号码、机主身份证号码等。

5) 确定对象间的关系。关系的表述方式为“从对象到对象”，在某特定时段内某手机号码给

另一手机号码拨号则称两手机的机主间存在通话关系, 如两人在 n 个特定时段内通话则他们存在 n 次通话关系。

2.2.2 模型实现算法 基于图论, 以无向带权图 $G = (V, E)$ 实现模型 $M = (O, K, P, A)$, 将对象视为无向图的结点, 以对象的主键标识, 对象与对象之间的关系用结点间的边来表示, 权重为对象关联关系的次数。

符号定义: 以无向带权图 $G = (V, E)$ 实现模型 $M = (O, K, P, A)$, 将对象视为无向图的结点, 用对象的主键进行标识, 对象与对象之间的关系用结点间的权重用无向图的边来表示。 DS 为预处理后的通话数据集, $D(d)$ 为特定通话时段; $\Theta_{PID} = \{x | x \text{ 为数据集 } DS \text{ 中主叫号码的身份证号码}\}; n = |\Theta_{PID}|$ 为集合 Θ_{PID} 的元素个数; $\Theta_{HMD} = \{y | y \text{ 为数据集 } DS \text{ 中被叫号码的身份证号码}\}; m = |\Theta_{HMD}|$ 为集合 Θ_{HMD} 的元素个数; $V_i \in \Theta_{PID} (1 \leq i \leq n)$, $V_j \in \Theta_{HMD} (1 \leq j \leq m)$; $G = (V, E)$ 为由结点 V_i 、 V_j 生成的无向带权图; 边 E_{ij} 表示对象 V_i 与对象 V_j 间是否存在关联关系, E_{ij} 的取值公式为

$$E_{ij} = \begin{cases} 0, & \text{对象 } V_i \text{ 与 } V_j \text{ 之间不存在关系} \\ t, & \text{对象 } V_i \text{ 与 } V_j \text{ 之间存在 } t \text{ 次关系} \end{cases} \quad (1)$$

基于通话数据 DS (数据格式如表 3 所示) 实现“通话”模型, 生成犯罪网络 (无向带权图 G) 的算法伪码描述如下:

```

sortDSbyD(d) and V_i asc
initial, set E_ij = 0 (1 ≤ i ≤ n, 1 ≤ j ≤ m)
foreach D(d)
for i = 1 to n
begin
if V_i ∉ G then add V_i to G
foreach V_j
begin
if V_j ∉ G then
begin
add V_j to G
set E_ij = E_ij + 1
end
if V_j ∈ G then set E_ij = E_ij + 1
end
end
generate G

```

2.3 模型应用

关系数据模型基于特定的数据源, 可单独用于已知嫌疑人和未知嫌疑人的犯罪团伙识别, 也可多

个模型组合进行犯罪网络挖掘。

2.3.1 单模型网络挖掘 基于单个模型, 可以生成特定时间内模型中对象与对象之间的关系。如利用“通话”模型可生成人与人的关联关系, 结点之间的联系表示对象之间至少存在一次直接的关系。

2.3.2 组合模型网络挖掘 实际应用中, 基于单个模型有时候不足以深度挖掘犯罪网络结构。如仅仅基于“通话”或“转账”模型, 只能挖掘出特定时间内的通话或经济关系, 存在一定的偶然性, 并不足以证明对象之间存在“事实上”的团伙关系。如果某些人在特定时间内既多次通话, 又多次进行转账, 则其属于同一团伙的可能性就更高。利用本文提出的关系数据模型, 可实现组合挖掘, 不同模型只要彼此存在一个共有对象, 并且对象的主键相同, 就可以组合进行网络挖掘。

3 实验与分析

3.1 实验环境与实验场景

本文采用为期半年的通话与转账数据作为数据源, 其中通话数据约 5 000 万条记录, 转账数据约 1 000 万条记录。实验前已对数据进行预处理。

本文分别利用单个模型和多个模型组合的方式, 挖掘已知嫌疑人和未知嫌疑人情况下的犯罪网络, 实验模型选用一定特定时段内的“通话”和“转账”模型。由于通话的频率一般都高于转账, 本实验对通话次数选取较高度量值。

3.2 实验结果与分析

3.2.1 已知嫌疑人挖掘其他团伙成员 已知嫌疑人“王 X”, 身份证号为“41012219871006XXXX”, 挖掘其可能的团伙成员。分别利用“通话”和“转账”模型挖掘与嫌疑人“王 X”半年内有通话或转账关系的可能团伙成员, 进一步利用“通话”及“转账”模型组合挖掘与嫌疑人“王 X”既通话又有转账往来的成员。实验结果如图 1-7 所示。



图 1 与已知嫌疑人有 1 次通话及转账关系的成员

Fig. 1 Members who have only one association with Wang

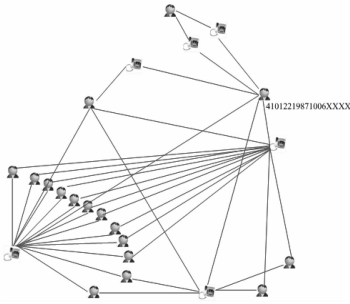


图 2 与已知嫌疑人有 5 次通话关系的成员
Fig. 2 Members who have 5 times mobile communications with Wang

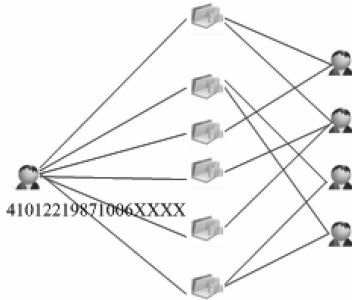


图 3 与已知嫌疑人有 2 次转账关系的成员
Fig. 3 Members who have twice transfer associations with Wang

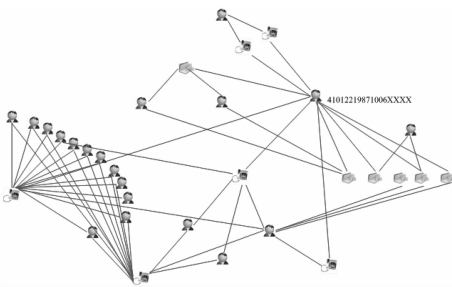


图 4 与已知嫌疑人有 5 次通话且 2 次转账关系的成员
Fig. 4 Members who have 5 times communication and twice transfer associations with Wang

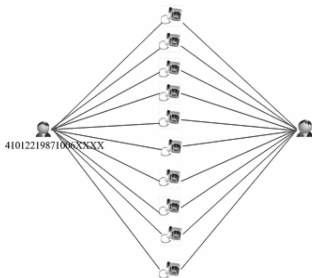


图 5 与已知嫌疑人有 10 次通话关系的成员
Fig. 5 Members who have 10 times associations with Wang



图 6 与已知嫌疑人有 3 次转账关系的成员
Fig. 6 Members who have 3 times transfer associations with Wang

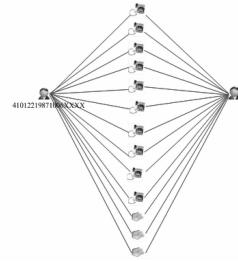


图 7 与已知嫌疑人有 10 次通话及 3 次转账关系的成员
Fig. 7 Members who have 10 times communications and 3 times transfer associations with Wang

实验结果表明：① 本文提出的关系数据模型可以挖掘并直观展示与已知嫌疑人王 X 存在通话转账关系的可能的团伙成员。② 特定时间内有 1 次通话转账关系的结点数量庞大，大多数并无必然联系。如表 5 所示，与王 X 有通话关系的人达到 465 人，转账的人数达到 121 人，其中绝大多数是单次通话或转账，无法确定是否属于同一团伙。③ 特定时间内有多次通话、转账、通话且转账关系的人，其属于某一团伙的可能性较高，需要重点关注。

表 5 与嫌疑人有 N 次关联关系的人数
Table 5 number of members who have N times associations with Wang

N	通话	转账	通话及转账
> = 1 (1)	465	121	65
> = 5 (2)	38	17	4
> = 10 (3)	5	2	1

3.2.2 未知嫌疑人挖掘所有潜在的犯罪网络 未知嫌疑人时，利用单个“通话”、“转账”模型可生成所有具有通话、转账关系的犯罪网络；进一步组合模型，可生成既通话且转账的犯罪网络。经验表明，通话或转账 1 次的关系中有大量的冗余信息，无法准确展示潜在的犯罪网络，故本实验首先对数据进行过滤，只考虑 2 次及以上的关联关系，

另外, 对象间存在的仅为1次的偶然的关联关系对犯罪团伙分析的作用并不大, 本文的算法并未对这部分数据做深入的处理。后续研究将对数据进行深入的预处理, 提高模型分析效率; 同时, 考虑对不同数据源(如通话记录、转账记录、定位记录等)合并分析, 构建如“通话-转账-位置”的复杂模型, 从数据挖掘角度就发现团伙、研究团伙内部角色等问题进行深入广泛的分析。

参考文献:

- [1] FREEMAN L. Centrality in social network: conceptual clarification [J]. *Social Networks*, 1979, 1: 215 - 239.
- [2] LIU Xiaoming, BOLLEN J, ENSON M L. Co-authorship networks in the digital library research community [J]. *Information Project & Management*, 2005, 41: 1462 - 1480.
- [3] XU Jennifer, CHEN Hsinchun. CrimeNet Explorer: A framework for criminal network knowledge discovery [J]. *ACM Transactions on Information Systems*, 2005, 23 (2): 201 - 226.
- [4] KREBS V. Mapping networks of terrorist cells [J]. *Connections*, 2002, 24(3): 43 - 52.
- [5] SPARROW M K. The application of network analysis to criminal intelligence: An assessment of the prospects [J]. *Social Networks*, 1991, 13(3): 251 - 274.
- [6] COLES N. It's not what you know- It's who you know that counts. Analyzing serious crime groups as social networks [J]. *British Journal of Criminology*, 2001, 41(4): 580 - 594.
- [7] KLERKS P. The network paradigm applied to criminal organizations [J]. *Connection*, 2001, 24(3): 53 - 65.
- [8] WILLIAMS P. Transnational criminal networks [C] // J Arquilla and D. Ronfeldt (eds), *Networks and Netwars: The Future of Terror, Crime and Militancy*. Santa Monica: and Corporation, 2001: 61 - 97.
- [9] 高建强, 谭强, 崔永发. 一种基于通讯痕迹的社会网络团伙分析模型[J]. *计算机应用与软件*, 2012, 29(3): 206 - 208.
- [10] 黄慧霞. 跨境毒品犯罪组织结构的社会网络分析[J]. *中国人民公安大学学报: 社会科学版*, 2010, 143: 29 - 40.
- [11] 张德丰, 马子龙, 梁忠宏. 基于数据挖掘技术的算法研究[J]. *中山大学学报: 自然科学版*, 2004, 43(3): 36 - 39.
- [12] 乔少杰, 唐常杰, 彭京. 基于个性特征仿真邮件分析系统挖掘犯罪网络核心[J]. *计算机学报*, 2008, 31(10): 1795 - 1802.
- [13] 周利娟, 林鸿飞, 罗文化. 基于实体关系的犯罪网络识别机制[J]. *计算机应用研究*, 2011, 28(3): 998 - 1002.
- [14] 林聚任. *社会网络分析: 理论、方法与应用* [M]. 北京: 北京师范大学出版社, 2009.
- [15] 约翰斯科特. *社会网络分析法* [M]. 刘军, 译. 重庆: 重庆大学出版社, 2007: 33 - 53.
- [16] 罗兆波. 面向社会网络分析的数据挖掘方法研究 [D]. 杭州: 浙江大学, 2010.
- [17] 李亮. 基于社会网络分析的犯罪团伙识别系统 [M]. 上海: 上海交通大学, 2008.
- [18] 柳国华, 谢璨, 英春. 基于短信的社会网络行为分析 [J]. *计算机应用与软件*, 2011, 28(6): 220 - 22.
- [19] 陈鹏, 袁宏永. 犯罪组织结构的社会网络分析 [J]. *清华大学学报: 自然科学版*, 2011, 51(8): 1097 - 1101.
- [20] 马方. 犯罪网络分析: 社会网络分析在有组织犯罪研究中的应用 [J]. *西南政法大学学报*, 2012, 14(2): 35 - 43.
- [21] 潘芳, 张自力. 对犯罪网络的模糊层次聚类分析 [J]. *西南师范大学学报: 自然科学版*, 2009, 34(3): 200 - 214.
- [22] 张昕, 关志超, 杨东援. 基于多目标数据挖掘的城市交通仿真算法研究 [J]. *中山大学学报: 自然科学版*, 2007 (S2): 210 - 214.