



一种 IP 地址地理定位方法的模型

倪晶晶¹, 龚俭¹

(1. 东南大学计算机科学与工程学院, 南京, 211189)

摘要: 为满足网络行为监控和管理系统中活跃 IP 地址以及安全状态 IP 地址的跟踪定位服务, 需要提出有效的地理定位算法获取 IP 地址实际使用位置, 供用户查询。国内外已经有一些地理定位数据库存在, 但是它们提供的信息可能不全。针对那些不能直接从现有数据库查询的 IP 地址, 本文提出一种基于网络拓扑数据的定位方法模型: 利用 Traceroute 路径信息, 通过关联规则, 对 IP 地址进行关联分析, 找到那些可能在同一地理位置的 IP 地址, 形成 IP 地址组, 同组内, 从已定位 IP 地址推断未定位 IP 地址的实际使用地理位置。

关键词: 网络管理; 网络测量; IP 地址地理定位; traceroute; IP 地理位置关联分析

An Model of Geolocation Method for IP Address

Ni Jingjing¹, Gong jian¹

(1. School of Mechanical Engineering, Southeast University, Nanjing, 211189)

Abstract: To meet the need of interest IP's geolocation service in network behavior monitoring and management system, need to come up with effective geographic positioning algorithm for the user's query. There are already some geographical location databases at both home and abroad, but the information they provide may be incomplete. For those IP addresses who can not directly query from an existing database, we propose a geolocation method model based on network topology data: do association analysis for IP addresses through association rules using traceroute information, to find those IPs who may be in the same geographic location, as a group of IP addresses. In a group, we can infer the actual geographic location of a IP which is not located from IPs whose location is already known.

Key words: Network management; Network measurement; IP Geolocation; traceroute; geolocation Co-relation analysis

互联网地理定位是确定一个互联网用户地理位置的问题, 这通常被称为网络实体地理定位。IP 地址的地理定位对于很多针对客户端位置的网络应用是很重要的, 如, 根据地理信息显示当地新闻以及天气等相关信息的网络服务, 从 web 日志分析出客户的地理位置分布, 区域性的 P2P 建设, 自动选择语言显示网站内容, 根据地区政策控制内容发布, 类似内容分发网络的多媒体传输系统根据客户端的位置选择多媒体内容或最近的服务器等等。因此, 越来越多的地理定位服务也随之出现, 如 Maxmind, IP2Location, Quova 等皆提供收费的定位服务^[1]。

本文的内容属于“211 工程三期公共服务体系建设项目“中国教育和科研计算机网主干网和重点学科信

息服务体系升级扩容工程”, 为网络行为监控和管理系统提供 IP 地址拥有者信息查询, 以满足网络服务质量管理和网络有害行为分析追踪的需要; 同时, 为网络行为监控和管理系统中活跃 IP 地址与安全状态 IP 地址提供实际使用位置查询, 支持网管系统对流量数据的定位分析, 安全监测, 威胁源分析等。系统预计将安装在 3 个点, 分别为 CERNET 主节点中的北京大学, 上海交通大学和东南理工大学, 由统一的维护端——总控系统进行 IP 归属表维护。

综上所述, 需要获得 IP 地址两个部分的信息: (1) 管理者信息, 包括 IP 地址所属的网络运营商, 自治域号以及所属单位信息; (2) 实际使用位置信息, 精度要求到城市。对于 IP 地址管理归属信息, 可以通过查询 5 个地域性的 IP 地址管理机构 RIR, 即 ARIN (北美地区), LACNIC (拉丁美洲), RIPENCC (欧洲地区), APNIC (亚太地区), AFRINIC (非洲地区), 提供的 whois 数据库得到。本文着重介绍基于 Traceroute 信息的 IP 地址地理定位方法模型, 针对那些不能直接

1 作者简介: 倪晶晶, (1989-), 女, 硕士研究生, E-mail: jjni@njnet.edu.cn; 龚俭, (1957-), 男, 教授、博士生导师, E-mail: jgong@njnet.edu.cn.



2 Trace 中连续多跳 IP 地址关联分析

2.1 原理介绍

子序列 $[IP_1, IP_2, \dots, IP_n]$ ($n \geq 3$) 为一条 Trace 下相同 AS 的连续多跳 IP 地址。

如果该子序列的首尾 IP 地址 (IP_1, IP_n) 地理位置信息相同, 那么该子序列中所有 IP 地址均有相同地理位置信息^[8]。如图 2, $[IP_{A1}, IP_{A2}, IP_{A3} \dots IP_{An}]$ 为一个 IP 地址组, IP_{A1} 、 IP_{An} 使用位置均为城市 A, 则地址组中其它 IP 地址使用位置也为城市 A。

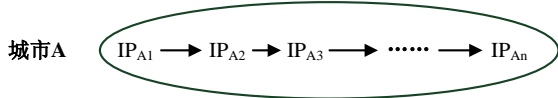


图 2 连续多跳属于相同 AS 的 IP 地址位置相同

假设存在一条 trace 下相同 AS 的连续多跳 IP 中首尾 IP 在同一城市, 而中间 IP 在另一城市, 如图 3, 即 $[IP_{A1}, IP_{B2}, IP_{A3} \dots IP_{An}]$ 为一个 IP 地址组。这表明一个自治域系统下, 对同一城市的两个 IP (IP_{A1}, IP_{An}) 路由选择时, IP_{A1} 放弃了最优路由 IP_{A2} , 而选择了与自己不在同一城市的 IP_{B2} , 不符合路由选择策略, 假设不成立。

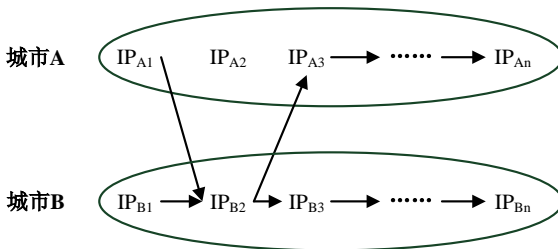


图 3 连续多跳属于相同 AS 的 IP 地址位置不相同

取 CAIDA 中实际测量的一条源宿 IP 地址分别为 84.88.81.122、178.134.97.85 的 Trace 路径进行多跳关联分析, 图 4 描述了路径的详细信息。

```
T 84.88.81.122 178.134.97.85 7 1249 1293660497 E 184.111.14.55 S 0 C 84.88.81.121,0.316,1 84.88.19.149,0.949,1 130.206.202.29,0.377,1 130.206.250.25,15.103,1 193.149.1.36,15.807,1 129.250.6.102,55.433,1 129.250.6.53,39.010,1 129.250.4.85,39.015,1 129.250.2.227,77.084,1 129.250.2.17,51.037,1 213.198.82.166,51.917,1 188.128.105.22,133.370,1 87.226.222.46,149.734,1 85.132.80.29,157.18 178.134.240.2,156.872,1
```

图 4 单个 Trace 多跳分析实例

整个 Trace 中有 16 跳 IP 地址信息, 分别获取其所对应的 AS 号, 结果如图 5。自治域号为 AS2914 时的子序列 $[129.250.6.102, 129.250.6.53, \dots,$

$129.250.2.17]$, 满足相同 AS 下, 至少有连续 3 跳 IP 地址。

84.88.81.121	AS13041
84.88.19.149	AS13041
130.206.202.29	AS766
130.206.250.25	AS766
193.149.1.36	AS6895
129.250.6.102	AS2914
129.250.6.53	AS2914
129.250.4.85	AS2914
129.250.2.227	AS2914
129.250.2.17	AS2914
213.198.82.166	AS2914
188.128.105.22	AS12389
87.226.222.46	AS12389
85.132.80.29	AS29049
178.134.240.2	AS35805
178.134.97.85	AS35805

相同AS子序列

图 5 单个 Trace 扩展 AS 图

通过 Maxmind 寻找子序列首尾 IP 地址 ($129.250.6.102, 129.250.2.17$) 的地理位置信息, 查询结果如图 6。因为位置信息相同, 根据关联规则, 可以推测子序列中所有 IP 地址均有相同地理位置信息。

GeoIP城市/ISP/机构查询结果						
IP 地址	国家	代码	地点	邮政	坐标	网络服务
129.250.6.102	US		Englewood, Colorado, United States, North America	80111	39.6237, -104.8738	NTT America
129.250.2.17	US		Englewood, Colorado, United States, North America	80111	39.6237, -104.8738	NTT America

图 6 子序列首尾 IP 位置信息

为了验证推测的正确性, 对子序列中的中间 IP 进行地理位置查询, 结果如图 7。可以发现, 中间 IP 地址地理位置信息均相同, 且和首尾 IP 在同一城市, 结果与建立的关联规则一致, 即如果存在子序列中所有 IP 属于同一 AS, 且首尾 IP 位置相同, 那么子序列中所有 IP 位置均相同。

IP 地址	国家	代码	地点	邮政	坐标	网络服务
129.250.6.53	US		Englewood, Colorado, United States, North America	80111	39.6237, -104.8738	NTT America
129.250.4.85	US		Englewood, Colorado, United States, North America	80111	39.6237, -104.8738	NTT America
129.250.2.227	US		Englewood, Colorado, United States, North America	80111	39.6237, -104.8738	NTT America

图 7 子序列中间 IP 地址位置信息



2.2 关联分析

多跳关联分析主要分为 3 个部分，结构如图 8:

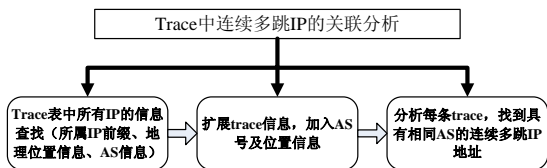


图 8 多跳 IP 地址关联分析结构图

- 1) 首先将 trace 路径文件中所有的 IP 整理出来(去掉重复), 并且按 IP 地址的大小排序, 并在中间数据库中找到其对应的地理位置信息, 在 BGP 路由信息中获取 IP 所属的 AS 信息^[9]。
- 2) 根据 1)中 IP 信息, 扩展 Trace 信息, 加入 IP 地址的 AS 信息以及地理位置信息, 得到格式如图 9 的文件。

```

Trace1
IP1 | country | region | city | AS1 |
IP2 | country | region | city | AS2 |
.....
IPn | country | region | city | ASn |

Trace2
.....
TraceN(最后一条 trace)
  
```

图 9 Trace 扩展图

- 3) 分析扩展后的 Trace1, Trace2,TraceN, 找到每条 trace 中具有相同 AS 号的连续多跳 IP(至少三跳), 写入文件, 供关联分析。存储格式如图 10:

```

IP1 | country | region | city | AS1 |
IP2 | country | region | city | AS1 |
.....
IPn | country | region | city | AS1 |

空行(以此为不同AS的分界)

IP1 | country | region | city | AS2 |
IP2 | country | region | city | AS2 |
.....
IPn | country | region | city | AS2 |
.....
  
```

图 10 关联 IP 组存储图

- 4) 至此, 已经找到 Trace 文件中所有适用此关联规则的 IP 地址组, 验证已存在的位置信息, 并填充未定位 IP 的位置信息。

3 相同接入路由的 IP 地址关联分析

3.1 原理介绍

设三元组 $[IP_{n-1}, IP_n, T]$ 表示一条 Trace 路径中 [倒数第二跳 IP, 目的 IP, 往返时延 T], IP_{n-1} 为 IP_n 的接入路由, $T/2$ 为 IP_{n-1} 到 IP_n 的单向时延。

现有三元组 $[IP_{1(n-1)}, IP_{1n}, T_1]$ 、 $[IP_{2(n-1)}, IP_{2n}, T_2]$, 若满足公式 1 中的条件, 表明 IP_{1n} 、 IP_{2n} 有相同接入路由, 且到该接入路由的时延信息 T 具有相似性: 水平方向上, 往返时延均小于 10ms, 即单向时延均小于 5ms; 垂直方向上, 单向时延差小于 2.5ms。

参数 10ms 与 2.5ms 均为学习经验值, 通过对实际运行时结果的学习, 调整参数值大小。例如, 几乎所有相同路由下, 往返时延约束为 20ms 的那些 IP 地址经验证都不在同一地理位置, 就会想到需要减小往返时延约束值, 使得约束更强, 地理位置在一起的可能性更大。

$$\left\{ \begin{array}{l} IP_{1(n-1)} = IP_{2(n-1)} \\ IP_{1n} \neq IP_{2n} \\ T_1, T_2 < 10 \text{ ms}, \text{ and } \left| \frac{T_1}{2} - \frac{T_2}{2} \right| < 2.5 \text{ ms} \end{array} \right\} \quad (1)$$

定义那些具有相同接入路由, 且到接入路由的时延两两相似的那些 IP 地址为一个 IP 地址组, 认为该 IP 地址组内的 IP 地址在同一地理位置。

根据公式 1, 分析 Traceroute 路径信息, 得到具有相同接入路由, 且满足时延约束的 IP 地址信息。具体实例如图 11:

59.51.186.146		
59.51.199.188	2.560000	
59.51.210.129	1.408000	
59.51.230.205	5.991000	
59.106.245.182		
49.212.161.150	0.044000	
49.212.163.84	0.542000	
49.212.169.240	0.866000	
49.212.182.55	0.601000	
49.212.189.138	0.066000	

图 11 相同接入路由实例

下面简单说明第一组数据, 第二组数据情况相同, 分析类似:

- 1) 59.51.199.188, 59.51.210.129, 59.51.230.205



均接在 59.51.186.146 这个路由 IP 下。

- 2) 往返时延 2.56ms、1.408ms、5.991ms，均小于 10ms。
- 3) 单向往返时延差最大值 $(5.991-1.408)/2=2.2915$ 小于 2.5ms。
- 4) 定义具有相同接入路由，且往返时延相似的 3 个 IP 地址：59.51.199.188、59.51.210.129、59.51.230.205 为一个 IP 地址组，根据关联规则，可以推测组内 IP 地址均在同一地理位置。

通过查询商业数据库 Maxmind，来验证推测的正确性。查询结果如图 12，可知满足时延限制，具有相同接入路由的 3 个 IP 地址使用位置的确相同，结果与定义的关联分析规则一致。

IP 地址	国家代码	地区	邮政编码	坐标	网络服务提供商	机构
59.51.199.188	CN	Guiyang, Guizhou Sheng, People's Republic of China	26.5833, 106.7167	26.5833, 106.7167	Chinanet Guizhou Province Network	Chinanet Guizhou Province Network
59.51.210.129	CN	Guiyang, Guizhou Sheng, People's Republic of China	26.5833, 106.7167	26.5833, 106.7167	Chinanet Guizhou Province Network	Chinanet Guizhou Province Network
59.51.230.205	CN	Guiyang, Guizhou Sheng, People's Republic of China	26.5833, 106.7167	26.5833, 106.7167	Chinanet Guizhou Province Network	Chinanet Guizhou Province Network

图 12 Maxmind 查询结果

3.2 关联分析

对 CAIDA 提供的 Traceroute 信息进行相同路由下的分析，结构如图 13：

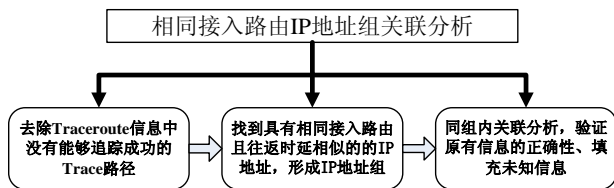


图 13 相同接入路由关联分析结构图

- 1) 首先去除 CAIDA 文件中路由追踪没有成功的那些 Trace 路径。
- 2) 分析只包含追踪成功的 Trace 路径文件，找到那些具有相同接入路由且往返时延相似的 IP 地址，形成 IP 地址组，写入文件，实例如图 11。
- 3) 对同组中的 IP 地址进行关联分析，验证已存在的位置信息，并填充未定位 IP 的位置信息。例如，

分析 59.51.199.188，59.51.210.129，59.51.230.205 这 3 个 IP 地址的位置相同，如果 59.51.210.129 的使用位置未知，就可以用 59.51.199.188 或 59.51.230.205 的位置信息关联。

4 总结与展望

本文介绍了一种 IP 地址地理定位模型：基于 Traceroute 信息的关联分析规则，利用现有的地理定位信息推测库中未知的定位信息。文中共提出了两种关联分析规则，分别为：同一 Trace 中连续多跳 IP 地址关联分析，相同接入路由且满足一定时延限制的 IP 地址关联分析，并对各个规则的原理与实现做了详细的介绍。

下一步工作就是运用此地理定位模型，为网络行为监控和管理系统中活跃 IP 地址以及安全状态 IP 地址提供实际使用位置，满足其跟踪定位等服务。同时，还需要对各个关联规则进行总结、效果分析，讨论关联规则的有效性，在研究中不断改进现有规则，发现其它新的规则，为地理定位提供更有效的方法模型。

参考文献

- [1] Venkata N. Padmanabhan, Lakshminarayanan Subramanian, "An Investigation of Geographic Mapping Techniques for Internet Hosts", SIGCOMM2001.
- [2] Maxmind 数据库: www.maxmind.com.
- [3] IP2location 数据库: www.IP2location.com.
- [4] QQ 纯真数据库: www.cz88.net.
- [5] CERNET NIC 数据库: www.nic.edu.cn.
- [6] W. Richard Stevens, TCP/IP 详解, 机械工业出版社, 2011.1 第 33 次印刷.
- [7] CAIDA: www.caida.org.
- [8] Yu Jiang, Binxing Fang, Mingzeng Hu, Xiang Cui, "Techniques for Determining the Geographic Location of IP Addresses in ISP Topology Measurement", J.Comput.Sci. & Technol: Vol.20:N0.5, Sept. 2005.
- [9] BGP: www.routeviews.org.