

# 基于 twitter 数据的网络关系提取及分析实验报告

## 一、实验目的

掌握从社交网络数据抽取用户关系并进行数据分析的方法。

## 二、实验环境

Windows 10 操作系统、python3.6、Gephi 0.9.2

## 三、实验内容

### 1.地震前后网络构建

根据附件 1 和附件 2 给出的推文数据，提取出用户发文时间、用户昵称、发文内容等关键字段，并根据这些信息构建地震前后的关系网络。

### 2.对构建好的网络进行分析

为了对比地震对用户关系带来的影响，分别对取得的英文使用者网络和日文使用者网络进行累积度分布变化统计、单独节点层面上的度变化统计并做对比分析。

### 3.网络可视化

使用复杂网络分析软件 Gephi 可以方便的对网络节点进行可视化。

## 四、实验步骤

### 1. 网络构建

#### (1) 分析 json 文件数据格式

```
response:
  window: "custom"
  page: 1
  total: 113
  perpage: 100
  last_offset: 102
  hidden: 1
  list:
    0:
      traceback_permalink: "http://twitter.com/there...status/43542423091150848"
      traceback_author_url: "http://twitter.com/therealjuliann"
      content: "Never be ashamed of who ... who you pretend to be."
      traceback_date: 1299216297
      topsy_author_img: "http://a0.twimg.com/prof...765va9rbw8e0t_normal.png"
      hits: 1
      topsy_traceback_url: "http://topsy.com/twitter...1150848?utm_source=otter"
      firstpost_date: 1299216309
      url: "http://twitter.com/there...status/43542423091150848"
      traceback_author_nick: "therealjuliann"
      highlight: "Never <span class=\"high-ight-term\">be</span>."
      topsy_author_url: "http://topsy.com/twitter...juliann?utm_source=otter"
      traceback_author_name: "Julian Alexander."
      mytype: "tweet"
      score: 10.4014
      traceback_total: 147
      title: "Never be ashamed of who ... who you pretend to be."
```

分析 json 文件中所需字段所处位置结构。目标字段分别为“response”字段下的“content”、“traceback\_date”、“traceback\_author\_nick”、“traceback\_author\_name”字段。

其中 nick 在 twitter 中唯一且只能由英文构成，name 可用多种语言。

#### (2) 提取关键字段（python3.6 代码示例）

```
1. fp = open(file_name, encoding='utf-8')
2. temp = json.loads(fp.read())
```

```

3.     page = temp['response']['list']
4.     key_record = defaultdict(int)
5.     for post in page:
6.         nick = post['trackback_author_nick']
7.         name = post['trackback_author_name']
8.         date = post['trackback_date']
9.         content = post['content']
10.    fp.close()

```

分别提取 json 文件中英文使用者和日文使用者的发帖信息。

### (3) 数据清洗

将数据写入 Excel 文档或数据库中，并对数据进行去重处理。

按照地震时间（2011 年 3 月 11 日 13:46）将数据分段。时间戳可以转化为年月日格式。

为方便后期网络对比实验，只考虑地震前后共同用户变化。所以可以在这一步去除非公共用户信息，也可以在构建网络时去重非公共用户。

分别统计地震前后英文网络和日文网络的用户数量、用户发文数量、@关系数量。

### (4) 网络构建

根据用户转发、评论关系构建网络。当用户 A 在其 content 字段中@用户 B，我们认为用户 A 与用户 B 之间存在联系。

分别对日语使用者地震前、日语使用者地震后、英文使用者地震前、英文使用者地震后各建网络。

NetworkX 是一个用 Python 语言开发的图论与复杂网络建模工具，内置了常用的图与复杂网络分析算法，可以方便的进行复杂网络数据分析、仿真建模等工作。networkx 支持创建简单无向图、有向图和多重图（multigraph）；内置许多标准的图论算法，节点可为任意数据；支持任意的边值维度，功能丰富，简单易用。

（网络构建部分代码 python3.6 版本示例：）

```

1. import networkx as nx
2. G = nx.Graph()
3. # 创建一个权重表，字典形式存储
4. w_nodes = defaultdict(defaultdict)
5. # 然后按照时间一个一个的构建
6. for row in range(1, len(data)):
7.     record = data[row]
8.     # check 网络中是不是已经有了这个节点
9.     c_nick = record['nick']
10.    if c_nick in G:
11.        pass
12.    # 如果没有 生成这个点
13.    else:
14.        G.add_node(c_nick)
15.    # check 这个帖子中 有没有 at 别人
16.    content = str(record['content'])

```

```

17.     if content:
18.         while find_string(content, '@'):
19.             # check 安特的这个人在不在网络
20.             ind = content.index('@')
21.             n_nick = ''
22.             while content[ind] != ' ' and content[ind] != ':' and ind
                < len(content) - 1:
23.                 ind = ind + 1
24.                 if content[ind] != ':':
25.                     n_nick += content[ind]
26.                 if n_nick in G:
27.                     pass
28.                 # 如果不在生成这个点
29.                 else:
30.                     G.add_node(n_nick)
31.                 # check 这两个人之间有没有边
32.                 # #如果有
33.                 if n_nick in neighbor_nodes(G,c_nick):
34.                     # 权重+1
35.                     G.add_edge(c_nick, n_nick)
36.                     if c_nick in w_nodes:
37.                         if n_nick in w_nodes[c_nick]:
38.                             w_nodes[c_nick][n_nick] += 1;
39.                 # 如果没有
40.                 else:
41.                     # 在这个节点和新的节点之间生成一条边
42.                     G.add_edge(c_nick, n_nick)
43.                     # 而且权重设为 1
44.                     w_nodes[c_nick][n_nick] = 1;
45.                 content = content[ind+1:]

```

(5) 去重非公共节点

```

1. def drop_diff_point(G1,G2):
2.     for node in G1.nodes():
3.         if not(node in G2):
4.             G1.remove_node(node)
5.     for node in G2.nodes():
6.         if not(node in G1):
7.             G2.remove_node(node)
8.     return G1,G2

```

(6) 将构建好的网络存入文件中

```

nx.write_gml(G,"Filename.gml")

```

## 2. 网络分析

(1) 使用 python 中 networkx 库对网络进行基本特性（平均度、最大联通片、群聚系数）分析

```

1. #输出网络平均度
2. def average_deg(G,name):
3.     d = nx.degree(G)
4.     print(name + "平均度")
5.     # print(d)
6.     print(np.array(list(d.values())).mean())

```

```

1. #输出最大联通片
2. def largest_com(G,name):
3.     largest_components=max(nx.connected_components(G),key=len)
4.     print(name + "的最大联通成分的大小为",end = '')
5.     print(len(largest_components))

```

```

1. #输出平均群聚系数
2. def average_clu(G,name):
3.     c = nx.average_clustering(G)
4.     print(name + "的平均群聚系数:",end = '')
5.     print(c)

```

## (2) 个人层面度分析

```

1. def individualdegree(G1,G2,name):
2.     nodes1 = G1.nodes()
3.     nodes2 = G2.nodes()
4.     degree1 = []
5.     degree2 = []
6.     for node in nodes1:
7.         if node in nodes2:
8.             degree1.append(G1.degree(node))
9.             degree2.append(G2.degree(node))
10.    # plt.scatter(degree1,degree2)#在双对坐标轴上绘制度分布曲线
11.    # plt.subplot(121)
12.    plt.title(name)
13.    plt.xlabel('before')
14.    plt.ylabel('after')
15.    plt.loglog(degree1,degree2,'o')#在双对坐标轴上绘制度分布曲线
16.    plt.show()

```

## (3) 累积度分析

```

1. #接下来做累积度分布
2. def cumlutive_degree_distribution(G):
3.     degree = []
4.     list = G.degree()
5.     for each_node in list:
6.         degree.append(list[each_node])
7.     xs = degree
8.     distKeys = range(min(xs), max(xs) + 1)

```

```

9.     pdf = dict([(k, 0) for k in distKeys])
10.    for x in xs:
11.        pdf[x] += 1
12.    pdf_temp=pdf
13.    scope = range(min(pdf),max(pdf)+1)
14.    for degree in scope:
15.        k=degree+1
16.        while k<=max(pdf):
17.            pdf[degree]+=pdf_temp[k]
18.            k+=1
19.    return pdf
20. #根据（图，名称，度列表，返回一个度分布图）
21. def draw_degree_chart(G,name,distribution):
22.     # degree=nx.degree_histogram(G)#返回图中所有节点的度分布序列
23.     degree = distribution
24.     # print(degree)
25.     y = np.array(list(degree.values()))
26.     # print(y)
27.     # y=[z/float(sum(degree))for z in degree]#将频次转化为频率，利用列表内
    涵
28.     y = y/y[0]#将频次转化为频率，利用列表内涵
29.     x=range(len(degree))#生成X轴序列，从1到最大度
30.     # y = degree
31.     if 'b' in name:
32.         color = 'r^'
33.         # 'r--':红色的需要;'bs':蓝色方块;'g^':绿色三角
34.         line = plt.loglog(x, y, color) # 在双对坐标轴上绘制度分布曲线
35.         plt.legend(line, ("before"))
36.     else:
37.         color = 'g^'
38.         # 'r--':红色的需要;'bs':蓝色方块;'g^':绿色三角 o 圆
39.         line = plt.loglog(x, y, color) # 在双对坐标轴上绘制度分布曲线
40.         plt.legend(line, ("a"))
41.     plt.title( name)
42.     plt.show()#显示图表

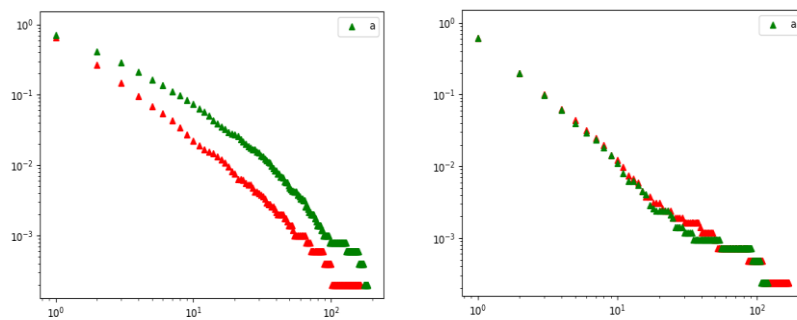
```

## 五、实验结果

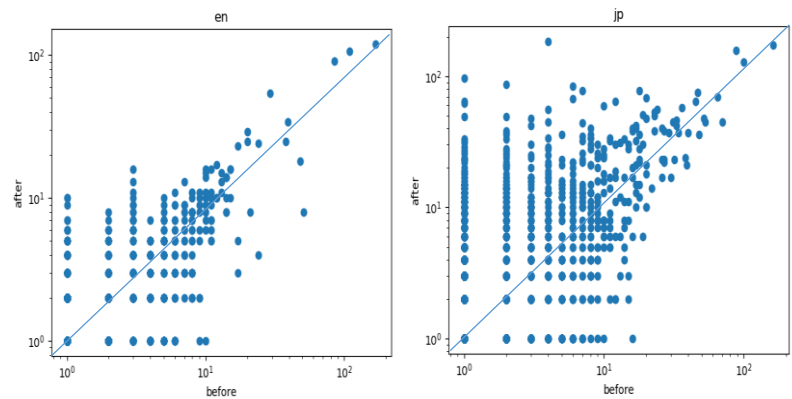
### 1. 网络基本特性信息统计对比

Tabel 1 Data description and basic network characteristics								
data	period	raw data			Extracted network			
		#users	#tweets	#links	#nodes	D	s	c
TP-JP	before	12863	40801	12956	5077	1.689	2591	0.000107
	after	32032	105668	52139	5077	3.214	3440	0.000208
TP-EN	before	17537	44993	15671	4200	1.274	1630	0
	after	17841	58429	15951	4200	1.235	1626	0.000847
D:average degree;s:size of giant connected component or giant strongly connected network;c:clustering coefficient.								

## 2. 累积度对比



## 3. 个人层面度对比



## 六、使用 Gephi 软件进行网络可视化和分析

### 1.打开之前保存的 gml 文件。



## 2.网络信息统计

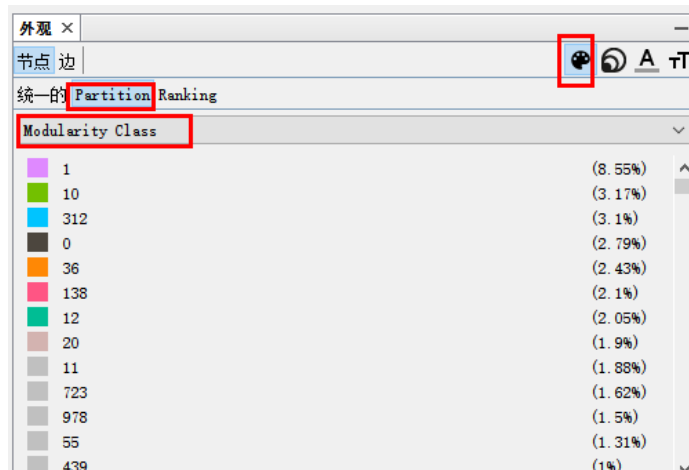


在统计板块可以对网络进行基本信息统计。

## 3.社区划分与渲染

另外，Gephi 内置快速模块化算法，可以根据网络结构对网络进行简单的社区划分。点击上图网络概述板块中的“模块化”，运行 **community detection** 算法。

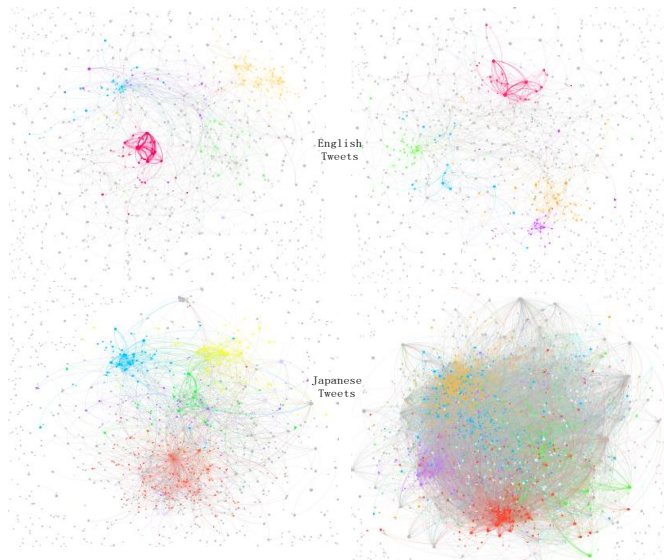
社区划分后，左上角外观部分可以对检测出的社区进行颜色渲染操作。



在左下角有布局选项，选择一个算法，点运行展开(lay out) 。

布局 ×	
OpenOrd	
<div>📄 运行</div>	
Stages	
液体 (%)	25
扩张 (%)	25
冷却 (%)	25
紧缩 (%)	10
Simmer (%)	15
OpenOrd	
Edge Cut	0.8
多线程	3
多次迭代	750
固定时间	0.2
随机种子	-1505014327479844142

结果：



参考视频教程：<https://www.udemy.com/gephi/>

参考案例：<https://blog.csdn.net/laoyang360/article/details/53616720>