# Large Graph Mining - Patterns, Tools and Cascade Analysis
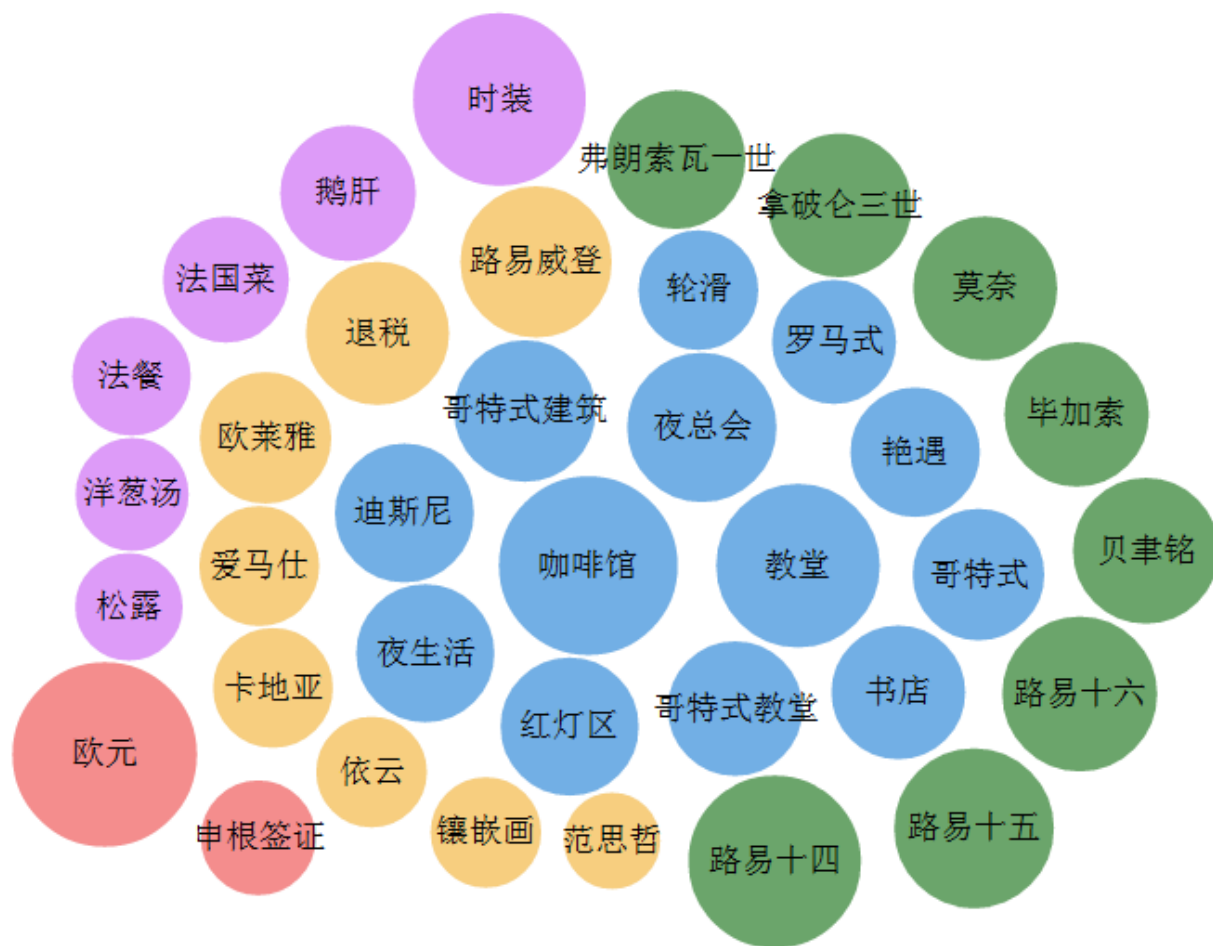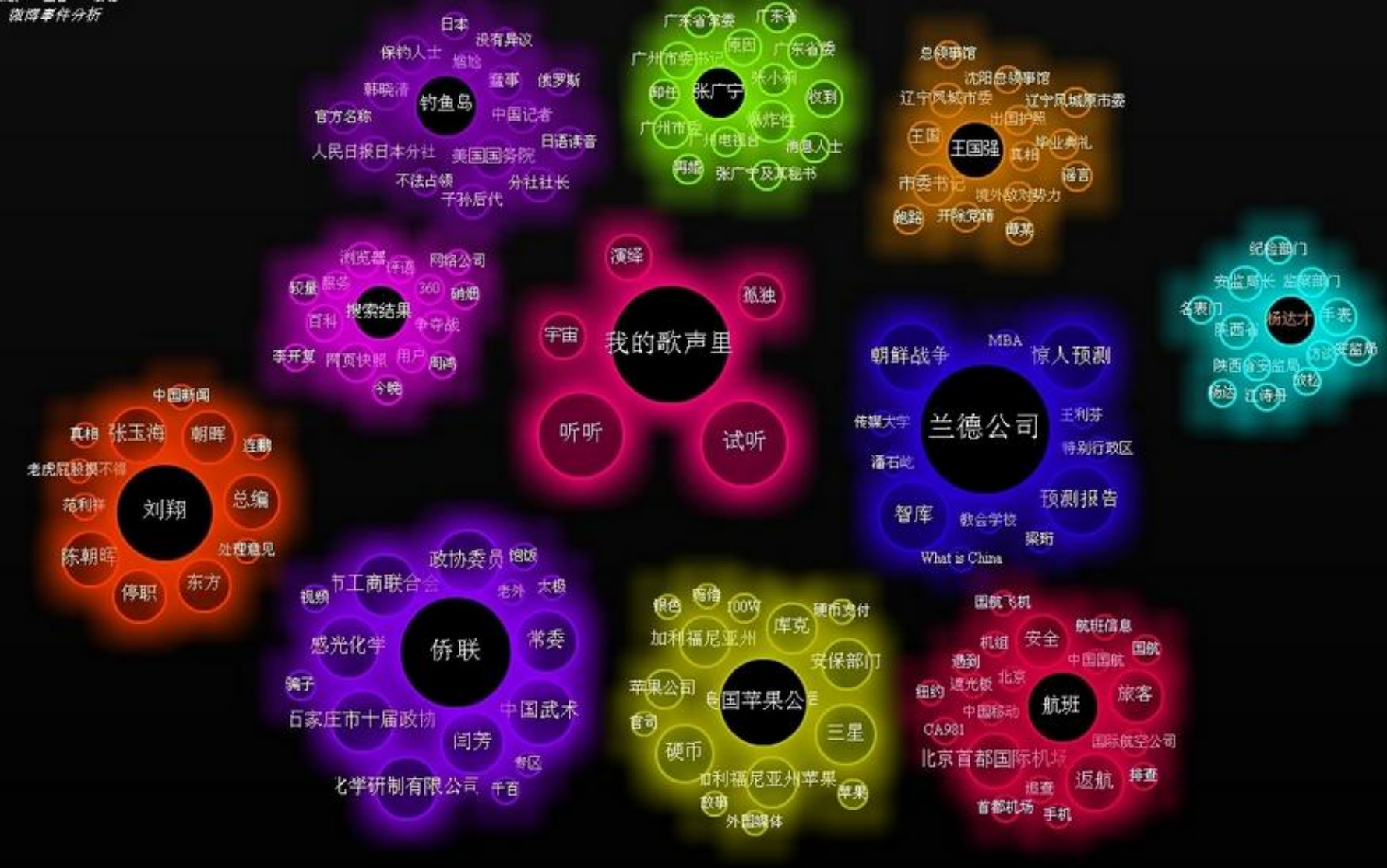
# outline

- simple tag cloud
- K-Means
- some python tools

# simple tag cloud

$$freq_{norm} = \frac{word_{freq} - min}{\max - min}$$

$$freq_{norm} = \frac{(word_{freq} - min)(upper_{bound} - lower_{bound})}{\max - min} + lower_{bound}$$

时装

弗朗索瓦一世

拿破仑三世

鹅肝

路易威登

莫奈

法国菜

轮滑

罗马式

法餐

退税

哥特式建筑

夜总会

毕加索

欧莱雅

艳遇

洋葱汤

迪斯尼

咖啡馆

教堂

哥特式

贝聿铭

松露

爱马仕

欧元

夜生活

书店

路易十六

卡地亚

红灯区

哥特式教堂

依云

路易十五

申根签证

镶嵌画

范思哲

路易十四

微博事件分析

# K-Means Algorithm
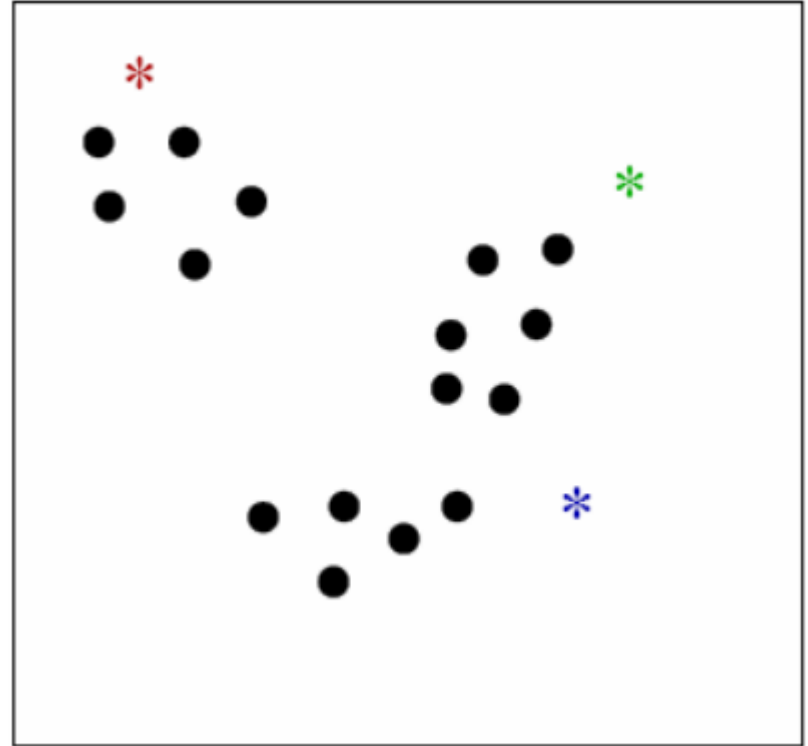
K = # of clusters (given);

One "mean" per cluster

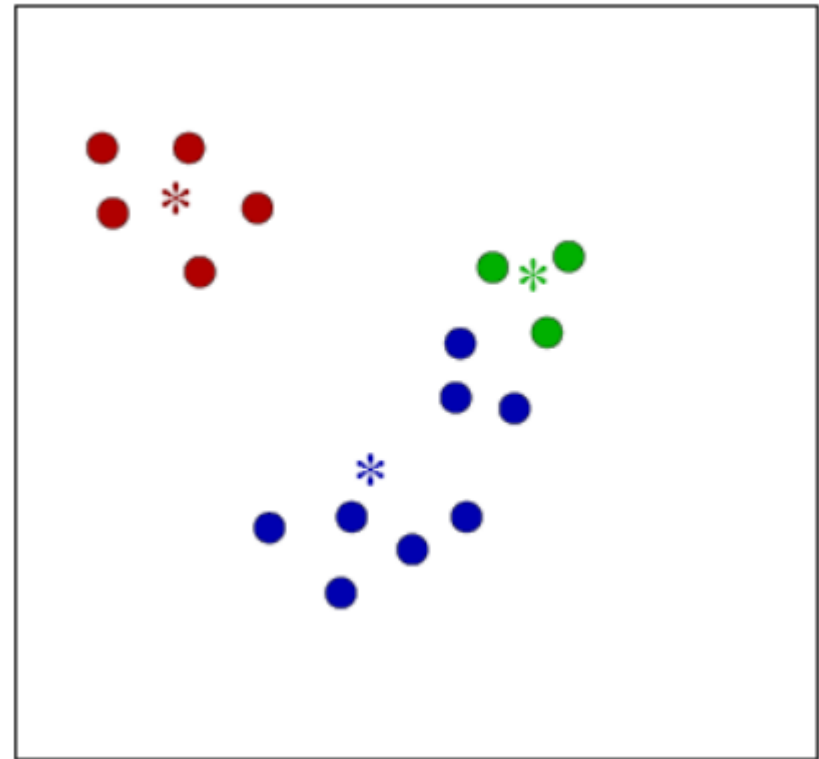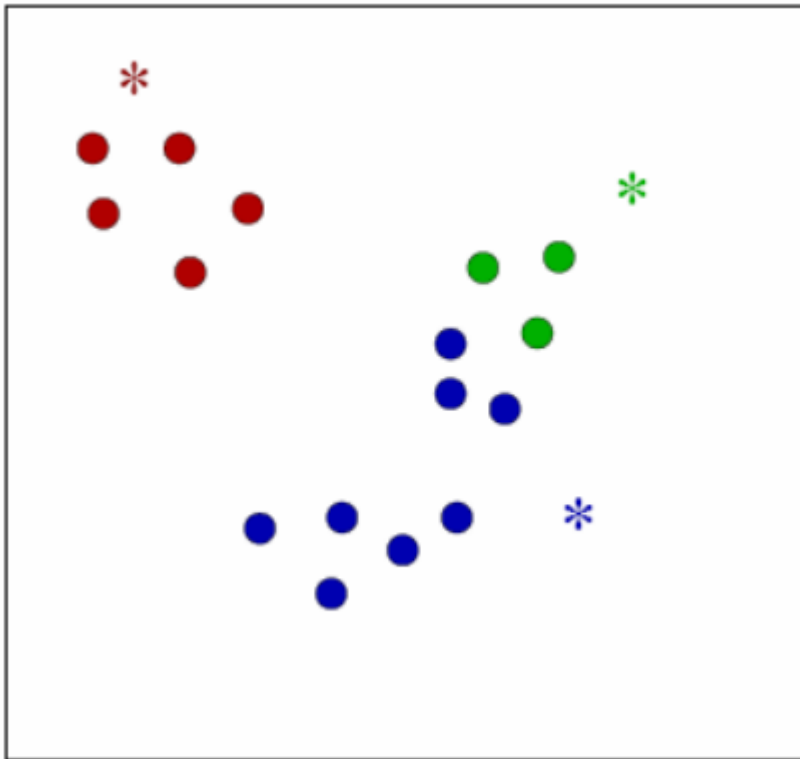• Initialize means

(e.g. by picking k samples at random)

• Iterate:

(1) assign each point to nearest mean
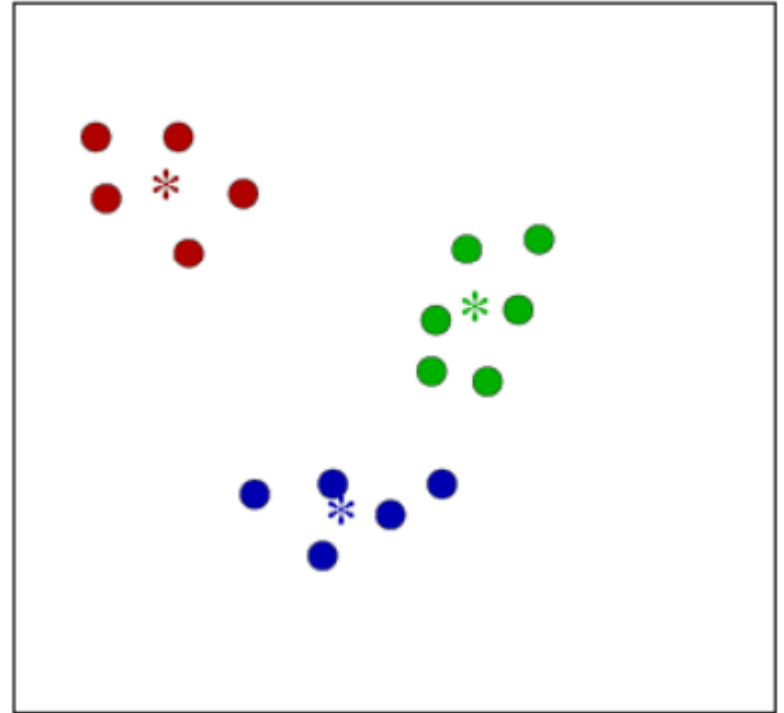
(2) move "mean" to center of its cluster

# Means Update

# K-Means Algorithm

- Complexity:

O(kn # of iterations)

- The object function is:



$$J = \min_{\{\mu_1, \cdots, \mu_k\}} \sum_{h=1} \sum_{x \in X_h} \|x - \mu_h\|^2$$

# K-Means Algorithm

- **Initialize** $\mu_1, \cdots, \mu_k$
- **do** classify n samples according to nearist $\mu_h$
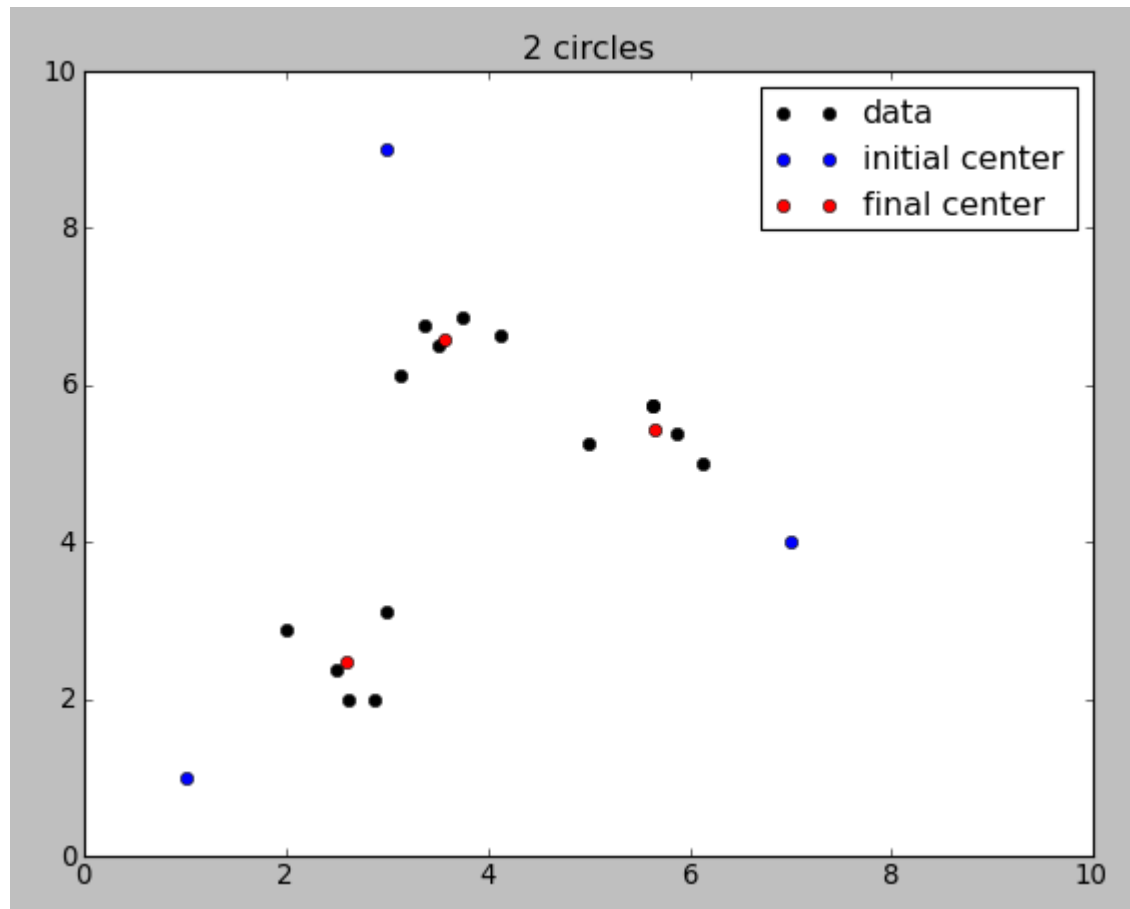
  recompute $\mu_h$
- **until** no change in $\mu_h$

# K-Means Algorithm

$$\frac{\partial}{\partial \mu_h} J = \frac{\partial}{\partial \mu_h} \sum_{i=1}^{k} \sum_{x \in X_h} \|x - \mu_h\|^2 = \sum_{i=1}^{k} \sum_{x \in X_h} \frac{\partial}{\partial \mu_h} \|x - \mu_h\|^2 = 0$$

$$\sum_{x \in X_h} 2(x - \mu_h) = 0 \rightarrow \mu_h = \frac{1}{n_h} \sum_{x \in X_h} x_h$$

# K-Means in Python

Simple k-means code using numpy and matplotlib

# Find related words

- 9.10 weibo corpus sample

Example:

- # 钓鱼岛 是 中国 的 #~ 我是 热血 爱国 好青年
- 这样的 第一个 教师节 也算是 难忘 了 。

100 cases

50：钓鱼岛是中国的

50：教师节

# Find related words

- Use K-Means to cluster into two classed

Feature size : 1338

Feature ∈{0.0, 1.0}

Two classes

$$\text{similarity} = \cos(\theta) = \frac{A \cdot B}{\|A\|\|B\|} = \frac{\sum\limits_{i=1}^{n} A_i \times B_i}{\sqrt{\sum\limits_{i=1}^{n} (A_i)^2} \times \sqrt{\sum\limits_{i=1}^{n} (B_i)^2}}$$

# Find related words

Pycluster used

labels, error, nfound = Pycluster.kcluster(weibo, 2)
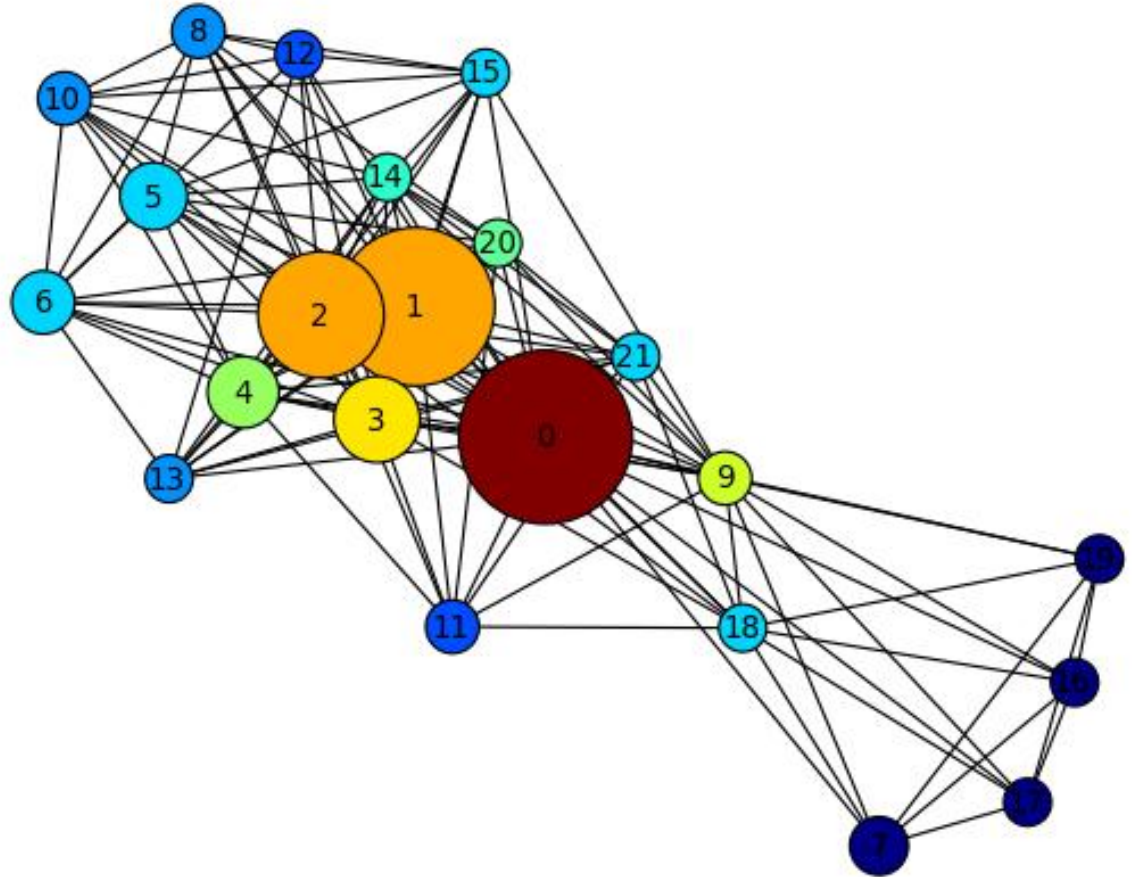
Average accuracy: 0.933

Bad examples:

教师节 向 德艺双馨 的 人民 艺术家 ** 致以 最真诚的 祝福**
钓鱼岛 是 中国 的 ， ** 是 世界 的

# Draw the result

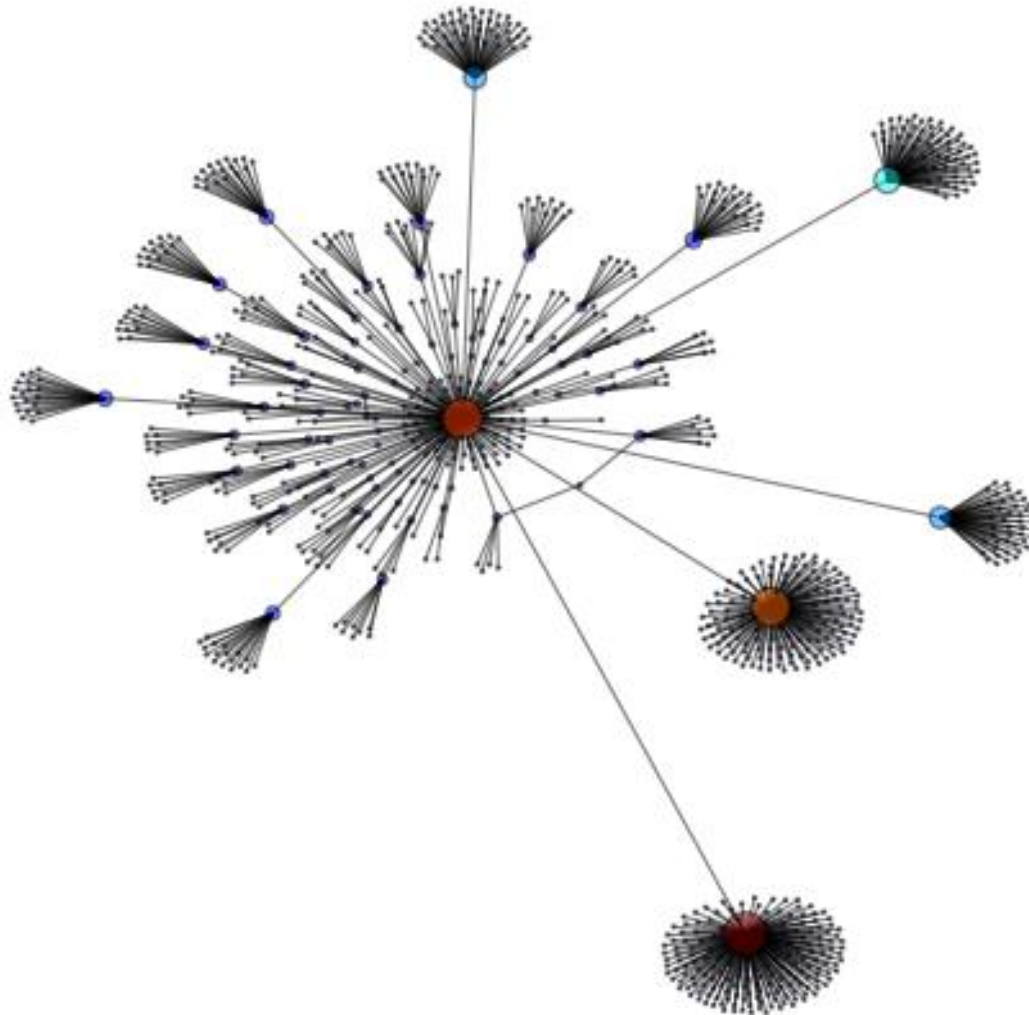0 钓鱼岛
1 中国
2 日本
3 领土
4 政府
5 我们
6 就是
7 有本事
8 日本人
9 小日本
10 明天

11 国有化
12 谴责
13 问题
14 退让
15 起来
16 垃圾
17 滚蛋
18 历史
19 破烂
20 固有
21 主权

# networkx

- import networkx as nx
- G = nx.Graph()
- G.add_node(*)
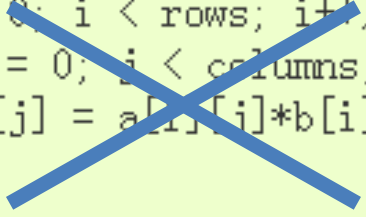- G.add_edge(*)
- nx.draw(G)

# use networkx

# Numpy and matplotlib

- NumPy is the fundamental package for scientific computing in Python.

- C = A*B



```
for (i = 0; i < rows; i++): {
    for (j = 0; j < columns; j++): {
        c[i][j] = a[i][j]*b[i][j];
    }
}
```

# Install them!

1、安装ipython
sudo apt-get install ipython

2、安装matplotlib
sudo apt-get install python-matplotlib

3、启动绘图环境
ipython -pylab

4、安装pycluster
sudo pip install pycluster

5、安装networkx
sudo easy_install networkx