15CS065 Users' Public Privacy Mining Software

LIU Jiacheng (53057374) Supervisor: Dr. LI Shuaicheng

Department of Computer Science City University of Hong Kong



PARTI Introduction

User Needs

How can I know an online user which is totally strange to me?

How can I know my friend more comprehensively through his/her online activity?

How can I know this guy/girl who is actively interacting with me?

I need a user-study application and it must be:

1. Automatic

Least human labor

2. Objective

Voice of data

3. Easy to perceive

Better with graphs

4. Comprehensive

Various perspectives

Project Overview

This project is

A user study application, focusing on Sina's Weibo.com, that satisfies the described user requirements:

1. Automation Only seed information needed

2. Objectiveness Models and statistics

3. Ease to perceive 14 graphs, 6 tables

4. Comprehensiveness 7 different aspects

general information, friends interaction, interested topics and timeline, etc.

Statistics

Line Count:

• ~25000 lines

Component Count:

• Live modules: 9

• Graphs: 14

• Tables: 6

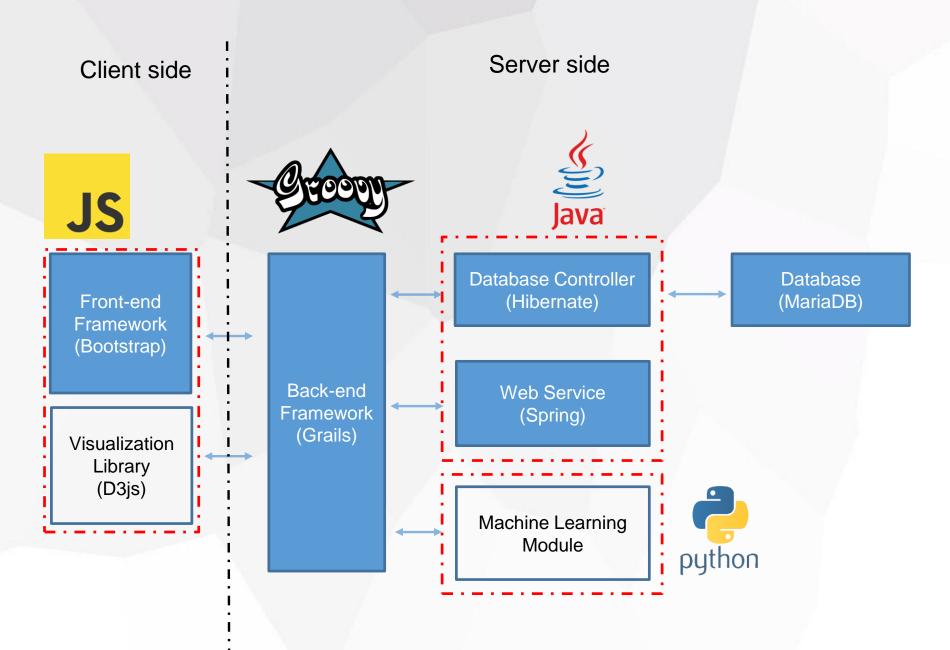
Outline

- I. Introduction
- II. System Design
- III. A Use Case
- IV. Algorithms Applied
- V. Conclusion



PARTII System Design

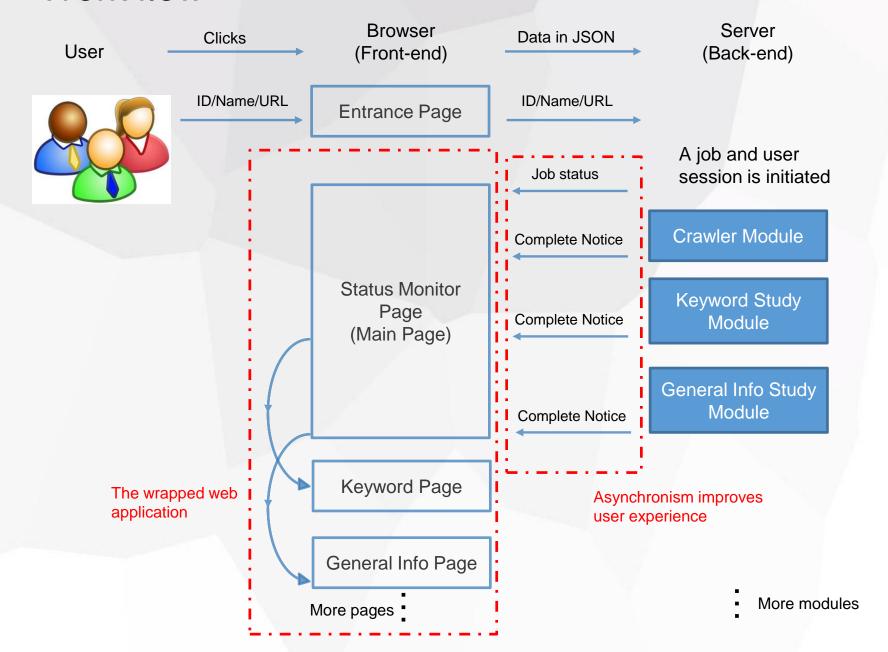
Structural Overview



Advantages

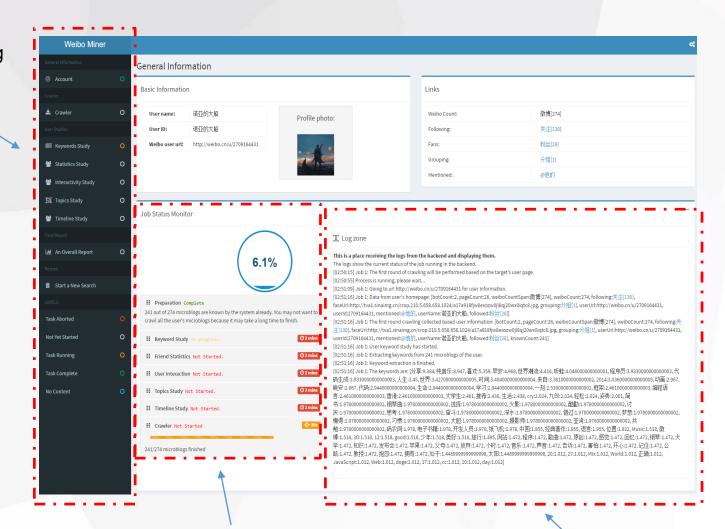
- 1. No dependency on hardware/software.
- 2. Stronger self-cohesion in front-end and back-end.
- 3. Exploiting the strengths of each language.

Work flow



Sample Main Page

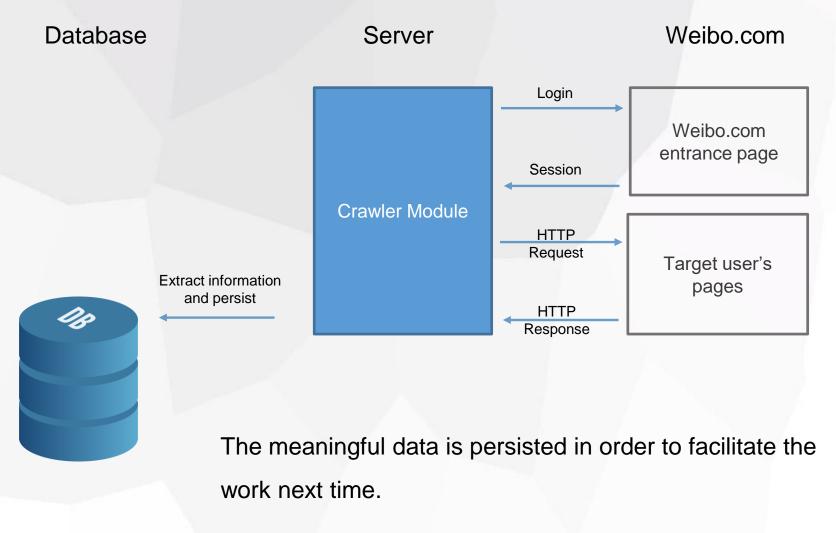
Status Monitoring Sidebar



Detailed module monitoring and time estimation

Logging

Crawler Mechanism

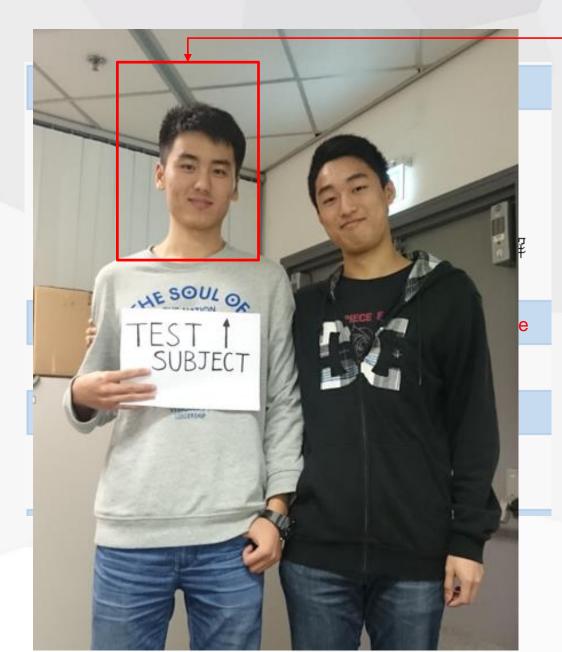


A trade-off between speed and storage.



PART III A Use Case

The person of test subject



Name: Mr. FANG

Gender: Male

Age: 22

Occupation:

Research Assistant in CityU

Hobbies:

Music, comics, programming, hiking

Others:

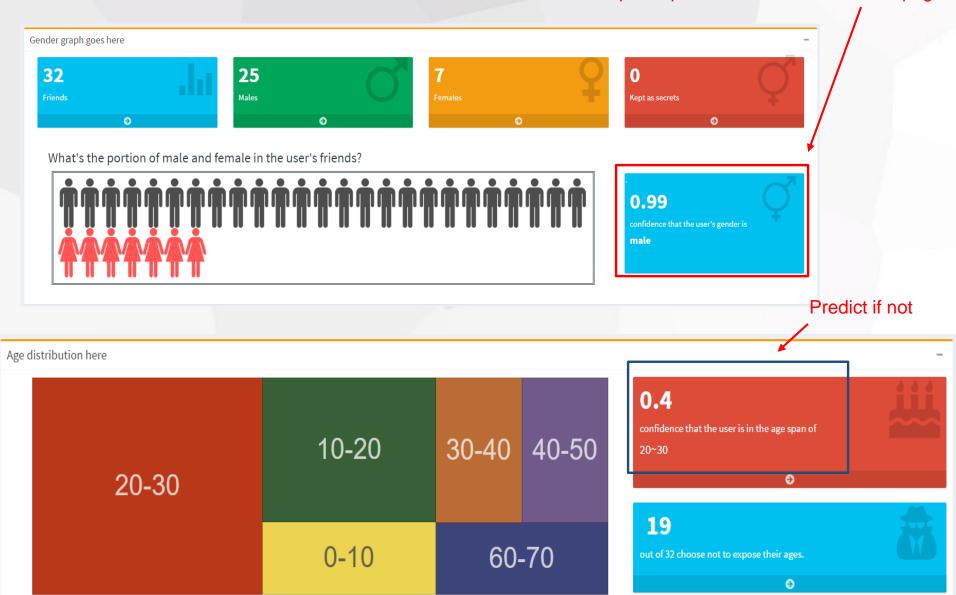
Placement in HSBC in 2014-2015

From Anhui province of China

Special thanks to Mr. Fang for his generosity and pertinent feedback!

Basic Information Study

Report if public on user's information page





friends are from **Beijing** province, accounting for the largest portion among all places.

"I registered with wrong geographical information to protect myself.
Plus I have many friends in Beijing."

confident that the user is from Beijing province.

0.99

Basic Information Study



Same report/predict flow on user's work/education experience.

User Interest Study

""My friends on Weibo.com are mainly university peer students. And we are all CS students."

How do the user's friends tag themselves?

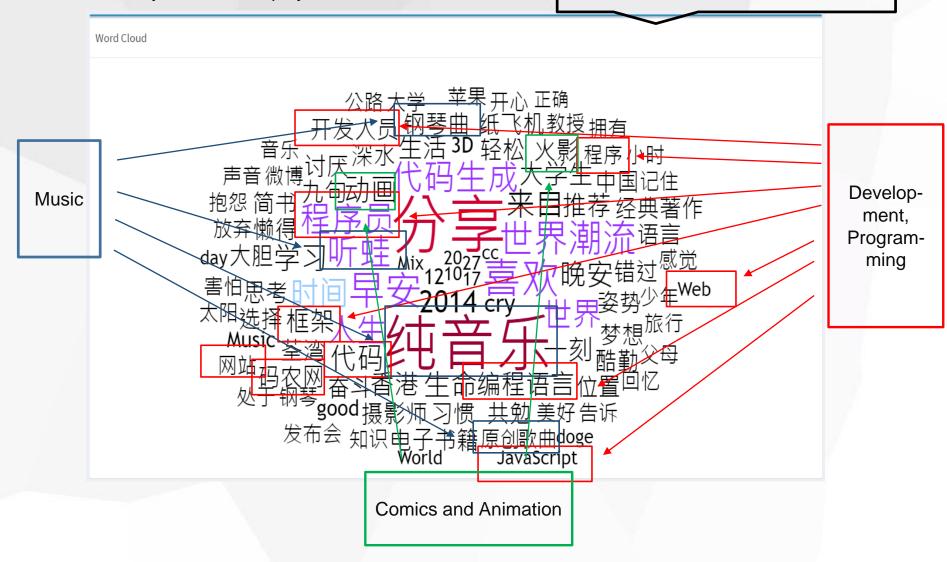


 User's interests can also be inferred from his/her friends' interests.

User Interest Study

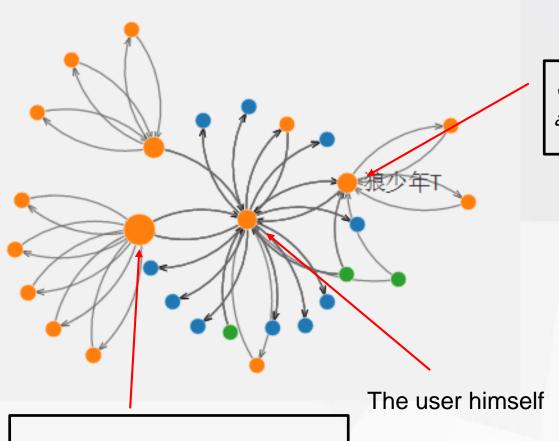
Keywords are displayed in a word cloud.

"I repost a lot from official accounts on music and programmer weeklies."



User Interactivity Study

"This is my best friend in childhood."



"This is my best friend, peer student and co-worker."

- The users that have interacted with the target form a directed network.
- Only the important nodes are displayed.
- Nodes are displayed in different colors to distinguish.

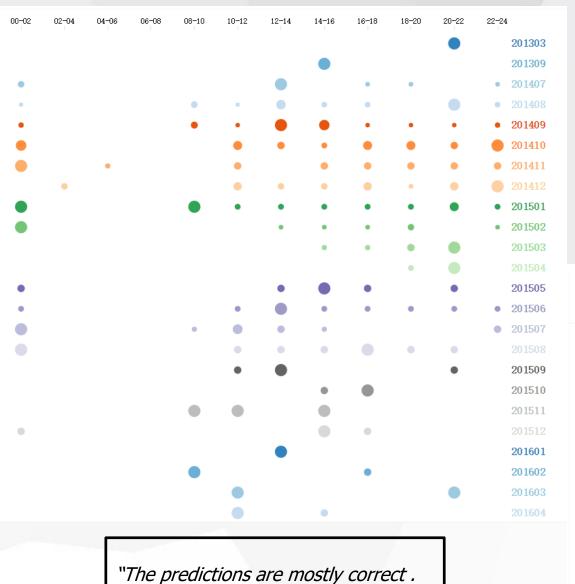
User Interactivity Study

		Ctaa		
They have c	ommented user's weibo			-
User	The torget bimself	Count	Ratio	
诺亚的大船	The target himself	24	9.95%	
狼少年T	Best friend, co-worker	8	3.31%	
棟暳紫	Girl friend	7	2.9%	
-我要成为神-	Roommate	3	1.24%	
hackeryoung	Peer student, friend	3	1.24%	
公馆教皇撒加1984		1	0.41%	
They have r	eposted user's weibo			-
User		Count	Ratio	
User FOREVER全力	인화 Childhood bost friend	Count	Ratio	
FOREVER全力	以赴 Childhood best friend	1	0.41%	
	以赴 Childhood best friend			
FOREVER全力 Rinka喵 诺亚的大船	以赴 Childhood best friend	1	0.41%	-
FOREVER全力 Rinka喵 诺亚的大船		1	0.41%	-
FOREVER全力 Rinka喵 诺亚的大船 They have li	iked user's weibo	1 1 1	0.41% 0.41% 0.41%	_
FOREVER全力 Rinka喵 诺亚的大船 They have li	iked user's weibo	1 1 1 Count	0.41% 0.41% 0.41% Ratio	_
FOREVER全力 Rinka喵 诺亚的大船 They have li User FOREVER全力	iked user's weibo 以赴	1 1 1 Count	0.41% 0.41% 0.41% Ratio	-
FOREVER全力 Rinka喵 诺亚的大船 They have li User FOREVER全力 楝暳紫	iked user's weibo 以赴	1 1 1 Count 4	0.41% 0.41% 0.41% Ratio 1.65% 0.41%	-

 The users who have most interacted with the target user are sorted and listed.

""Those with top interaction are my the people that are closest to me."

User Life Style Study



Better if it's finer grinded in 24 hours."

- The post time of each microblog is mapped to a 24-hour table.
- Predict life schedule.

T Predictions on the user's daily schedule

The user is active in these time slots:

0:00~2:00 12:00~14:00 20:00~22:00

The user is inactive in these time slots:

2:00~4:00 4:00~6:00 8:00~10:00

This is the predicted bedtime based on the activity of user:

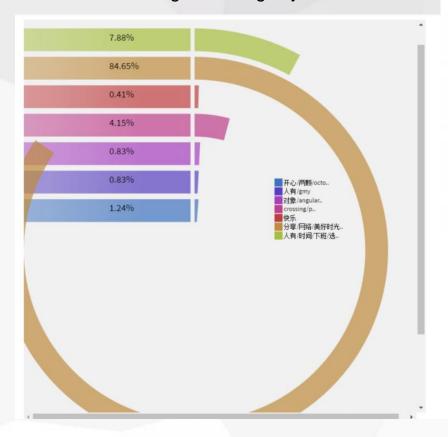
2:00~10:00

Topic Study

Reposts, music, moods, "chicken soup of the soul".

K-means++:

Clustering microblogs by content



ıag

开心/两颗/october

人有/gmy

对象/angularis

crossing/portraitpaint

快乐

Closest weibo

a little story . . . Farewell, 20. Tom Rosenthal - Go Solo.

http://t.cn/R7gkhNS //讨厌自己明明不甘平凡,却又不好好努力。 镜花水月 //往事浓淡,色如青,已轻。经年悲喜,净如镜,已静。 7小时22公里 边走边瞧@(̄- ̄)@ 新年加油,学业第二,身体第一! 李榮浩 李白

http://t.cn/RvYTDzn Stories 黑石瞳 鲁 鲁修 HITOMI,大三的记忆~ 空灵的 歌声, 激荡的心灵 http://t.cn/hbLvaf //【舌尖上的京朝[笑哈哈]】这个系列太 萌了! //我明白你会来,所以我等。——沈从文

9大船:3

9大船:2

9大船:2

9大船: 10

为大船:1

分享/网络/美好时光/中国/美好/集中式/扶好/余光中/胆小 **诺亚的大船**: 204 怕事/日常事务

人有/时间/下班/选择题/立即行动

诺亚的大船: 19

Total:

诺亚的大船: 241

"Yes it's okay to perceive my microblogs as more of reposts about music or moods. But the correlation between the tags or keywords is not significant."

Topic Study

Topic Modeling (Latent Dirichlet Allocation):

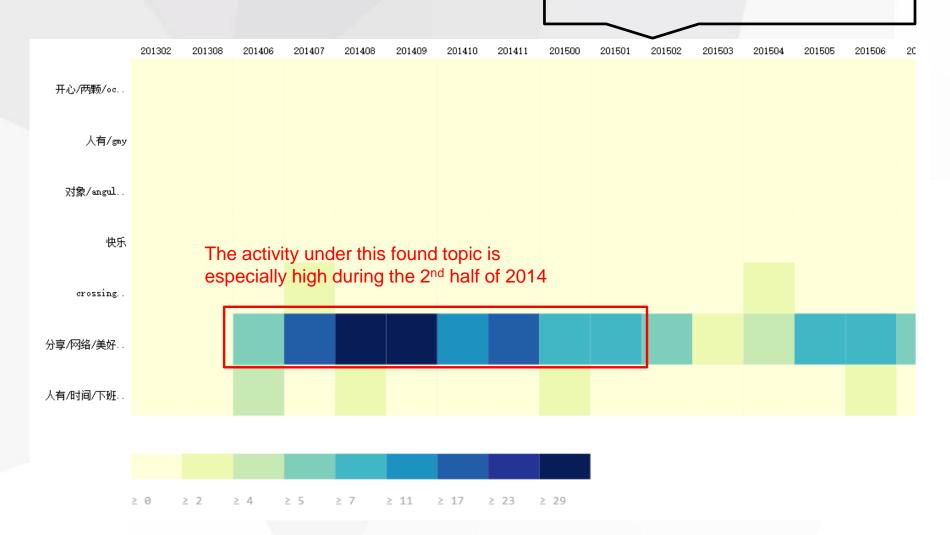
- Assume each piece of microblog is composed of words from multiple topics.
- Extract semantic clusters from the contents.

Topics			
Index	Programming Keywords		
1	来自 架构 css 代码 框架 html 程序员 位置选择 网站 香港 程序 人生 简单		
2	公路中国推荐大学互联网转发学习微博流星雨前端浪漫这是值得		
3	cry 父母 一刻 轻松 手机 世界 思考 哈哈哈 good 网友		
4	分享纯音乐 听蛙 喜欢一首 music 音乐 hitomi 一生 文章 蜡烛 Music		
5	时间生活世界努力习惯生命早安晚安教授九句		

"The semantic clusters are more obvious and indicative."

Timeline Study

"Yeah...I didn't have much to do at work...And I mostly interacted with girl friend so I was active in that period."





Algorithms Applied

Algorithms Applied

Network

PageRank

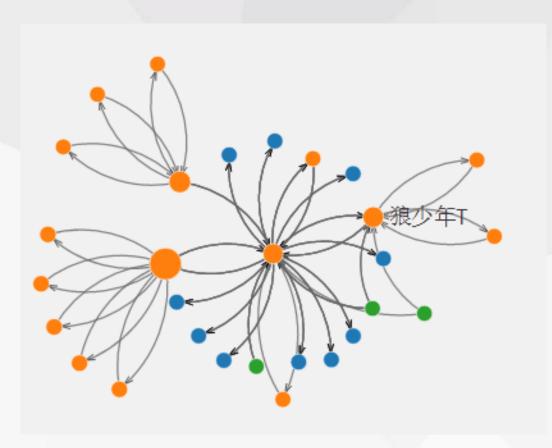
Classification

Logistic Regression, Random Forest, Gradient Boosting Trees (GBT), etc.

Clustering

K-means++, Latent Dirichlet Allocation (LDA)

Graph Definition



Node:

User

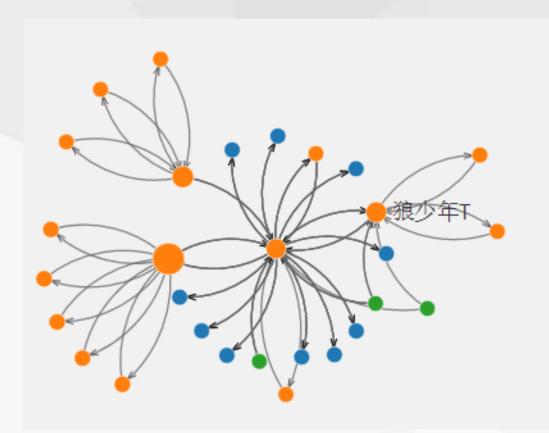
Link:

 There is a directed link from user A to B if A has commented / reposted / liked / mentioned B.

Weight:

The count of interaction

Characteristic of Weibo.com



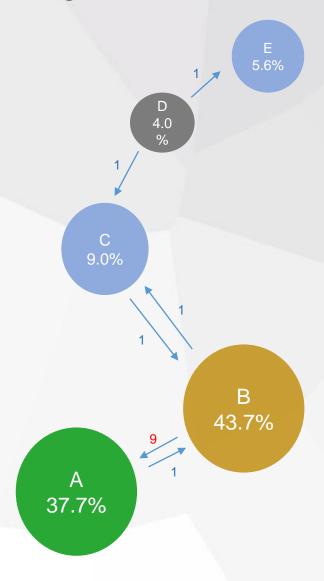
Celebrity Effect [1]

The correlation between inlinks and outlinks is low.

Base Assumption:

Outlinks outweigh inlinks.

PageRank



Step 1: Construct transition matrix [2,3]

$$H = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 0.9 & 0 & 0.1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0.5 & 0 & 0.5 \\ 0.2 & 0.2 & 0.2 & 0.2 & 0.2 \end{bmatrix}$$

Step 2: Initialize by evenly distribute weight to all the nodes:

$$\pi = egin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$
 π_A denotes importance of node A

Step 3: "Google Matrix" [4]

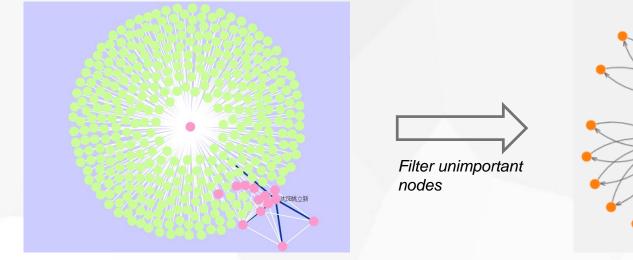
$$G = \theta H + (1 - \theta) \frac{1}{N} 11$$
 $\frac{\theta}{\theta}$ is the dampening factor.

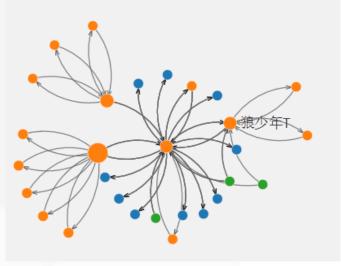
Step 4: Iteratively distribute the importance, until it converges:

$$\pi^T[k] = \pi^T[k-1]G$$

Benefits

- 1. Meaningful result
- 2. Better visualization.







PART 4 Conclusion

The Value

Data Mining:

- 1. Models and algorithms
- 2. Automated process
- 3. Providing prior knowledge about the user

Data Visualization:

- 1. Raw data are extracted and structured.
- 2. Inspiration from graphs

Software Engineering:

- 1. Integration of frameworks
- 2. Strengths of languages
- 3. Flexibility

Reference

- 1. Fan, P., Li, P., Jiang, Z., Li, W., Wang, H. Measurement and analysis of topology and information propagation on Sina-Microblog, *Intelligence and Security Informatics (ISI)*, 2011 IEEE International Conference
- 2. Page, L., Brin, S., Motwani, R. and Winograd, T. The PageRank citation ranking: bringing order to the web. 1999
- Csendes, T., Antal, E. PageRank Based Network Algorithms for Weighted Graphs with Applications to Wine Tasting and Scientometrics, *Proceedings of the 8th International Conference on Applied Informatics, Vol. 2. pp. 209–216.*
- 4. Chiang, M. Networked Life 20 Questions And Answers. *Cambridge University Press; 1 edition*, Sep. 2012, p. 71.

Q&A