# Social Big Data - Applications
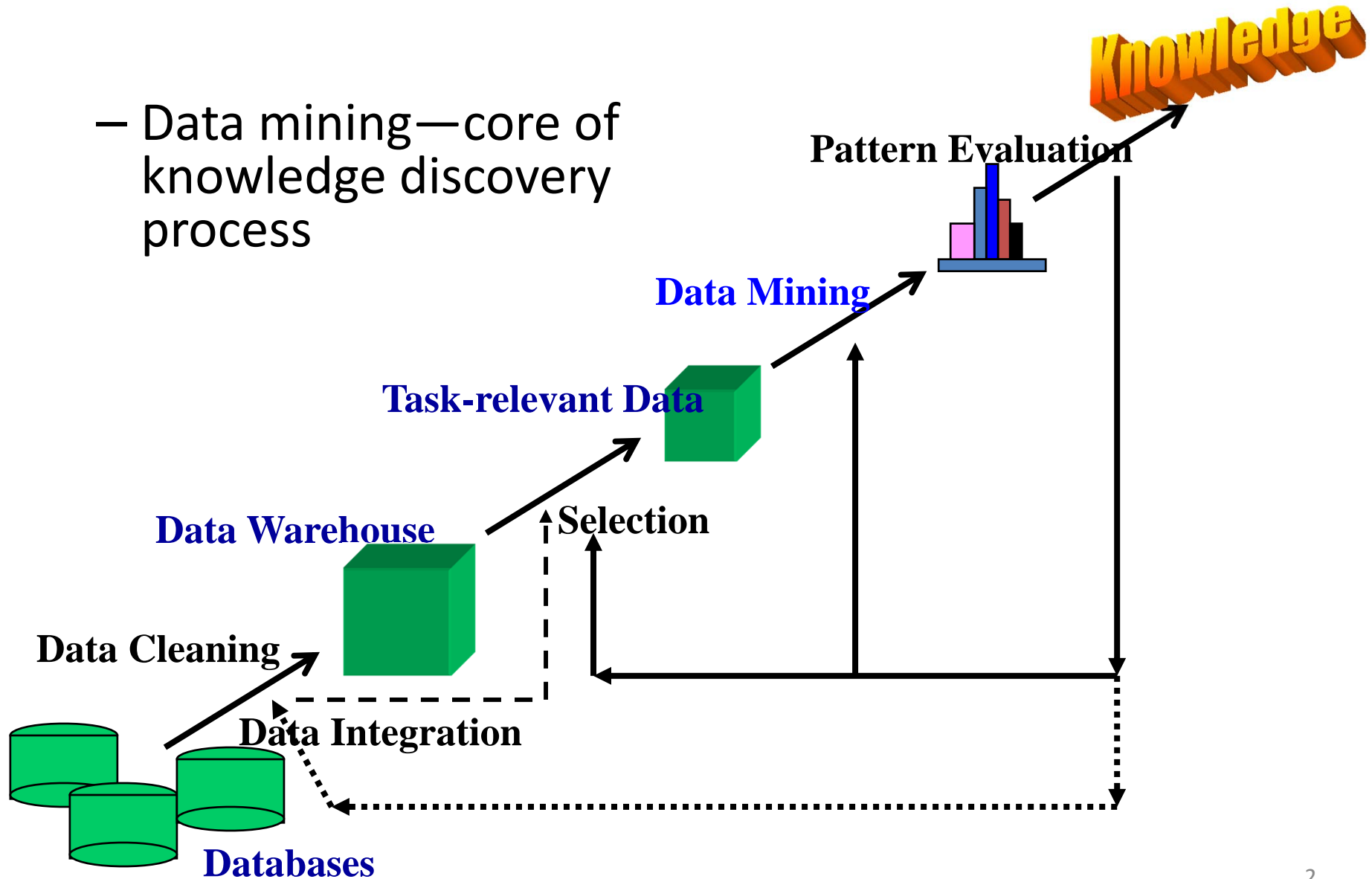
Dr. Hong Huang
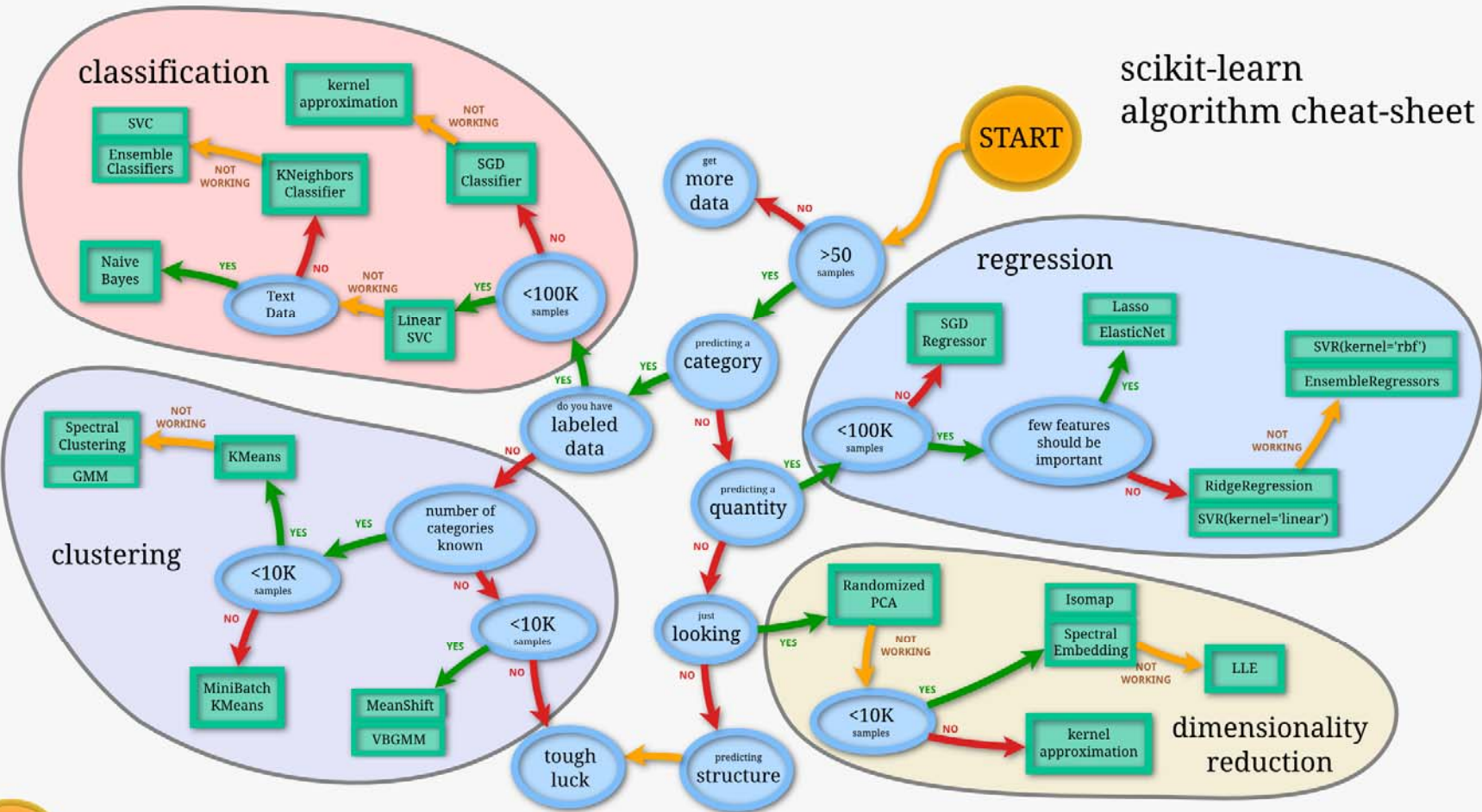
# Knowledge Discovery (KDD) Process

– Data mining—core of knowledge discovery process



**Pattern Evaluation**

**Knowledge**

**Data Mining**

**Task-relevant Data**

**Data Warehouse**

**Selection**

**Data Cleaning**

**Data Integration**

**Databases**

# Data mining algorithms



scikit-learn algorithm cheat-sheet

# What can we do in social network?

- Community identification
- Influential user identification
- Link prediction
- Point of interest recommendation
- Disease prediction
- Crime prediction
- Event monitoring
- …

# Opinion Leader Mining in CQA

# Community question answering (CQA) sites

- What is CQA site
  - Allow users to answer the questions posted by other users
  - Give positive or negative judgments to answers provided by others via voting
  - Popular QA portals: Yahoo! Answers, Stack Overflow, Quora, Zhihu

# An innovative CQA -- Zhihu

- ## What is Zhihu (Chinese Quora)
  - Traditional QA functions
    - Ask & answer questions
    - Vote answers
  - Social functions
    - Follow users
    - Follow topics and questions

user profile

followees

followers

topics the user follows

follow him

topic tags

question

follow the question

vote

answer

# Opinion leader identification

- **What is opinion leader**
  - Give their influential comments and opinions, put forward guiding ideas, agitate and guide the public to understand social problems[1]

- **Topical opinion leader in Zhihu**
  - Give authoritative and influential answers, comments and other activities in some topic area
  - Play an important role in promoting formation and management of online public opinion and knowledge base

[1] Lazarsfeld, P.F., Berelson, B., Gaudet, H.: The People's Choice: How the Voter Makes up His Mind in a Presidential Campaign. New York: Columbia University Press, (1948)
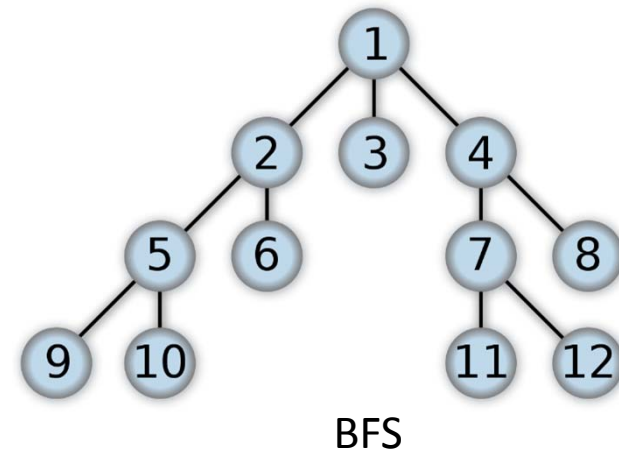
# Benefits of opinion leader identification

– Users: realize public opinion and authoritative knowledge, get specific answers efficiently

– Zhihu: invite them to attend public activities(e.g., editing, publication) to attract more users

– Marketing: influence customer opinions on products and services

– Government: realize, guide and interfere public opinion on the internet

# Dataset collection

- We gathered the Zhihu dataset through a web-based crawler from March to June in 2016
  - Start the crawler using a set of 10 popular Zhihu users from a webpage
  - The crawler follows a **Breadth First Search** (BFS) pattern through the follower and followee links of each user.



BFS

# Dataset

- Each user data contains
  - user ID, the lists of the user's followers and followees, the user's answers and questions posted

- Each answer/question data contains
  - its topic and the number of received votes

| Total number of users | 1,411,669 |
|---|---|
| Total number of questions | 701,982 |
| Total number of answers | 4,047,183 |
| Total number of topics | 160,664 |

# Data analysis tools

- Depending on needs, some data analysis tools are as follows
  - MATLAB, or its open-source alternatives, Scilab and GNU Octave (great at dealing with numbers)
  - Python with libraries like Numpy, Scipy and Matplotlib (great for general purpose data analysis- particularly good at interacting with other tools)
  - R (Great for statistics)
- Use Python to process the Zhihu dataset

# Python libraries for data science

- **NumPy**
  - the foundational library for scientific computing in Python, and many of the libraries on this list use NumPy arrays as their basic inputs and outputs. In short, NumPy introduces objects for multidimensional arrays and matrices, as well as routines that allow developers to perform advanced mathematical and statistical functions on those arrays with as little code as possible.

- **SciPy**
  - builds on NumPy by adding a collection of algorithms and high-level commands for manipulating and visualizing data. This package includes functions for computing integrals numerically, solving differential equations, optimization, and more.

- **Pandas**
  - adds data structures and tools that are designed for practical data analysis in finance, statistics, social sciences, and engineering. Pandas works well with incomplete, messy, and unlabeled data (i.e., the kind of data you're likely to encounter in the real world), and provides tools for shaping, merging, reshaping, and slicing datasets.

- **scikit-learn**
  - builds on NumPy and SciPy by adding a set of algorithms for common machine learning and data mining tasks, including clustering, regression, and classification. As a library, scikit-learn has a lot going for it. Its tools are well-documented and its contributors include many machine learning experts. What's more, it's a very curated library, meaning developers won't have to choose between different versions of the same algorithm. Its power and ease of use make it popular with a lot of data-heavy startups, including Evernote, OKCupid, Spotify, and Birchbox.

- **NetworkX**
  - allows you to create and analyze graphs and networks. It's designed to work with both standard and nonstandard data formats, which makes it especially efficient and scalable. All this makes NetworkX especially well suited to analyzing complex social networks.

# Data preprocessing

- Data cleaning
  - Remove incomplete user data
- Data transformation



Relationship matrix



Topic vector

# Initial analysis (1/2)



- Power law distribution
  - It has been identified in social science
  - It means that a small portion have extremely high degree while most have low degree

# Initial analysis (2/2)



- CDF (Cumulative distribution function) & CCDF (Complementary Cumulative Distribution Function)
  - 81% of the users did not ask any question and 72% of users did not publish any answer
  - about 38% of users have no follower and more than 99% of users have followees

# Methods on opinion leader identification



- Opinion leader identification is a ranking problem
  - PageRank (Larry Page), HITS (Jon Kleinberg)

- A typical ranking method – PageRank
  - An algorithm used by Google Search to rank websites in their search engine results
  - It works by counting the number and quality of links to a page to determine a rough estimate of how important the website is

18

# PageRank

- used by Google Search to rank websites in their search engine results (Larry Page)

PageRank works by counting the number and quality of links to a page to determine a rough estimate of how important the website is. The underlying assumption is that more important websites are likely to receive more links from other websites

# Original PageRank algorithm

- PR(A) = (1-d) + d (PR(T1)/C(T1) + ... + +PR(Tn)/C(Tn))
- Where:
  - PR(A) is the PageRank of page A
  - PR(Ti) is the PageRank of pages Ti which link to page A
  - C(Ti) is the number of outbound links on page Ti
  - d is a damping factor which can be set between 0 and 1

# A simple example of PageRank

- We regard a small web consisting of three pages A, B and C
  - Page A links to the pages B and C, page B links to page C and page C links to page A
  - The damping factor d is usually set to 0.85, but to keep the calculation simple we set it to 0.5
- Calculating PageRank

PR(A) = 0.5 + 0.5 PR(C)

PR(B) = 0.5 + 0.5 (PR(A) / 2)

PR(C) = 0.5 + 0.5 (PR(A) / 2 + PR(B))

We get the following PageRank values for the single pages:

PR(A) = 14/13 = 1.07692308

PR(B) = 10/13 = 0.76923077

PR(C) = 15/13 = 1.15384615

The sum of all pages' PageRanks is 3 and thus equals the total number of web pages

# Topical opinion leader identification in Zhihu

- Our method considers three aspects
  - Social network structure
    - Based on PageRank
  - Topical interest similarity
    - A user's influence on each follower depends on the topic interest similarity between them
  - Knowledge authority
    - The higher authority (the number of votes) a user has, the higher probability he has to impact his followers

- Set link weight between users according to the topical interest similarity and knowledge authority

# Results and evaluation

| Topic | Top 5 opinion leaders in each topic |
|---|---|
| Movie | xiepanda, liuniandate, vikinglau, WxzxzW, chen-yao-39-75 |
| Psychology | xiepanda, liuniandate, WxzxzW, zhang-xiao-wei-23, yezhuang |
| Travel | WxzxzW, chico-62, xu-wen-39, li-zhi-qiang-peter, qiu-shi-19-94 |
| Food | xiepanda, anshi, weijiali, ji-li-ji-li, liuniandate |
| Fitness | WxzxzW, chico-62, xiepanda, summer.li, guo-fu-lin |
| Internet | xiepanda, WxzxzW, liuniandate, big_caaat, 8king |
| Fashion | WxzxzW, 8king, sickberry, liuniandate, xiepanda |
| Pioneer | wangxing, zhou-kui, xiepanda, liuniandate, dreamcog |
| Design | WxzxzW, 8king, xiepanda, soulchef, xiaoxiao |
| Finance | xiepanda, liuniandate, WxzxzW, ji-li-ji-li, big_caaat |

- Overall evaluation
  - They always published lots of topic-related posts and received a great number of votes, and have a great number of followers including some important followers

# Topical opinion leader examples

- "xiepanda", "liuniandate"and "WxzxzW"are identified among the top 5 topical opinion leaders in most topics
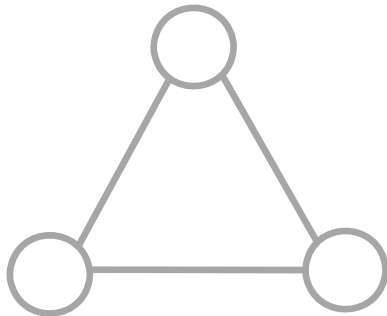  - This kind of users are so-called **cewebrity**, who often posts a large number of content about various topics to acquire fame on the Internet so that they have more than 200K followers including many important ones.
- "yezhuang" is identified as a opinion leader in psychology
  - He is a psychology trainer, whose posts are all related to psychology
  - His answers received 458 of average vote count. He is also followed by more than 40K users including some top-ranked ones
- Most of pioneer-related top 5 opinion leaders are successful company founders in real life
  - "wangxing" posted mainly about Pioneer and has 61,268 followers including a few of influential users "yuyue-51", "zhou-kui", and "GavinQi". He founded many popular websites such as Meituan, Fanfou and Renren

# Triadic Closure and Its Influence in Social Networks
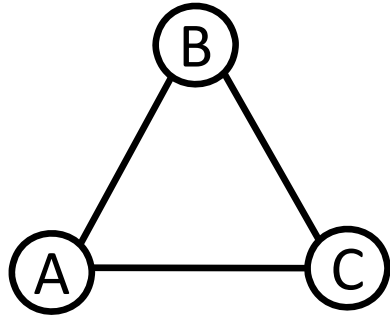
*– A Microscopic View of Network Structural Dynamics*

# Two Issues in Triadic Closure Process
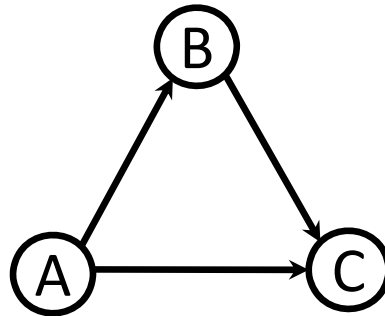


**Problem 1**

- Problems:
  - 1: will open triad be closed at time t?
  - 2: at $t + \Delta t$, will the tie AB become stronger or weaker? Also BC?
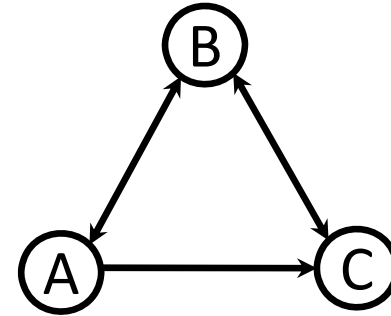
# Related Work



Undirected [1]          Directed – Case 1 [2]          Directed – Case 2 [3]

[1] Zignani, M., et.al. Link and triadic closure delay: Temporal metrics for social network dynamics. In ICWSM'14.
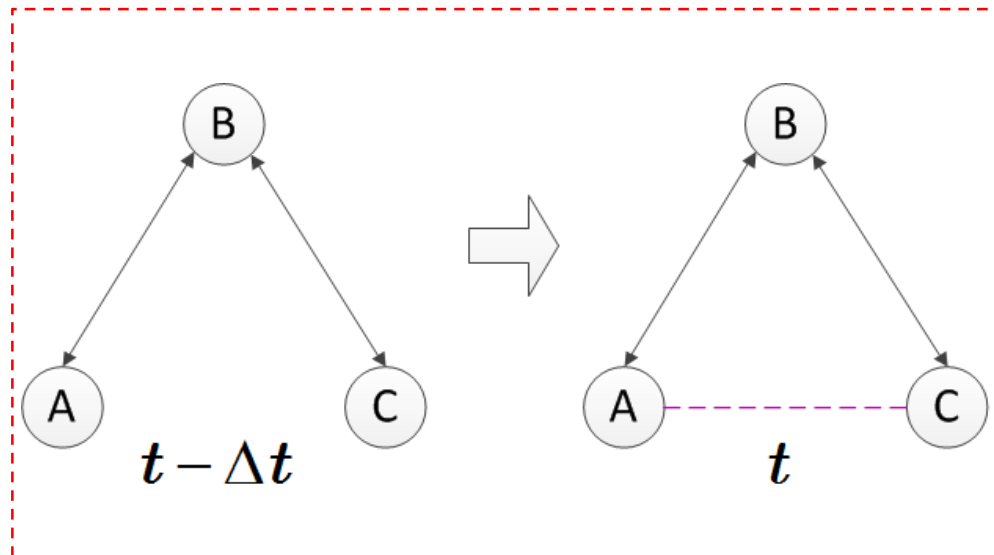[2] Romero, D. M. and Kleinberg, J. The directed closure process in hybrid social-information networks, with an analysis of link formation on twitter. Stat, 2010.
[3] Lou, T et.al. Learning to predict reciprocity and triadic closure in social networks. TKDD, 2013.
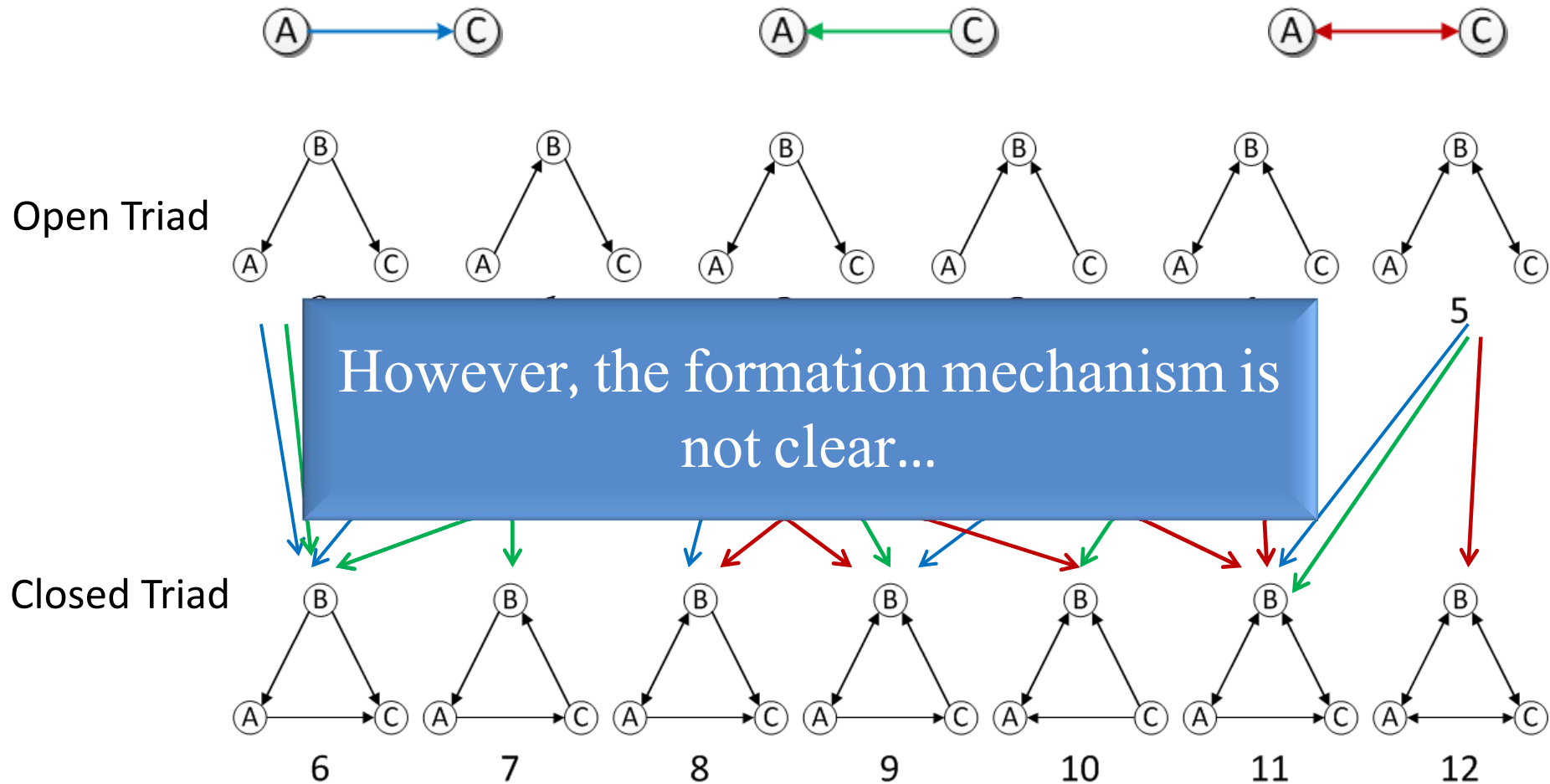
Triad: A group of three people

# WHY TRIADS? TRIADIC CLOSURES?

# What are underlying factors that trigger triadic closure?

# Open Triad to Triadic Closure



Open Triad

Closed Triad

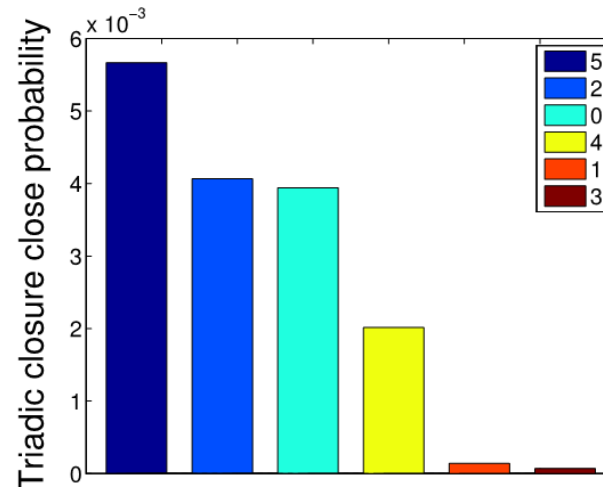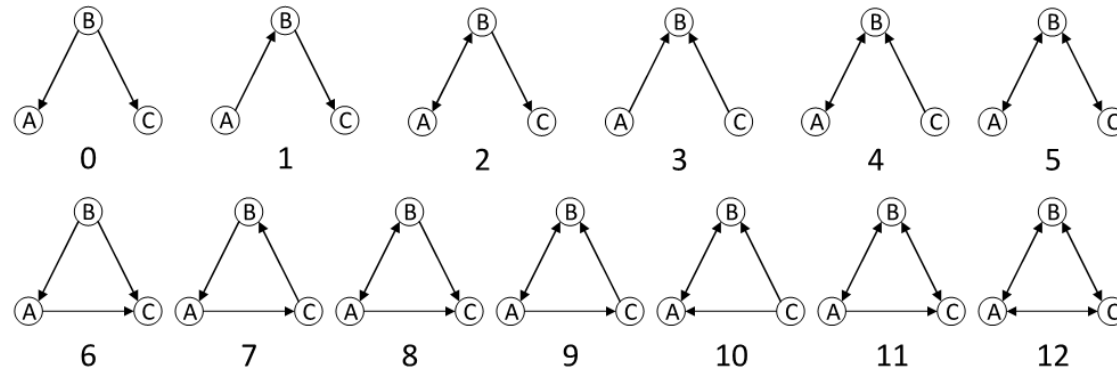However, the formation mechanism is not clear...

6  7  8  9  10  11  12

Milo R, Itzkovitz S, Kashtan N, et al.. Superfamilies of evolved and designed networks. Science, 2004

# Look into one social network...

- Time span: Sep 29[th], 2012 – Oct 29[th], 2012
- 1.7 million nodes
- 400 million following links
- 200 out-degree per user
- 360 thousand new links
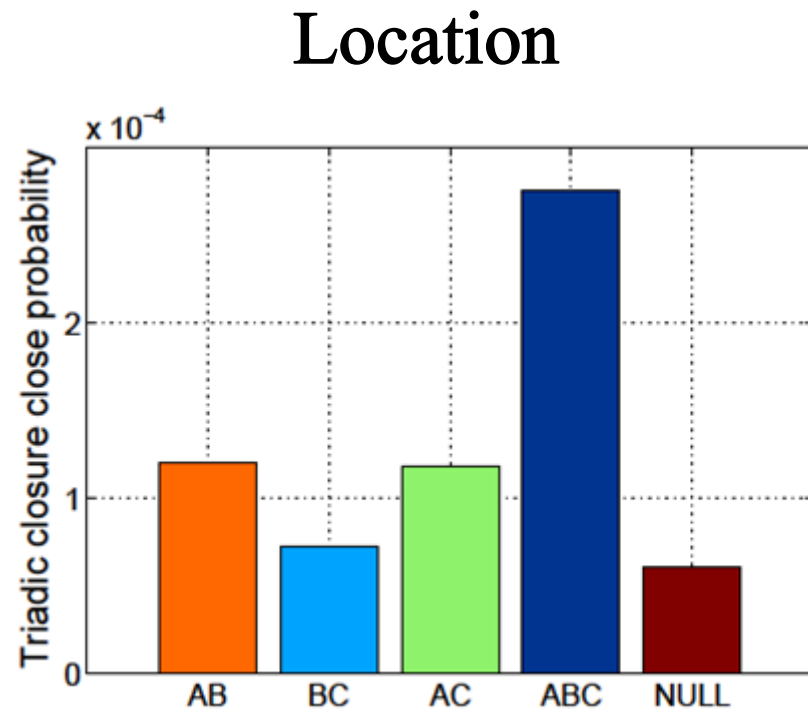- 746 thousand newly formed closed triads per day

# Observation - Network Topology
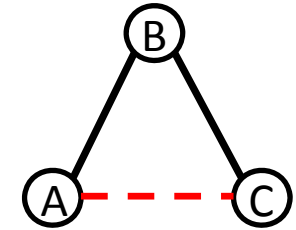


Y-axis: probability that each open triad forms triadic closures
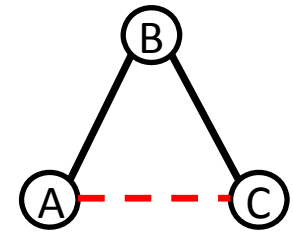
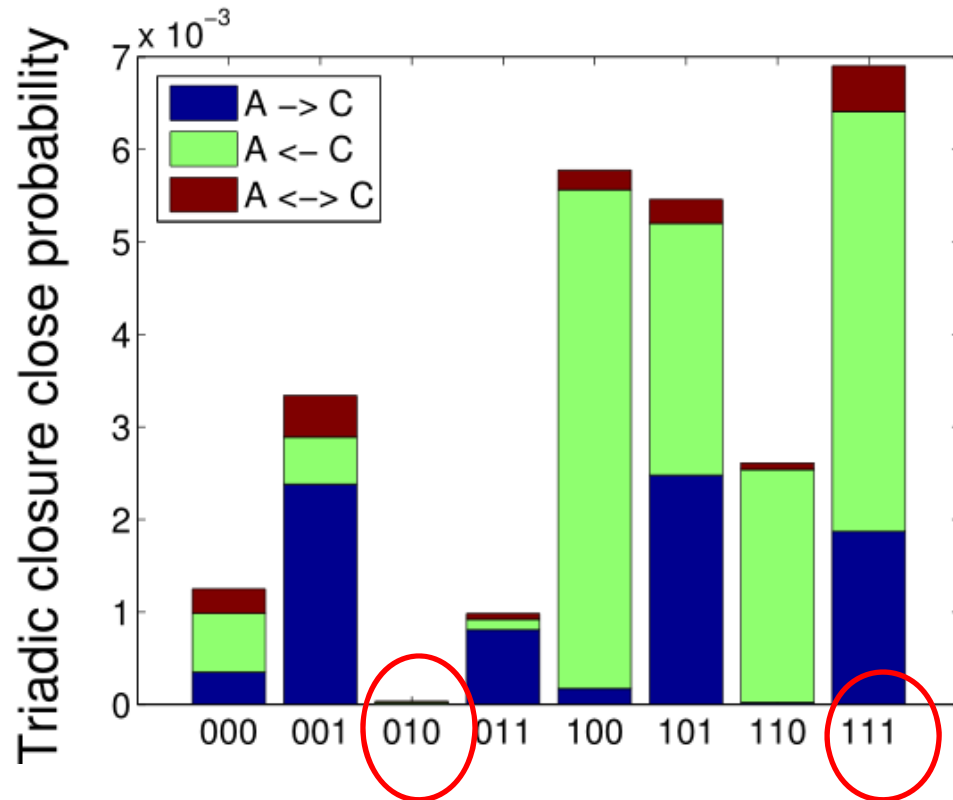# Observation - Demography



## Location
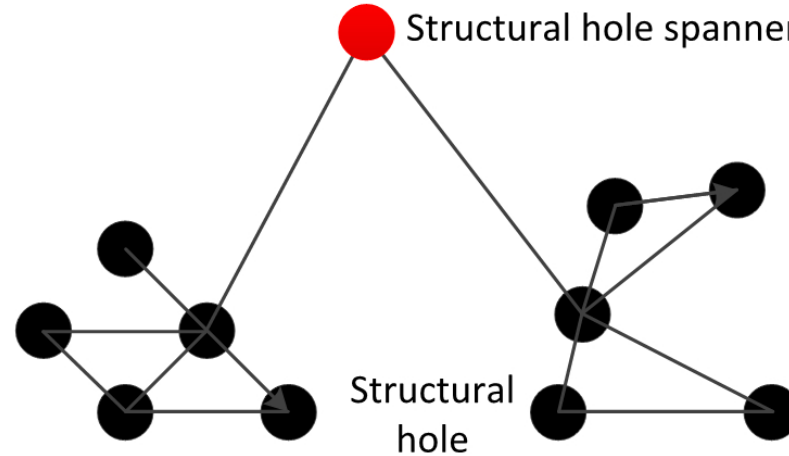


## Gender

AB means A and B are from the same city
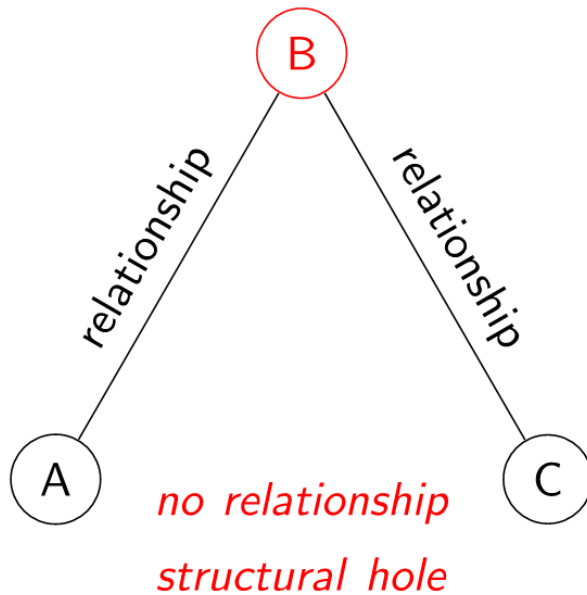
# Observation – Opinion Leader



0—ordinary user
1—opinion leader

# Structural Hole

- **Structural hole:** two separate clusters possess non-redundant information.



relationship

relationship

*no relationship*

*structural hole*

Structural hole spanner

Structural hole

Lou T, Tang J. Mining structural hole spanners through information diffusion in social networks, www2013
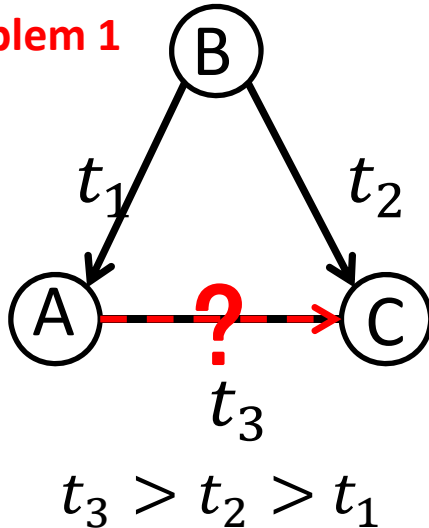
# Observation - Structural Hole



0—ordinary user
1—structural hole spanner

# Problem Formalization

$t_3 > t_2 > t_1$

- Given a network $G^t = (V, E, X, Y)$,
  - $X$ features defined for candidate triads
  - $Y$ whether an open triad become closed or not (Problem 1)

- **Goal**: Predict the formation of triadic closure given learned features and the network
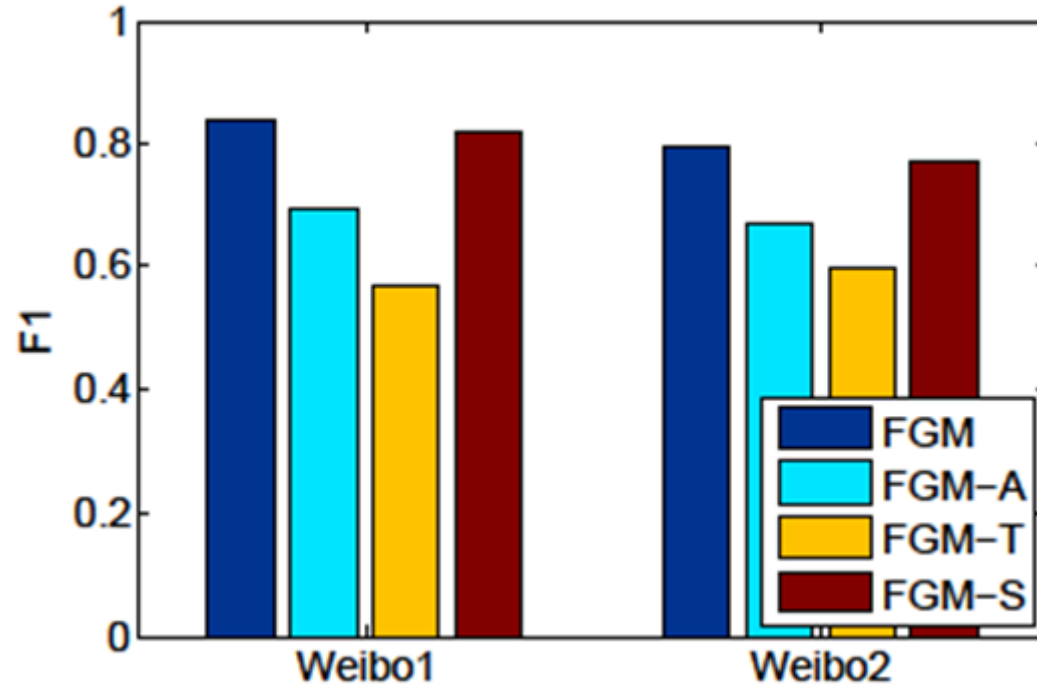
$$\varphi = logP(Y|X, G)$$

# Experiments & Results

- Datasets: Weibo
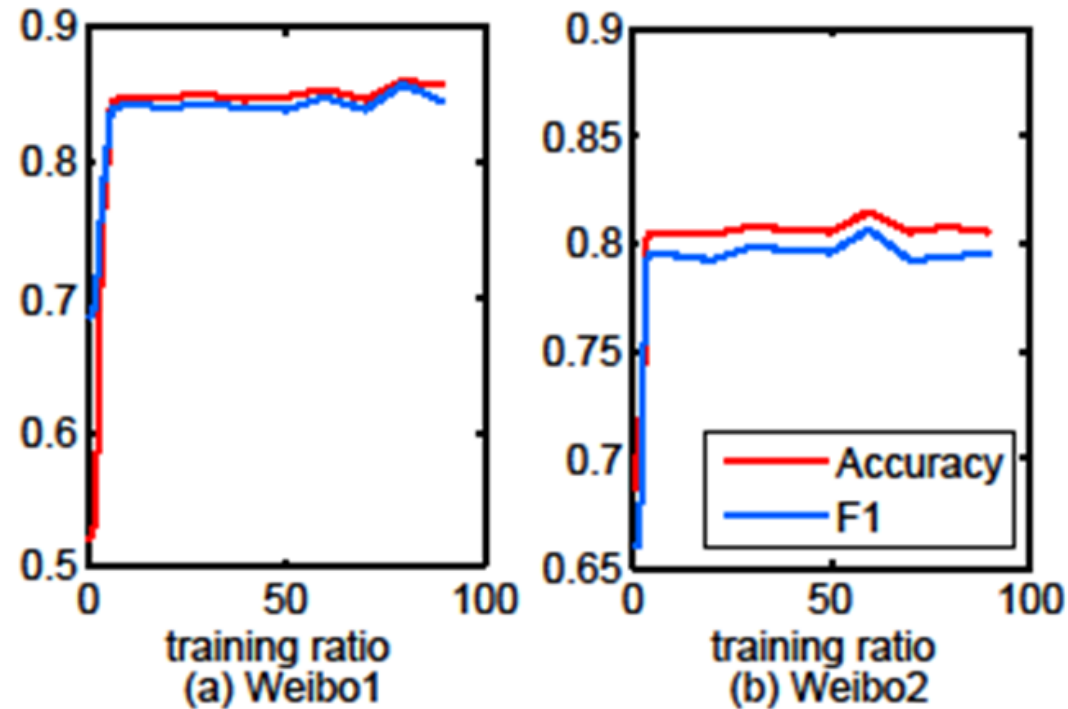- 50% as a training set, 50% as a testing set

Problem 1

| Algorithm | Precision | Recall | F1 | Accuracy |
|-----------|-----------|--------|------|----------|
| SVM | 0.890 | 0.844 | 0.866 | 0.882 |
| Logistic | 0.882 | 0.913 | 0.897 | 0.885 |
| Our | **0.901** | **0.953** | **0.926** | **0.931** |

# Factor Contributions



Temporal factor > Attribute factor > Social factor

# Performance - Training Data Ratio



(a) Weibo1 — (b) Weibo2

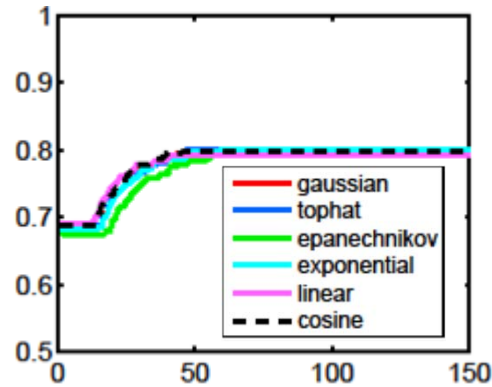Robust to the size of training dataset

# Performance - Convergence



(a) Weibo1 — # Iterations

(b) Weibo2 — # Iterations

(c) Weibo1 — # Iterations

(d) Weibo2 — # Iterations

Legend (b): gaussian, tophat, epanechnikov, exponential, linear, cosine

Legend (d): 10%, 30%, 50%, 70%, 90%

# iterations < 50

# Conclusion



Problem 1

$t - \Delta t$   $t$

- Studying triadic closure process
- Uncovering underlying factors that trigger triadic closure
- Proposing  efficient models to predict triadic

# Announcement

- Next class on 8th June will start at 10:30am and end at 12:00.