

Community and Role Model

Zhongjing Yu

Outline

Role model: Role is a important factor to model for **detect behavior** or **detect community**.

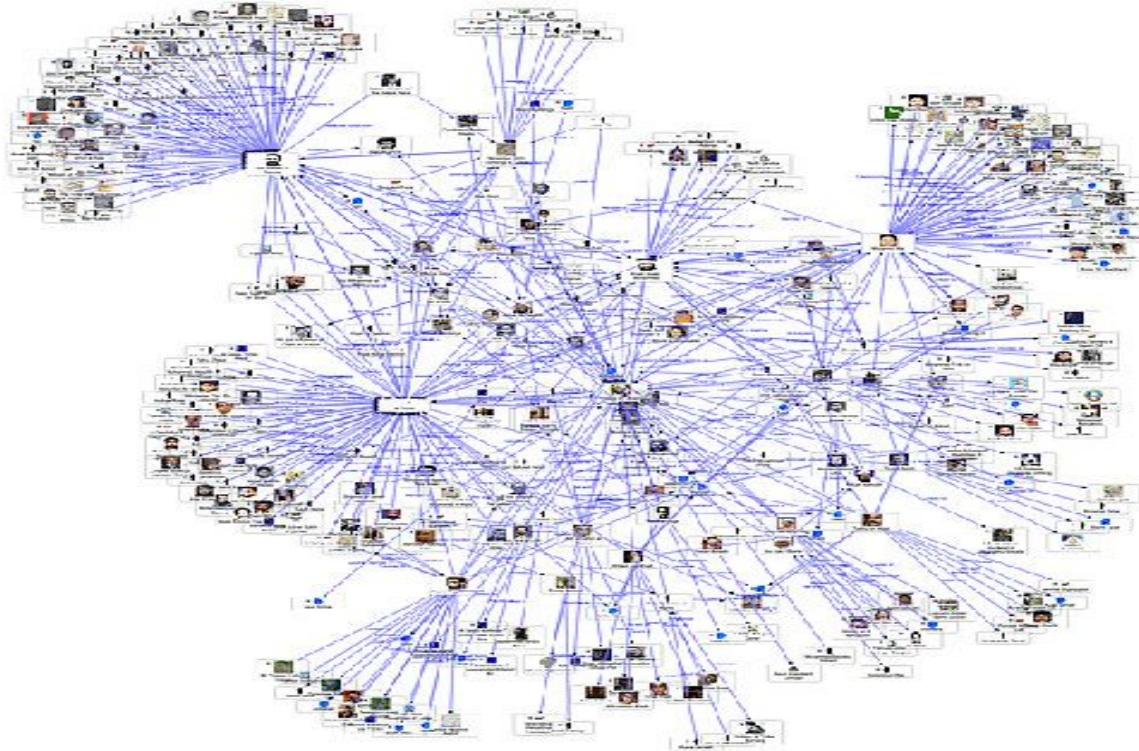
- Firstly , introduce a model to **identify roles** and application of the model.
- Secondly , introduce a Probabilistic Community and Role Model ,which consider **role** and other factors are significant to model.

(SSRM: Structural Social Role Mining for Dynamic Social Networks,
Probabilistic Community and Role Model for Social Networks)

SSRM: Structural Social Role Mining for Dynamic Social Networks

Afra Abnar, Mansoureh Takaffoli,
Reihaneh Rabbany, Osmar R. Zaiane

Background:



Structure of role model exists in a social network and each individuals play various of roles

Outline

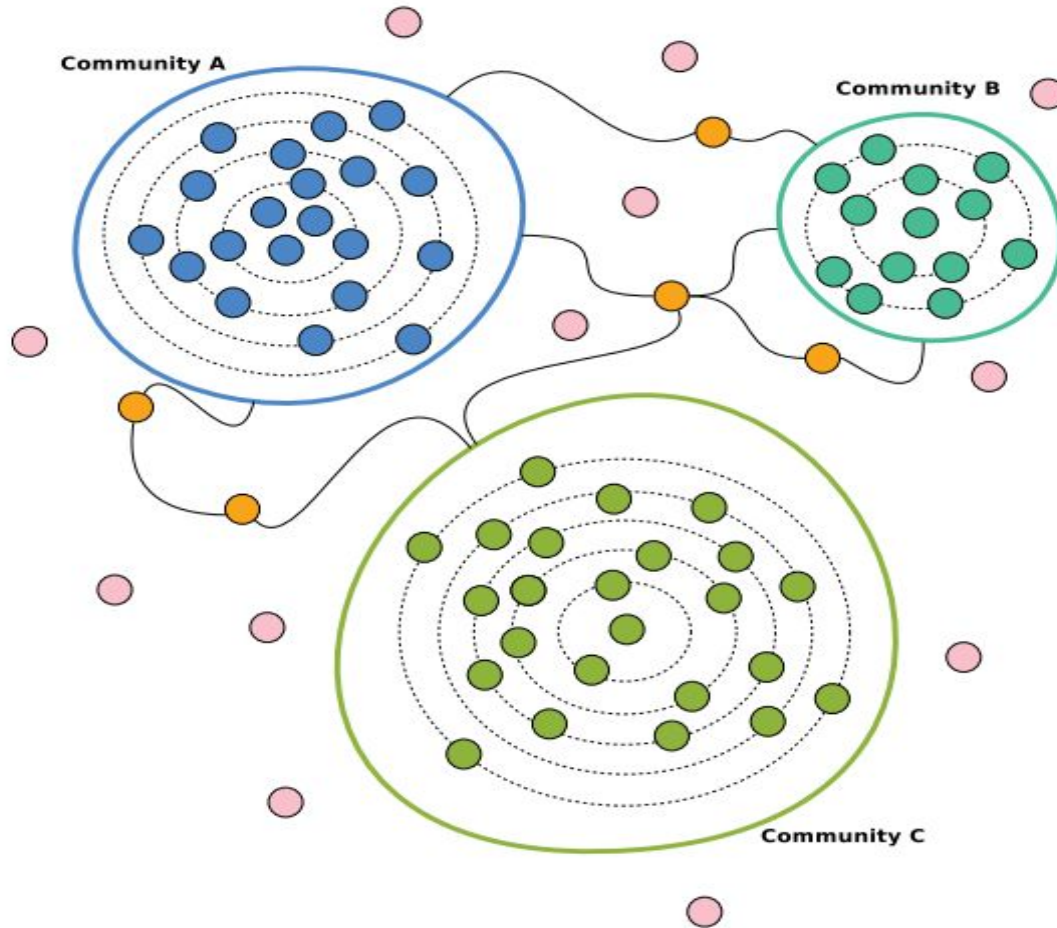
- **Target:** identify roles in social networks
- **Analysis** some kinds of roles
- How to **identify**
 - which measure adapted
 - C-Betweenness**
 - L-Betweenness**
 - MedExtractor** algorithm
- **Experiment**

Target: identify roles in social networks

condition:

- The SSRM framework is built up **two characteristic** of human societies, given structural properties:
 - ❖ role-taking behavior of the individuals with each other (**edges**)
 - ❖ A social network is considered only structural properties of nodes(**nodes**)

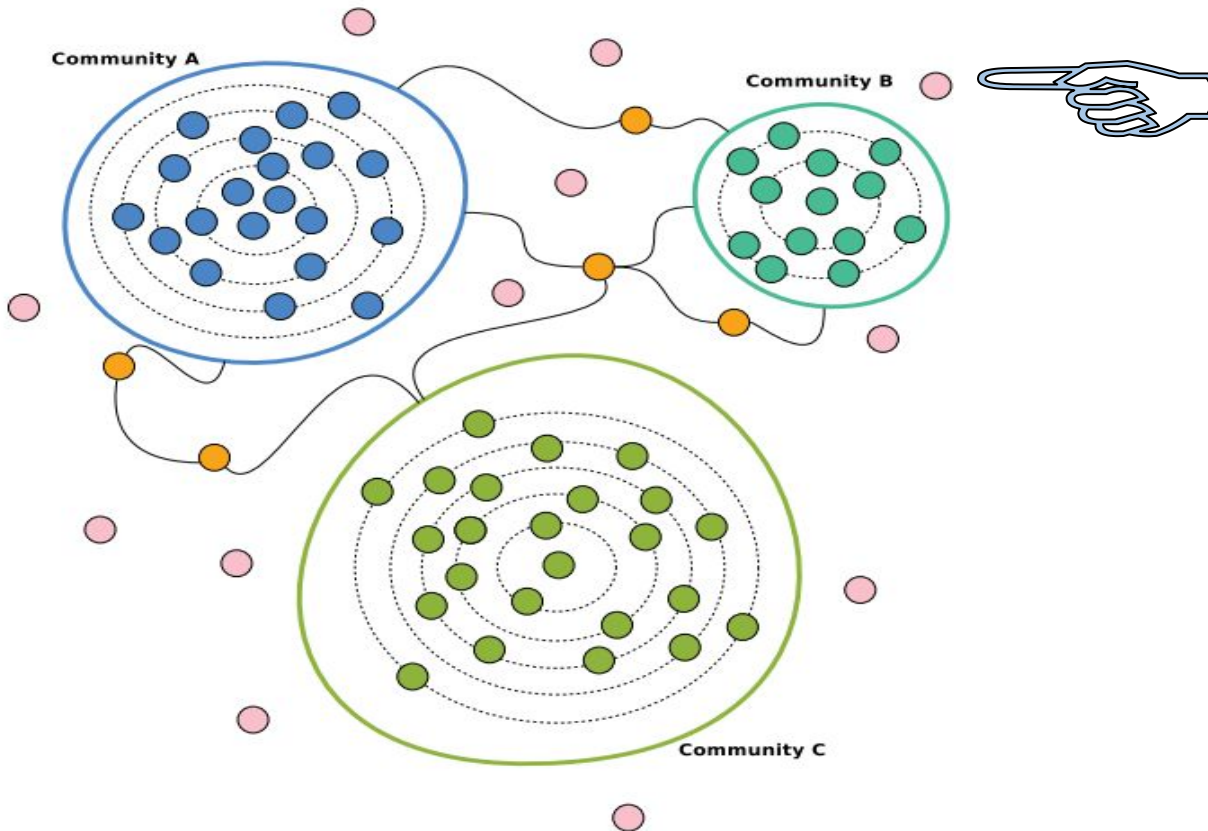
Analysis



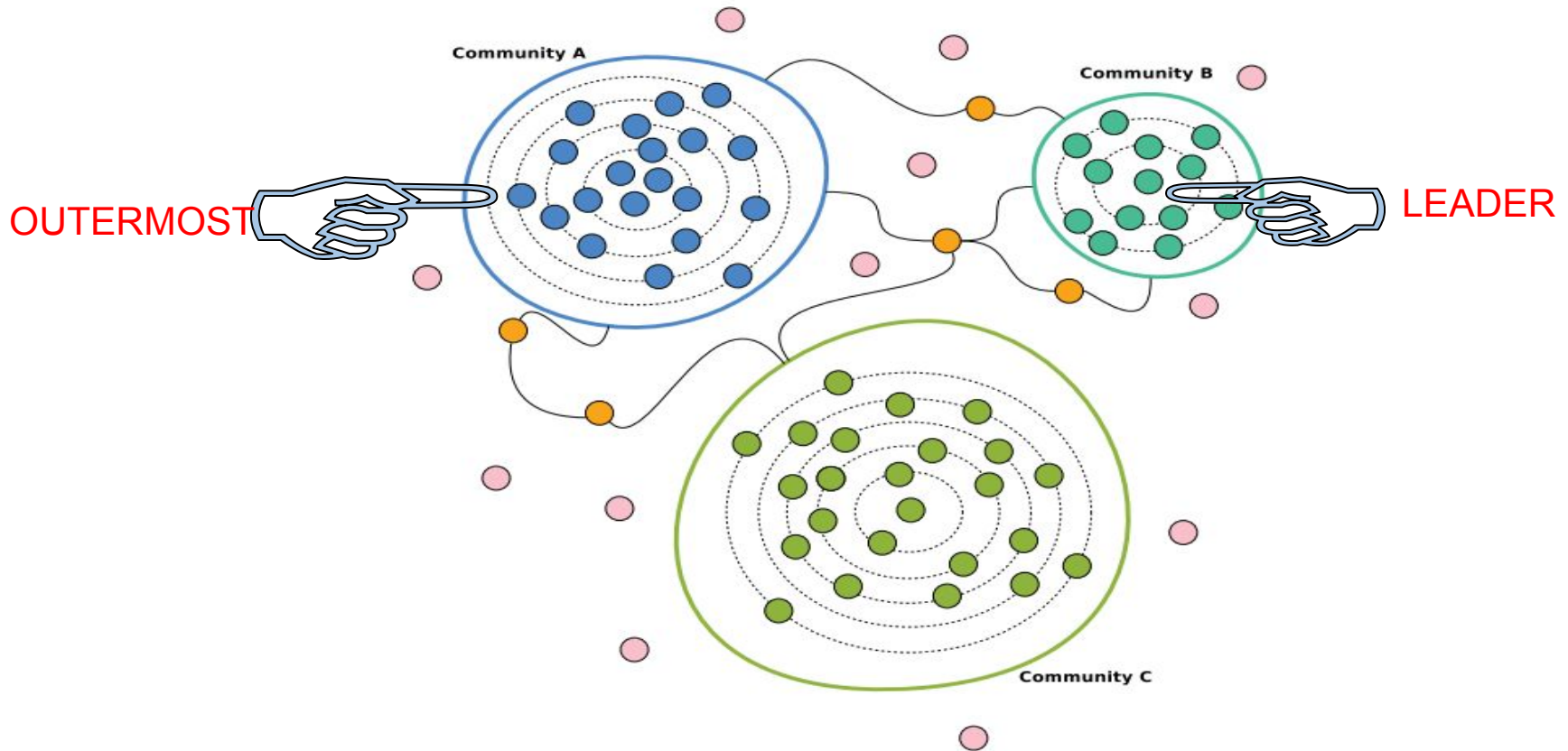
- ❖ **LEADERS**: important nodes in a community
- ❖ **OUTERMOSTS**: least significant individuals
- ❖ **MEDIATORS**: connect different communities
- ❖ **OUTSIDERS**: no affiliate to any one community

How to identify

- ❖ **OUTSIDER** members not belong to **any community**



- ❖ **LEADER** (a) adapts an appropriate measure M is used to score the members of the community;(b) the probability distribution function(pdf) for the importance scores is estimated.
- ❖ **OUTERMOST** members are identified in contrast to the leaders.



Measure

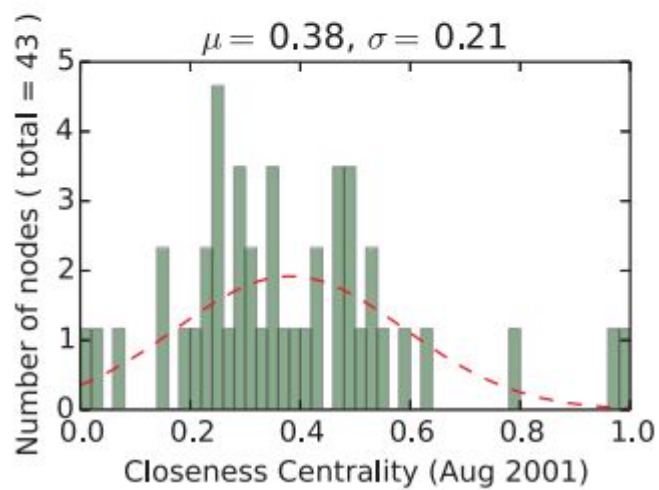
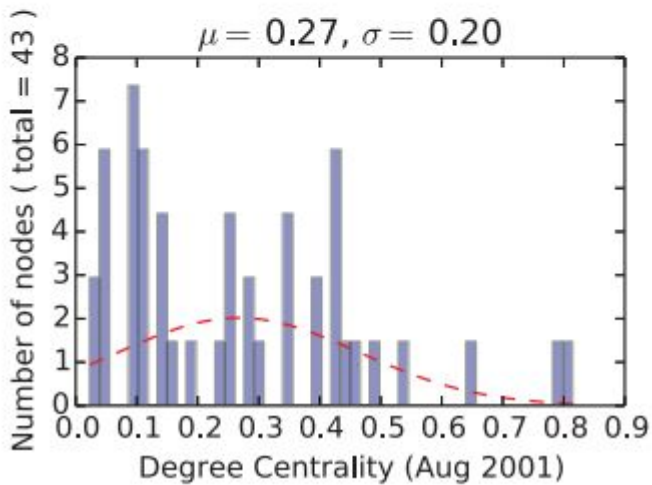
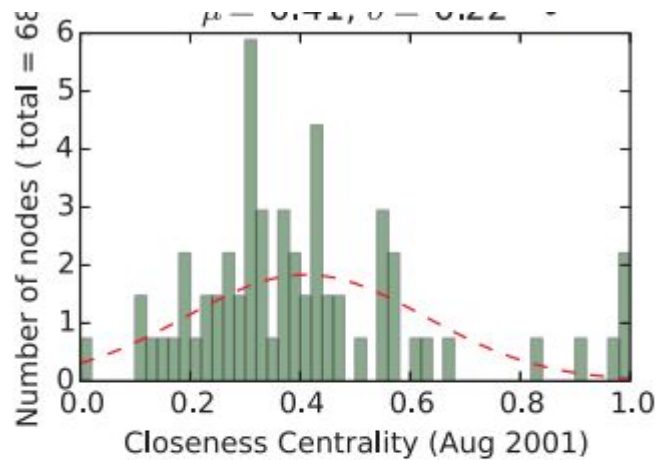
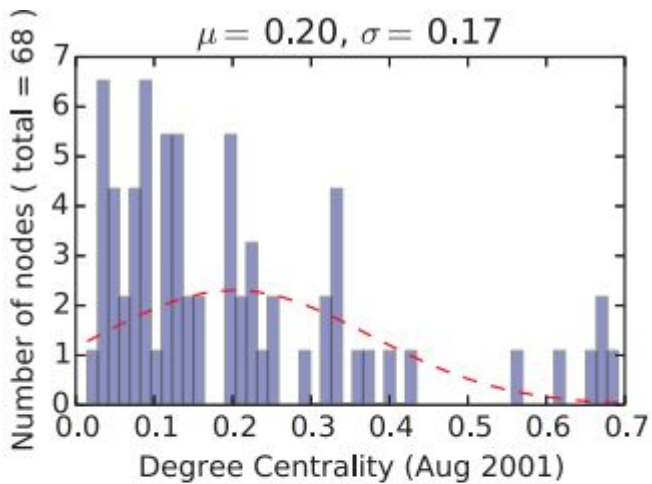
- degree centrality:
- closeness centrality:
 - Closeness centrality of a node u is the reciprocal of the sum of the **shortest path distances** from u to all $n-1$ other nodes. Since the sum of distances depends on the number of nodes in the graph, closeness is **normalized** by the sum of minimum possible distances $n-1$

$$C(u) = \frac{n - 1}{\sum_{v=1}^{n-1} d(v, u)},$$

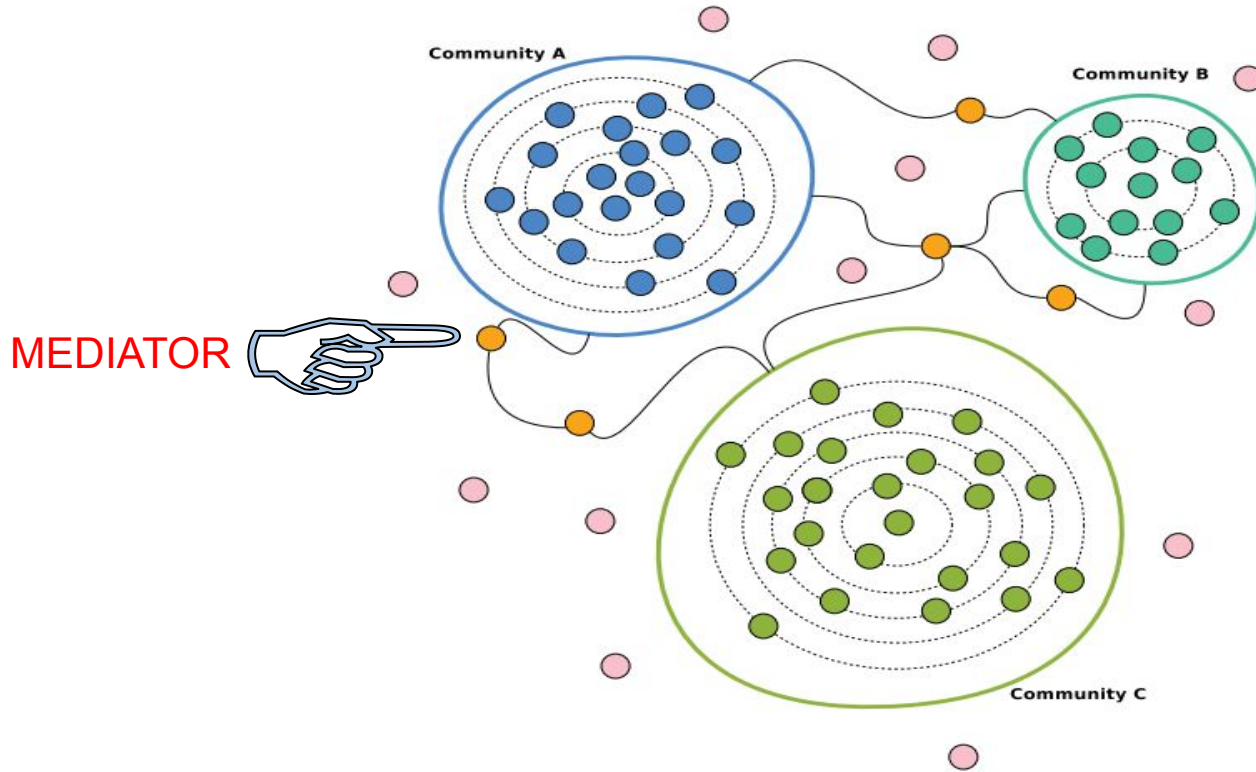
where $d(v,u)$ is the shortest-path distance between v and u , and n is the number of nodes in the graph.

-----NetworkX

»



- **MEDIATOR**: Commonly-used **betweenness centrality** ranks nodes based to **the number of shortest paths** that pass through the nodes.



Betweenness centrality

C-Betweenness counts the numbers of **shortest paths** between *different communities* that pass through a node.

- ◆ s_p and e_p denote the **start point** and **end point** of the shortest path p
- ◆ c_v return **community** that node v **belongs to**.
- ◆ the **set of all shortest paths** that *connect different communities* as

$$CPaths = \{p \mid c_{s_p} \neq c_{e_p}\}$$

- ◆ $I_p(p, v)$ return 1 if node v resides on path p , and 0 otherwise.
- ◆ C-Betweenness of node v is defined as:

$$CBC(v) = \frac{1}{2} \sum_{p \in CPaths} I_p(p, v)$$

L-Betweenness denotes not only the shortest paths between leaders of *different communities*.

◆ **leaderSet(c)** : the set of leaders of community c , then consider $LPath$ as:

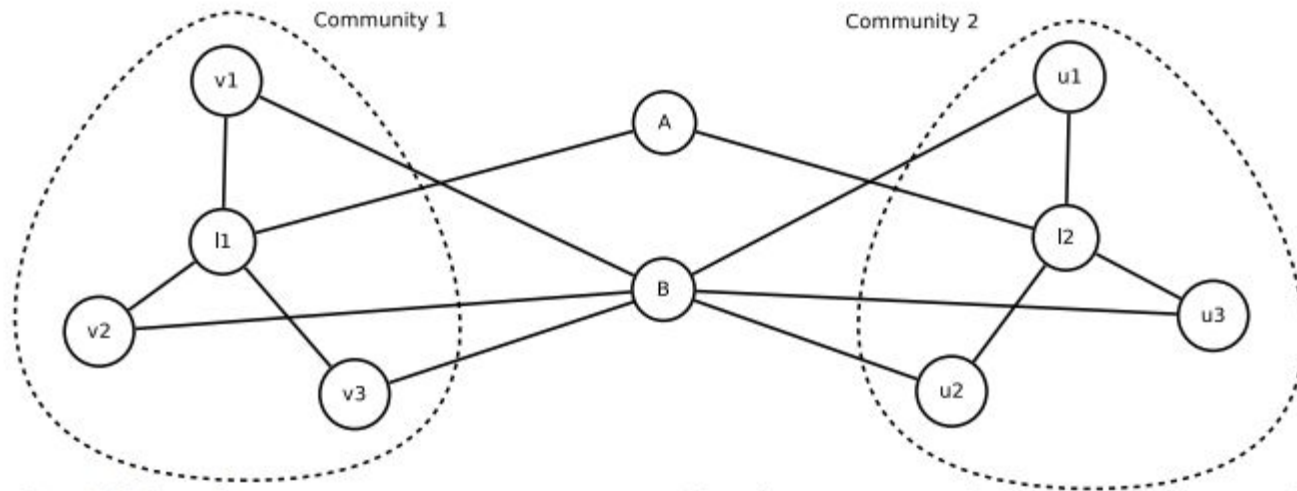
$$LPaths = \{p \in CPaths \mid \exists c_i, c_j : \\ s_p \in leaderSet(c_i) \wedge e_p \in leaderSet(c_j)\}$$

◆ **L-Betweenness** of a node v , $LBC(v)$ is defined as:

$$LBC(v) = \frac{1}{2} \sum_{p \in LPath} I_p(p, v)$$

(2 is omitted for directed graph)

C-Betweenness score, L-Betweenness score



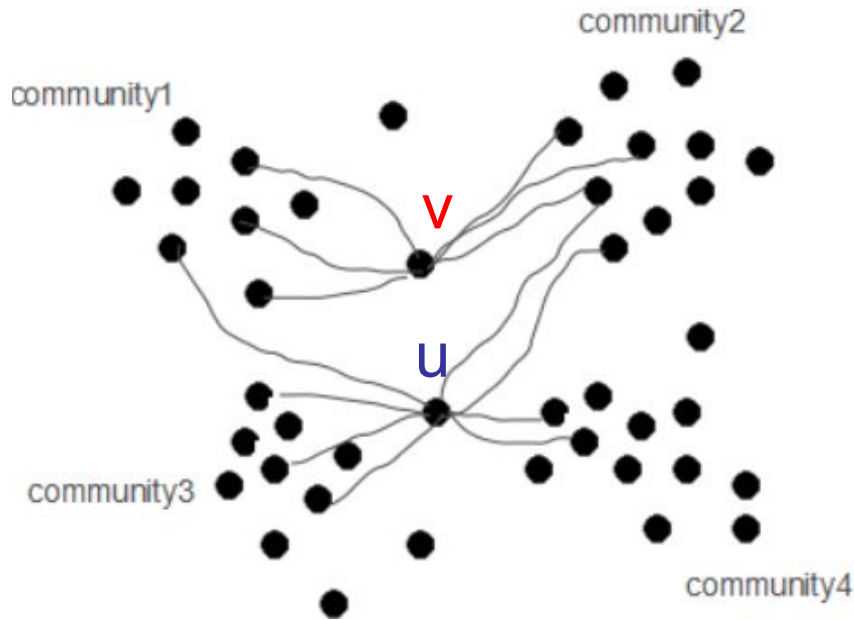
This figure presents a synthetic network consisting of **two communities**. Leaders of two communities (l_1 and l_2) are connected to the node A, while other nodes are all connected to node B. Computing LBC and CBC for all nodes of the graph, the results are as follows: $LBC(A) = 1$, $LBC(B) = 0$, $CBC(A) = 7$, $CBC(B) = 12$, $\forall i : CBC(v_i) = CBC(u_i) = 6$, and $CBC(l_1) = CBC(l_2) = 3$

which is more important, for A, B ?

the measure is not sufficient.

Question and new notions

the CBC values of two nodes are same .



considering the number of distinct communities related with the mediator, there come up with the notion of **diversity score**.

Define two various for the diversity score: DS_{count} , DS_{pair}

DS_{count} is defined as *the number of distinct communities* connected through a node. let $I_d(c_i, v)$ return 1 if $\exists p \in CPaths : s_p \in c_i \wedge v \in p$. and DS_{count} as follows:

$$DS_{count}(v) = \frac{1}{2} \sum_{c_i} I_d(c_i, v) \quad \text{for undirected networks.}$$

(division by 2 is omitted for directed graphs.)

DS_{pair} count **pairs of communities** that have **at least one shortest path** between their numbers passing through node v . $I_d(c_i, c_j, v)$, return 1 if

$$\exists p \in CPaths : s_p \in c_i \wedge e_p \in c_j \wedge v \in p.$$

$$DS_{pair}(v) = \frac{1}{2} \sum_{c_i} \sum_{c_j \neq c_i} I_d(c_i, c_j, v)$$

in this papper, it propose **MedExtractor** algorithm to **identify ranked nodes** connecting the max number of communities as mediators

Algorithm 1 MedExtractor: Find Mediators from SortedList based on their Mediacy Score

```
1: procedure ExtractMediators (Graph  $G$ , OrderedList  $L$ )
2:    $\triangleright G$  is the graph associated with a network
3:    $\triangleright L$  is descending OrderedList containing nodes of the network sorted based on their mediacy score.
4:    $mediatorSet = \{\}$   $\triangleright$  set of selected nodes as mediators
5:    $connectedComs = \{\}$   $\triangleright$  set of communities connected to each other by nodes in mediatorSet
6:   while  $connectedComs.size < G.CommunityCount$  do
7:      $n \leftarrow L.top()$ 
8:     for all Community  $c \in n.incedentCommunities()$  do
9:       if  $c \notin connectedComs$  then
10:        Add  $n$  to  $mediatorSet$ 
11:        Add  $c$  to  $connectedComs$ 
12:       end if
13:     end for
14:      $L.remove(n)$ 
15:   end while
16: end procedure
```

Choosing a Centrality Measure:

the probability of finding more prominent mediators between larger communities is **higher** in comparison to the smaller communities.

➤ **normalized CBC as follows:**

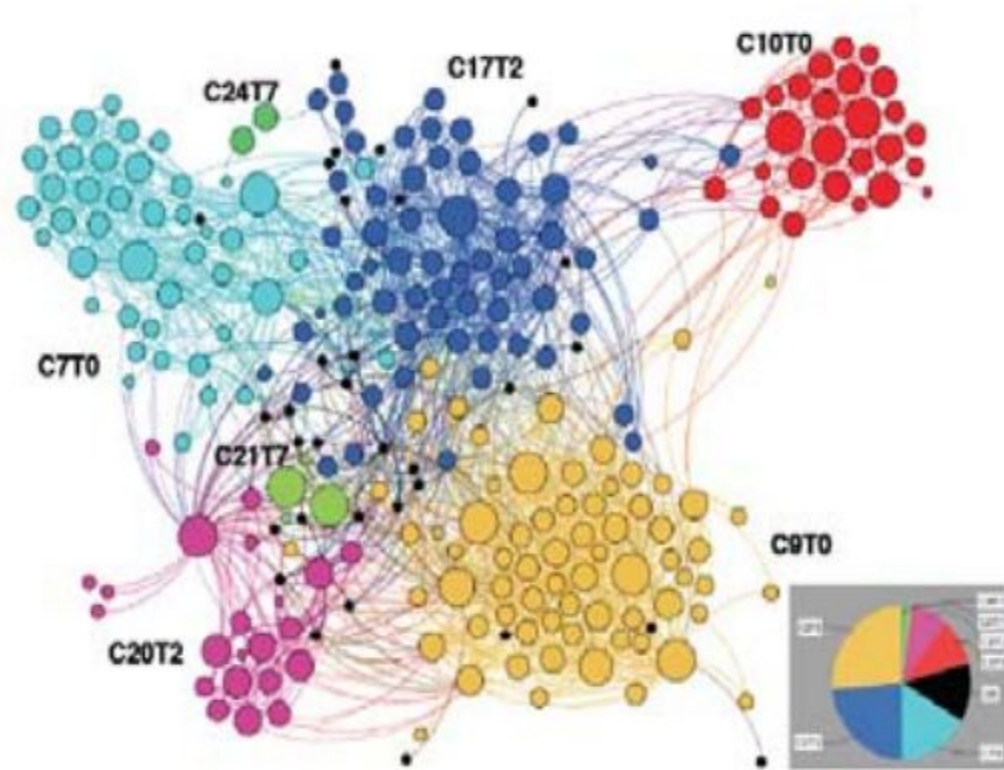
$$NBC(v) = \frac{1}{2} \sum_{p \in CPaths} \frac{I_p(p, v)}{\min(|c_{s_p}|, |c_{e_p}|)}$$

➤ **define mediator score:**

$$MS(v) = NBC(v) \times DS_{count}(v).$$

Experiment

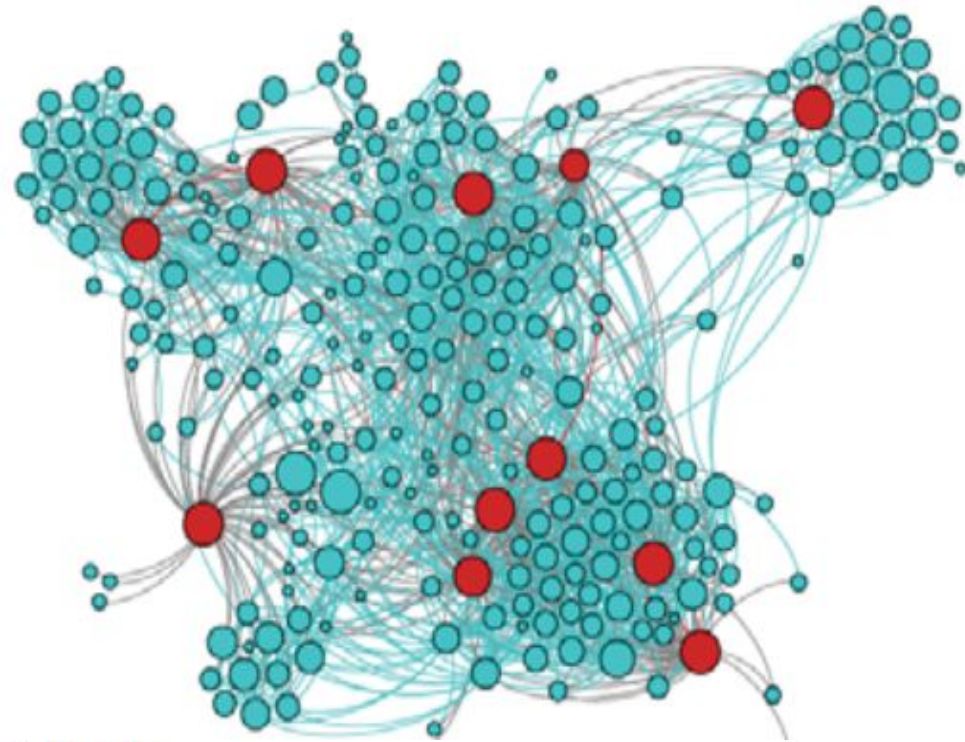
Outsiders



(a)

communities within the Enron email network in August 2001. Colors represent communities except for black that represents outsiders. Size of a node shows its centrality in its community

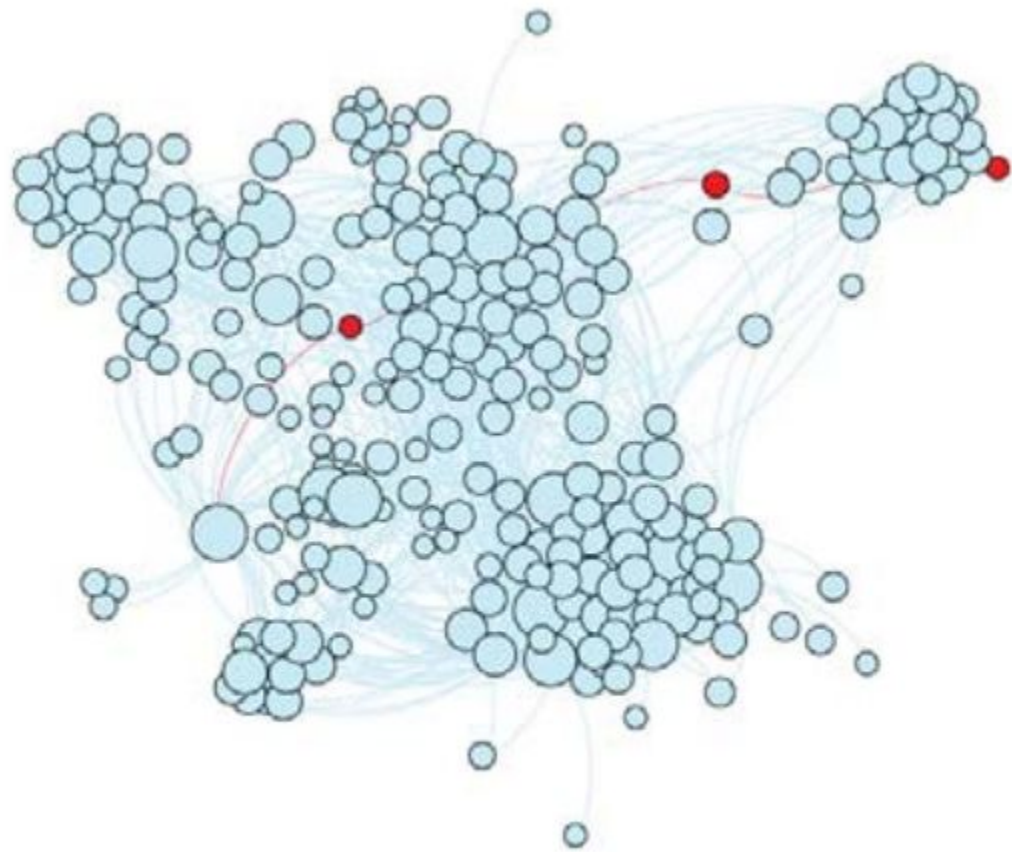
Leaders:



| Email | Community | Position |
|---------------------------|-----------|---|
| jeff.dasovich@enron.com | C9T0 | Executive/Director for State Government Affairs |
| giger.dernehl@enron.com | C9T0 | |
| richard.shapiro@enron.com | C9T0 | VP regulatory affairs (Enron's top lobbyist) |
| kimberly.watson@enron.com | C10T0 | Director |
| d.steffes@enron.com | C9T0 | Vice President |
| kenneth.lay@enron.com | C17T2 | CEO, chairman, and chief executive officer |
| e.haedicke@enron.com | C7T0 | Managing director |
| susan.mara@enron.com | C9T0 | California director of Regulatory Affairs |
| billy.lemmons@enron.com | C17T2 | Vice President |
| becky.spencer@enron.com | C7T0 | |
| l.denton@enron.com | C20T2 | Lawyer |

TABLE I: Leaders of the network shown in Figure 5b, their community affiliation and position in the Enron organization

Outermost



(c)

outermosts showed in red, just two communities have nodes serving as outermosts. From right to left, outermosts are Ava Garcia (probably an [assistant](#) according to the body of some emails), Shirley Crenshaw (probably an [assistant](#)), and Leslie Reves a module [manager](#).

Role change:

| | Jan-01 | Feb-01 | Mar-01 | Apr-01 | May-01 | Jun-01 | Jul-01 | Aug-01 | Sep-01 | Oct-01 | Nov-01 | Dec-01 | |
|----------------------------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|---|
| jeff.dasovich@enron.com | 0 | 0.91 | 0.97 | 1 | 0.91 | 0.86 | 1 | 1 | 0.94 | 0 | 0 | 0 | 8 |
| tana.jones@enron.com | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0.89 | 1 | 0 | 0 | 7 |
| james.steffes@enron.com | 0.82 | 1 | 0.85 | 0.73 | 1 | 0.79 | 0 | 0 | | | 0 | | 6 |
| richard.shapiro@enron.com | 0 | | | | | | 0.82 | 1 | 0.99 | 1 | 0.86 | 0.87 | 6 |
| d.steffes@enron.com | 0 | | | | | | | 0.98 | 0.94 | 0.96 | 1 | 0.85 | 5 |
| marie.heard@enron.com | 0 | | | | | | 1 | 0 | 1 | 0.94 | 1 | 0.78 | 5 |
| becky.spencer@enron.com | 0 | | | 0.98 | 0 | 1 | 1 | 0.96 | 0 | 0 | 0 | 0 | 4 |
| louise.kitchen@enron.com | 0 | | | | | | | | 1 | 0.73 | 1 | 1 | 4 |
| ginger.dernehl@enron.com | 1 | 0 | 0 | 0 | 0.86 | 0 | 0 | 1 | 0 | 0 | 0 | 0.87 | 4 |
| susan.mara@enron.com | 0 | | | | | 0.9 | 0.79 | 0.91 | 0.94 | 0 | 0 | 0 | 4 |
| kathryn.sheppard@enron.com | 0 | | | | | | | | 1 | 1 | 1 | 1 | 4 |
| alan.comnes@enron.com | 0 | | 0.84 | 0 | 0 | 1 | 0 | 0 | 0.92 | 0 | 0 | 0 | 3 |
| mary.hain@enron.com | 0.84 | 0.92 | 1 | 0 | 0 | | | | | | | | 3 |
| janel.guerrero@enron.com | 0 | | | 0.73 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 3 |
| john.lavorato@enron.com | 0 | | | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.91 | 0 | 3 |

- **Changes** of nodes serving as **leaders** in the Enron dataset

| | Jan-01 | Feb-01 | Mar-01 | Apr-01 | May-01 | Jun-01 | Jul-01 | Aug-01 | Sep-01 | Oct-01 | Nov-01 | Dec-01 | |
|---------------------------------------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|---|
| jeff.dasovich@enron.com | 0.043 | 0.069 | 0 | 0.218 | 0.252 | 0.264 | 0.791 | 0 | 0 | 0.114 | 0 | 0 | 7 |
| L.denton@enron.com | | | | | | | 0.547 | 0 | 0.083 | 0.144 | 0.139 | 0 | 4 |
| alan.comnes@enron.com | 0 | | | | | | | | 0.095 | 0 | 0 | 0.114 | 2 |
| rhonda.denton@enron.com | 0 | 0.113 | 0.634 | 0 | 0 | 0 | | | | | | | 2 |
| janet.butler@enron.com | 0 | | | | 0.588 | 0.272 | 0 | 0 | 0 | 0 | 0 | 0 | 2 |
| shelley.corman@enron.com | 0 | | | | | 0.236 | 0 | 0 | 0 | 0 | 0.097 | 0 | 2 |
| cheryl.johnson@enron.com | 0 | | | | | | 0.335 | 0 | 0 | 0 | 0 | 0 | 1 |
| susan.scott@enron.com | 0 | | | 0.285 | 0 | 0 | 0 | 0 | 0 | | | | 1 |
| kam.keiser@enron.com | 0 | | | | | | | 0.486 | 0 | 0 | 0 | 0 | 1 |
| L.nicolay@enron.com | 0 | | | | | | | | 0.046 | 0 | 0 | 0 | 1 |
| stanley.horton@enron.com | 0 | | | | 0.144 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| mark.frevert@enron.com | 0 | | 0.135 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| kenneth.lay@enron.com | 0 | | | | | | | 1 | 0 | 0 | 0 | 0 | 1 |
| tim.belden@enron.com | 0 | | | | | | | | | | | 0.126 | 1 |
| mary.hain@enron.com | 0 | | 0.512 | 0 | 0 | | | | | | | | 1 |
| d.steffes@enron.com | 0 | | | | | | | | | | 0.091 | 0 | 1 |
| janel.guerrero@enron.com | 0 | | | | | | | | | | | 0.095 | 1 |
| deshonda.hamilton@enron.com | 0 | | 0.463 | 0 | 0 | | | | | | | | 1 |
| outlook.team@enron.com | 0 | | | 0.197 | 0 | 0 | 0 | 0 | | | | | 1 |
| stephanie.miller@enron.com | 0 | | | | | | 0.16 | 0 | 0 | 0 | 0 | 0 | 1 |
| k.allen@enron.com | 0 | | | | | | 0.424 | 0 | 0 | 0 | 0 | 0 | 1 |
| bob.ambrocik@enron.com | | | | | | | | | | 0.073 | | | 1 |
| Number of mediators in each timeframe | 1 | 2 | 4 | 3 | 3 | 3 | 5 | 2 | 3 | 3 | 3 | 3 | |

Changes of nodes serving as **mediators** in Ernon dataset

conclusion

In this paper, SSRM was proposed to analyze social networks, based on taking behavior of people, considering the **existence of community structures**.and it indeed identified **outsiders ,outermosts, mediators, and leaders** and found interesting information about the people associated with nodes having the role of a leader or a mediator.it applied to **detect the behavior** of node in social networks.

However , we usually don't know the structure of social network, but we want to model and consider influence of roles, attributes and so on.

Probabilistic Community and Role Model for Social Networks

Yu Han and Jie Tang

Outline

- Introduction
- Main idea
- Analysis
- Formulation and Define
- Model Description
- Experiments

Introduction

- There are not only **visible elements** , but also **invisible elements** to affect the structure of social network
- **Recovery** the structure of social network through many samples

Question:

People's behaviors not only depend on their **own attributes**, but also on their **neighbors** and **communities**.

How to model to capture the **intrinsic relations** between all these element?

How to use a social network model to handle issues such as **community detection** and **behavior prediction**?

Main Idea

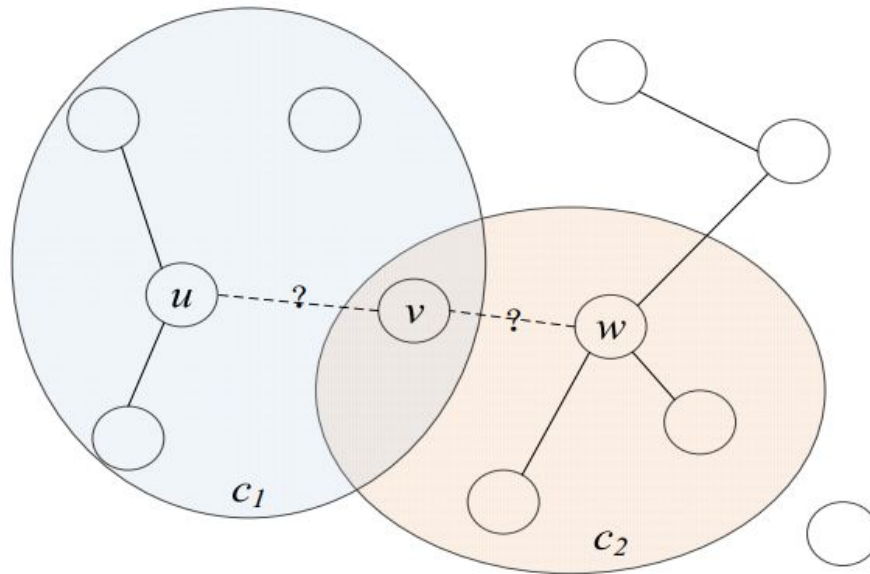
- Propose a **unified probabilistic framework**, the Community Role Model (CRM), to model a social network. CRM incorporates **all the information** of nodes and edges that form a social network. the methods based on **Gibbs sampling** and an **EM algorithm**.

Gibbs sampling or a Gibbs sampler is a **Markov chain Monte Carlo (MCMC) algorithm** for **obtaining a sequence of observations** which are **approximated** from a specified **multivariate probability distribution** (i.e. from the joint probability distribution of two or more random variables), when direct **sampling is difficult**.
-----wikipedia

- CRM can be used not only to **represent** a social network, but also to **handle various application problems** with better performance than a baseline model, without any modification to model.

Analysis

- Community



which probability is higher ?

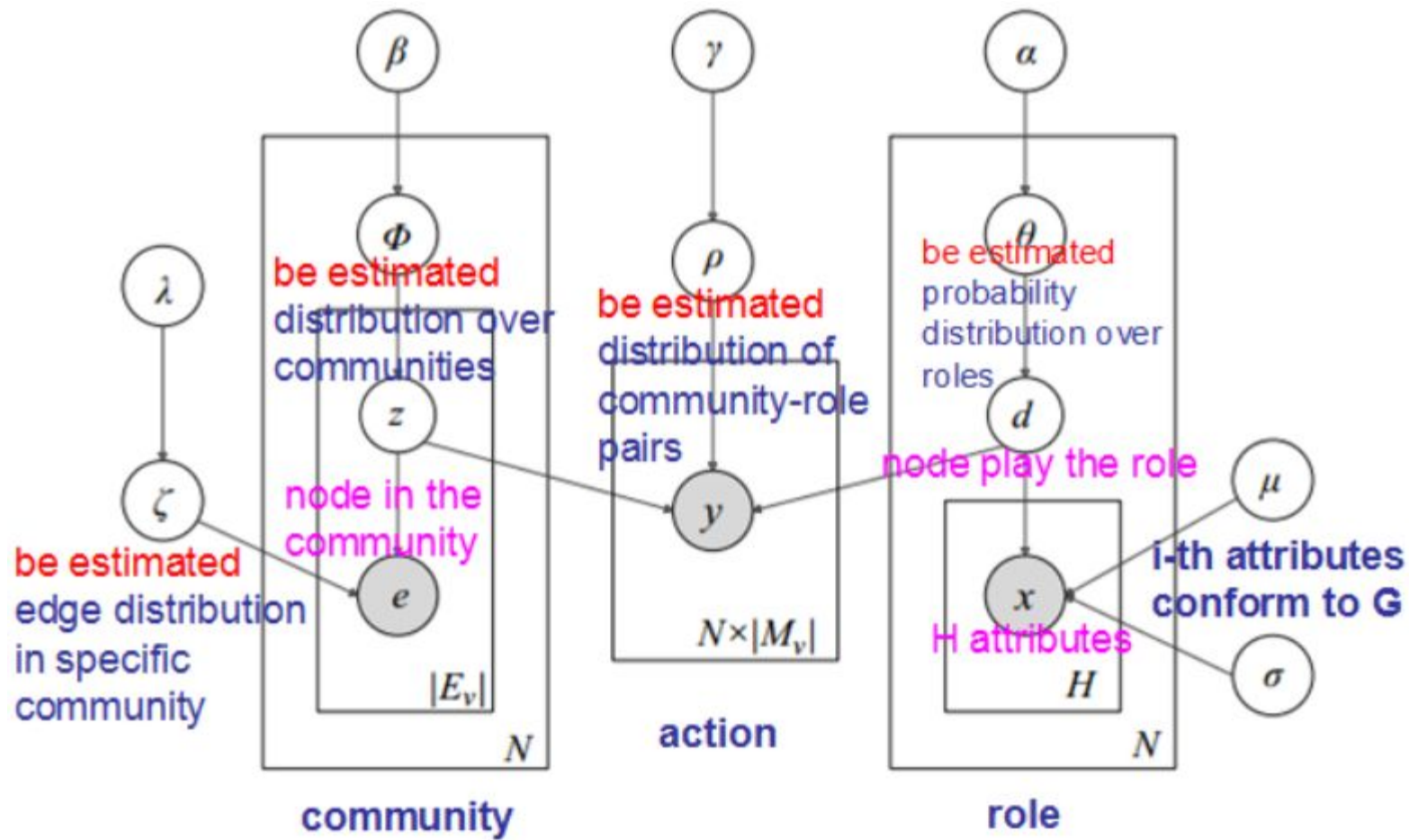
Each node may belong to **several communities**, and whether it has a **link to other nodes** might depend on the **communities to which it belongs**. Thus we assume that each node has **a distribution over the communities**

- Role

- Each node has many **attributes** . we can classify the node into **clusters**,and each cluster can be regarded as a **role that nodes play**.
- The **attributes of each role** satisfy a specific distribution such as Gaussian distribution.Each node has **a distribution over roles**

- Action

- Most nodes tend to take similar actions with nodes in the **same community**.
- Moreover, whether a node takes an action may also depend on the **role** it plays.
- Consider the distribution that the node has **over both communities and roles**



Formulation and Define

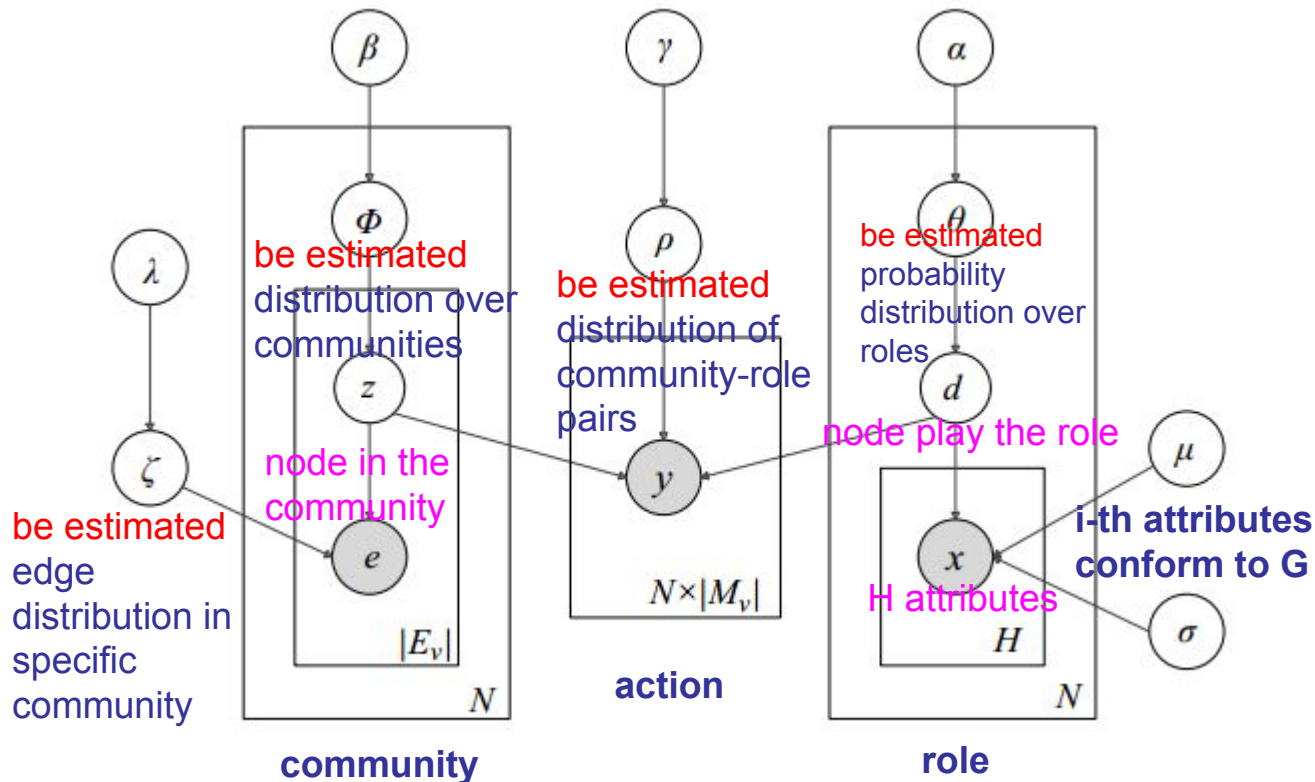
- Structure of a **social network**: $G=(V,E,X)$
- Users(**nodes**): V , $|V|=N$
- **Edges**: E ,an $N*N$ matrix,with each element $e_{v,u} = 0$ or 1 indicating whether user v has a link to user u .
- The set of edges that **associate** with v :
- Notation **X** : $N * H$, **H** is the number of all attributes and element $x_h^{(v)} \in X$ denotes the h -th attribute of user v

- Community:** A social network consists of multiple communities, denoted as $c = [1, 2, \dots, C]$. Each community has a **multinomial distribution** over all pair (u, v) , denoted as ζ . in community c , subject to
$$\sum_{v,u} \zeta_{u,v}^{(c)} = 1$$
- Node Distribution over Communities:** $\phi_c^{(v)}$ denotes the probability for v to be **located in c** . and is subject to
$$\sum_c \phi_c^{(v)} = 1$$
- Role:** a node may play multiple different roles, denoted as $r = [1, 2, \dots, R]$, Each role has a **set of parameters** for the distribution the attributes conform to. Here we use Gaussian distribution. if a node plays role r , its h -th attribute conforms to $N(u_{r,h}, \sigma_{r,h}^2)$

- **Nodes Distribution over Roles**: each node has a **multinomial distribution over roles**, which is denoted as $\theta_r^{(v)}$, the probability for v to play role r and is subject to $\sum_r \theta_r^{(v)} = 1$
- **Action**: For different kinds of social networks, actions take different forms. $y_m^{(v)}$ denote **a repost action** of user v . and set $t=0$ as the start point. During time period $[0, T]$, there are M messages posted by the user that v follows. $y_m^{(v)} = 0$ or 1 to denote whether v reposts the m -th message during a reasonable time period $[0, T]$.
- **Community-Role Pair**: whether a node would take an action depends on the **communities** and the **role** it play.
 - ρ : the distribution of community-role pairs over action.
 - $\rho^{\tau, r}$: the probability for $y_m = 1$, where $\tau = 1(c_u \neq c_v)$.
 that the "community" in "community-role pair" represents whether the node and its target belong to the **same community**

Model Description

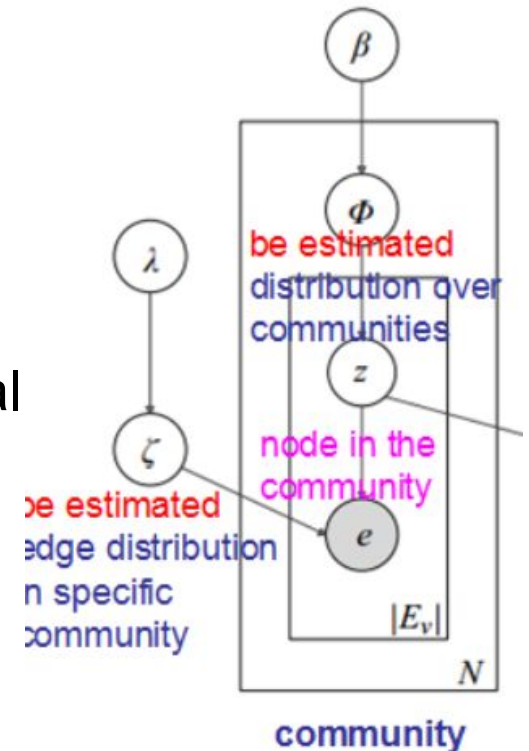
- ❖ **Goal**: devise a **probabilistic generative model**, CRM, to represent a social network
- ❖ **CRM** assumes that a social network can be generated through **three processes**, every process based on **edges, node attributes, and actions**



Edge

For each node v in graph:

1. Draw ζ from $\text{Dirichlet}(\lambda)$;
2. Draw a ϕ_v from $\text{Dirichlet}(\beta)$ prior;
3. For **each edge** $e_{v,i}$
 - Draw a community $z_{v,i} = c$ from multinomial distribution
 - Draw an edge $e_{v,i}$ from a multinomial $\zeta^{(c)}$ specific to community c .



The distribution of the edge E is as:

$$\begin{aligned}
 p(E|\beta, \lambda) = & \int p(\zeta|\lambda) \prod_v \int p(\phi_v|\beta) \\
 & \cdot \prod_{|E_v|} \sum_{z_{v,i}} p(z_{v,i}|\phi_v) p(e_v|z_{v,i}, \zeta) d\phi_v d\zeta.
 \end{aligned} \tag{1}$$

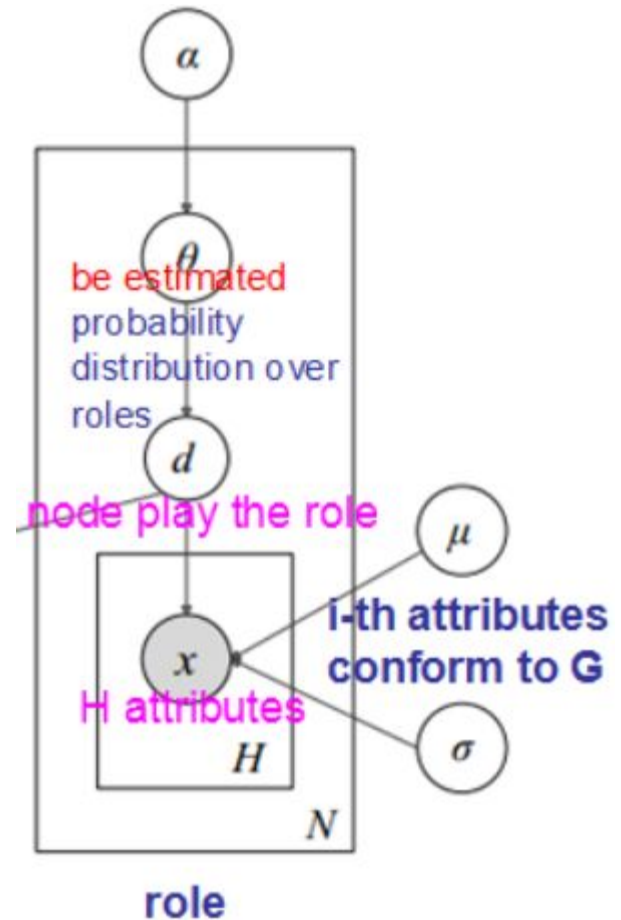
Node

For **each node** v in graph:

1. Draw θ_v from Dirichlet(α) prior;
2. Draw a role $d_v = r$ from multinomial distribution θ_v
3. For **each** attribute of v , draw a value $x_h^{(r)} \sim G(u_{r,h}, \sigma_{r,h}^2)$

The joint distribution of attributes X is defined as :

$$p(X|\alpha, \mu, \sigma) = \prod_v \int p(\theta_v|\alpha) \cdot \sum_{d_v} p(d_v|\theta_v) \prod_h p(x_h^{(v)}|d_v, \mu_{r,k}, \sigma_{r,k}) d\theta_v. \quad (2)$$

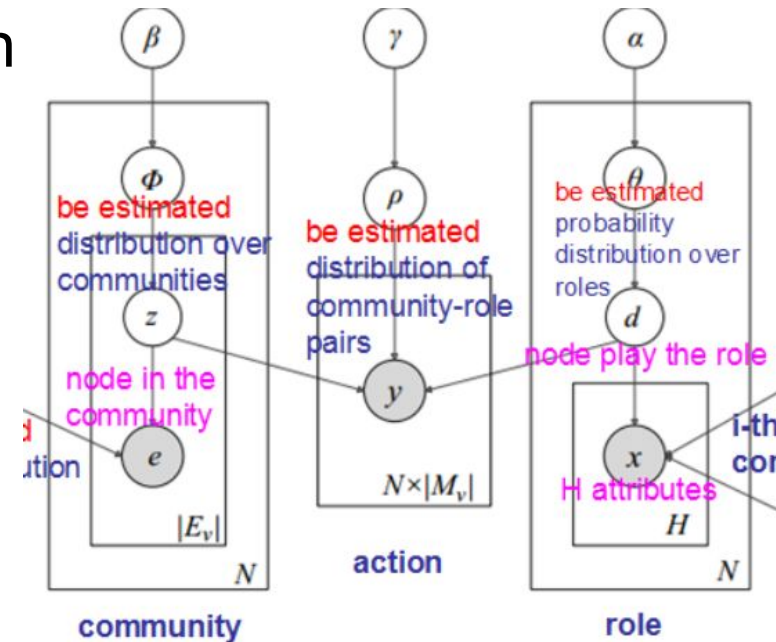


Actions

1. Draw ρ from Dirichlet(γ) prior;
2. Draw a community c_v for v from ϕ_v ;
3. Draw a community c_u for u , which post the message m , from ϕ_u
4. Draw a role r from θ_v ;
5. Draw $y_m \sim \text{Bernoulli}(\rho_{c_v, r})$

The joint distribution fo actions Y is defined as:

$$p(Y|\gamma, \phi, \theta) = \int p(\rho_{\tau, r}) \prod_v \sum_{\tau} \sum_r p(r|\theta_v) p(\tau|\phi_v) p(y_m^{(v)} | \rho_{\tau, r}) d\rho_{\tau, r}. \quad (3)$$



Inference and Parameters Estimation

- Using Gibbs sampling to estimate ζ and ϕ
- The **posterior probability** of $z_{v,i}$ is calculated by

$$p(z_{v,i} = c | \mathbf{z}_{-v,-i}, E) \propto \frac{n_{-v,-i,c}^{(v)} + \beta}{|E_v| + |C|\beta} \frac{n_{-v,-i,c}^{(e)} + \lambda}{n_{-v,-i,\cdot}^{(e)} + |E|\lambda}. \quad (4)$$

- Parameters ζ and ϕ can be estimated by:

$$\phi_{v,c} = \frac{n_{v,c} + \beta}{|E_v| + |C|\beta}, \quad (5)$$

$$\zeta_{c,e} = \frac{n_{c,e} + \lambda}{n_c + |E|\lambda}. \quad (6)$$

- The likelihood of X can be written as:

$$\mathcal{L} = \prod_v \prod_h \sum_{d_v} \frac{\theta_{v,r}}{\sqrt{2\pi}\sigma_{r,h}} e^{-\frac{(x_{v,h} - \mu_{r,h})^2}{2\sigma_{r,h}^2}}. \quad (7)$$

- E-step, estimate the h-th item of θ given the current parameters by:

$$\theta_{v,r} = \frac{\prod_h (2\pi)^{-\frac{1}{2}} \sigma_{r,h}^{-1} e^{-\frac{(x_{v,h} - \mu_{r,h})^2}{2\sigma_{r,h}^2}}}{\sum_{d_v} \prod_h (2\pi)^{-\frac{1}{2}} \sigma_{r,h}^{-1} e^{-\frac{(x_{v,h} - \mu_{r,h})^2}{2\sigma_{r,h}^2}}}. \quad (8)$$

- M-step, update parameters u and σ

$$\mu_{r,h} = \frac{\sum_v \theta_{v,r} x_{v,h}}{\sum_v \theta_{v,r}},$$

$$\sigma_{r,h} = \sqrt{\frac{\sum_v \theta_{v,r} (x_{v,h} - \mu_{r,h})^2}{\sum_v \theta_{v,r}}}.$$

- Only need to estimate ρ , because θ and ϕ have been estimated.

$$p(a_v = \tau, d_v = r | a_{-v}, r_{-v}, \mathbf{y}) \propto (\phi_v \phi_v^T) \theta_v \frac{n_{-v, -m, \tau, r} + \gamma}{|M| + 2|H|\gamma}. \quad (11)$$

- And parameters can be estimated by :

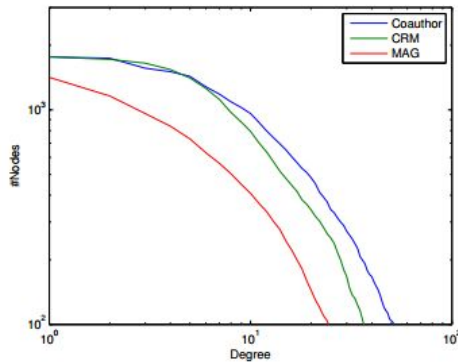
$$\rho = \frac{n_{v,m,\tau,r} + \gamma}{|M| + 2|H|\gamma}. \quad (12)$$

Experiments

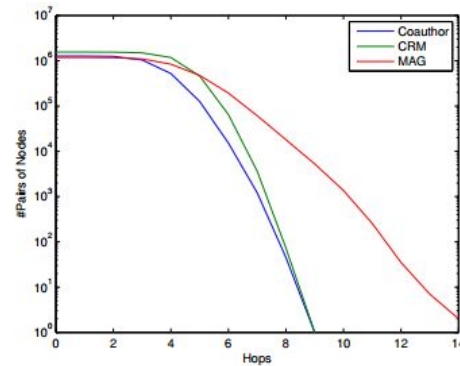
- ❖ Parameters ζ represent communities discovered by CRM
- ❖ Parameters ρ can be used to predict user's actions
- ❖ Structure recovery.

Structure recovery

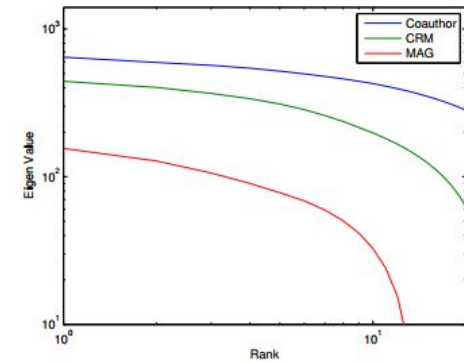
- Comparing with real data, MAG¹.



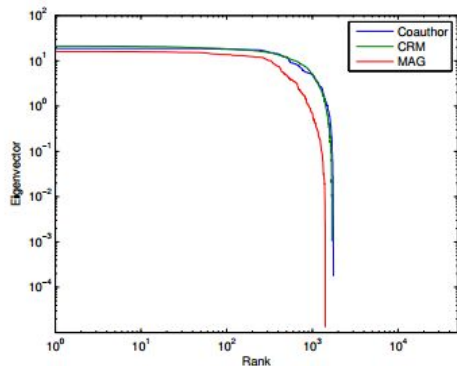
(a) Degree



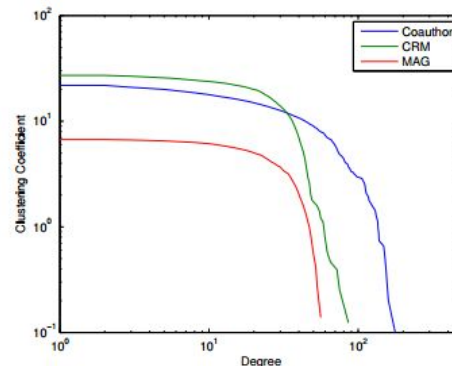
(b) Pairs of Nodes



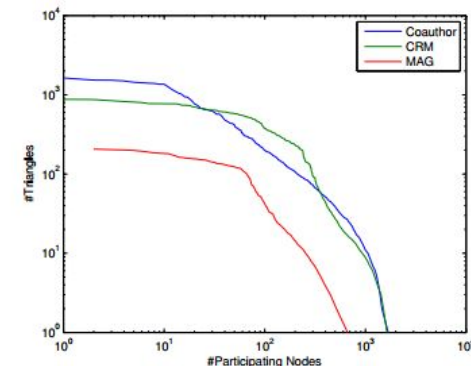
(c) Eigenvalues



(d) Eigenvector



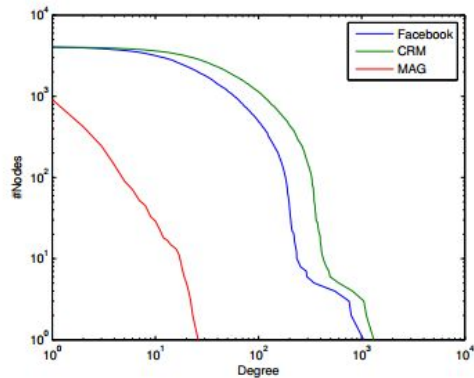
(e) Clustering Coefficient



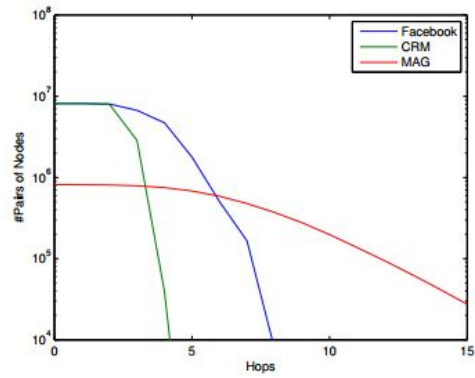
(f) Triangle Participation Ratio

Metric values of the Coauthor network and the two networks generated by CRM and MAG. CRM **outperforms** MAG for every metric

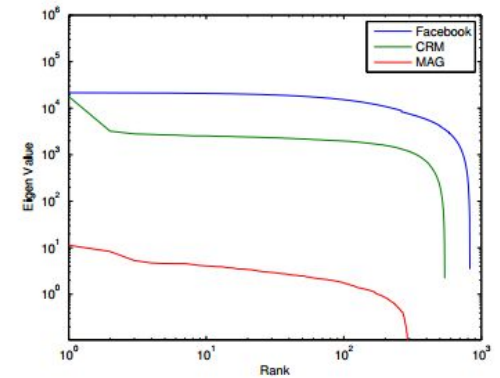
1.M. Kim and J. Leskovec. Modeling social networks with node attributes using the multiplicative attribute graph model. arXiv preprint arXiv:1106.5053, 2011



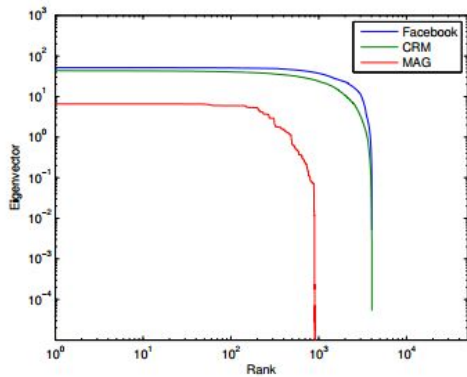
(a) Degree



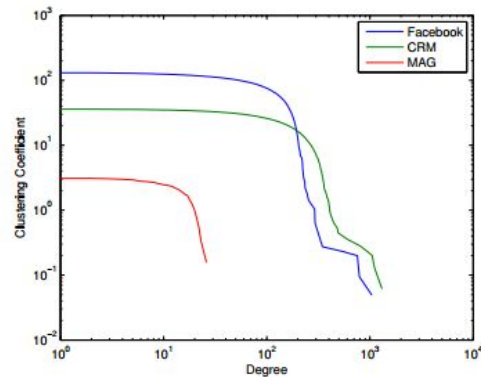
(b) Pairs of Nodes



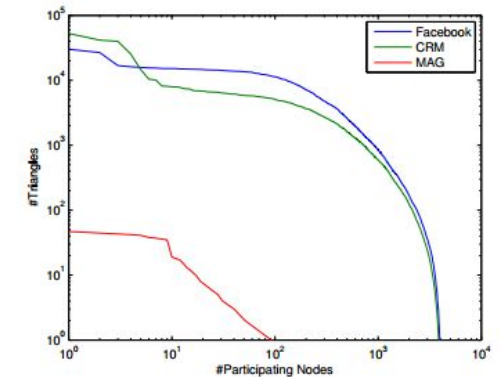
(c) Eigenvalues



(d) Eigenvector



(e) Clustering Coefficient



(f) Triangle Participation Ratio

Metric values of the Facebook network and the two networks generated by CRM and MAG. CRM **outperforms** MAG for every metric

Behavior Prediction

Classifies users into three roles: **opinion leader**, **structural hole spanner**, and **ordinary users**. Whether a user **reposts a message** greatly depends on the role it plays (ρ : a node would take an action depends on the **communities** and the **role** it play. action mean that a **reposted action** of a node).

| Date set | Method | Precision | Recall | F1-measure | AUC |
|----------|--------|-----------------------|-----------------------|-----------------------|-----------------------|
| Coauthor | SVM | 0.8838(0.1725) | 0.5562(0.3183) | 0.6827(0.2054) | 0.7360(0.1111) |
| | SMO | 0.8647(0.1218) | 0.8142(0.1260) | 0.8387(0.1138) | 0.9218(0.0366) |
| | LR | 0.8668(0.1242) | 0.8292(0.1022) | 0.8476(0.1016) | 0.9642(0.0196) |
| | NB | 0.8183(0.1830) | 0.8115(0.1444) | 0.8149(0.1549) | 0.9417(0.0335) |
| | RBF | 0.8552(0.1058) | 0.8353(0.1165) | 0.8451(0.1081) | 0.9477(0.0271) |
| | C4.5 | 0.8328(0.0518) | 0.8015(0.1286) | 0.8169(0.1478) | 0.9065(0.1165) |
| | CRM | 0.8562(0.1490) | 0.8630(0.0598) | 0.8596(0.1013) | 0.9800(0.0199) |
| Weibo | SVM | 0.5067(0.1405) | 0.5027(0.1185) | 0.5047(0.1150) | 0.6068(0.1113) |
| | SMO | 0.5074(0.1464) | 0.5209(0.1099) | 0.5141(0.1271) | 0.6145(0.0363) |
| | LR | 0.5199(0.1306) | 0.5469(0.1073) | 0.5331(0.1157) | 0.6330(0.0377) |
| | NB | 0.5112(0.1245) | 0.5692(0.1083) | 0.5386(0.1172) | 0.6397(0.0394) |
| | RBF | 0.5225(0.1361) | 0.4679(0.1117) | 0.4937(0.1217) | 0.5945(0.0085) |
| | C4.5 | 0.5237(0.1367) | 0.5322(0.1114) | 0.5279(0.1211) | 0.6271(0.1083) |
| | CRM | 0.7017(0.1300) | 0.7305(0.1079) | 0.7158(0.1149) | 0.8174(0.0233) |

Average prediction performance of different methods on the Coauthor and Weibo datasets. The numbers **enclosed in brackets** are standard deviations

Behavior Prediction

CRM achieves much **better performance** than other methods.

| Data Sets | Precision | Recall | F1-measure | AUC |
|-----------|-----------|--------|------------|--------|
| Coauthor | 0.37% | 13.76% | 7.04% | 9.45% |
| Weibo | 36.22% | 40.14% | 38.14% | 32.08% |

Improvement shown by CRM over SVM,SMO,LR,NB,RBF,and C4.5 in terms of precision,recall, F1-measure,and AUC

community detection

detect community with parameters ζ

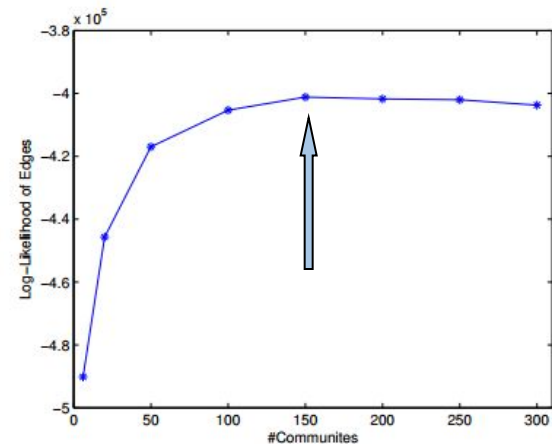
- We must decide the number of communities C before detecting communities with CRM.
- The probability ζ of a edge in different communities is different .compute the sums of log-likelihood for edges and action with :

$$\mathcal{L}(\text{edges}) = \sum_{i=1}^{|E|} \ln p(e_i), \quad (13)$$

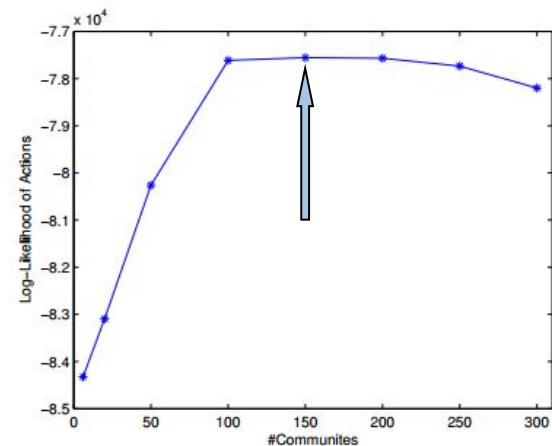
$$\mathcal{L}(\text{actions}) = \sum_{i=1}^{|Y|} \ln p(y_i). \quad (14)$$

❖ $c=150$ may be the best choice .

❖ Through the training of the model, we obtain the community distribution over node



(a) Sum of log-likelihood of edges changes with C



(b) Sum of log-likelihood of actions changes with C

conclusion

- Know how to model a social network through many samples, capturing its information, including structure recovery , behavior prediction
- Applying CRM to real-world datasets, and obtain better performance .

the end , thanks