

E6885 Network Science Lecture 8: ***Dynamic Probabilistic Complex Networks and Graph Database (I)***

Ching-Yung Lin, Dept. of Electrical Engineering, Columbia University

October 28th, 2013



Course Structure

Class Date	Lecture	Topics Covered
09/09/13	1	Overview of Network Science
09/16/13	2	Network Representation and Feature Extraction
09/23/13	3	Network Partitioning, Clustering and Visualization
09/30/13	4	Network Analysis Use Case
10/07/13	5	Network Sampling, Estimation, and Modeling
10/14/13	6	Network Topology Inference
10/21/13	7	Network Information Flow
10/28/13	8	Dynamic & Probabilistic Networks and Graph Database
11/11/13	9	Final Project Proposal Presentation
11/18/13	10	Graph Databases II
11/25/13	11	Information Diffusion in Networks
12/02/13	12	Impact of Network Analysis
12/09/13	13	Large-Scale Network Processing System
12/16/13	14	Final Project Presentation

Dynamic Probabilistic Complex Network Model

The Most Difficult Challenge: State-of-the-Arts?

→ **Our Objectives: Find important people, community structures, or information flow in a network, which is *dynamic*, *probabilistic* and *complex*, in order allocate resources in a large-scale mining system.**

- Social Networks in sociological and statistic fields: focus on (1) overall network characteristics, (2) dynamic random graphs, (3) binary edges, etc. → Not consider probabilistic nodes/edges or individual nodes/edges.
- Epidemic Networks & Computer Virus Network: focus on (1) overall network characteristics – when will an outbreak occurs, (2) regular / random graphs. → Not focus on individual nodes/edges.
- (Computer) Communication Networks: focus on (1) packet transmission – information is not duplicated, or (2) broadcasting – not considering individual nodes/edges or complex network topology.
- WWW: focus on (1) topology description, (2) binary edges and ranked nodes (e.g., Google PageRank) → Not consider probabilistic edges

What is a Dynamic Probabilistic Complex Network?

- Example: <http://smallblue.research.ibm.com>

<http://smallblue.research.ibm.com/publications/netsci2007.pdf>



Modeling a Dynamic Probabilistic Complex Network

- [Assumption] A DPCN can be represented by a Dynamic Transition Matrix $\mathbf{P}(t)$, a Dynamic Vertex Status Random Vector $\mathbf{Q}(t)$, and two dependency functions f_M and g_M .

$$\mathbf{P}(t) @ \begin{bmatrix} \mathbf{p}_{1,1}(t) & \mathbf{p}_{2,1}(t) & \cdots & \cdots & \mathbf{p}_{N,1}(t) \\ \mathbf{p}_{1,2}(t) & \mathbf{p}_{2,2}(t) & & & \mathbf{p}_{N,2}(t) \\ \vdots & \vdots & \ddots & & \vdots \\ \vdots & \vdots & & \ddots & \vdots \\ \mathbf{p}_{1,N}(t) & \mathbf{p}_{2,N}(t) & \cdots & \cdots & \mathbf{p}_{N,N}(t) \end{bmatrix}, \quad \mathbf{Q}(t) @ \begin{bmatrix} \mathbf{q}_1(t) \\ \mathbf{q}_2(t) \\ \vdots \\ \vdots \\ \mathbf{q}_N(t) \end{bmatrix},$$

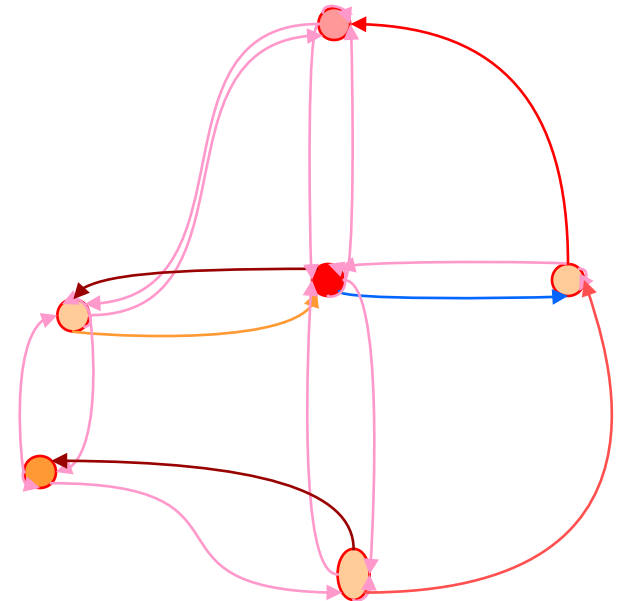
$\mathbf{P}(t + \delta t) @ f_M(\mathbf{Q}(t), \mathbf{P}(t)),$
 $\mathbf{Q}(t + \delta t)$
 $@ g_M(\mathbf{P}(t + \delta t), \mathbf{Q}(t), \mathbf{P}(t)),$

$$\mathbf{p}_{i,j}(t) @ \begin{bmatrix} \Pr(y_{i,j}(t) = SE_1) \\ \Pr(y_{i,j}(t) = SE_2) \\ \vdots \\ \Pr(y_{i,j}(t) = SE_{\Omega_E}) \end{bmatrix}, \quad \mathbf{q}_i(t) @ \begin{bmatrix} \Pr(x_i(t) = SV_1) \\ \Pr(x_i(t) = SV_2) \\ \vdots \\ \Pr(x_i(t) = SV_{\Omega_V}) \end{bmatrix},$$

$$\sum_{\omega \in \Omega_E} \Pr(y_{i,j}(t) = SE_{\omega}) = 1, \quad \sum_{\omega \in \Omega_V} \Pr(x_i(t) = SV_{\omega}) = 1,$$

$x_i(t)$: the status value of vertex i at time t .

$y_{i,j}(t)$: the status value of edge $i \rightarrow j$ at time t .



Modeling a Dynamic Probabilistic Complex Network – cont'd

- Also the Network Topology should follow the characteristics of complex network:

Network topology follows power-law:

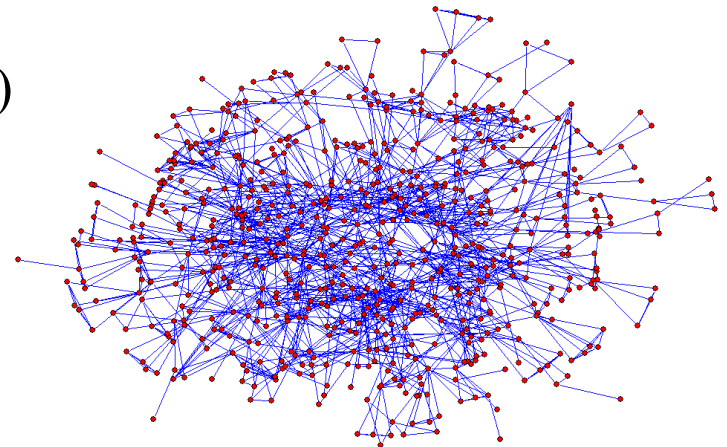
$$\Pr(\sum_i u(p_{i,j}) = l) : S \cdot l^{-d}$$

$$u(p_{i,j}) = \begin{cases} 1, & \text{if } \exists t, \Pr(y_{i,j}(t) \neq \text{null}) > 0 \\ 0, & \text{else} \end{cases}$$

d is typically in the range of $2 \sim 2.5$.

and the clustering coefficient C is typically > 0.2 .

$$C = \Pr(u(p_{j,k}) = 1 \mid u(p_{i,j}) = 1, u(p_{i,k}) = 1)$$



Modeling a Binary DPCN of binary nodes and edges

- A Binary DPCN can be represented by a Dynamic Transition Matrix $P(t)$, a Dynamic Vertex Status Random Vector $Q(t)$, and two dependency functions f_M and g_M .

$$\mathbf{P}(t) @ \begin{bmatrix} p_{1,1}(t) & p_{2,1}(t) & \cdots & \cdots & p_{N,1}(t) \\ p_{1,2}(t) & p_{2,2}(t) & & & p_{N,2}(t) \\ \vdots & \vdots & \ddots & & \vdots \\ \vdots & \vdots & & \ddots & \vdots \\ p_{1,N}(t) & p_{2,N}(t) & \cdots & \cdots & p_{N,N}(t) \end{bmatrix}, \quad \mathbf{Q}(t) @ \begin{bmatrix} q_1(t) \\ q_2(t) \\ \vdots \\ \vdots \\ q_N(t) \end{bmatrix},$$

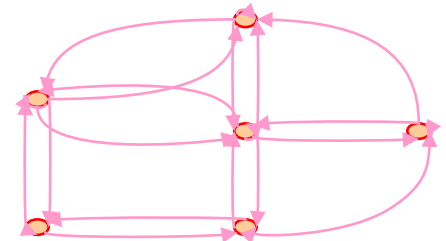
$\mathbf{P}(t + \delta t)$
 $@ f_M(\mathbf{Q}(t), \mathbf{P}(t)),$
 $\mathbf{Q}(t + \delta t)$
 $@ g_M(\mathbf{P}(t + \delta t), \mathbf{Q}(t), \mathbf{P}(t)),$

$$p_{i,j}(t) @ \Pr(\text{edge}_{i \rightarrow j}(t) = 1), \quad q_i(t) @ \Pr(x_i(t) = 1),$$

$$\Pr(\text{edge}_{i \rightarrow j}(t) = 1) + \Pr(\text{edge}_{i \rightarrow j}(t) = 0) = 1,$$

$$\Pr(x_i(t) = 1) + \Pr(x_i(t) = 0) = 1,$$

$x_i(t)$: the status value of vertex i at time t .



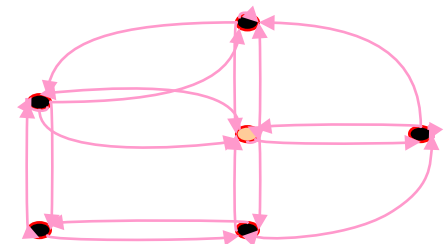
Markov Model is a special case of Binary DPCN

- Markov Model

$$\mathbf{P} @ \begin{bmatrix} p_{1,1} & p_{2,1} & \cdots & \cdots & p_{N,1} \\ p_{1,2} & p_{2,2} & & & p_{N,2} \\ \vdots & \vdots & \ddots & & \vdots \\ \vdots & \vdots & & \ddots & \vdots \\ p_{1,N} & p_{2,N} & \cdots & \cdots & p_{N,N} \end{bmatrix}, \quad \mathbf{Q}(t) @ \begin{bmatrix} q_1(t) \\ q_2(t) \\ \vdots \\ \vdots \\ q_N(t) \end{bmatrix}, \quad \begin{aligned} &\mathbf{Q}(t + \delta t) \\ &@g(\mathbf{P}, \mathbf{Q}(t)) \\ &= \mathbf{P} \cdot \mathbf{Q}(t) \end{aligned}$$

$$\sum_{j=1}^N p_{i,j} = 1 \quad q_i(t) = \Pr(x_i(t) = 1)$$

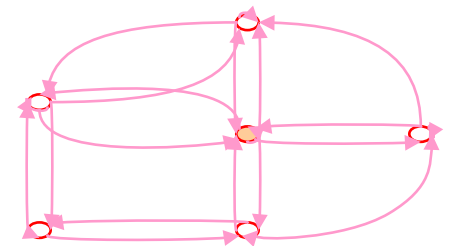
$x_i(t)$: the status value of vertex i at time t .



Many Prior Researches are based on Markov Models

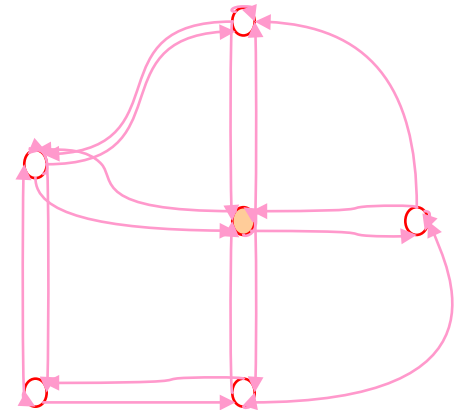
- Random Walks:

$$\lim_{t \rightarrow \infty} Q_V(t) = 1 - (1 - P_E Q_V(0)) \cdot (1 - P_E^{(2)} Q_V(0)) \cdot \dots \cdot (1 - P_E^{(\infty)} Q_V(0))$$



Markov Model is not appropriate to model information flow

- Random Walks assume the existence of a token → unique existence.
- However, information can be duplicated at nodes.
 - New models are needed.



Outline

- Complex Network: Characteristics and Examples
- Dynamic Probabilistic Complex Network
- **Information Flow in Dynamic Probabilistic Complex Network**
- Summary and Conclusion

Information Flow in Dynamic Probabilistic Complex Network (*Let's call it: Behavioral Information Flow (BIF) Model*)

- [Assumption] Edge can be represented by a four-state S-D-A-R (Susceptible-Dormant-Active-Removed) Markov Model. Nodes can be represented by three states S-A-I (Susceptible-Active-Informed) Model.

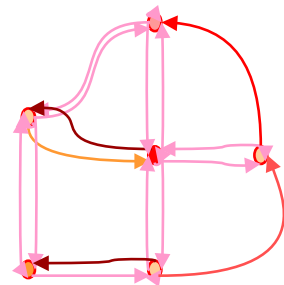
$$\mathbf{P}(t) @ \begin{bmatrix} \mathbf{p}_{1,1}(t) & \mathbf{p}_{2,1}(t) & \cdots & \cdots & \mathbf{p}_{N,1}(t) \\ \mathbf{p}_{1,2}(t) & \mathbf{p}_{2,2}(t) & & & \mathbf{p}_{N,2}(t) \\ \vdots & \vdots & \ddots & & \vdots \\ \vdots & \vdots & & \ddots & \vdots \\ \mathbf{p}_{1,N}(t) & \mathbf{p}_{2,N}(t) & \cdots & \cdots & \mathbf{p}_{N,N}(t) \end{bmatrix}, \quad \mathbf{Q}(t) @ \begin{bmatrix} \mathbf{q}_1(t) \\ \mathbf{q}_2(t) \\ \vdots \\ \mathbf{q}_N(t) \end{bmatrix}, \quad \begin{array}{l} \mathbf{P}(t + \delta t) \\ @f(\mathbf{M}, \mathbf{Q}(t), \mathbf{P}(t)), \\ \\ \mathbf{Q}(t + \delta t) \\ @g(\mathbf{P}(t + \delta t), \mathbf{Q}(t), \mathbf{P}(t)), \end{array}$$

$$\mathbf{p}_{i,j}(t) = \begin{bmatrix} \Pr(y_{i,j}(t) = S) \\ \Pr(y_{i,j}(t) = D) \\ \Pr(y_{i,j}(t) = A) \\ \Pr(y_{i,j}(t) = R) \end{bmatrix} @ \begin{bmatrix} \sigma_{i,j} \\ \psi_{i,j} \\ \mu_{i,j} \\ \rho_{i,j} \end{bmatrix},$$

$$\sigma_{i,j} + \psi_{i,j} + \mu_{i,j} + \rho_{i,j} = 1$$

$$\mathbf{q}_i(t) = \begin{bmatrix} \Pr(x_i(t) = S) \\ \Pr(x_i(t) = A) \\ \Pr(x_i(t) = I) \end{bmatrix} @ \begin{bmatrix} \lambda_i \\ \eta_i \\ \nu_i \end{bmatrix},$$

$$\lambda_i + \eta_i + \nu_i = 1$$

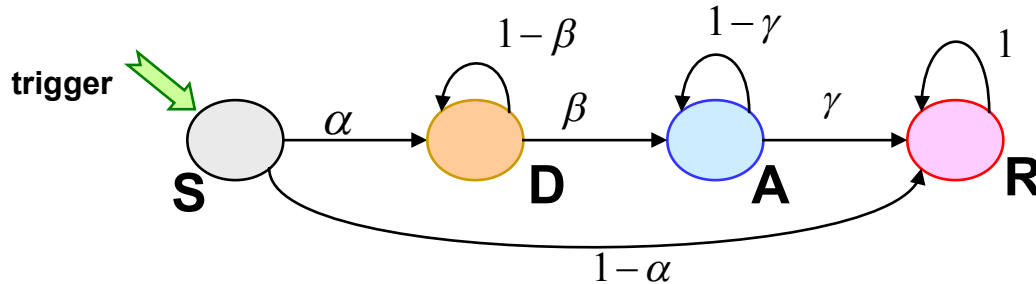


Major Difference between BIF and Prior Modeling Methods in Epidemic Research and Computer Virus Fields

- Model Human Nodes as S-I-R (Susceptible, Infected, and Removed).
- Did not consider individual node's behavior distinctly in network structure/topology → did not consider edge status.
- We propose to model edge status as (autonomous) S-D-A-R Markov Model (Susceptible, Dormant, Active, Removed)
- We propose to model human node behavior as S-A-I (Susceptible, Active, and Informed).

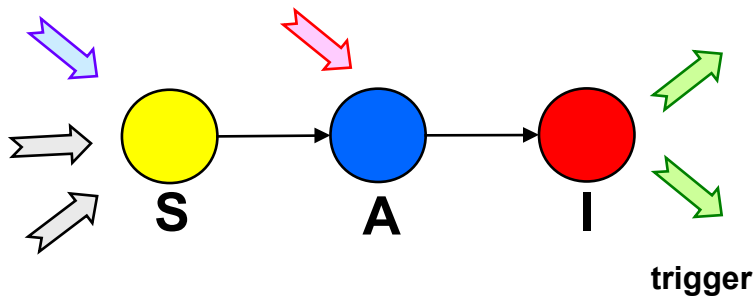
Edges are Markov State Machines, Nodes are not

- State transitions of edges: S-D-A-R model. (Susceptible, Dormant, Active, and Removed) This indicates the time-aspect changes of the state of edges.

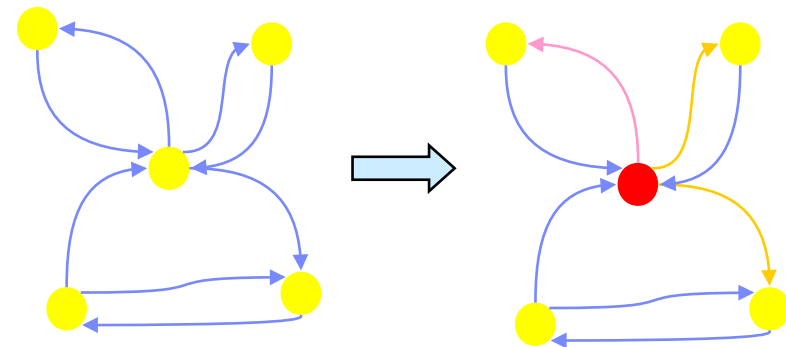


Edge view

- States of nodes: S-A-I model. (Susceptible, Active, and Informed) Trigger occurs when the start node of the edge changes from state S to state I :



Node view



Network view

Edge State Probability and Network Configuration Model

- Nodes and Edges

$$\mathbf{P}(t + \delta t) = f(\mathbf{M}, \mathbf{Q}(t), \mathbf{P}(t)),$$

- Network Configuration Model (which is learned by training). It includes the network topology information, long-term edge probability, and delay parameter).

$$\mathbf{M} @ \begin{bmatrix} (\alpha_{1,1}, \beta_{1,1}, \gamma_{1,1}) & (\alpha_{2,1}, \beta_{2,1}, \gamma_{2,1}) & \cdots & \cdots & (\alpha_{N,1}, \beta_{N,1}, \gamma_{N,1}) \\ (\alpha_{1,2}, \beta_{1,2}, \gamma_{1,2}) & (\alpha_{2,2}, \beta_{2,2}, \gamma_{2,2}) & & & (\alpha_{N,2}, \beta_{N,2}, \gamma_{N,2}) \\ \vdots & \vdots & \ddots & & \vdots \\ \vdots & \vdots & & \ddots & \vdots \\ (\alpha_{1,N}, \beta_{1,N}, \gamma_{1,N}) & (\alpha_{2,N}, \beta_{2,N}, \gamma_{2,N}) & \cdots & \cdots & (\alpha_{N,N}, \beta_{N,N}, \gamma_{N,N}) \end{bmatrix},$$

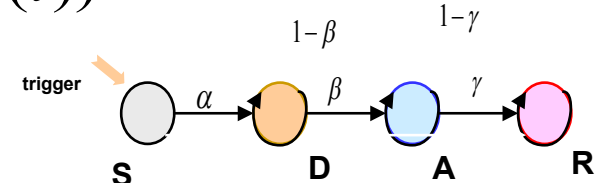
- $\alpha_{ij} = 0 \rightarrow$ No Edge between i and j
- Our KDD 2005 paper is a special case that $\alpha_{ij} = 1$ or 0 , and did not model $(\beta_{ij}, \gamma_{ij})$

Edge State Probability Update function *s.t*

$$\mathbf{P}(t + \delta t) = f(\mathbf{M}, \mathbf{Q}(t), \mathbf{P}(t))$$

- Given three different cases:

1. On trigger: $x_i(t - \delta t) \neq I, x_i(t) = I$



$$\mathbf{p}_{i,j}(t + \delta t) = \begin{bmatrix} \sigma'_{i,j} \\ \psi'_{i,j} \\ \mu'_{i,j} \\ \rho'_{i,j} \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 & 0 \\ \alpha_{i,j} & 1 - \beta_{i,j} & 0 & 0 \\ 0 & \beta_{i,j} & 1 - \gamma_{i,j} & 0 \\ 1 - \alpha_{i,j} & 0 & \gamma_{i,j} & 1 \end{bmatrix} \begin{bmatrix} \sigma_{i,j} \\ \psi_{i,j} \\ \mu_{i,j} \\ \rho_{i,j} \end{bmatrix} @\mathbf{F} \cdot \mathbf{p}_{i,j}(t),$$

2. No trigger – node not informed yet: $x_i(t - \delta t) \neq I, x_i(t) \neq I$

$$\mathbf{p}_{i,j}(t + \delta t) = \mathbf{p}_{i,j}(t),$$

3. No trigger – node has been informed: $x_i(t - \delta t) = I, x_i(t) = I$

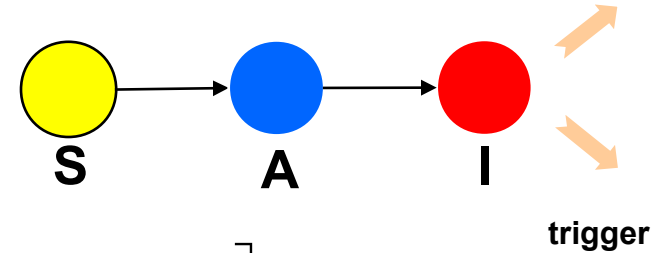
$$\mathbf{p}_{i,j}(t + \delta t) = \mathbf{F} \cdot \mathbf{p}_{i,j}(t),$$

- Therefore, consider the probabilities of node states, then we get $f(\cdot)$:

$$\mathbf{p}_{i,j}(t + \delta t) = v_i \cdot \mathbf{F} \cdot \mathbf{p}_{i,j}(t) + (1 - v_i) \cdot \mathbf{p}_{i,j}(t)$$

Nodes: State Transitions Determined by Incoming Edges

$$\mathbf{Q}(t + \delta t) = g(\mathbf{P}(t + \delta t), \mathbf{Q}(t), \mathbf{P}(t)),$$

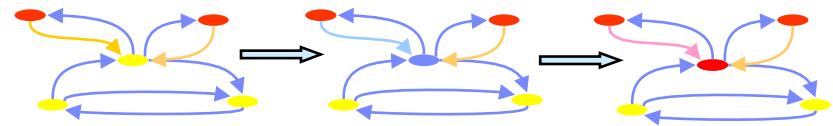


- Node State Probability Update Function $g(\cdot)$:

$$\mathbf{q}_i(t + \delta t) = \begin{bmatrix} \lambda'_i \\ \eta'_i \\ \nu'_i \end{bmatrix} = \begin{bmatrix} \prod_{n \in \Omega_{V,i}} (1 - \mu'_{n,i}) & 0 & 0 \\ 1 - \prod_{n \in \Omega_{V,i}} (1 - \mu'_{n,i}) & \prod_{n \in \Omega_{V,i}} (1 - \gamma_{n,i}) \mu_{n,i} & 0 \\ 0 & 1 - \prod_{n \in \Omega_{V,i}} (1 - \gamma_{n,i}) \mu_{n,i} & 1 \end{bmatrix} \begin{bmatrix} \lambda_i \\ \eta_i \\ \nu_i \end{bmatrix} @ \mathbf{Q} \cdot \mathbf{q}_i(t),$$

$$\Pr(\exists n \in \{1 \dots N\}, y_{n,i}(t + \delta t) = R, y_{n,i}(t) = A)$$

$$= 1 - \prod_{n \in \Omega_{V,i}} (1 - \gamma_{n,i}) \mu_{n,i}$$



and $\Omega_{V,i}$ is the set of all source nodes of the incoming edges of Node i : $\Omega_{V,i} = \{n \mid \forall n \in \{1 \dots N\}, \alpha_{n,i} > 0\}$

Network view

Two special considerations for information propagation behavior

- No Reverse Propagation:
 - Add an update criteria to $f(\cdot)$:

$$if(y_{i,j}(t + \delta t) = R, y_{i,j}(t) = A) \Rightarrow y_{j,i}(t + \delta t) = R$$

- This constraint does not affect $Q(t)$.
- It makes the probabilities change to: $\tilde{\psi}'_{j,i} = \tilde{\mu}'_{j,i} = 0$
 and $\tilde{\rho}'_{j,i} = \rho'_{j,i} + \psi'_{j,i} + \mu'_{j,i}$

- No Simultaneously Communication from one person (e.g., phone calls):
 - Add a constraint criteria to $f(\cdot)$:

$$if(y_{i,j}(t) = A) \Rightarrow \forall m \in \Omega_{U,i} \neq j, y_{i,m}(t) \neq A$$

- And, also: $\Omega_{U,i} = \{m \mid \forall m \in \{1 \dots N\}, \alpha_{i,m} > 0\}$

where $\Omega_{U,i}$ is the set of all end nodes of the outgoing edges of Node i :

- The probabilities should be:

$$if(y_{k,i}(t) = A) \Rightarrow \forall n \in \Omega_{V,i} \neq k, y_{n,i}(t) \neq A$$

$$\mathbf{p}_{i,m}(t + \delta t) = \mathbf{p}_{i,m}(t) \quad \text{and} \quad \mathbf{p}_{n,i}(t + \delta t) = \mathbf{p}_{n,i}(t)$$

An Application of Information Flow Prediction – find important people

- Who are the most likely people to talk about this information at a specific time given the current observation?

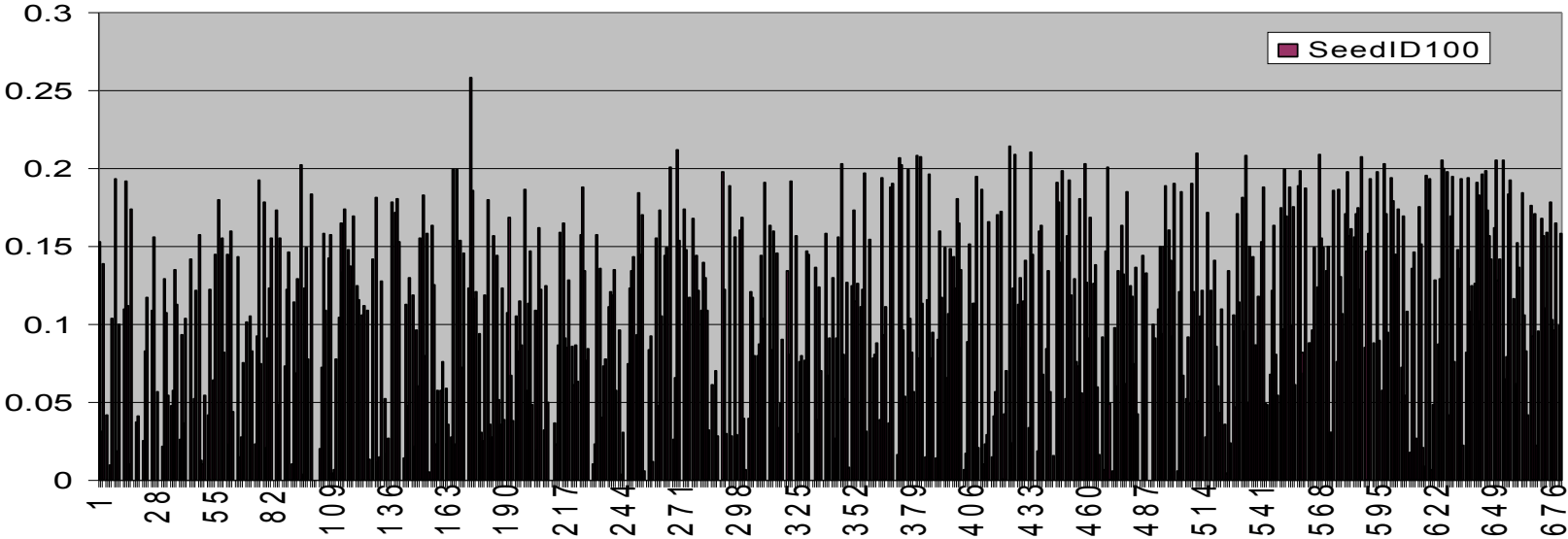
$$(m, n) = \arg \max_{m, n \in \{1 \dots N\}} (\mu_{m, n}(t + \tau)) \text{ given } \mathbf{Q}(t) \text{ or } (\mathbf{P}(t), \mathbf{Q}(t))$$

- For a given concrete observation, the values in the given priors $\mathbf{P}(t), \mathbf{Q}(t)$ are either 0 or 1.
- For speaker recognition results, the priors can be confidence values between $0 \sim 1$.

Predicting behavioral information flow – Algorithm I

- Monte Carlo Method: Simulate each DPCN information flow for 1000 times.
- It takes 12 seconds to use MC simulation to predict the process. (For a given model and test 679 nodes, it takes a PC 130 mins for calculate the probabilities if the information flow starts from different 679 seeds).

The Probabilities of the Nodes Receives Information



Outline

- Complex Network

- Dynamic Probabilistic Complex Network

- Information Flow in Dynamic Probabilistic Complex Network
 - Who should we monitor? Where to put the sniffers?
 - Training, Effect of Noises, and Dynamic Model Updates in Dynamic Probabilistic Complex Network

- Demo

- Next Steps
 - Communities in Dynamic Probabilistic Complex Network

Training the Network Configuration

- Train the Network Configuration Model \mathbf{M} using the observation data in Time $0 - T$:

$$\alpha_{i,j} = \frac{L}{K}$$

$$\beta_{i,j} = \frac{1}{E[w]}$$

$$\gamma_{i,j} = \frac{1}{E[d]}$$

where K is the number of times that node x_i becomes active during time period $0-T$. L is the count of the number of times that edge y_{ij} becomes active.

In the Markov Model of SDAR, the duration of the information staying in the D state is a Poisson distribution with mean value = $1/\beta_{ij}$. We can then estimate the parameter β_{ij} based on the mean waiting time $E[w]$ of training data. Similarly, we can get γ_{ij} based on the mean active duration $E[d]$ of training data.

Impact of Classification Error on BIF Model

- Consider two types of errors:
 - Speaker Recognition Error
 - E.g. DIG scenario → nodes may have miss and false alarm. This types of error would cause edge error.
 - Topic Detection Error
 - E.g. classification of email content → nodes are correct. Edges may have miss and false alarm.
 - Combination of both errors
 - E.g., if we are doing both topic detection and speaker recognition in DIG, then the above two types of error will be combined.

An Application of Information Flow Prediction – enhance speaker recognition accuracy

- The probability that a given pair talks about this information at a specific time given the current observation?

$$\Pr(y_{i,j}(t + \tau) = A) = \mu_{i,j}(t + \tau) \quad \text{given } \mathbf{Q}(t) \text{ or } (\mathbf{P}(t), \mathbf{Q}(t))$$

- These probabilities can serve as prior confidence for speaker recognition.

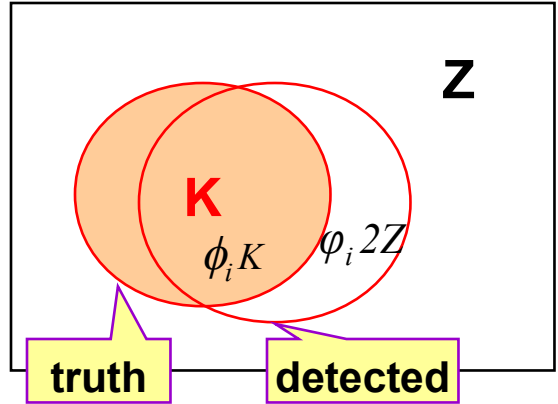
- Assume the classification precision rate on the speaker (node) i is ϕ_i and the false alarm rate on the speaker i is φ_i .

- Then the expected number of times that the node is counted is:

$$L = \phi_i \phi_j L + \varphi_i \varphi_j Z$$

- And the link is counted is: $\tilde{K} = \phi_i K + \varphi_i \cdot 2Z$

- Therefore,
$$\tilde{\alpha}_{i,j} = \frac{\tilde{L}}{\tilde{K}} = \frac{\phi_i \phi_j L + \varphi_i \varphi_j Z}{\phi_i K + \varphi_i \cdot 2Z}$$



- If we assume a universal precision and false alarm rate at all speakers, then:
$$\tilde{\alpha}_{i,j} = \frac{L}{\tilde{K}} = \frac{\phi^2 L + \varphi^2 Z}{\phi_i K + \varphi_i \cdot 2Z}$$

Assume the average waiting time of links and the average transmission duration of links are the same regardless of the links observed, then:

$$\beta_{i,j} = \beta_{i,j} \quad \text{and} \quad \gamma_{i,j} = \gamma_{i,j}$$

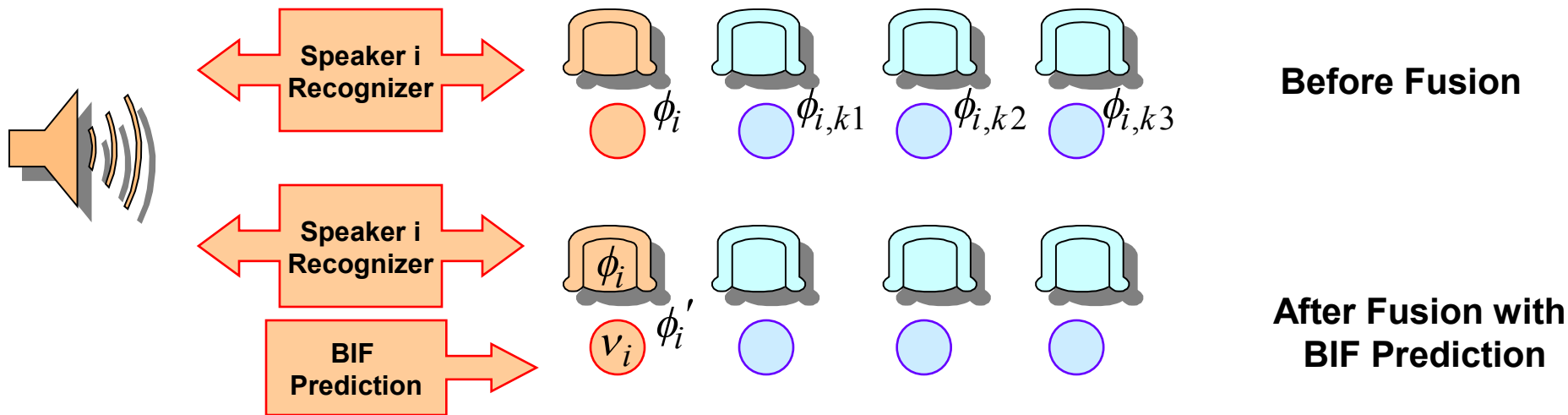
- If we assume the false alarm rate is small and can be neglected when the number of nodes is large, then
$$\alpha_{i,j} \approx \phi \cdot \alpha_{i,j}$$

Speaker Recognition Accuracy can be Improved by Fusion of Original Speaker Recognition and Predicted Node Probability

- We can use this fusion method to combine both speaker recognition result and the estimated node probability:

$$\phi'_i = \frac{\phi_i \cdot v_i}{\phi_i \cdot v_i + \sum_k \phi_{i,k} \cdot v_k}$$

which is guaranteed to be increasing when $v_i > \forall v_k$



Dynamic Updates on Speaker Recognition Result and the BIF Model

- Assume a special case that a speaker is usually mistakenly classified as the other speaker. E.g., given a true Speaker i speaking, her voice is sometimes classified as Speaker k . (but not the reverse direction)
- Let (n) represents the n -th dynamic update of the model. Each time the model is updated based on Slide ‘Impact of Classification Error on Model’
- Based on the previous slide, we shall get:

$$\phi_{(n)} = \frac{\phi \cdot v_{i(n)}}{\phi \cdot v_{i(n)} + (1 - \phi) \cdot v_{k(n)}} \quad \text{where} \quad \kappa_{(n)} = \frac{\phi}{\phi + (1 - \phi) \cdot \kappa_{(n)}} \quad \text{where} \quad \kappa_{(n)} = \frac{v_{k(n)}}{v_{i(n)}}$$

- Based on the previous two slides, we can get: $\kappa_{(n)} = \frac{1 - \phi_{(n-1)}}{\phi_{(n-1)}} \cdot \kappa_{(n-1)}$

- This value can be calculated as: $\kappa_{(n)} = \left[\frac{1 - \phi}{\phi} \right]^n \cdot \kappa_{(0)}$

which quickly converges to zero. $\rightarrow \lim_{n \rightarrow \infty} \phi_{(n)} = 1$

Dynamic Updates on Speaker Recognition Result and the Model – cont'd

- Assume another special case that a speaker is usually mistakenly classified as the other speaker. E.g., given a true Speaker i speaking, her voice is sometimes classified as Speaker k . And, also, Speaker k 's voice can be confused as Speaker i .
- Following similar steps as in the previous slide, we shall get:

$$\phi_{(n)} = \frac{\phi \cdot v_{i(n)}}{\phi \cdot v_{i(n)} + (1 - \phi) \cdot v_{k(n)}} @ \frac{\phi}{\phi + (1 - \phi) \cdot \kappa} \quad \text{where } \kappa_{(n)} = \kappa$$

- If we assume the confusion error is not uniformly the same, i.e., asymmetric error between speakers, then:

$$\phi_{i,(n)} = \frac{\phi_i \cdot v_{i(n)}}{\phi_i \cdot v_{i(n)} + (1 - \phi_i) \cdot v_{k(n)}} @ \frac{\phi_i}{\phi_i + (1 - \phi_i) \cdot \kappa_{(n)}}$$

$\kappa_{(n)} @ \frac{v_{k(n)}}{v_{i(n)}}$ which depends on the network topology and does not have a closed-form solution.

Noise Factor II – Impact of Classification Error from Topic Classification

- Assume the classification precision rate on the edge $i \rightarrow j$ is $\theta_{i,j}$, and the false alarm rate on the edge is $\omega_{i,j}$.
- Then the expected number of times that the edge is counted is:

$$\hat{K}_i = \theta_{i,j}K + \omega_{i,j}Z$$

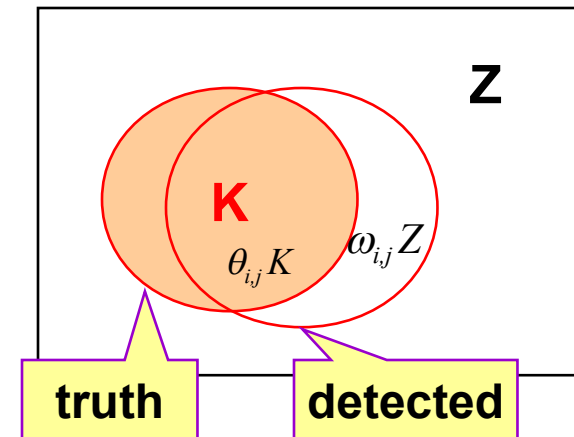
- And the link is counted is:

$$\hat{L} = \theta_{i,j}L + \omega_{i,j}Z$$

- Therefore,
$$\hat{\alpha}_{i,j} = \frac{\hat{L}}{\hat{K}} = \frac{\theta_{i,j}L + \omega_{i,j}Z}{\theta_{i,j}K + \omega_{i,j}Z}$$

- If the false alarm rate can be neglected:

$$\hat{\alpha}_{i,j} = \frac{\theta_{i,j}L}{\theta_{i,j}K} = \alpha_{i,j}$$



Noise Factor II – Impact of Classification Error from Topic Classification – cont'd

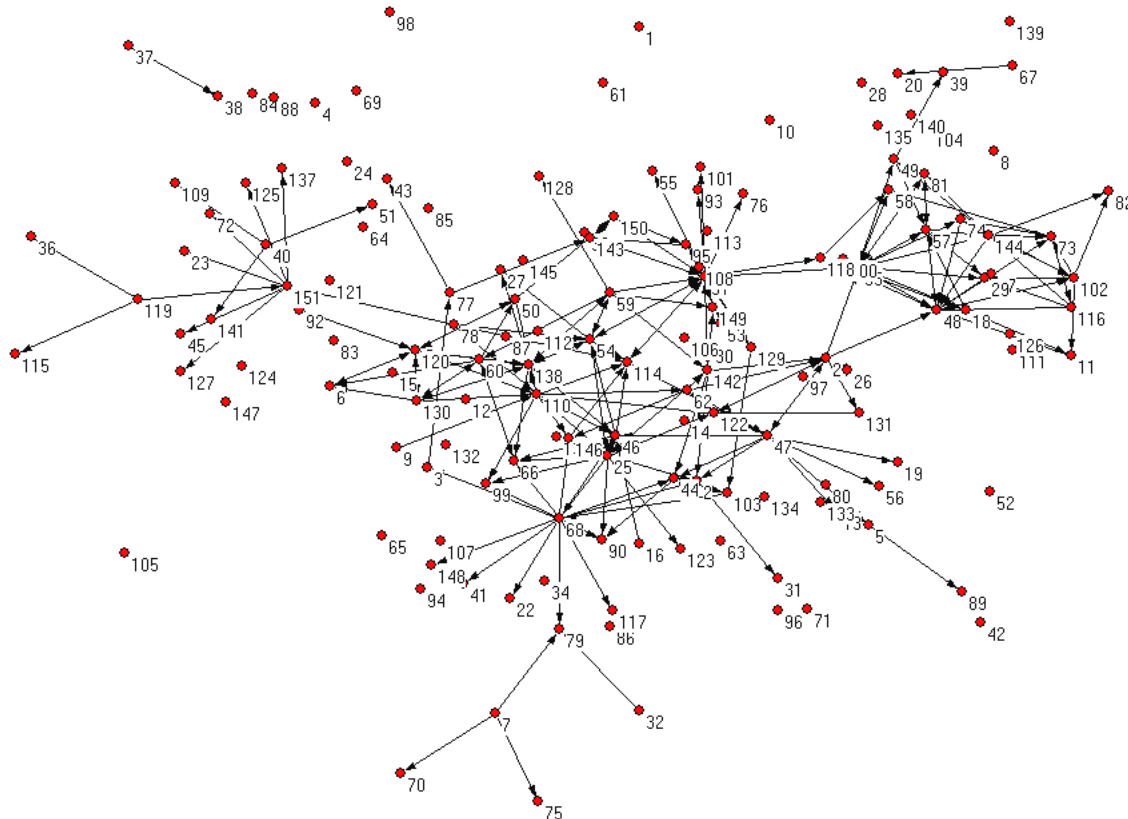
- If the false alarm rate can be neglected → The probability that a classification error occurs at an conversation record is equal to the probability that the nodes are detected.
- Therefore, since the propagation coefficient α 's are the same, there will be no effects on the information flow prediction.
- If we consider both the topic classification error and the speaker recognition error together, the information flow prediction will be the same as the case when there is only speaker recognition error.

Multiple Topics

- Each topic is one information flow model. Multiple Topics can be considered as a combination of these models.

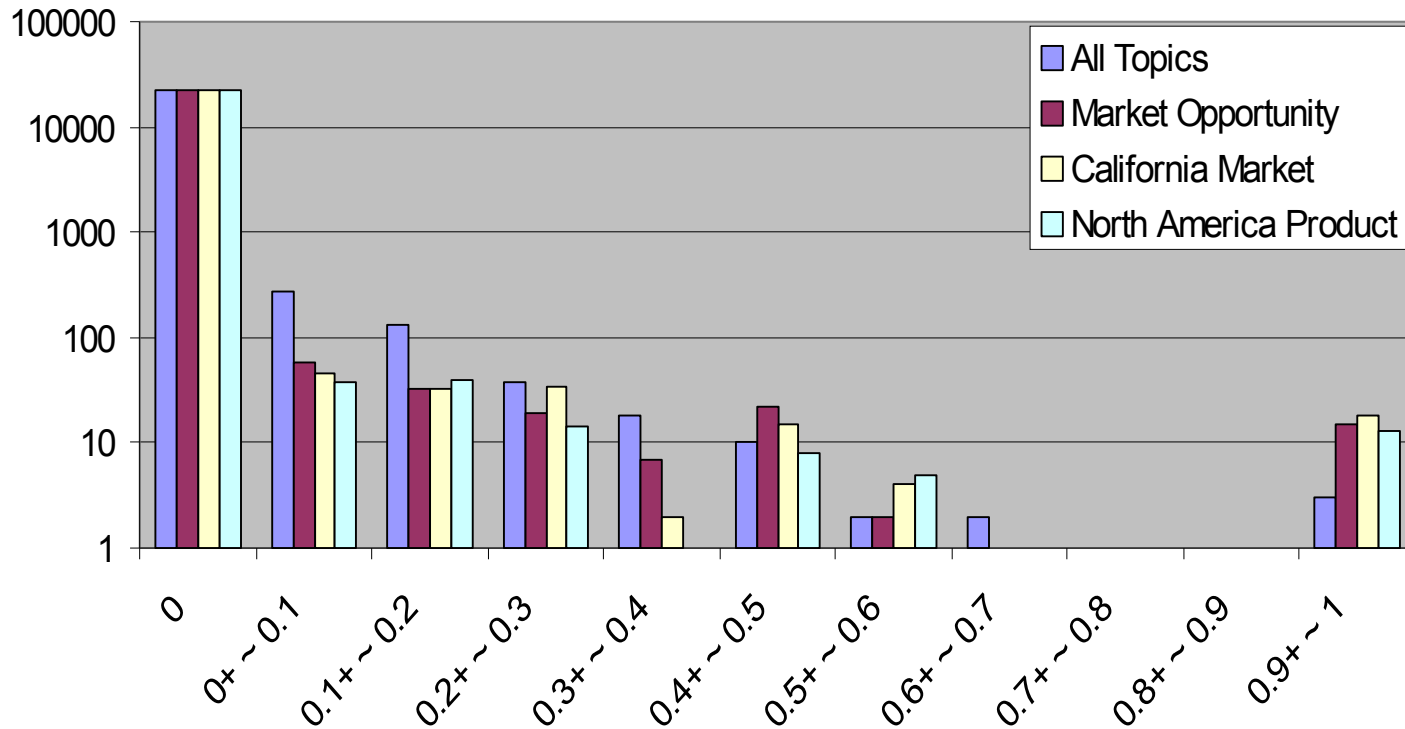
We can first find the experts and then see how this community works (II)

- Rosalee Fleming played an important role at “Market Opportunities.” She received info from Actor 119 (Mike Carson) and Actor 23 (James Steffes – VP of Gov. Affairs of Enron.)
- Actor 68 (Rod Havslett -- CFO) is also a major information spreader.



We can estimate the parameters in the DPCN model

- Example: a histogram of the alpha values by applying the DPCN model on Enron Dataset.



Graph Databases – See Presentations