

Scan Statistics on Graphs

# “Us” vs “Them”

Anomaly Detection in Streaming Graphs

Carey E. Priebe

Department of Applied Mathematics & Statistics  
Johns Hopkins University, Baltimore, MD, USA

June 2012  
Bristol, England

## Collaborators



Heng Wang



Youngser Park



Minh Tang

John Conroy, David Marchette, Andrey Ruhkin, Nam Lee, . . .

# Outline

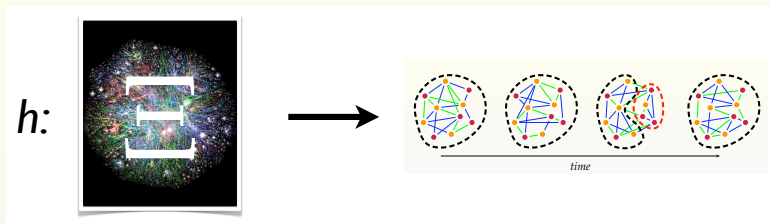
Introduction

Definitions

Theory & Simulation

Conclusions & Discussion

# Introduction



$$h_V : \Xi \rightarrow \mathcal{K}_V = \{\text{red}, \text{orange}\}$$

$$h_E : \Xi \rightarrow \mathcal{K}_E = \{\text{green}, \text{blue}\}$$

The map  $h$  provides a time series of (vertex- and edge-attributed) graphs

$$\{G_t\} = \{G(\Xi_t, h)\} = \{(V_t, E_t)\}$$

# Motivation

## Us:



*C.E. Priebe, J.M. Conroy, D.J. Marchette, Y. Park*  
"Scan Statistics on Enron Graphs", *Computational & Mathematical Organization Theory*, Vol 11, No 3, pp 229-247, 2005

## Them:

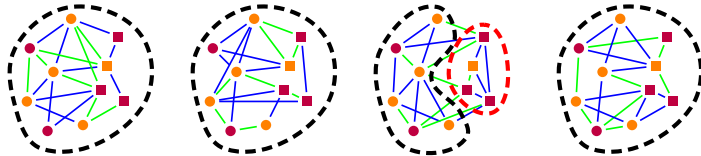


*X. Wan, J. Janssen, N. Kalyaniwalla and E. Milios,*  
"Statistical analysis of dynamic graphs", *Proceedings of AISB06: Adaptation in Artificial and Biological Systems*, v3, pp.176-179, 2006.

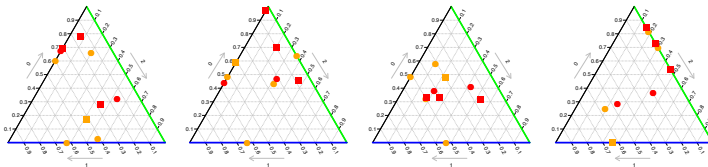


*X. Wan, N. Kalyaniwalla,*  
"Capturing causality in communications graphs", *DIMACS/DyDAn Workshop on Computational Methods for Dynamic Interaction Networks*, 2007.

# A Latent Process Model For Time Series of Attributed Graphs



time



time

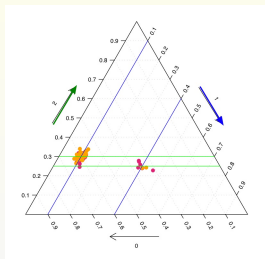
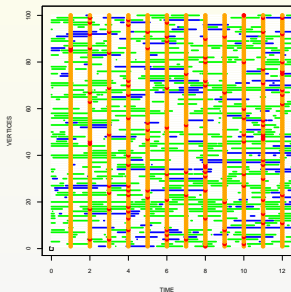


N.H. Lee and C.E. Priebe,

"A Latent Process Model for Time Series of Attributed Random Graphs", *Statistical Inference for Stochastic Processes*, Vol. 14, No. 3, pp. 231-253, 2011.

# A Latent Process Model For Time Series of Attributed Graphs

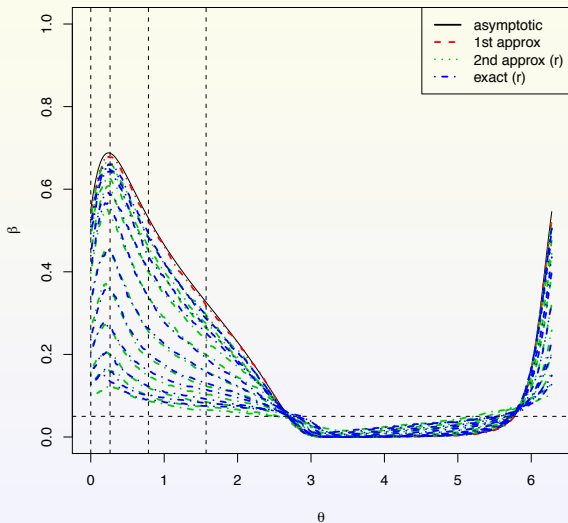
Lee & Priebe (SISP, 2011)



$$Q_i = \begin{pmatrix} -1 & & & 1 \\ & \ddots & & \vdots \\ & & -1 & 1 \\ \frac{\pi_{i,1}}{\pi_{i,d}} & \dots & \frac{\pi_{i,d-1}}{\pi_{i,d}} & -\frac{\sum_{k=1}^{d-1} \pi_{i,k}}{\pi_{i,d}} \end{pmatrix}$$

# A Latent Process Model For Time Series of Attributed Graphs

Lee & Priebe (SISP, 2011)





## Scan Statistics

“moving window analysis” [1922 R.A. Fisher, 1965 J. Naus]:

to scan a small “window” (*scan region*) over data, calculating some *locality statistic* for each window; e.g.,

- number of events for a point pattern,
- average pixel value for an image.

scan statistic  $\equiv$  maximum of locality statistics:

If maximum of observed locality statistics is large, then the inference can be made that

*there exists a subregion of excessive activity  $\implies$  detection!*

# Scan Statistics on Graphs

Let  $G = (V, E)$  be a graph.

$(k^{th})$  neighborhood of  $v$ :  $N_k[v; G] = \{w \in V : d(v, w) \leq k\}$ ,  $k \geq 0$ ,

$(k^{th})$  scan region of  $v$ :  $\Omega(N_k[v; G])$ ,

$(k^{th})$  locality statistic of  $v$ :  $\Psi_k(v; G) = |E(\Omega(N_k[v; G]))|$ ,

$(k^{th})$  scan statistic of  $G$ :  $M_k(G) = \max_{v \in V(G)} \Psi_k(v; G)$ .

“Maximum activity in  $k$ -neighborhood”

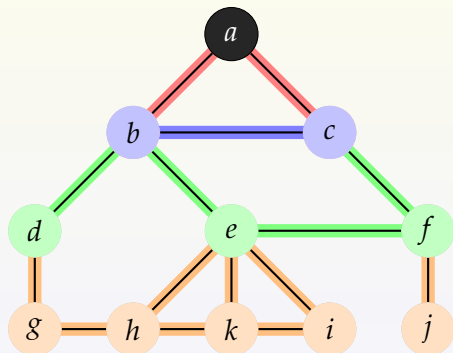


C.E. Priebe, J.M. Conroy, D.J. Marchette, Y. Park

“Scan Statistics on Enron Graphs”, Computational & Mathematical Organization Theory, Vol 11, No 3, pp 229-247, 2005

## Example of Scan Statistics

$$|V| = 11, |E| = 15, \Psi_k(v; G) = |E(\Omega(N_k[v; G]))|$$



$k$	$E(\Omega(N_k[\mathbf{a}; G]))$	$\Psi_k(\mathbf{a})$
0	$\ominus$	2
1	$\ominus + \ominus$	3
2	$\ominus + \ominus + \ominus$	7
3	$\ominus + \ominus + \ominus + \ominus$	15

	$d$	$\Psi_1$	$\Psi_2$	$\Psi_3$
a	2	3	7	15
b	4	5	14	15
c	3	4	8	15
d	2	2	8	14
e	5	7	15	15
f	3	3	12	15
g	2	2	7	14
h	3	4	8	15
i	2	3	7	15
j	1	1	3	12
k	3	5	10	15

# Scan Statistics and Time Series

Let  $\{G_t\}$ ,  $t = 1, \dots, t_{max}$ , be a time series of graphs.

$(k^{th})$  scan region:  $\Omega(N_k(v; G_t))$ .

$(k^{th})$  locality statistic:  $\Psi_k(v; G_t) = |E(\Omega(N_k(v; G_t)))|$ .

$(k^{th})$  scan statistic:  $M_k(G_t) = \max_{v \in V(G_t)} \Psi_k(v; G_t)$ .

# Us vs Them

**Us:**

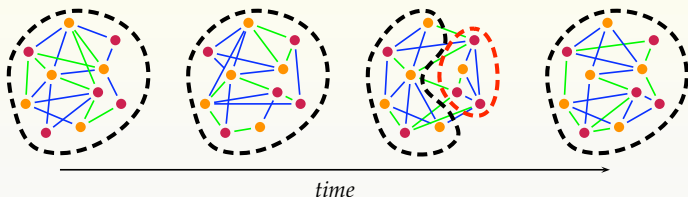
$$\Psi_{t;t'}^k(v) = |E(\Omega(N_k[v; G_{t'}]; G_{t'}))|$$

**Them:**

$$\Phi_{t;t'}^k(v) = |E(\Omega(N_k[v; G_t]; G_{t'}))|$$

$$t' \leq t$$

# Normalization



vertex normalization:

$$\tilde{J}_{t,\tau}^k(v) = \begin{cases} J_{t,t}^k(v) & \tau = 0 \\ \frac{J_{t,t}^k(v) - \hat{\mu}_{t,\tau}^k(v)}{\max(\hat{\sigma}_{t,\tau}^k(v), 1)} & \tau > 0 \end{cases}$$

$$\hat{\mu}_{t,\tau}^k(v) = \frac{1}{\tau} \sum_{t'=t-\tau}^{t-1} J_{t,t'}^k(v)$$

$$\hat{\sigma}_{t,\tau}^k(v) = \sqrt{\frac{1}{\tau-1} \sum_{t'=t-\tau}^{t-1} (J_{t,t'}^k(v) - \hat{\mu}_{t,\tau}^k(v))^2}$$

## Us vs Them

**Us:**

$$\Psi_{t;t'}^k(v) = |E(\Omega(N_k[v; G_{t'}]; G_{t'}))|$$

**Them:**

$$\Phi_{t;t'}^k(v) = |E(\Omega(N_k[v; G_t]; G_{t'}))|$$

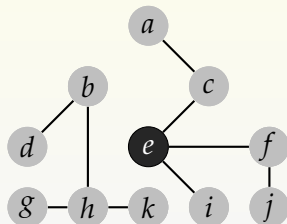
$$t' \leq t$$

$$\hat{\mu}_{t,\tau}^k(v) = \frac{1}{\tau} \sum_{t'=t-\tau}^{t'-1} J_{t,t'}^k(v)$$

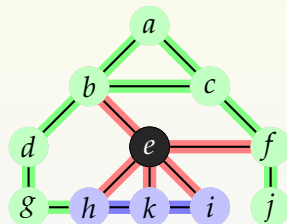
# Normalization

$$\tau = 1$$

$$t = t^* - 1$$



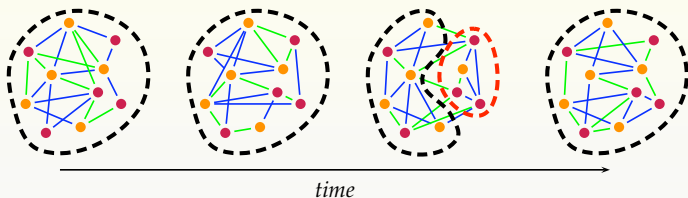
$$t = t^*$$



	$J_{t^*, t^*}^k(e)$		$\hat{\mu}_{t^*, \tau}^k(e)$		$\tilde{J}_{t^*, \tau}^k(e)$	
	$k=0$	$k=1$	$k=0$	$k=1$	$k=0$	$k=1$
us	5	7	3	3	2	4
them	5	7	2	4	3	3



# Normalization



temporal normalization:

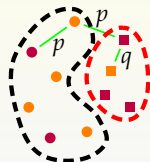
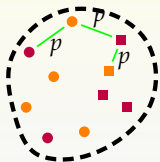
$$S_{t,\tau,\ell}^k = \begin{cases} \tilde{M}_{t,\tau}^k = \max_v(\tilde{J}_{t,\tau}^k(v)) & \ell = 0 \\ \frac{\tilde{M}_{t,\tau}^k - \tilde{\mu}_{t,\tau,\ell}^k}{\max(\tilde{\sigma}_{t,\tau,\ell}^k, 1)} & \ell > 0 \end{cases}$$

$$\tilde{\mu}_{t,\tau,\ell}^k = \frac{1}{l} \sum_{t'=t-l}^{t'-1} \tilde{M}_{t',\tau}^k$$

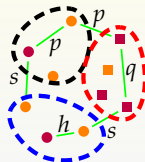
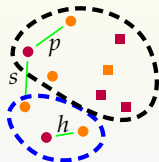
$$\tilde{\sigma}_{t,\tau,\ell}^k = \sqrt{\frac{1}{l-1} \sum_{t'=t-l}^{t'-1} (\tilde{M}_{t',\tau}^k - \tilde{\mu}_{t,\tau,\ell}^k)^2}$$

# Heterogeneous Null

homogeneous null:  $p < q$



heterogeneous null:  $p \leq s < \{h, q\}$



# Theory

$$H_0 : \{Q_v\} = \{Q_1^0(t), \dots, Q_n^0(t)\}, \forall t$$

$$H_A : \begin{cases} \{Q_v\} = \{Q_1^0(t), \dots, Q_n^0(t)\}, & t < t^* - 1 \\ \{Q_v\} = \{Q_1^A(t), \dots, Q_m^A(t), Q_{m+1}^0(t), \dots, Q_n^0(t)\}, & t \geq t^* - 1 \end{cases}$$

1st approximation  $\implies$  stochastic block model:

$$H_0 : G_t \stackrel{\text{iid}}{\sim} SBM(P_{B \times B}, n_{B \times 1}), \forall t$$

$$H_A : \begin{cases} G_t \stackrel{\text{iid}}{\sim} SBM(P_{B \times B}, n_{B \times 1}), & t \leq t^* - 1 \\ G_t \stackrel{\text{iid}}{\sim} SBM(P_{B \times B} + \text{diag}(0, \dots, 0, \delta), n_{B \times 1}), & t \geq t^* \end{cases}$$

# Theory

## maxdegree, 1st approximation

$$\beta(\text{us}) - \beta(\text{them})$$

$$(\ell = 0, \tau = 1)$$

### Theorem

$$\lim_{n \rightarrow \infty} S \stackrel{\mathcal{L}}{=} \sum_{c=1}^{C_{T,H}} \pi_{T,H,c} g(\cdot; \theta_{T,H,c}).$$

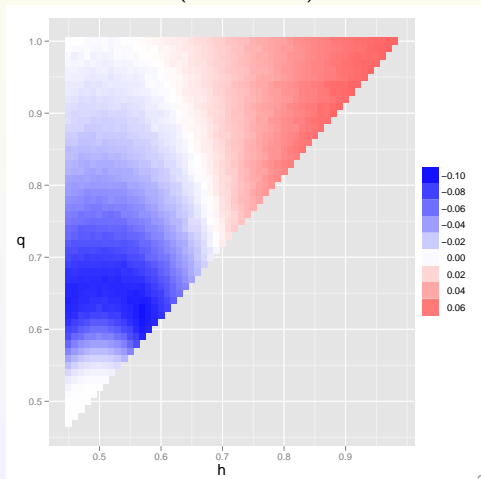
$g$ : Gumbel

$C$ : number of components

$\pi$ : mixture coefficients

$T \in \{\Psi, \Phi\}$

$H \in \{H_0, H_A\}$



# Theory

## maxdegree, 1st approximation

$$\beta(\text{us}) - \beta(\text{them})$$

( $\ell = 0, \tau = 1$ )

### Theorem

$$\lim_{n \rightarrow \infty} S \stackrel{\mathcal{L}}{=} \sum_{c=1}^{C_{T,H}} \pi_{T,H,c} g(\cdot; \theta_{T,H,c}).$$

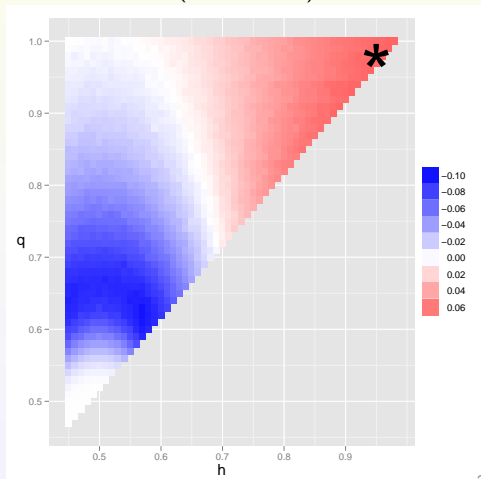
$g$ : Gumbel

$C$ : number of components

$\pi$ : mixture coefficients

$T \in \{\Psi, \Phi\}$

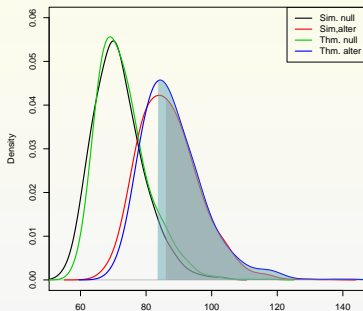
$H \in \{H_0, H_A\}$



# Theory & Simulation

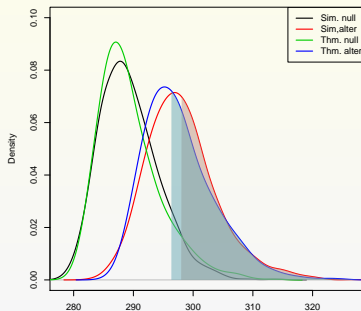
## maxdegree, 1st approximation

Thm vs Sim, ell=0,tau=1, us



N = 2000 Bandwidth = 2.203

Thm vs Sim, ell=0,tau=1, them



N = 2000 Bandwidth = 1.322

---

 $\beta(us)$ 
 $\beta(them)$ 

0.50

theory

0.42

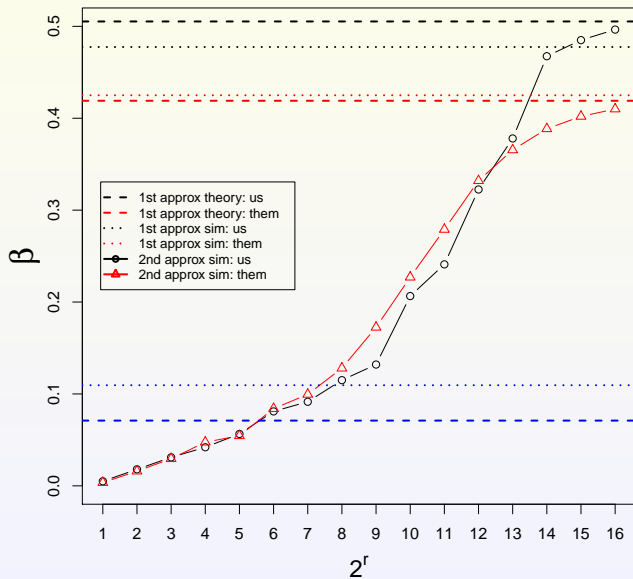
0.48

simulation

0.43

# Theory & Simulation

maxdegree

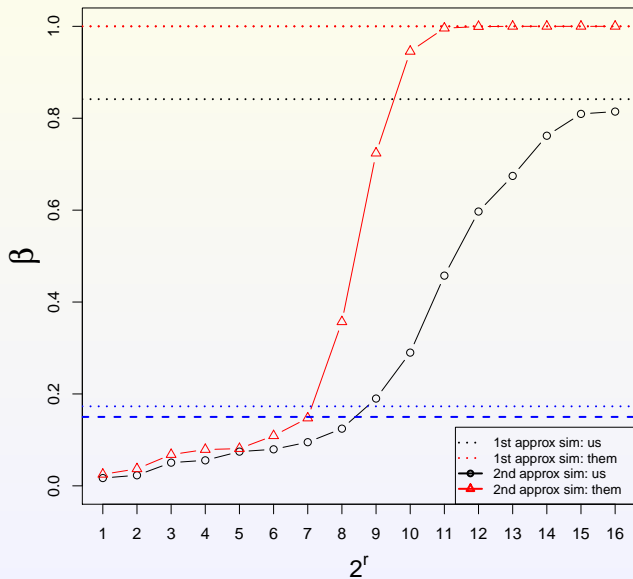


$\tau = 1$

$\tau = 0$   
(both  
us & them)

# Theory & Simulation

scan



$\tau = 1$

$\tau = 0$   
(both  
us & them)



# Conclusion







“them” is admissible!

## Discussion

- there remains theory, simulation, and experiments yet to be done . . .
- *power vs. computational complexity tradeoff* for scan statistics on streaming graphs!

## More References

<http://www.cis.jhu.edu/~parky/CEP-Publications>

-  A. Rukhin and C.E. Priebe, "On the Limiting Distribution of a Graph Scan Statistic," *Communications in Statistics - Theory and Methods*, Vol. 41, No. 7, pp. 1151-1170, 2012.
-  H. Pao, G.A. Coppersmith and C.E. Priebe, "Statistical Inference on Random Graphs: Comparative Power Analyses via Monte Carlo," *Journal of Computational and Graphical Statistics*, Vol. 20, No. 2, pp. 395-416, 2011.
-  A. Rukhin and C.E. Priebe, "A Comparative Power Analysis of the Maximum Degree and Size Invariants for Random Graph Inference," *Journal of Statistical Planning and Inference*, Vol. 141, pp. 1041-1046, 2011.
-  D.J. Marchette and C.E. Priebe, "Scan Statistics for Interstate Alliance Graphs," *Connections*, Volume 28, Issue 2, pp. 43-64, 2008.
-  M. Tang, Y. Park, N.H. Lee, and C.E. Priebe, "Attribute fusion in a latent process model for time series of graphs," submitted, 2011.
-  C.E. Priebe, G.A. Coppersmith, and A. Rukhin, "You say 'graph invariant,' I say 'test statistic' ", *ASA Sections on Statistical Computing Statistical Graphics, SCGN Newsletter*, 21, 2010.

## Leopold Kronecker to Hermann von Helmholtz:

*“The wealth of your practical experience with sane and interesting problems will give to mathematics a new direction and a new impetus.”*



Kronecker



Helmholtz