

Scan Statistics on Enron Graphs

Carey E. Priebe

cep@jhu.edu

Johns Hopkins University

Department of Applied Mathematics & Statistics

Department of Computer Science

Center for Imaging Science

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

— Leopold Kronecker to Hermann von Helmholtz

Citation

C.E. Priebe, J.M. Conroy, D.J. Marchette, and Y. Park,
“Scan Statistics on Enron Graphs,”
Computational and Mathematical Organization Theory,
to appear.

SIAM International Conference on Data Mining
Workshop on Link Analysis, Counterterrorism and Security
Newport Beach, California
April 23, 2005

<http://www.ams.jhu.edu/~priebe/sseg.html>

Outline

- Scan Statistics on Graphs, and on Time Series thereof
- Scan Statistics on Enron Graphs

Introduction

- **Objective:** *to develop and apply a theory of scan statistics on random graphs to perform change point / anomaly detection in graphs and in time series thereof.*
- **Hypotheses:**
 - H_0 : homogeneity
 - H_A : local subregion of excessive activity

Scan Statistics

- “moving window analysis”:

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,
- ...

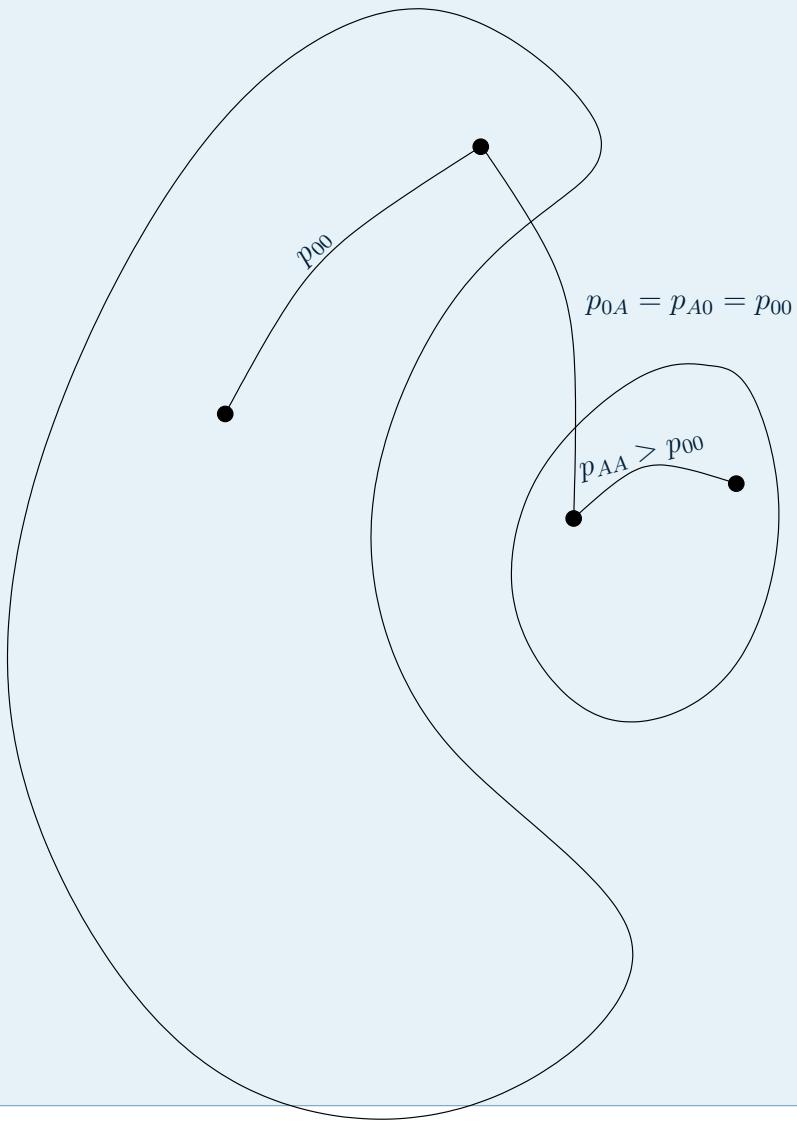
- **scan statistic \equiv maximum of locality statistic**:

If maximum of observed locality statistics is large,
then the inference can be made that
there exists a subregion of excessive activity.

Scan Statistics on Graphs

- directed graph (digraph): $D = (V, A)$
 - order: $|V(D)|$
 - size: $|A(D)|$
 - k -th order neighborhood of v :
$$N_k[v; D] = \{w \in V(D) : d(v, w) \leq k\}$$
 - scan region (example: induced subdigraph): $\Omega(N_k[v; D])$
 - locality statistic (example: size): $\Psi_k(v) = |A(\Omega(N_k[v; D]))|$
 - “scale-specific” scan statistic: $M_k(D) = \max_{v \in V(D)} \Psi_k(v)$

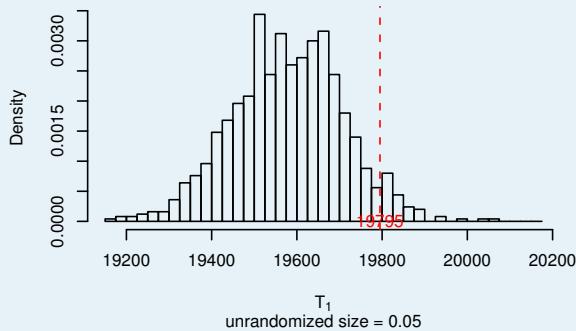
H_{A_1}



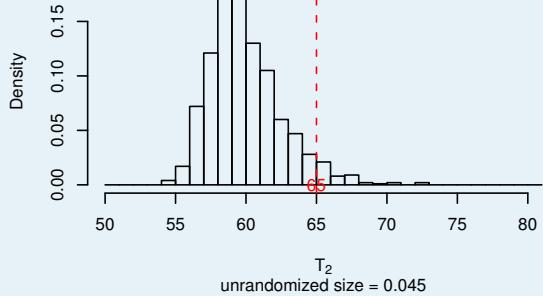
< > - +

Monte Carlo Simulation

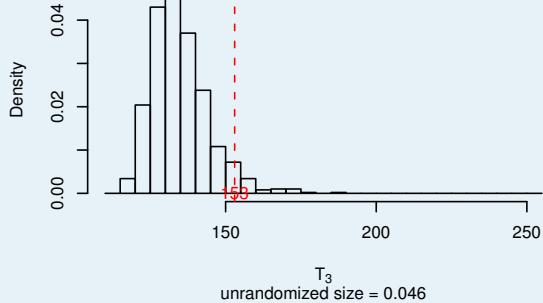
$T = \text{size}$



$T = \text{scan0}$



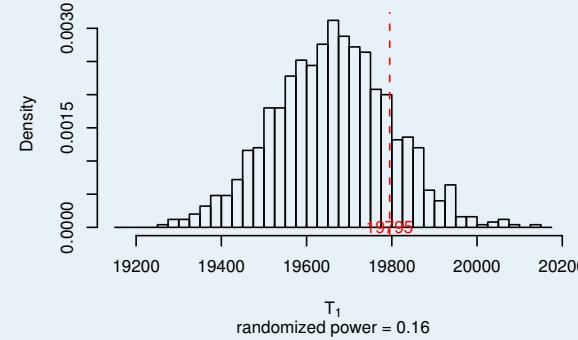
$T = \text{scan1}$



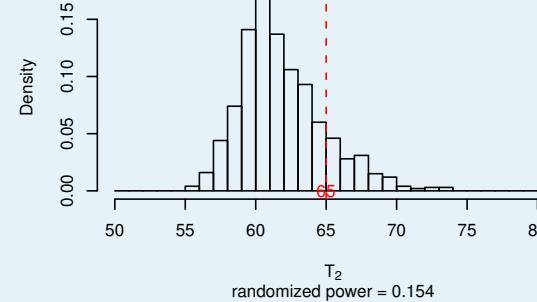
$n = 1000, n' = 13, \alpha = 0.05, MC = 1000$

< > - +

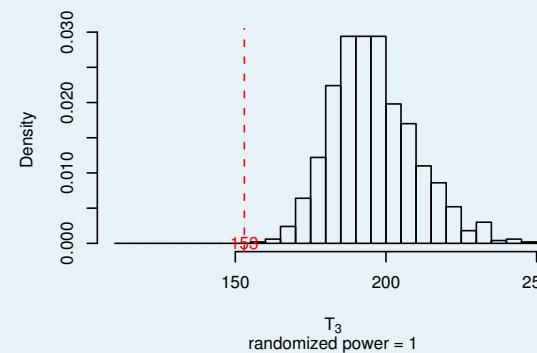
$\beta = 0.160$



$\beta = 0.154$

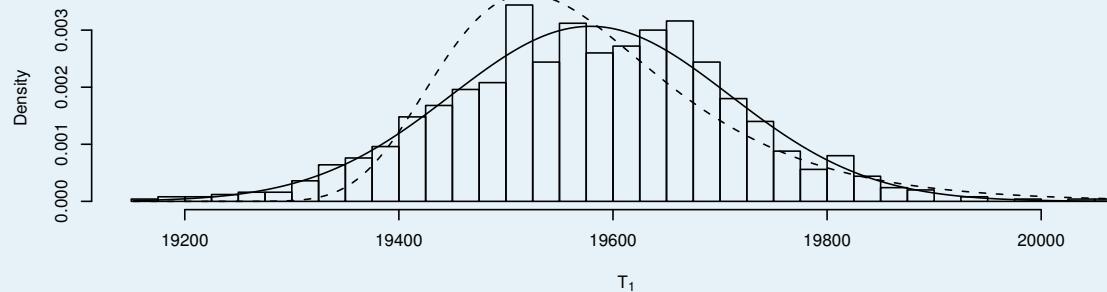


$\beta = 1.000$

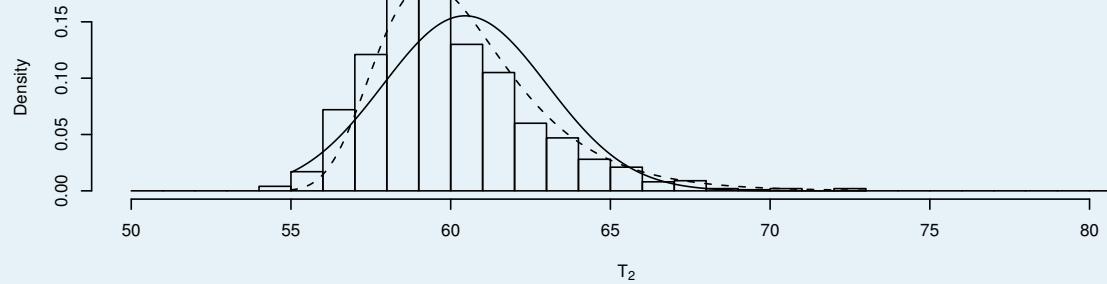


Gumbel Conjecture

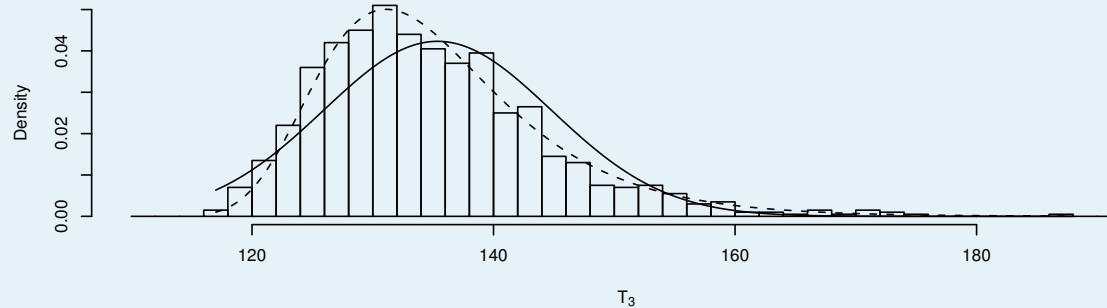
$T = \text{size}$



$T = \text{scan0}$

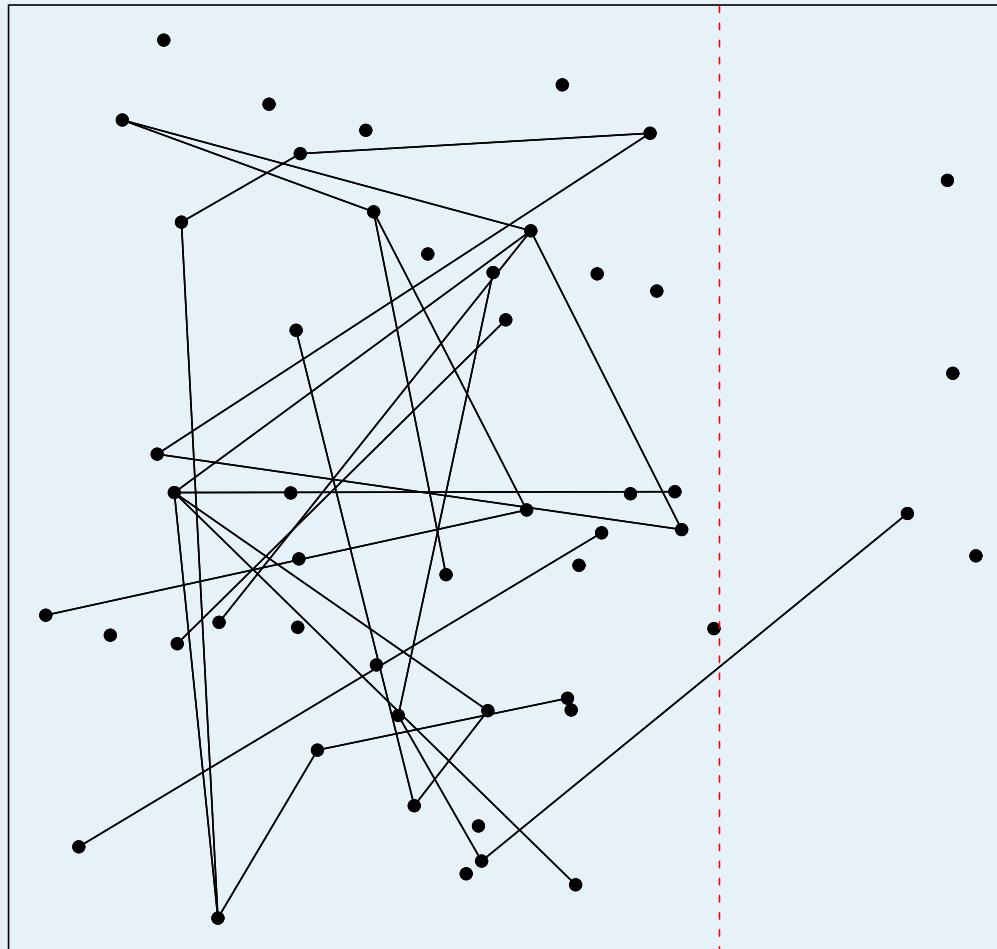


$T = \text{scan1}$



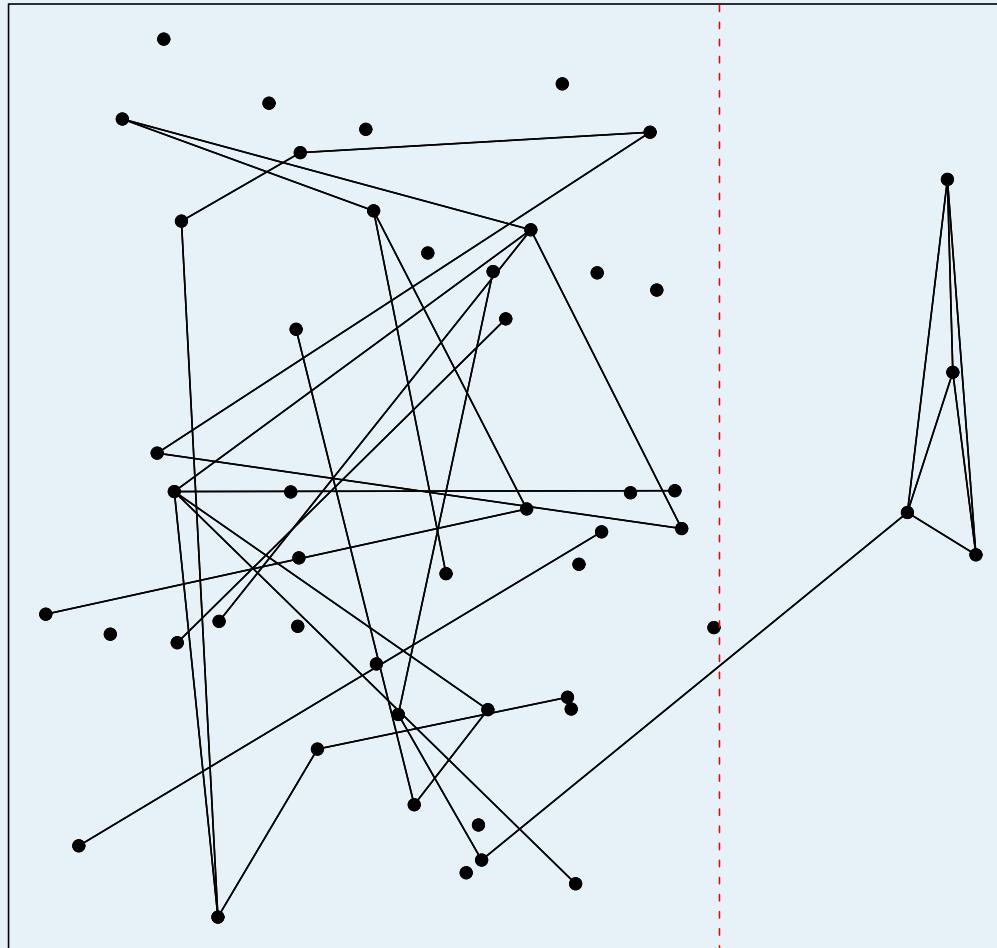
$n = 1000, MC = 1000$

Example: H_0



$size = 26, scan0 = 5, scan1 = 5, scan2 = 11.$

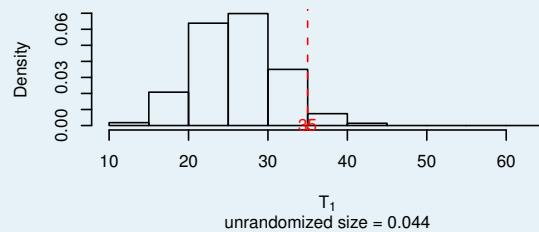
Example: H_{A_1}



$size = 32, scan0 = 5, scan1 = 7, scan2 = 11.$

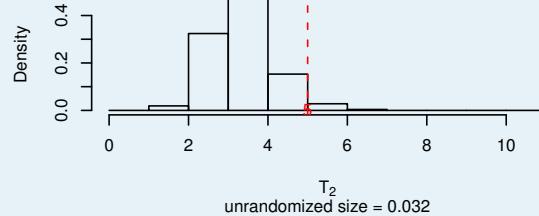
Example: Monte Carlo Simulation H_0 vs H_{A1}

$T = \text{size}$



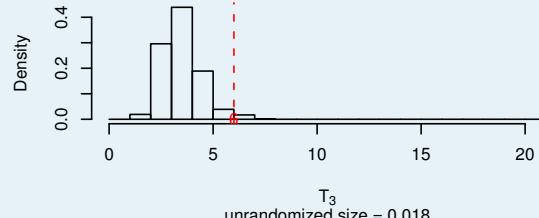
$$\beta = 0.273$$

$T = \text{scan0}$



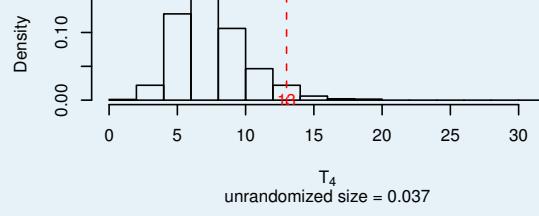
$$\beta = 0.350$$

$T = \text{scan1}$



$$\beta = 0.976$$

$T = \text{scan2}$



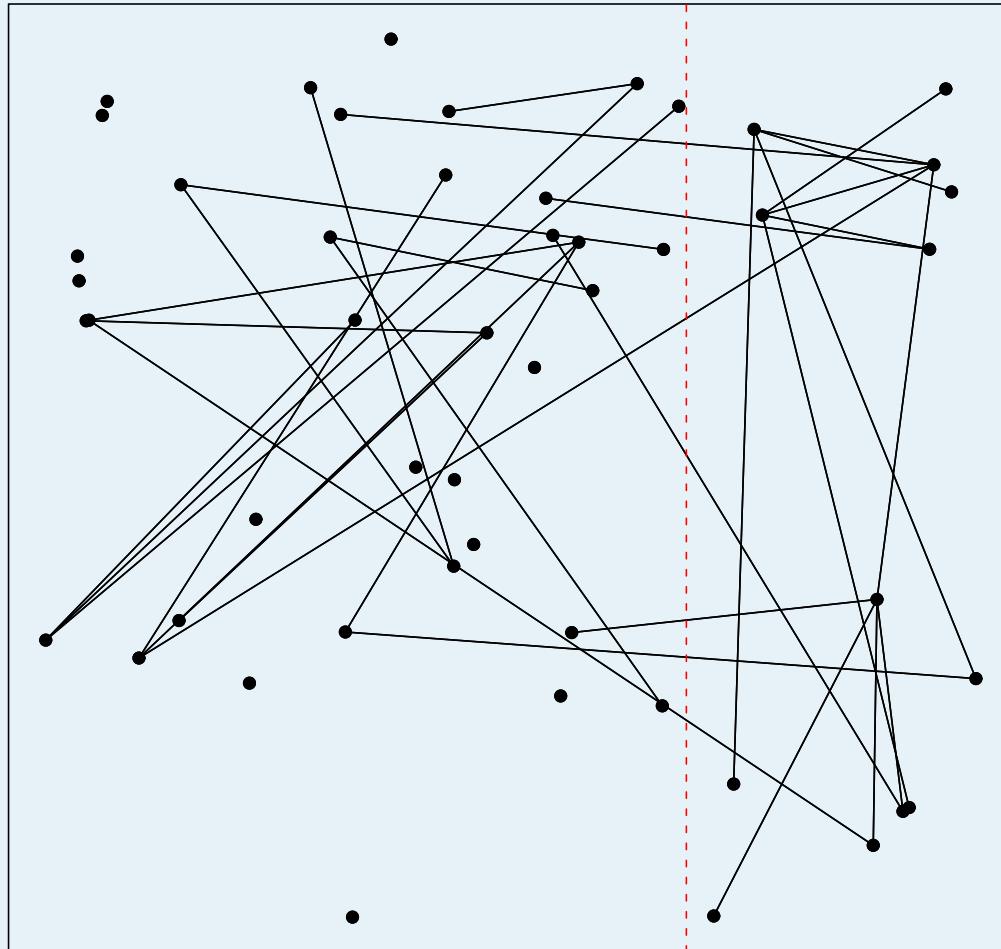
$$\beta = 0.451$$

$$n = 50, n' = 4, \alpha = 0.05, MC = 1000$$

$$H_{A_2}$$

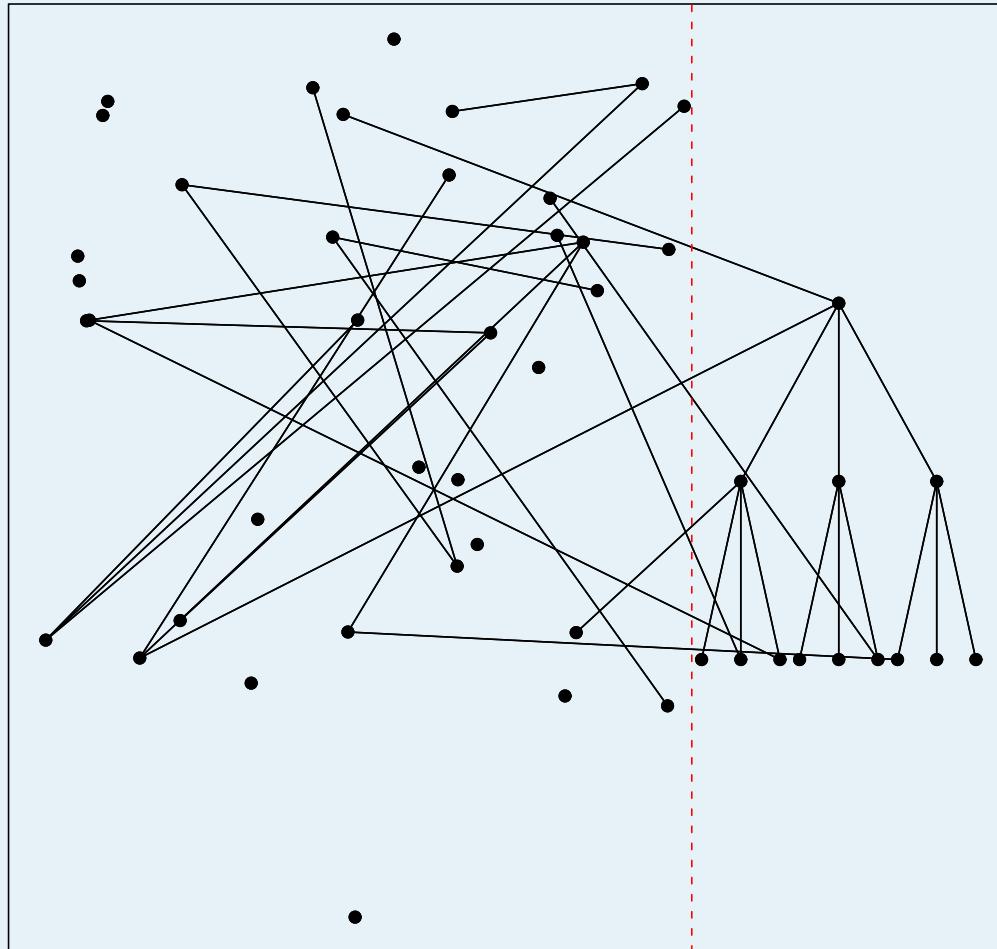
consider now a *structured* alternative . . .

Example: H_{A_2}



$size = 34, scan0 = 5, scan1 = 5, scan2 = 17.$

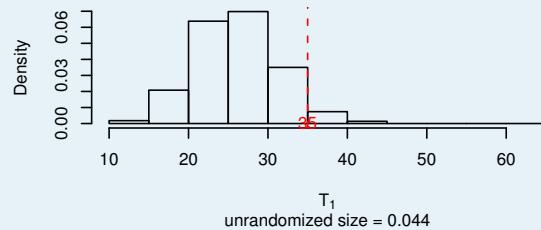
Example: H_{A_2}



$size = 34, scan0 = 5, scan1 = 5, scan2 = 17.$

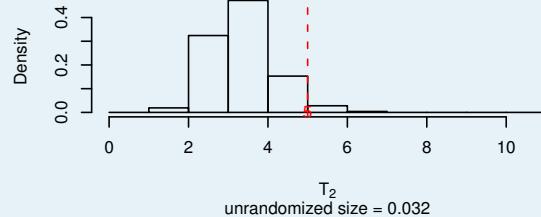
Example: Monte Carlo Simulation H_0 vs H_{A_2}

$T = \text{size}$



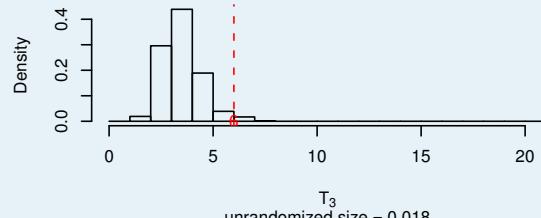
$$\beta = 0.633$$

$T = \text{scan0}$



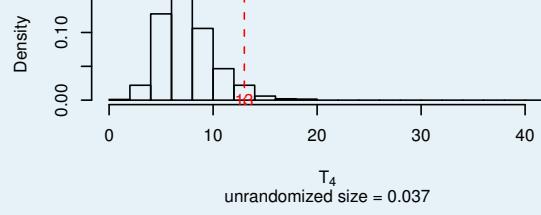
$$\beta = 0.545$$

$T = \text{scan1}$



$$\beta = 0.501$$

$T = \text{scan2}$

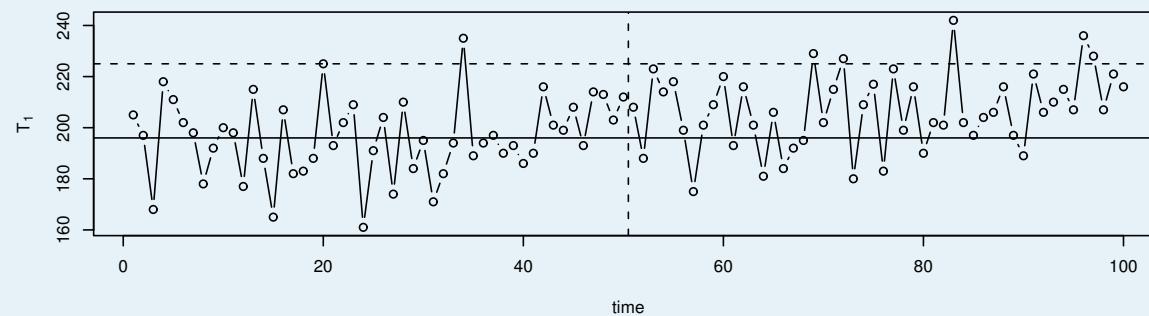


$$\beta = 0.915$$

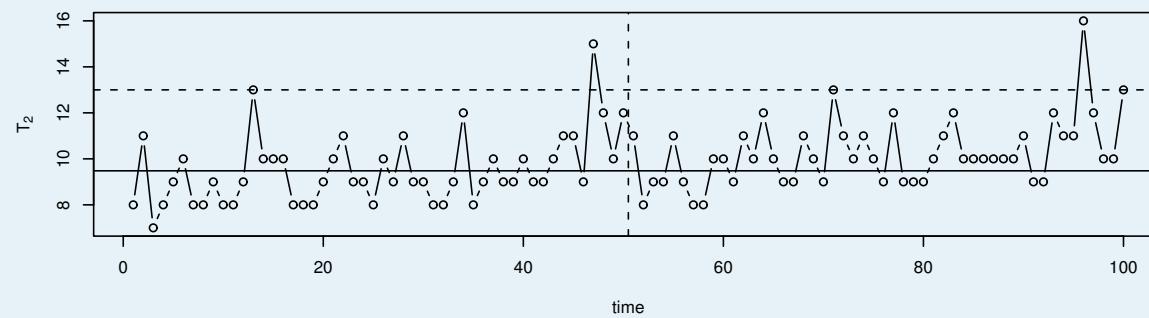
$$n = 50, n' = 13, \alpha = 0.05, MC = 1000$$

Time Series

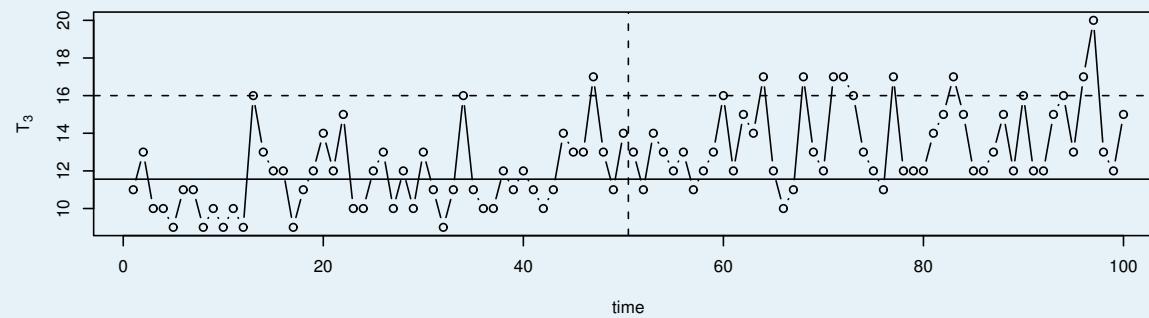
$T = \text{size}$



$T = \text{scan0}$



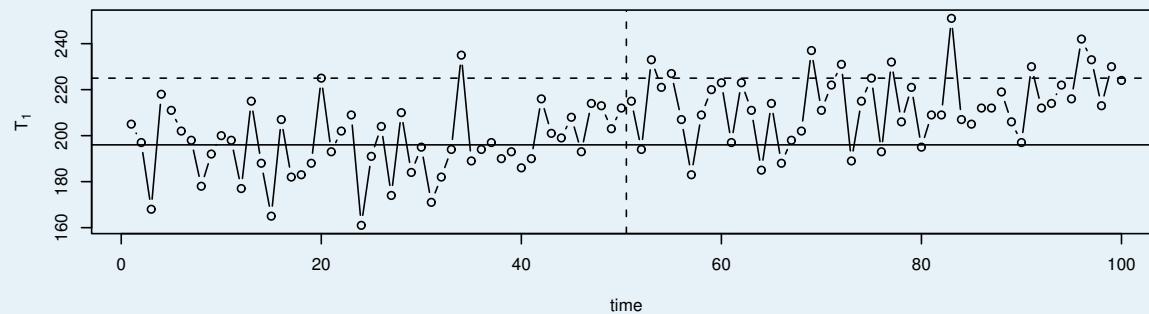
$T = \text{scan1}$



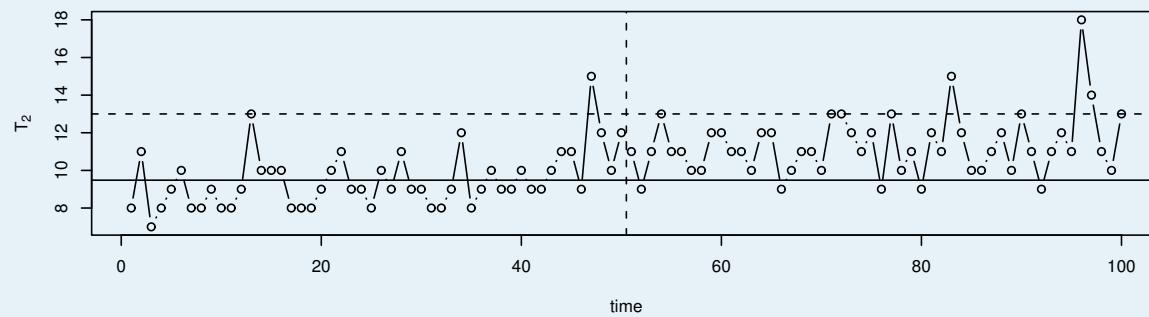
$$n = 100, n' = 4, \alpha = 0.05$$

Time Series

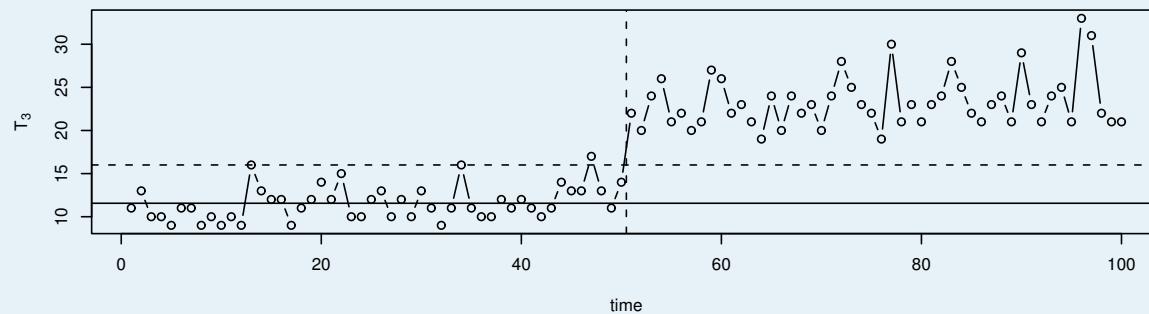
$T = \text{size}$



$T = \text{scan0}$



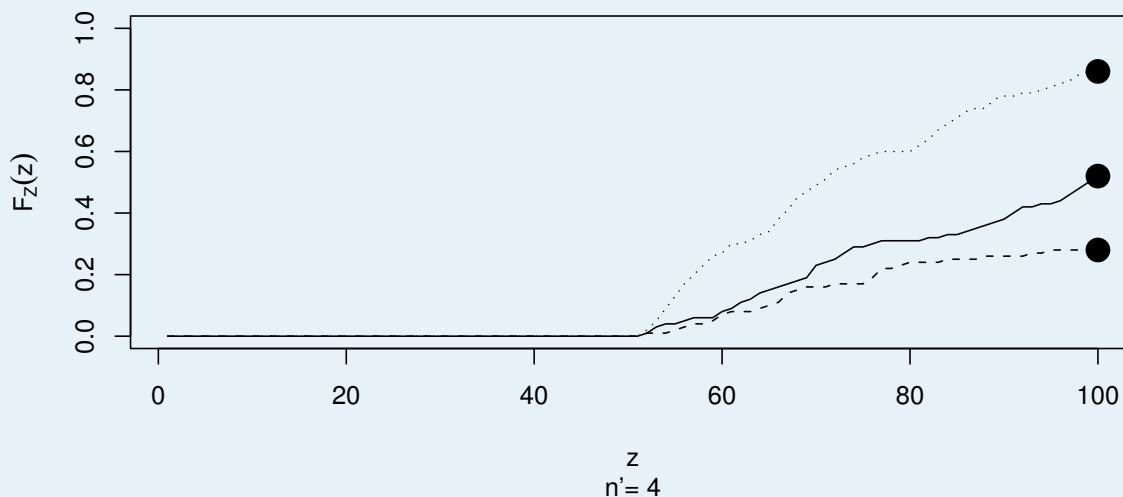
$T = \text{scan1}$



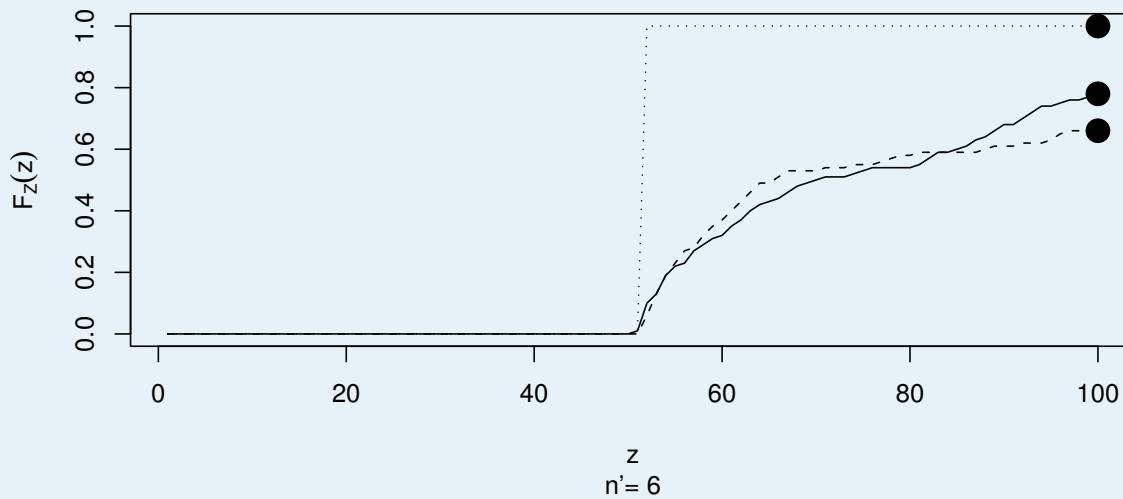
$$n = 100, n' = 6, \alpha = 0.05$$

Time Series

$n = 4$

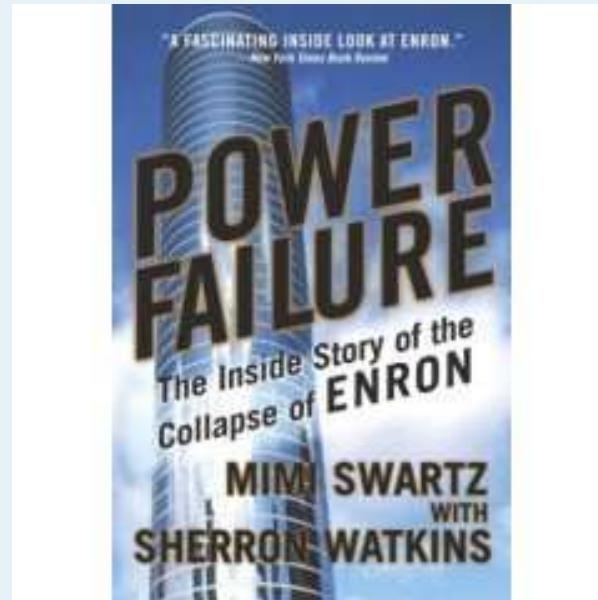


$n = 6$



$n = 100, n' = 4\&6, \alpha = 0.05$

Enron



< > - +

Enron Data

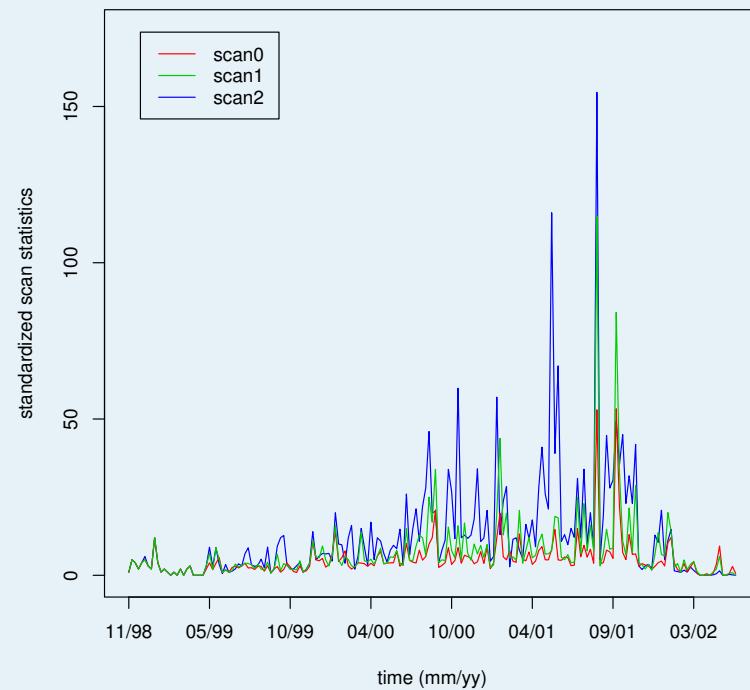
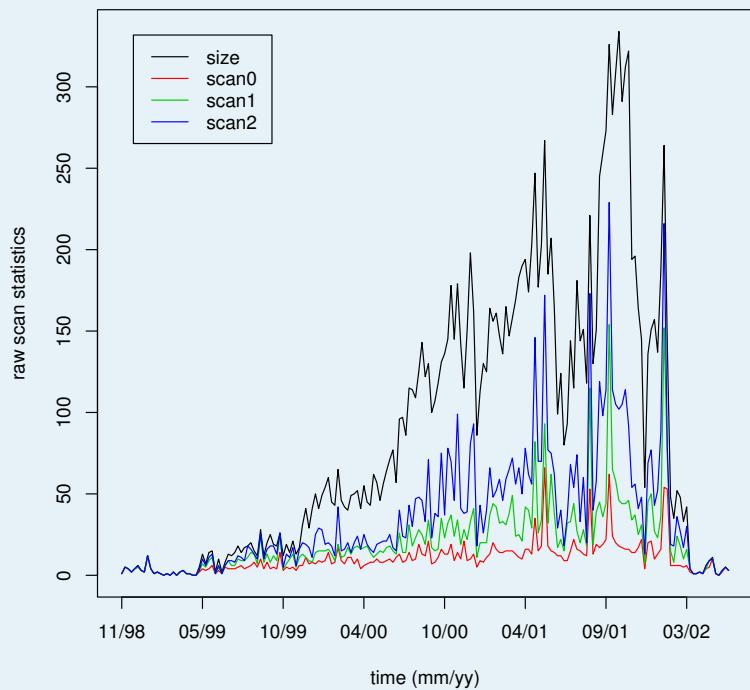
- 125,409 distinct messages from 184 unique “From” field, mostly Enron executives.
- 189 weeks, from 1988 through 2002.
- directed edges (arcs) $A_t = \{(v, w) : \text{vertex } v \text{ sends at least one email to vertex } w \text{ during the } t\text{-th week (“To”, “CC”, or “BCC”)}\}$

910948020 114 169	albert.meyers	Albert Meyers	Employee	Specialist
910948020 114 169	a..martin	Thomas Martin	Vice President	
911477940 114 123	andrea.ring	Andrea Ring	N/A	
911477940 114 123	andrew.lewis	Andrew Lewis	Director	
911481840 114 123	andy.zipper	Andy Zipper	Vice President	Enron Online
911481840 114 123	a..shankman	Jeffrey Shankman	President	Enron Global Mkt
911481840 114 123	barry.tycholiz	Barry Tycholiz	Vice President	
911481840 114 123	benjamin.rogers	Benjamin Rogers	Employee	Associate
911892180 114 38	bill.rapp	Bill Rapp	N/A	
911892180 114 38	bill.williams	xxx		

Statistics and Time Series

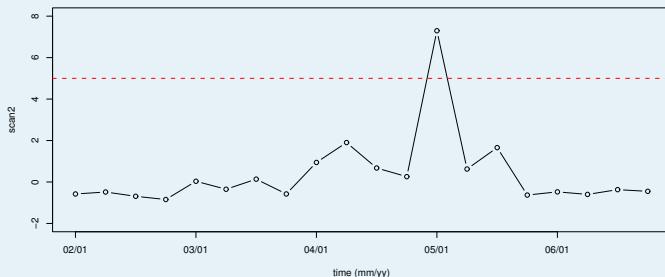
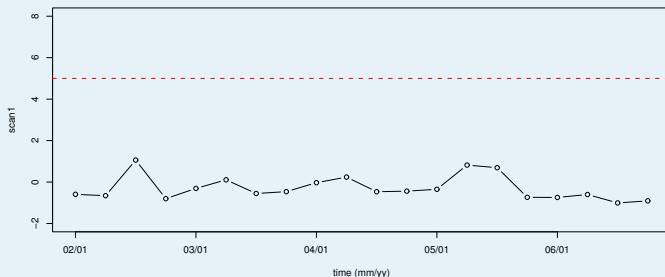
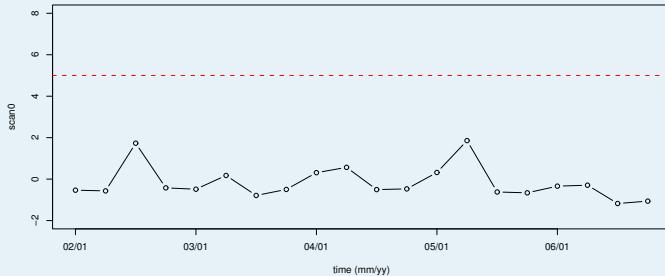
- **scale- k locality statistics:** $\Psi_{k,t}(v) = |A(\Omega(N_k[v; D_t]))|$
 - $k = 0$: $\Psi_{0,t}(v) = \text{outdegree}(v; D_t)$.
 - **scan statistic:** $M_{k,t} = \max_v \Psi_{k,t}(v); k = 0, 1, 2$
- **vertex-dependent standardized locality statistic:**
 - $\tilde{\Psi}_{k,t}(v) = (\Psi_{k,t}(v) - \hat{\mu}_{k,t,\tau}(v)) / \max(\hat{\sigma}_{k,t,\tau}(v), 1)$
 - $\hat{\mu}_{k,t,\tau}(v) = \frac{1}{\tau} \sum_{t'=t-\tau}^{t-1} \Psi_{k,t'}(v)$
 - $\hat{\sigma}_{k,t,\tau}^2(v) = \frac{1}{\tau-1} \sum_{t'=t-\tau}^{t-1} (\Psi_{k,t'}(v) - \hat{\mu}_{k,t,\tau}(v))^2$
 - **standardized scan statistic:** $\tilde{M}_{k,t} = \max_v \tilde{\Psi}_{k,t}(v)$

Statistics and Time Series



Anomaly Detection

- temporally-normalized scan statistics: $S_{k,t} = (\tilde{M}_{k,t} - \tilde{\mu}_{k,t,\ell}) / \max(\tilde{\sigma}_{k,t,\ell}, 1)$
- detection: time t such that $S_{k,t} > 5$
- $t^* = 132$ (May, 2001)



Detection Graph D_{132}

$$\arg \max_v \Psi_{0,132}(v) = john.lavorato$$

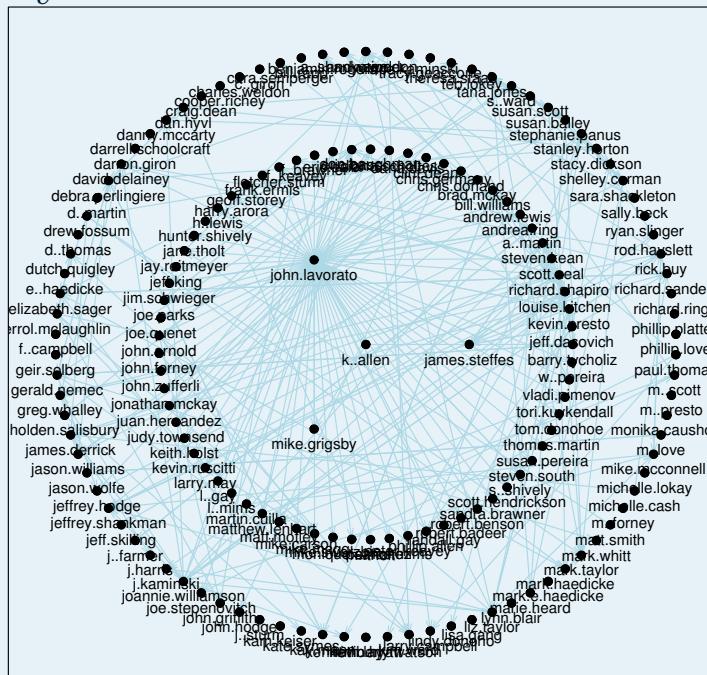
$$\arg \max_v \Psi_{1,132}(v) = john.lavorato$$

$$\arg \max_v \Psi_{2,132}(v) = richard.shapiro$$

$$\arg \max_v \tilde{\Psi}_{0,132}(v) = richard.shapiro$$

$$\arg \max_v \tilde{\Psi}_{1,132}(v) = joannie.williamson$$

$$\arg \max_v \tilde{\Psi}_{2,132}(v) = k..allen$$



Detection Graph D_{132}

Details for the ‘detection’ graph D_{132}

time t^*	132 (week of May 17, 2001)		
$\text{size}(D_{132})$	267		
scale k	$M_{k,132}$	$\widetilde{M}_{k,132}$	$S_{k,132}$
0	66	8.3	0.32
1	93	7.8	-0.35
2	172	116.0	7.30
3	219	174.0	5.20
number of isolates	50		

Anomaly Detection (*Aliasing*)

- $v^* = \arg \max_v \tilde{\Psi}_{2,132}(v) = k..allen$
- $k..allen == phillip.allen?$
 - $k..allen$ had no activity before $t^* = 132$.
 - At $t^* = 132$, $phillip.allen$ switched to $k..allen$.
- Matched Filter:
 - For each vertex $v \in V \setminus \{v^*\}$,

$$s_{t^*, \kappa}(v; v^*) = \sum_{t' = t^* - \kappa}^{t^* - 1} |N_1(v; D_{t'}) \cap N_1(v^*; D_{t^*})|$$

- Is this a detection we want?

New York Times (May 22, 2005)

The New York Times

Week in Review

NYTimes.com

Go to a Section

Site Search:

NYT Since 1996

Submit

Enron Offers an Unlikely Boost to E-Mail Surveillance

By GINA KOLATA

Published: May 22, 2005

AS an object of modern surveillance, e-mail is both reassuring and troubling. It is a potential treasure trove for investigators monitoring suspected terrorists and other criminals, but it also creates the potential for abuse, by giving businesses and government agencies an efficient means of monitoring the attitudes and activities of employees and citizens.

E-Mail This

Printer-Friendly
Reprints

Multimedia

► GRAPHIC



[Finding Patterns in Corporate Chatter](#)

Now the science of e-mail tracking and analysis has been given a unlikely boost by a bitter chapter in the history of corporate malfeasance - the Enron scandal.

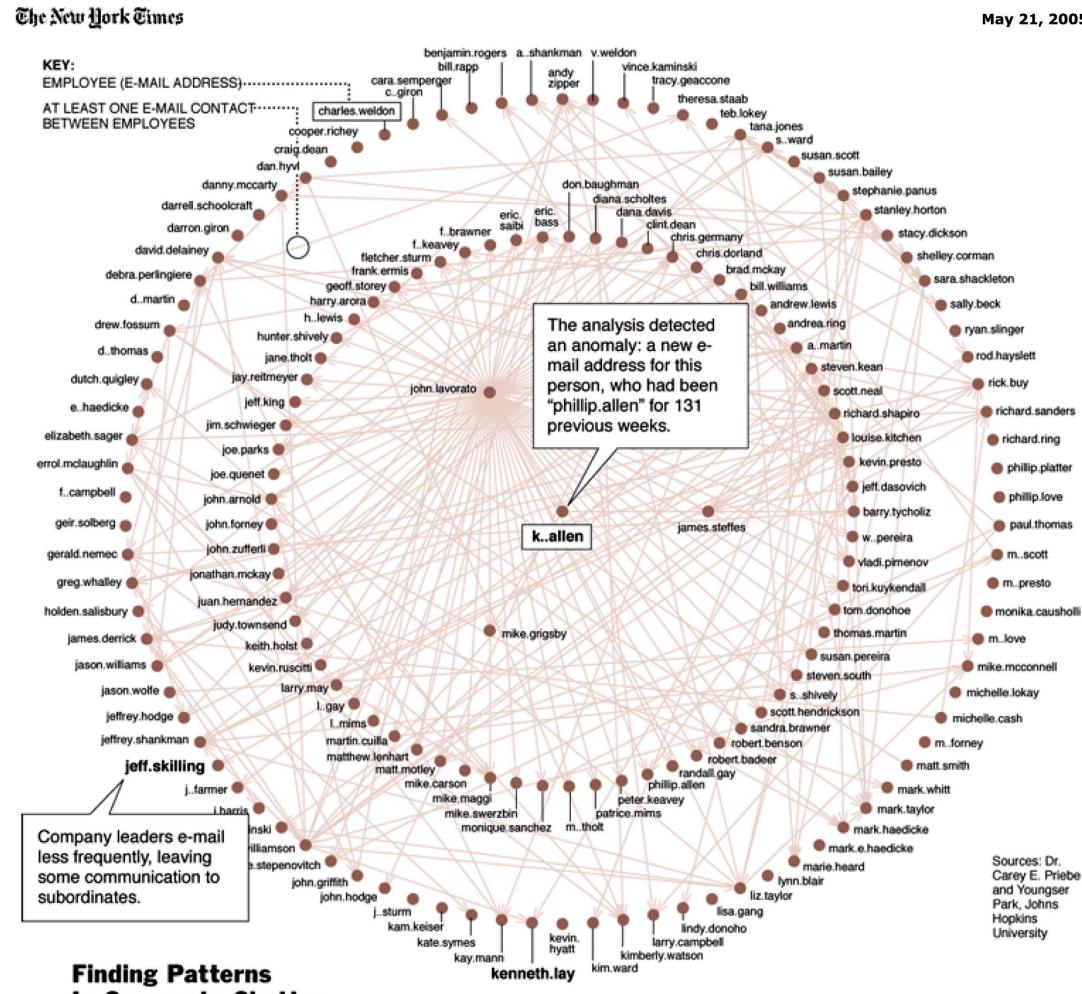
In 2003, the Federal Energy

< > - +

New York Times (May 22, 2005)

The New York Times > Week in Review > Image > Finding Patterns in Corporate Chatter

05/23/2005 08:32 AM



Finding Patterns In Corporate Chatter

Computer scientists are analyzing about a half million Enron e-mails. Here is a map of a week's e-mail patterns in May 2001, when a new name suddenly appeared. Scientists found that this week's pattern differed greatly from others, suggesting different conversations were taking place that might interest investigators. Next step: word analysis of these messages.

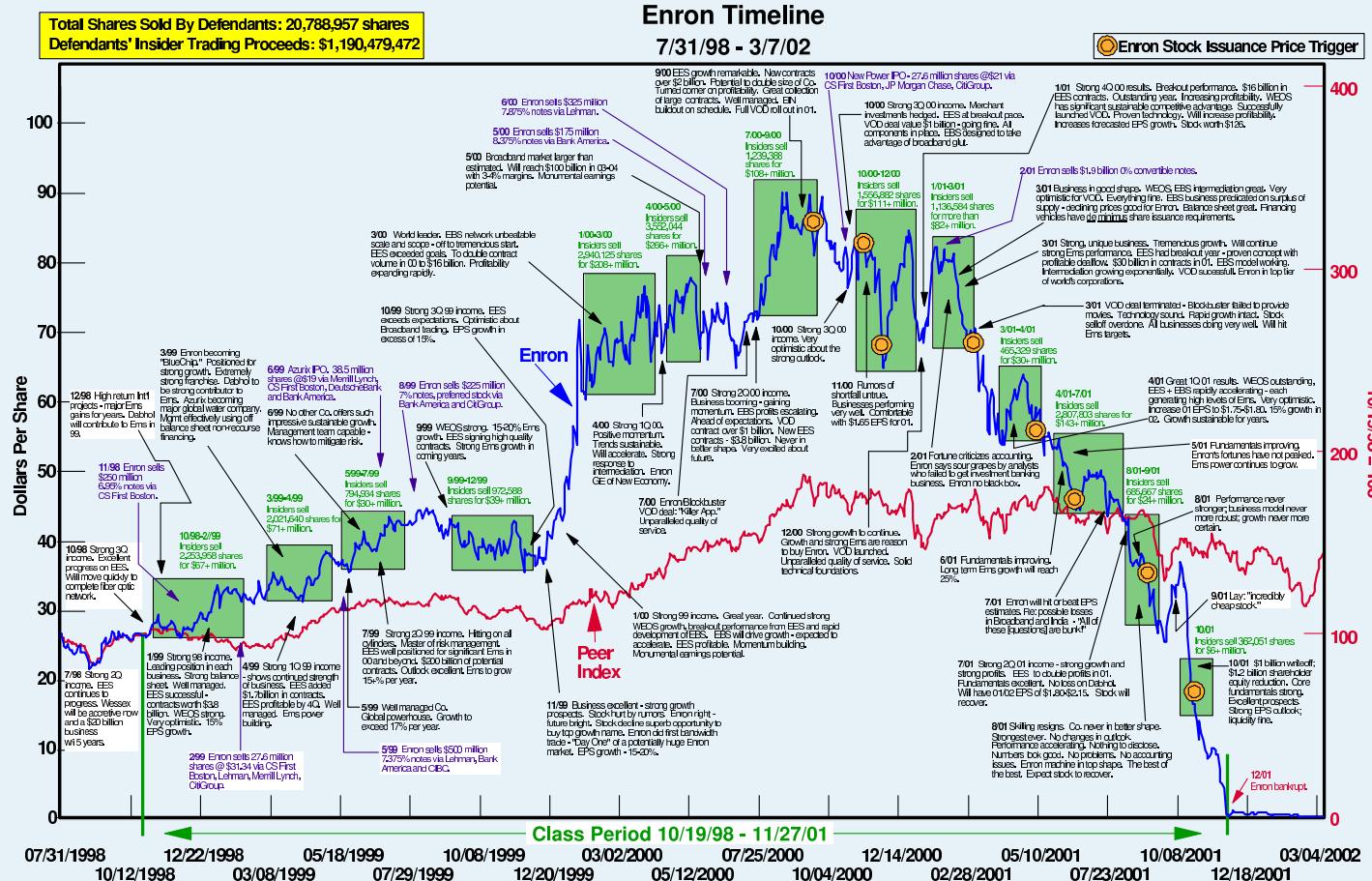
Bill Marsh/The New York Times

 Close Window

Copyright 2005 The New York Times Company

Scan Statistics on Enron Graphs – p.29/39

Enron Timeline



Anomaly Detection (*another*)

- Non-zero activity: $\tilde{\Psi}_{k,t}(v) \cdot I\{\hat{\mu}_{0,t,\tau}(v) > c\}$
 - For $c = 1$, $v^* = \text{roy.hayslett}$ at $t^* = 152$.

scale k	$\Psi_{k,t^*-5:t^*}(v^*)$
0	[1 , 2 , 1 , 3 , 1 , 2]
1	[1 , 2 , 2 , 9 , 2 , 4]
2	[1 , 2 , 2 , 19 , 4 , 175]
3	[1 , 2 , 2 , 58 , 6 , 268]

Anomaly Detection (*another*)

- *roy.hayslett* communicates with *sally.beck*, who is a $k = 0$ detection!

scale k	$\Psi_{k,t^*-5:t^*}(v)$
0	[3 , 2 , 0 , 2 , 3 , 62]
1	[3 , 3 , 0 , 3 , 6 , 154]
2	[4 , 3 , 0 , 37 , 11 , 229]
3	[4 , 3 , 0 , 98 , 16 , 267]

Anomaly Detection (*chatter*)

- Seek a detection in which the excess activity is due to **chatter** amongst the 2-neighbors!

$$\tilde{\Psi}'_t(v) = \left(\tilde{\Psi}_{2,t}(v) \cdot \mathcal{I}_{t,\tau}(v) \right) / \max(\gamma_t(v), 1)$$

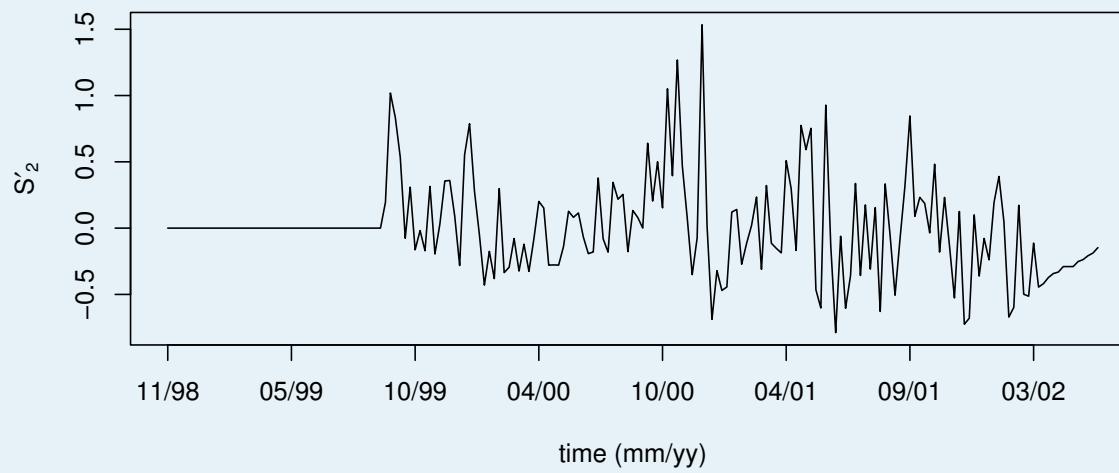
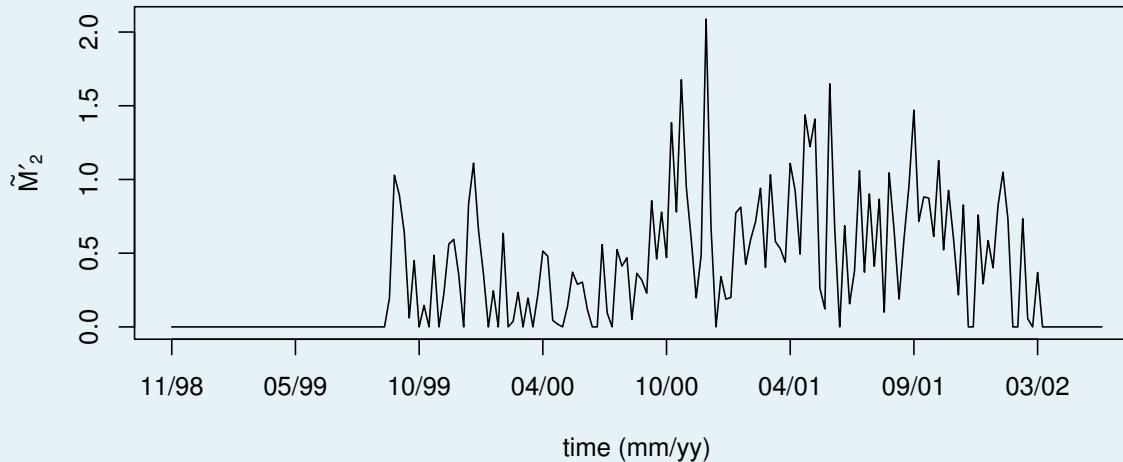
$$\mathcal{I}_{t,\tau}(v) = I_1 \times I_2 \times I_3$$

$$I_1 = I\{\hat{\mu}_{0,t,\tau} > c_1\},$$

$$I_2 = I\{\Psi_0(v) < \hat{\sigma}_{0,t,\tau}(v)c_2 + \hat{\mu}_{0,t,\tau}(v)\},$$

$$I_3 = I\{\Psi_1(v) < \hat{\sigma}_{1,t,\tau}(v)c_3 + \hat{\mu}_{1,t,\tau}(v)\}.$$

Anomaly Detection (*chatter*)

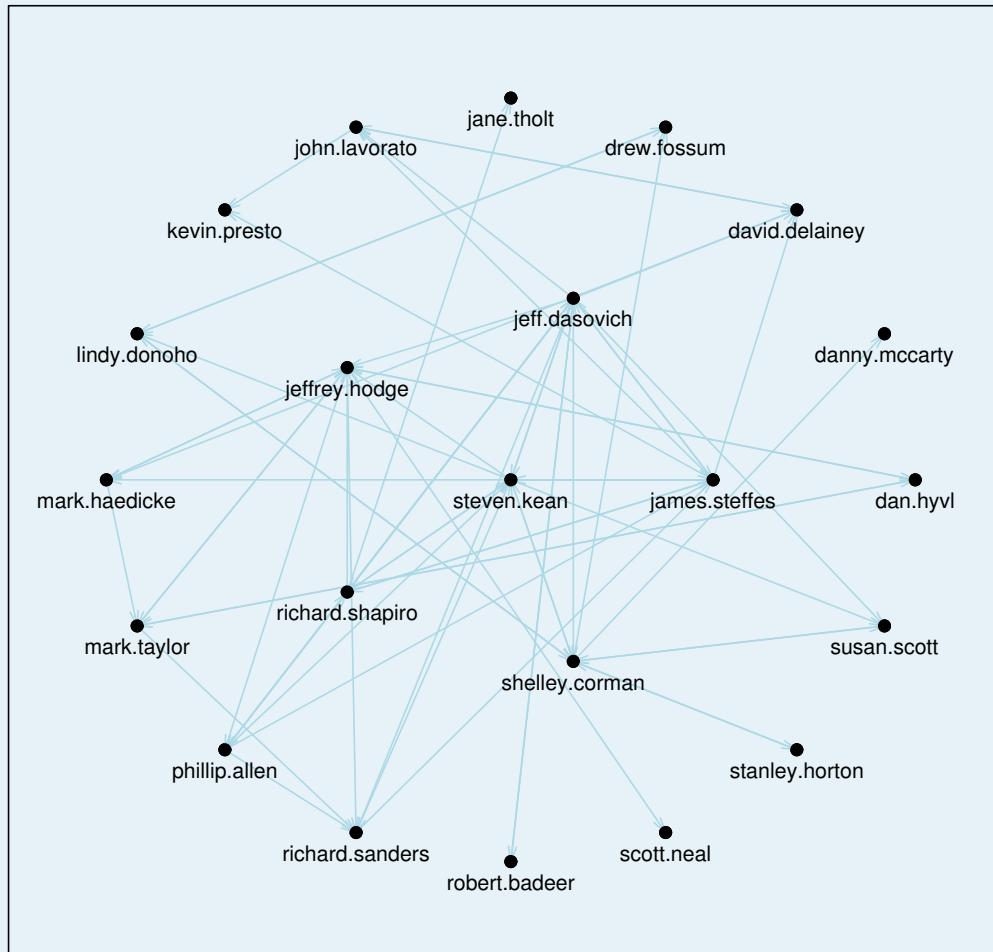


Anomaly Detection (*chatter*)

- $(v^*, t^*) = (\text{steven.kean}, 109)$

scale k	$\Psi_{k,t^*-5:t^*}(v^*)$
0	[3 , 5 , 4 , 5 , 4 , 5]
1	[11 , 13 , 10 , 10 , 11 , 18]
2	[14 , 35 , 21 , 38 , 13 , 65]

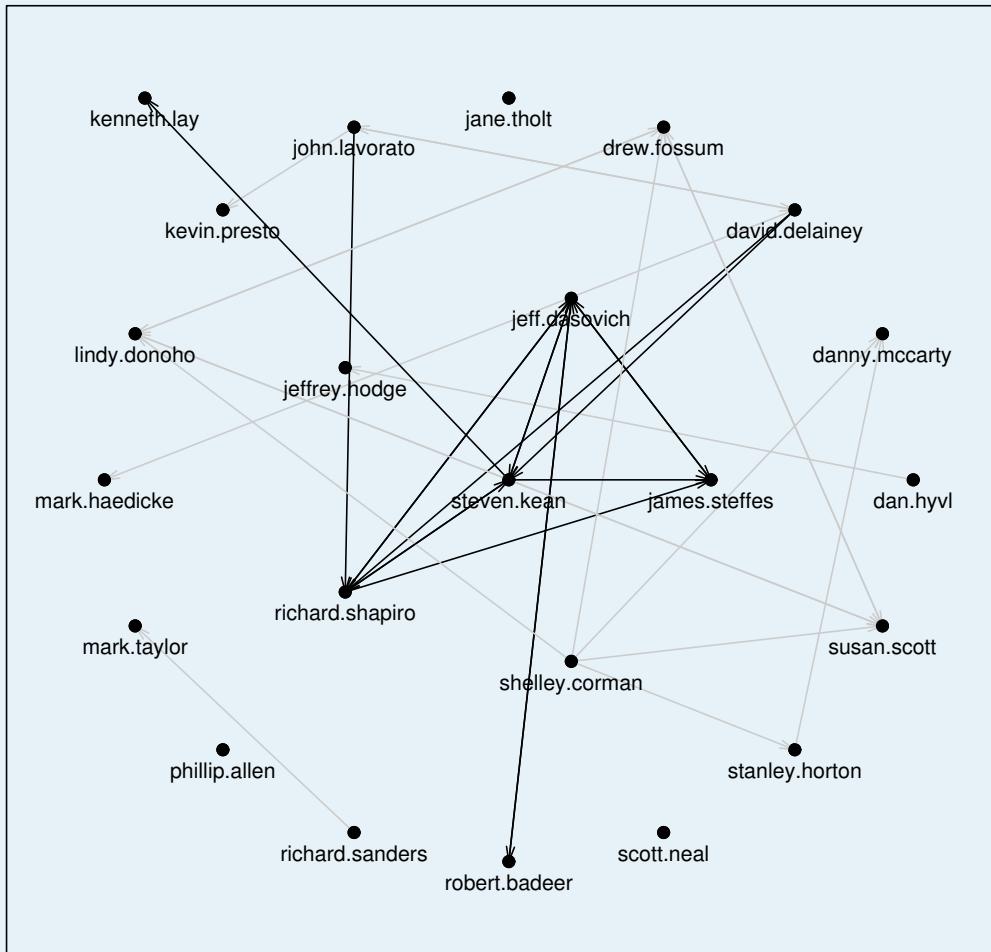
Anomaly Detection (*chatter*)



$$\Omega_{109}$$

< > - +

Anomaly Detection (*chatter*)



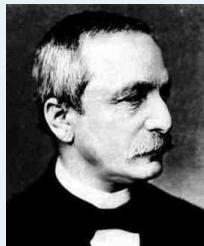
$$\Omega_{108}$$

Discussion

- scan statistics offers promise for detecting anomalies in time series of graphs.
- extensions:
 - weighted graphs (# of messages),
 - coloured hypergraphs (“To”, “CC”, or “BCC”),
 - ...
 - sliding window (online analysis),
 - ...
 - exponential smoothing, detrending, variance stabilization, ...
- [“Content and Scan Statistics for Enron” — John Conroy, Thursday](#)

Kronecker Quote

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



– Leopold Kronecker to Hermann von Helmholtz –

