

Talk 3: Graph Mining Tools – Tensors, communities, parallelism

Christos Faloutsos

CMU

Overall Outline

- Introduction – Motivation
- Talk#1: Patterns in graphs; generators
- Talk#2: Tools (Ranking, proximity)
- ➔ • Talk#3: Tools (Tensors, scalability)
- Conclusions

Outline

- ➔ • Task 4: time-evolving graphs – tensors
- Task 5: community detection
- Task 6: virus propagation
- Task 7: scalability, parallelism and hadoop
- Conclusions

Thanks to

- Tamara Kolda (Sandia)



for the foils on tensor
definitions, and on TOPHITS

Detailed outline

- Motivation
- Definitions: PARAFAC and Tucker
- Case study: web mining

Examples of Matrices:

Authors and terms

	data	mining	classif.	tree	...
John	13	11	22	55	...
Peter	5	4	6	7	...
Mary
Nick
...

Motivation: Why tensors?

- Q: what is a tensor?

Motivation: Why tensors?

- A: N-D generalization of matrix:

KDD'09	data	mining	classif.	tree	...
John	13	11	22	55	...
Peter	5	4	6	7	...
Mary
Nick
...

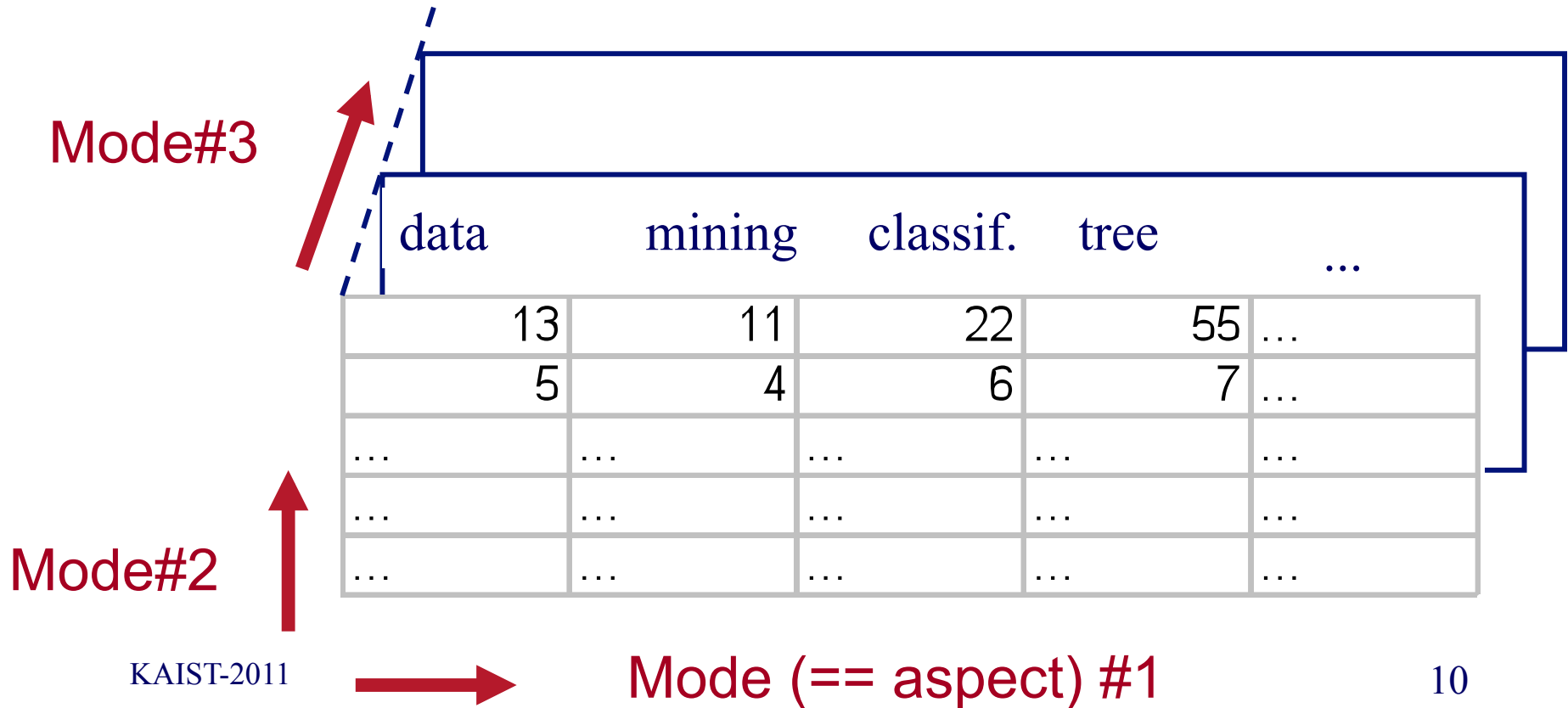
Motivation: Why tensors?

- A: N-D generalization of matrix:

	data	mining	classif.	tree	...
John	13	11	22	55	...
Peter	5	4	6	7	...
Mary
Nick
...

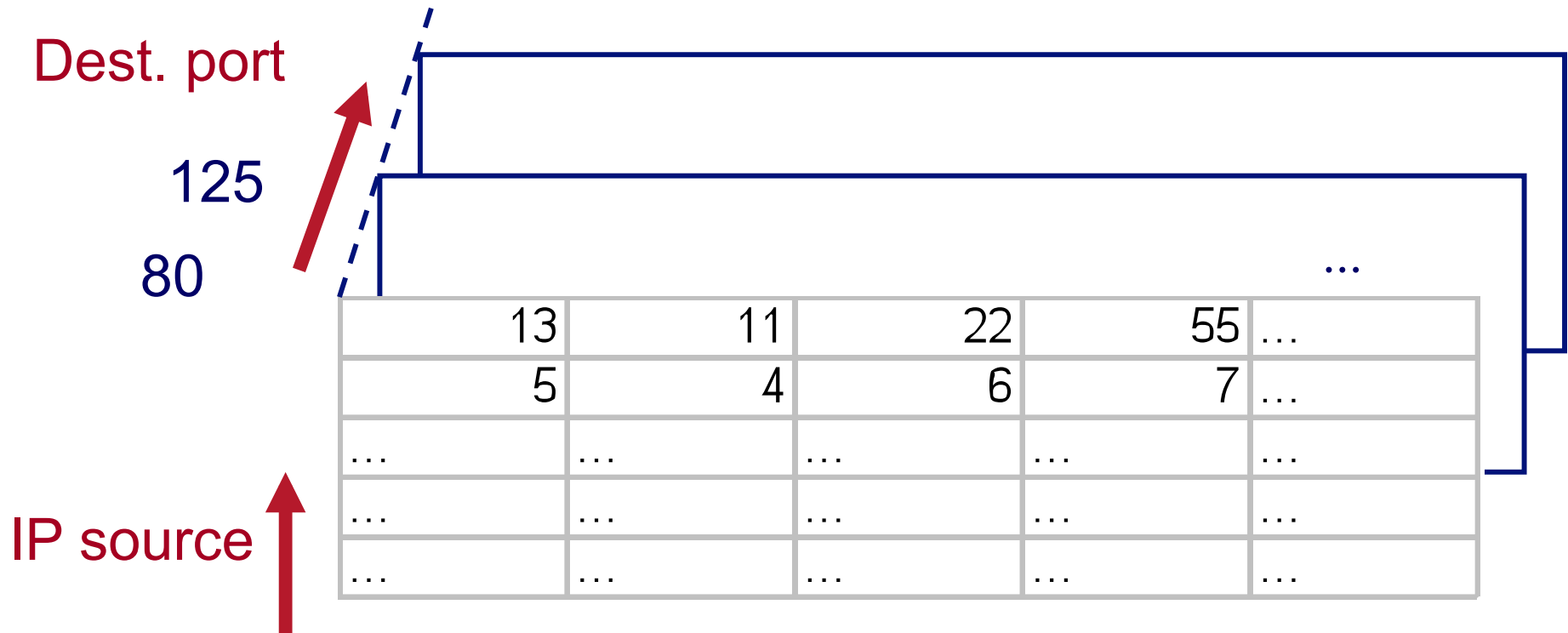
Tensors are useful for 3 or more modes

Terminology: ‘mode’ (or ‘aspect’):



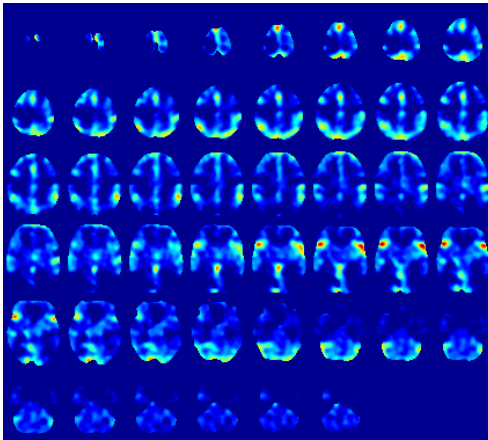
Notice

- 3rd mode does not need to be time
- we can have more than 3 modes



Notice

- 3rd mode does not need to be time
- we can have more than 3 modes
 - Eg, fMRI: x,y,z, time, person-id, task-id



From DENLAB, Temple U.
(Prof. V. Megalooikonomou +)

<http://denlab.temple.edu/bidms/cgi-bin/browse.cgi>

Motivating Applications

- Why tensors are useful?
 - web mining (TOPHITS)
 - environmental sensors
 - Intrusion detection (src, dst, time, dest-port)
 - Social networks (src, dst, time, type-of-contact)
 - face recognition
 - etc ...

Detailed outline

- Motivation
- ➔ • Definitions: PARAFAC and Tucker
- Case study: web mining

Tensor basics

- Multi-mode extensions of SVD – recall that:

Reminder: SVD

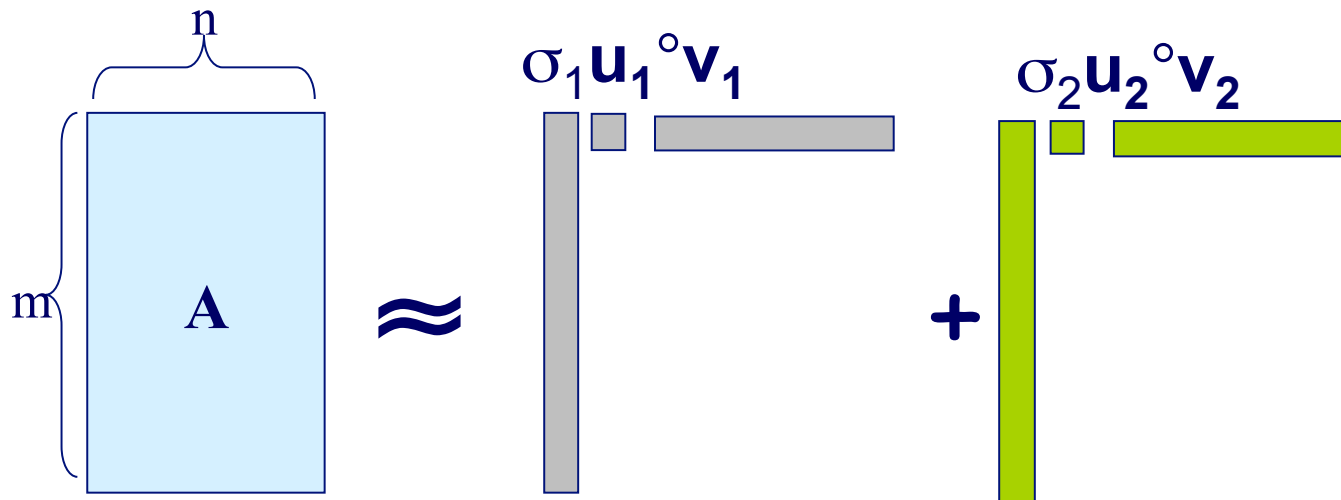
$$\mathbf{A} \approx \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \sum_i \sigma_i \mathbf{u}_i \circ \mathbf{v}_i$$

The diagram illustrates the SVD decomposition of matrix \mathbf{A} . Matrix \mathbf{A} is shown as a light blue rectangle with dimensions m by n . It is approximated by the product of three matrices: \mathbf{U} , $\mathbf{\Sigma}$, and \mathbf{V}^T . Matrix \mathbf{U} is a tall, narrow rectangle with dimensions m by k , shown with a vertical green stripe. Matrix $\mathbf{\Sigma}$ is a small square with dimensions k by k , shown with a 2x2 grid. Matrix \mathbf{V}^T is a wide, short rectangle with dimensions k by n , shown with a horizontal green stripe. The approximation symbol is a double wavy line.

– Best rank- k approximation in L2

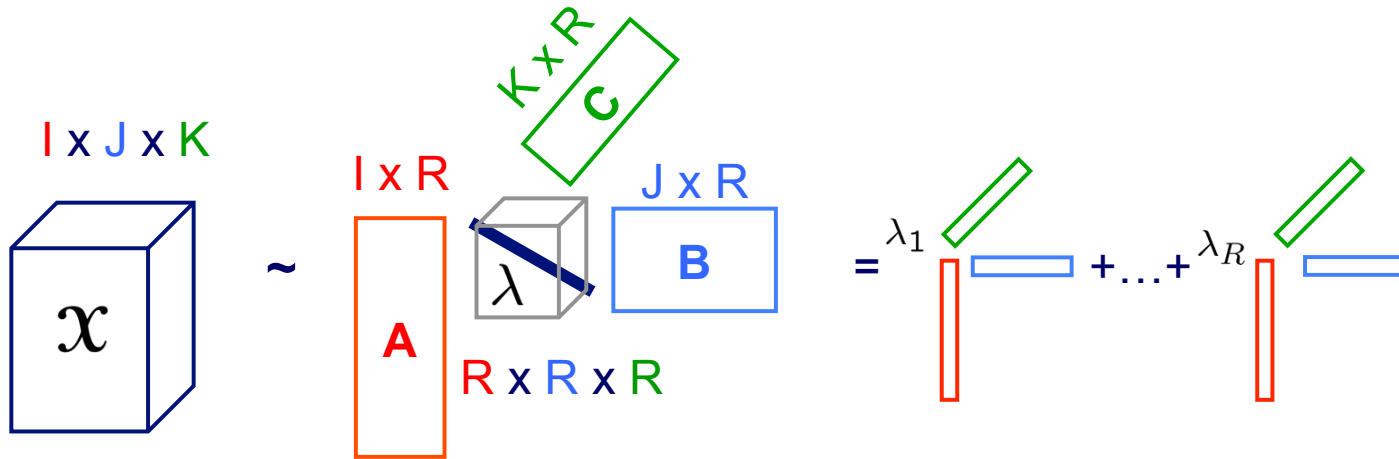
Reminder: SVD

$$\mathbf{A} \approx \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \sum_i \sigma_i \mathbf{u}_i \circ \mathbf{v}_i$$



– Best rank- k approximation in L2

Goal: extension to ≥ 3 modes



$$\mathcal{X} \approx [\lambda; \mathbf{A}, \mathbf{B}, \mathbf{C}] = \sum_r \lambda_r \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r$$

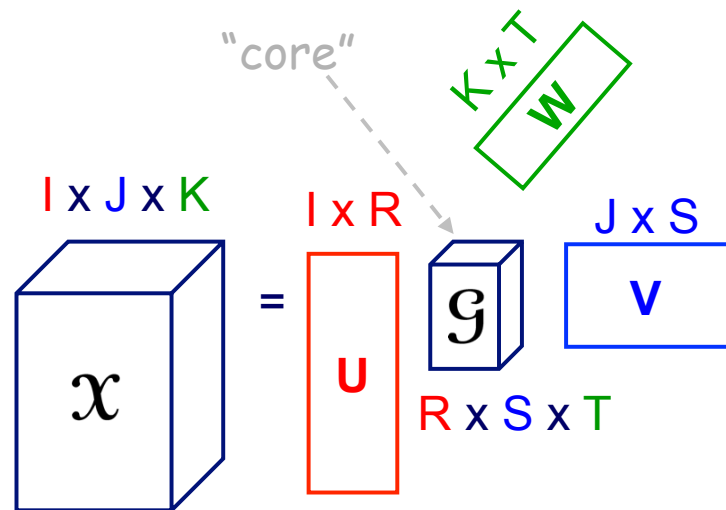
Main points:

- 2 major types of tensor decompositions: PARAFAC and Tucker
- both can be solved with “alternating least squares” (ALS)

Specially Structured Tensors

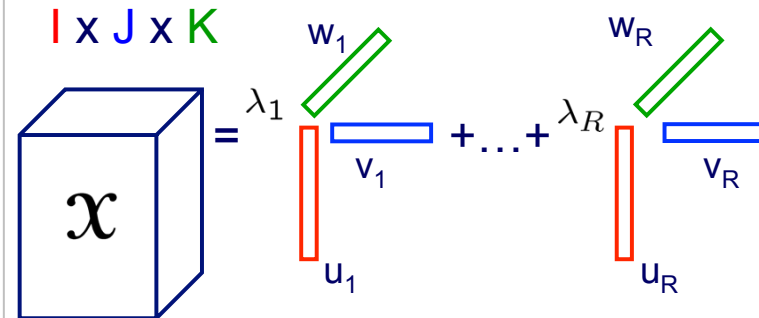
• Tucker Tensor

$$\begin{aligned}\mathcal{X} &= \mathcal{G} \times_1 \mathbf{U} \times_2 \mathbf{V} \times_3 \mathbf{W} \\ &= \sum_r \sum_s \sum_t g_{rst} \mathbf{u}_r \circ \mathbf{v}_s \circ \mathbf{w}_t \\ &\equiv [\mathcal{G}; \mathbf{U}, \mathbf{V}, \mathbf{W}] \end{aligned} \left. \vphantom{\begin{aligned}\mathcal{X} \\ &= \sum_r \sum_s \sum_t g_{rst} \mathbf{u}_r \circ \mathbf{v}_s \circ \mathbf{w}_t \\ &\equiv [\mathcal{G}; \mathbf{U}, \mathbf{V}, \mathbf{W}]\end{aligned}} \right\} \text{Our Notation}$$

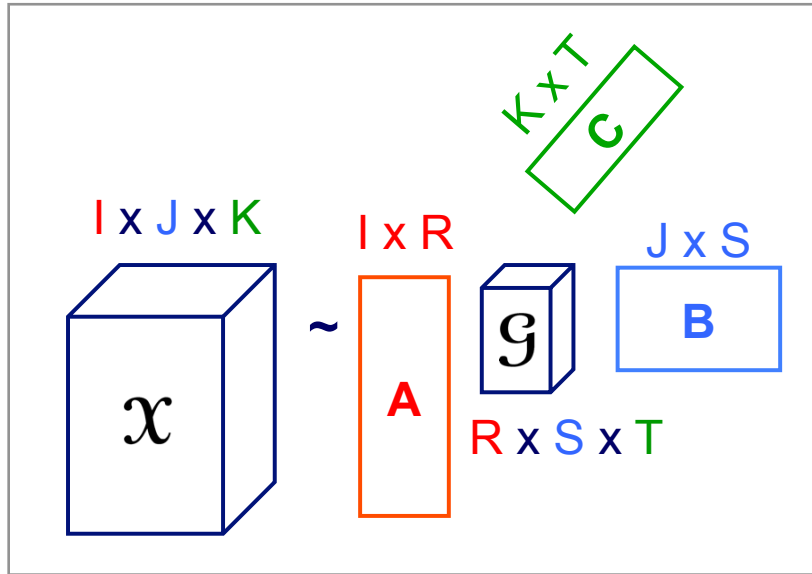


• Kruskal Tensor

$$\begin{aligned}\mathcal{X} &= \sum_r \lambda_r \mathbf{u}_r \circ \mathbf{v}_r \circ \mathbf{w}_r \\ &\equiv [\lambda; \mathbf{U}, \mathbf{V}, \mathbf{W}] \end{aligned} \left. \vphantom{\begin{aligned}\mathcal{X} \\ &= \sum_r \lambda_r \mathbf{u}_r \circ \mathbf{v}_r \circ \mathbf{w}_r \\ &\equiv [\lambda; \mathbf{U}, \mathbf{V}, \mathbf{W}]\end{aligned}} \right\} \text{Our Notation}$$



Tucker Decomposition - intuition



- author x keyword x conference
- A : author x author-group
- B : keyword x keyword-group
- C : conf. x conf-group
- \mathcal{G} : how groups relate to each other

Intuition behind core tensor

- 2-d case: co-clustering
- [Dhillon et al. Information-Theoretic Co-clustering, KDD'03]

med. doc cs doc

term group x
doc. group

$$\begin{bmatrix} .05 & .05 & .05 & 0 & 0 & 0 \\ .05 & .05 & .05 & 0 & 0 & 0 \\ 0 & 0 & 0 & .05 & .05 & .05 \\ 0 & 0 & 0 & .05 & .05 & .05 \\ .04 & .04 & 0 & .04 & .04 & .04 \\ .04 & .04 & .04 & 0 & .04 & .04 \end{bmatrix}$$

med. terms

cs terms

common terms

$$\begin{bmatrix} .5 & 0 & 0 \\ .5 & 0 & 0 \\ 0 & .5 & 0 \\ 0 & .5 & 0 \\ 0 & 0 & .5 \\ 0 & 0 & .5 \end{bmatrix}$$

$$\begin{bmatrix} .3 & 0 \\ 0 & .3 \\ .2 & .2 \end{bmatrix}$$

$$\begin{bmatrix} .36 & .36 & .28 & 0 & 0 & 0 \\ 0 & 0 & 0 & .28 & .36 & .36 \end{bmatrix} =$$

doc x
doc group


$$\begin{bmatrix} .054 & .054 & .042 & 0 & 0 & 0 \\ .054 & .054 & .042 & 0 & 0 & 0 \\ 0 & 0 & 0 & .042 & .054 & .054 \\ 0 & 0 & 0 & .042 & .054 & .054 \\ .036 & .036 & .028 & .028 & .036 & .036 \\ .036 & .036 & .028 & .028 & .036 & .036 \end{bmatrix}$$

term x
term-group

Tensor tools - summary

- Two main tools
 - PARAFAC
 - Tucker
- Both find row-, column-, tube-groups
 - but in PARAFAC the three groups are identical
- (To solve: Alternating Least Squares)

Detailed outline

- Motivation
- Definitions: PARAFAC and Tucker
-  • Case study: web mining

Web graph mining

- How to order the importance of web pages?
 - Kleinberg's algorithm HITS
 - PageRank
 - Tensor extension on HITS (**TOPHITS**)

Google Web Images Video News Maps more »

tensor Search Advanced Search Preferences

Turn OFF Personalized Search (Beta) for these results »

Web Personalized Results 1 - 10 of about 12,800,000 for tensor [definition]. (0.31 seconds)

Tensor - Wikipedia, the free encyclopedia
Examples of physical tensors are the energy-momentum tensor, the inertia tensor ... Tensorial 3.0 Tensorial is a general purpose tensor calculus package for ...
en.wikipedia.org/wiki/Tensor - 55k - Cached - Similar pages

Tensor product - Wikipedia, the free encyclopedia
There is a general formula for the product of two (or more) tensors, as ... The tensor product inherits all the indices of its factors. ...
en.wikipedia.org/wiki/Tensor_product - 41k - Cached - Similar pages

Tensor Trucks
Manufacturer of skateboard trucks. Check out team members, videos and apparel.
www.tensortrucks.com/ - 3k - Cached - Similar pages

Time and Attendance & Access Control through Smart Cards ...
Tensor manufacture and supply Smart Card and Biometric Time and Attendance & Access Control Software and Systems.
www.tensor.co.uk/ - 7k - Cached - Similar pages

Free Textbook Tensor Calculus and Continuum Mechanics
A free downloadable textbook on introductory tensor analysis and continuum

Sponsored Links

Tensor
Bargain Prices. Smart Deals. Save on Tensor!
Shopzilla.com

Wire Tensioners
For coil and motor winding machines Mechanical or electronic tensioners
www.digmotor.com

Tensor
Looking for Tensor? Find exactly what you want today
www.eBay.com

Tensor
Shop For Tensor Here With The Convenience Of OneCart™!
SHOP.COM

YAHOO! SEARCH

Web Images Video Local Shopping more »

tensor Search

Search Results 1 - 10 of about 2,870,000 for tensor - 0.74 sec. (About this page)

Also try [tensor lamps](#), [tensor lighting](#), [tensor corporation](#), [tensor product](#) More...

Sponsored Results

- Tensor Skateboard Trucks**
www.AllegroMedical.com - Great Selection and Fast Shipping Order Online Today and Save.
- Purchase Tensor Bandages at HCD**
www.homecaredelivered.com - Save on our full line of wound care supplies.

1. **Tensor - from MathWorld**
An n th-rank tensor in m -dimensional space is a mathematical object that has n ... Each index of a tensor ranges over the number of dimensions of space. ...
mathworld.wolfram.com/Tensor.html - More from this site

2. **Tensor - Wikipedia, the free encyclopedia**
The term 'tensor' has slightly different meanings in mathematics and physics. ... algebra and differential geometry, a tensor is a multilinear function. ...
Quick Links: [importance and applications](#) - [History](#) - [The choice of approach](#)
en.wikipedia.org/wiki/Tensor - 50k - Cached - More from this site

Sponsored Results

Tensor
Find Deals on Tensor and other Sporting Equipment at DealTime...
www.dealtime.com

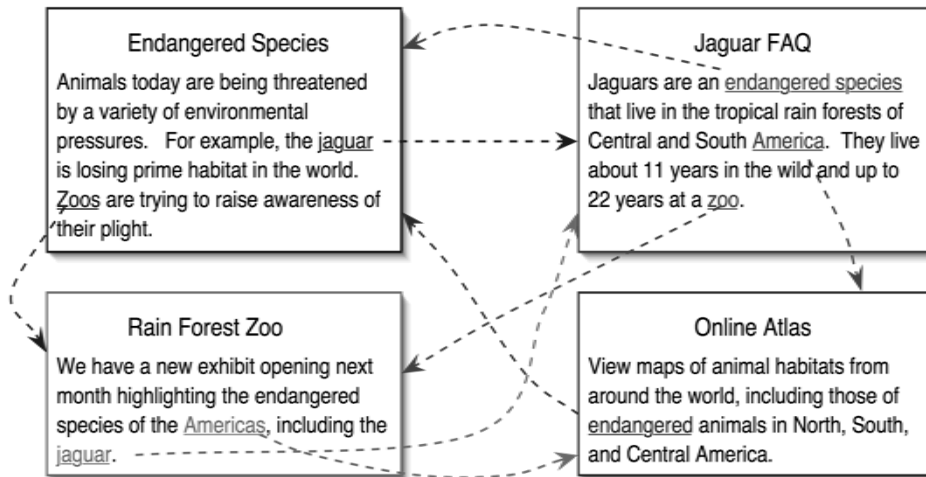
Tensor Compare Prices
Find Bargains on Tensor at thousands of trusted online stores. Get...
www.bizrate.com

Tensor
We are writing an on-line e-book with code: "Pseudocolor in Pure...
www.youvan.com

Tensor at Shopping.com
Find, compare and buy products in categories ranging from sports...
www.shopping.com

Tensor
Shop eBay for anything and

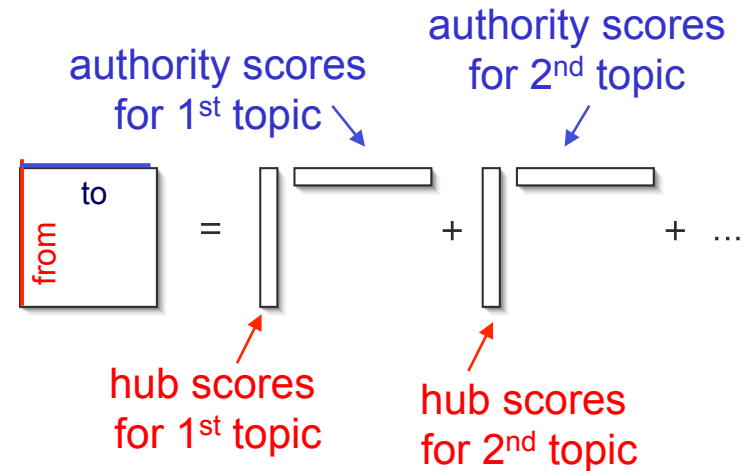
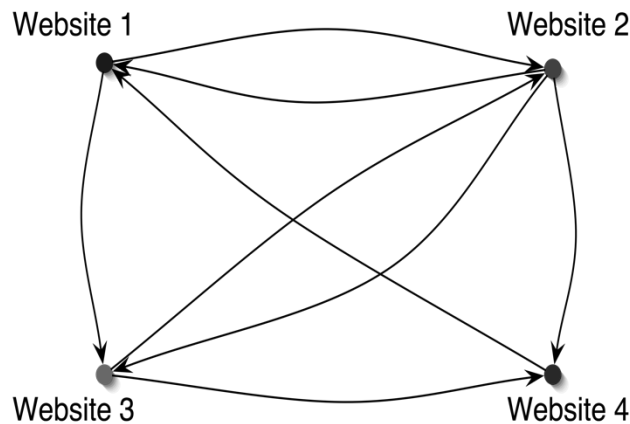
Kleinberg's Hubs and Authorities (the HITS method)



Sparse adjacency matrix and its SVD:

$$x_{ij} = \begin{cases} 1 & \text{if page } i \text{ links to page } j \\ 0 & \text{otherwise} \end{cases}$$

$$\mathbf{X} \approx \sum_r \sigma_r \mathbf{h}_r \circ \mathbf{a}_r$$



HITS Authorities on Sample Data

1st Principal Factor	
.97	www.ibm.com
.24	www.alphawo
.08	www-128.ibm
.05	www.develop
.02	www.research
.01	www.redbook
.01	news.com.cc

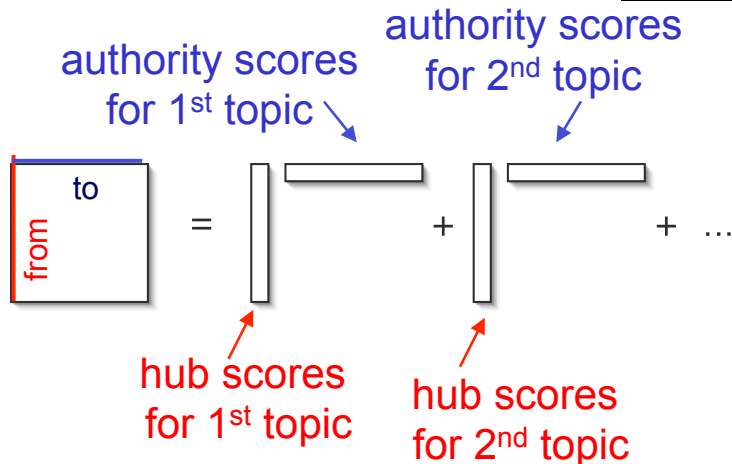
2nd Principal Factor	
.99	www.lehigh.edu
.11	www2.lehigh.edu
.06	www.lehigha
.06	www.lehighs
.02	www.bethleh
.02	www.adobe.c
.02	lewisweb.cc
.02	www.leo.lehi
.02	www.distanc
.02	fp1.cc.lehigh

3rd Principal Factor	
.75	java.sun.com
.38	www.sun.com
.36	developers.sun.
.24	see.sun.com
.16	www.samag.co
.13	docs.sun.com
.12	blogs.sun.com
.08	sunsolve.sun.c
.08	www.sun-catalo
.08	news.com.com

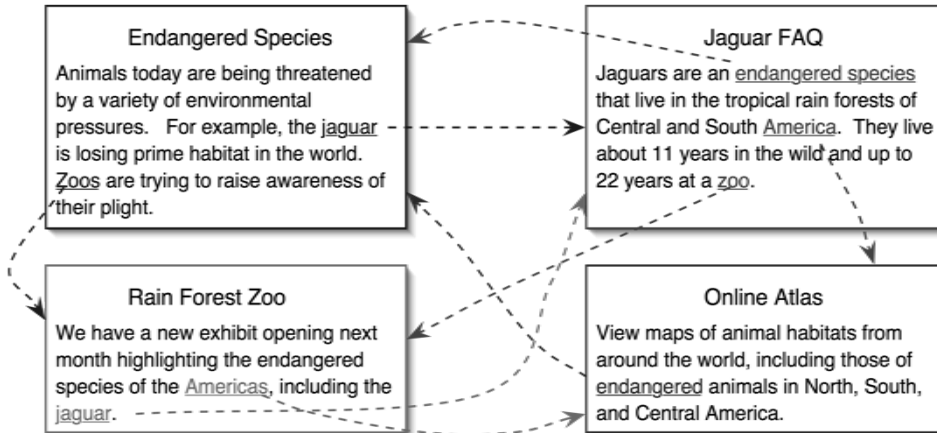
4th Principal Factor	
.60	www.pueblo.gsa.gov
.45	www.whitehouse.gov
.35	www.irs.gov
.31	travel.state
.22	www.gsa.g
.20	www.ssa.g
.16	www.censu
.14	www.govbe
.13	www.kids.g
.13	www.usdoj

We started our crawl from <http://www-neos.mcs.anl.gov/neos>, and crawled 4700 pages, resulting in 560 cross-linked hosts.

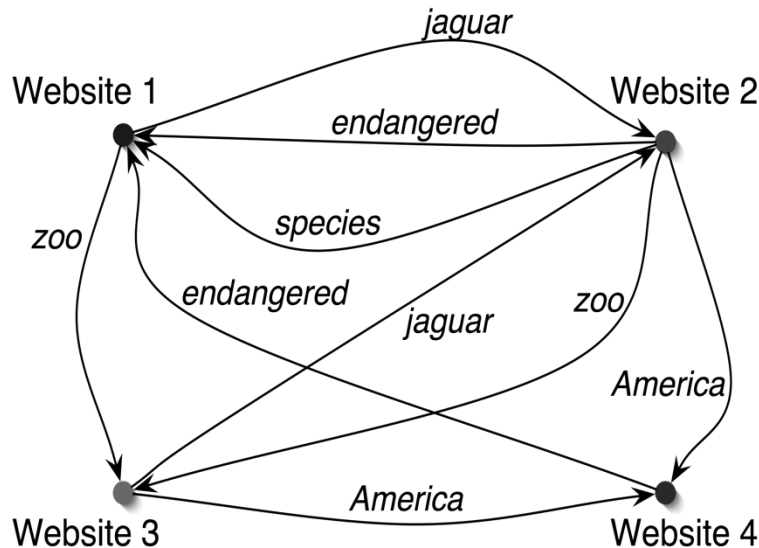
6th Principal Factor	
.97	mathpost.asu.edu
.18	math.la.asu.edu
.17	www.asu.edu
.04	www.act.org
.03	www.eas.asu.edu
.02	archives.math.utk.edu
.02	www.geom.uiuc.edu
.02	www.fulton.asu.edu
.02	www.amstat.org
.02	www.maa.org



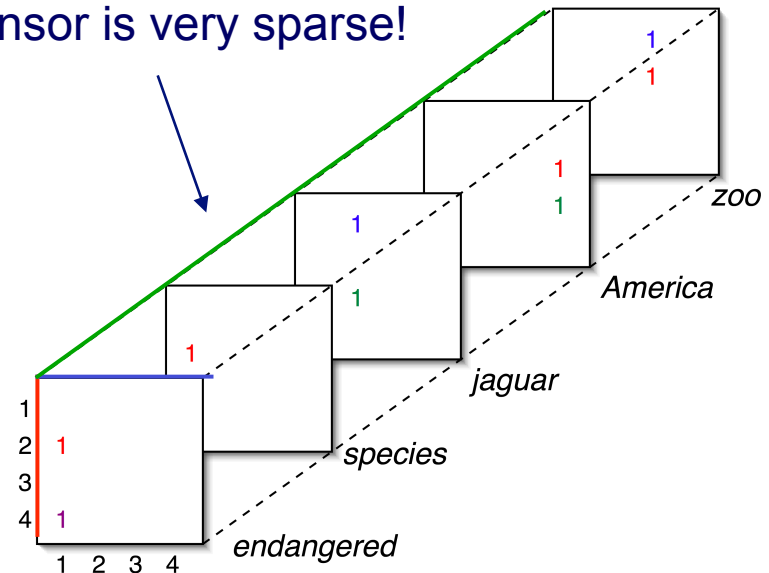
Three-Dimensional View of the Web



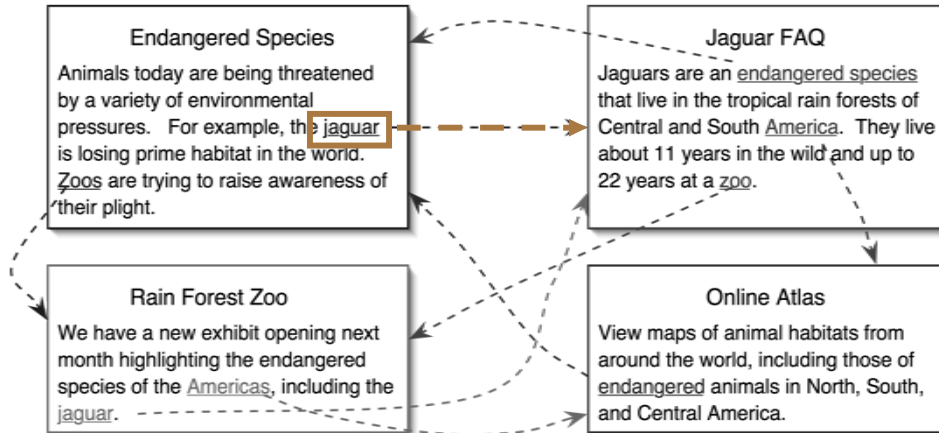
$$x_{ijk} = \begin{cases} 1 & \text{if page } i \rightarrow \text{page } j \\ & \text{with term } k \\ 0 & \text{otherwise} \end{cases}$$



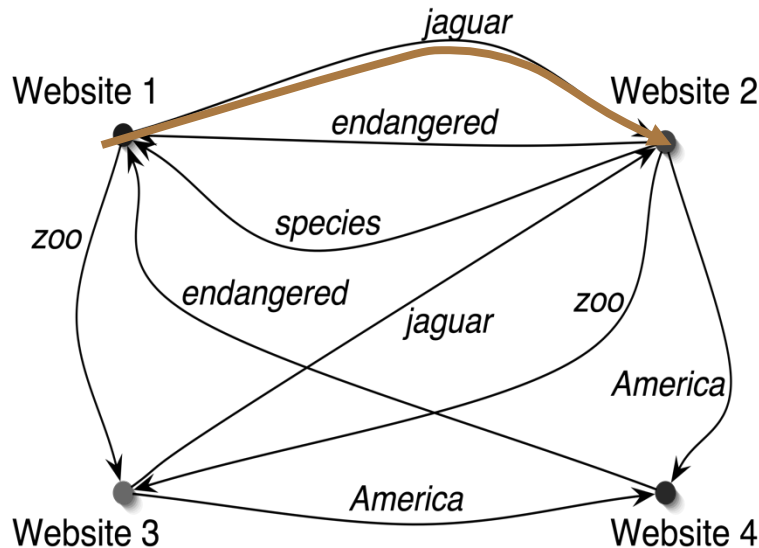
Observe that this tensor is very sparse!



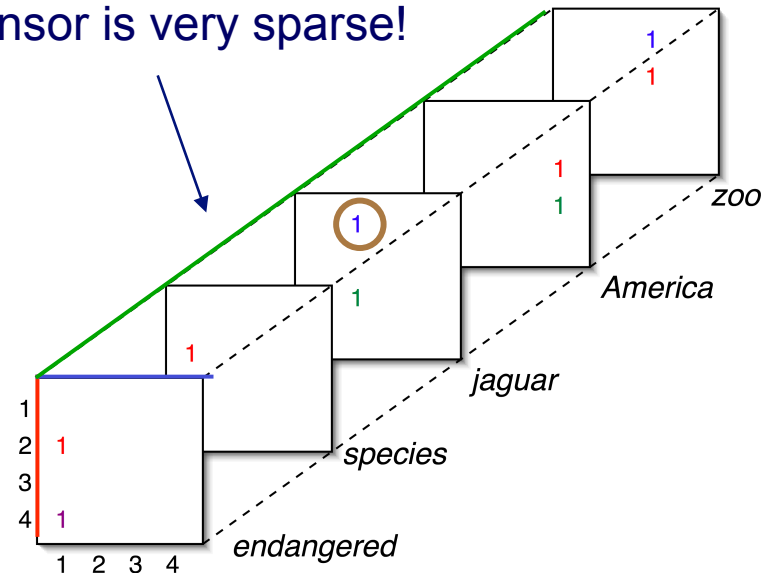
Three-Dimensional View of the Web



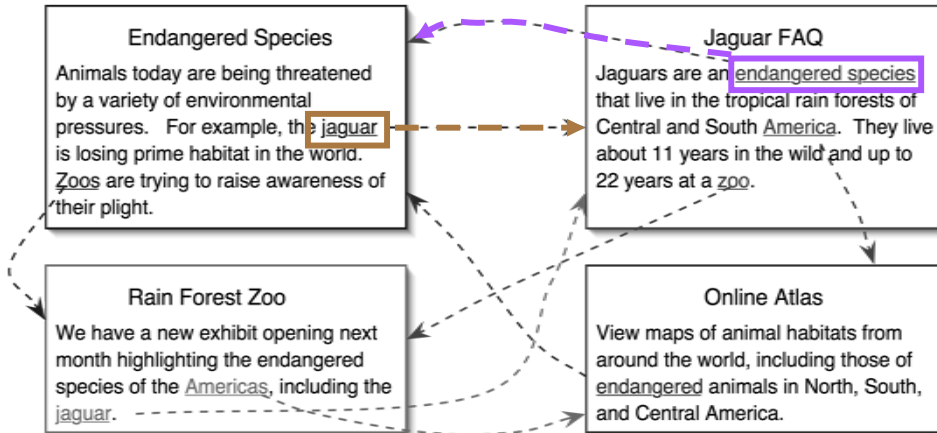
$$x_{ijk} = \begin{cases} 1 & \text{if page } i \rightarrow \text{page } j \\ & \text{with term } k \\ 0 & \text{otherwise} \end{cases}$$



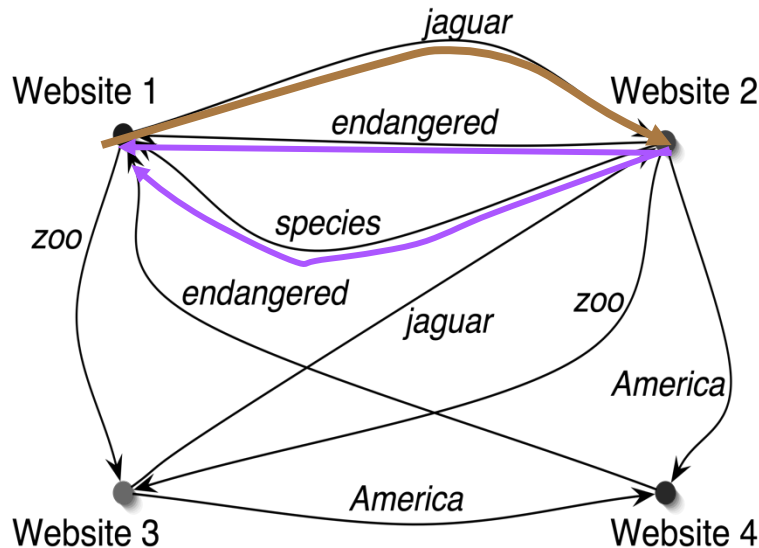
Observe that this tensor is very sparse!



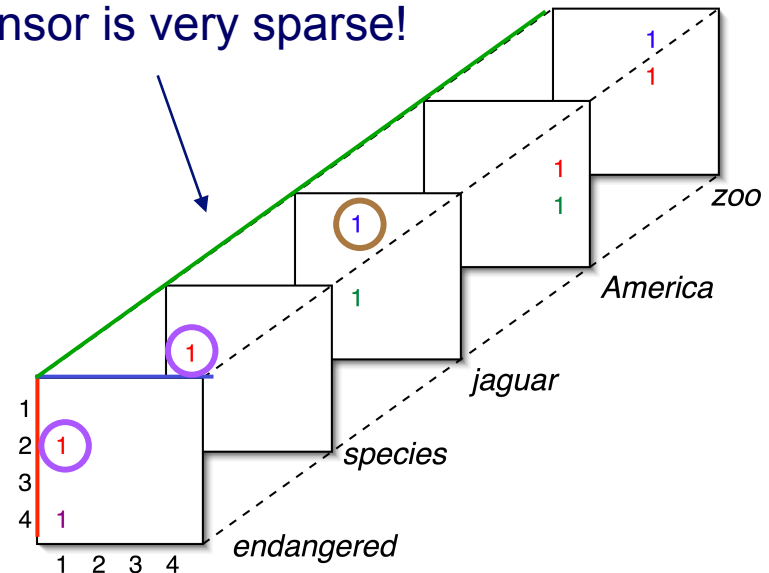
Three-Dimensional View of the Web



$$x_{ijk} = \begin{cases} 1 & \text{if page } i \rightarrow \text{page } j \\ & \text{with term } k \\ 0 & \text{otherwise} \end{cases}$$



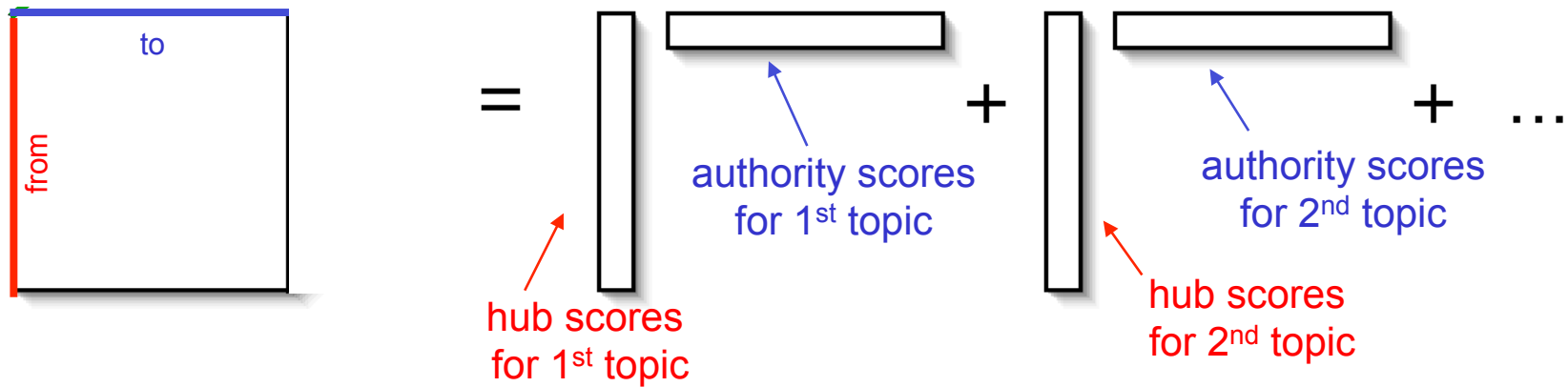
Observe that this tensor is very sparse!



Topical HITS (TOPHITS)

Main Idea: Extend the idea behind the HITS model to incorporate term (i.e., topical) information.

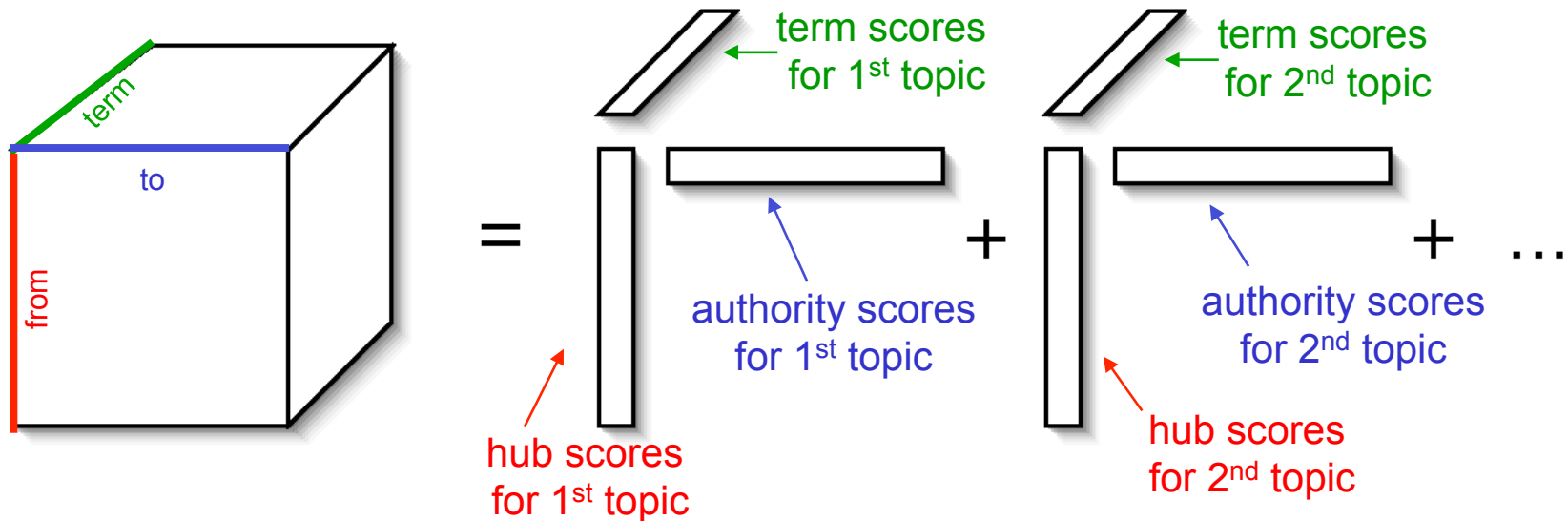
$$\mathbf{x} \approx \sum_{r=1}^R \lambda_r \mathbf{h}_r \circ \mathbf{a}_r$$



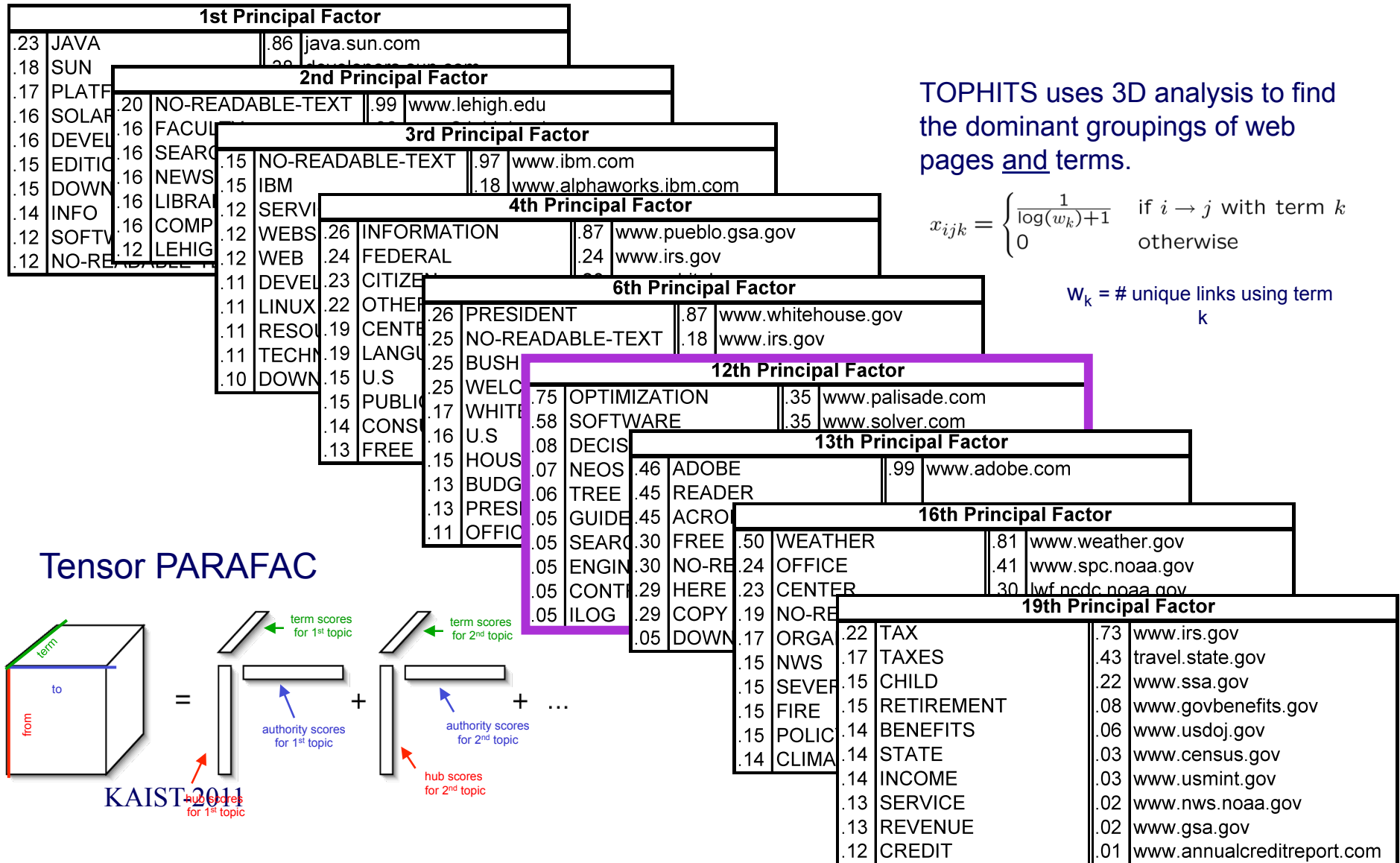
Topical HITS (TOPHITS)

Main Idea: Extend the idea behind the HITS model to incorporate term (i.e., topical) information.

$$\mathbf{x} \approx \sum_{r=1}^R \lambda_r \mathbf{h}_r \circ \mathbf{a}_r \circ \mathbf{t}_r$$



TOPHITS Terms & Authorities on Sample Data



TOPHITS uses 3D analysis to find the dominant groupings of web pages and terms.

$$x_{ijk} = \begin{cases} \frac{1}{\log(w_k)+1} & \text{if } i \rightarrow j \text{ with term } k \\ 0 & \text{otherwise} \end{cases}$$

$W_k = \frac{\# \text{ unique links using term } k}{k}$

Conclusions

- Real data are often in high dimensions with multiple aspects (modes)
- Tensors provide elegant theory and algorithms
 - PARAFAC and Tucker: discover groups

References

- T. G. Kolda, B. W. Bader and J. P. Kenny. *Higher-Order Web Link Analysis Using Multilinear Algebra*. In: ICDM 2005, Pages 242-249, November 2005.
- Jimeng Sun, Spiros Papadimitriou, Philip Yu. *Window-based Tensor Analysis on High-dimensional and Multi-aspect Streams*, Proc. of the Int. Conf. on Data Mining (ICDM), Hong Kong, China, Dec 2006

Resources

- See tutorial on tensors, KDD'07 (w/ Tamara Kolda and Jimeng Sun):

www.cs.cmu.edu/~christos/TALKS/KDD-07-tutorial

Tensor tools - resources



- Toolbox: from Tamara Kolda:
csmr.ca.sandia.gov/~tgkolda/TensorToolbox

• T. G. Kolda and B. W. Bader. ***Tensor Decompositions and Applications***. SIAM Review, Volume 51, Number 3, September 2009
csmr.ca.sandia.gov/~tgkolda/pubs/bibtgkfiles/TensorReview-preprint.pdf

• T. Kolda and J. Sun: Scalable Tensor Decomposition for Multi-Aspect Data Mining (ICDM 2008)

Outline

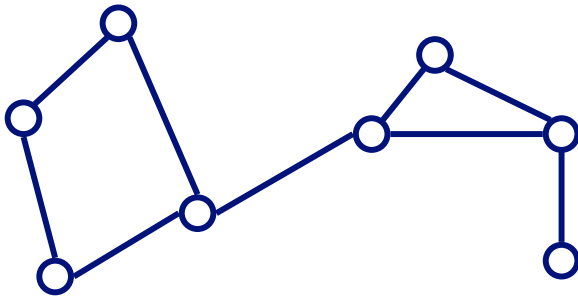
- Task 4: time-evolving graphs – tensors
- ➔ • Task 5: community detection
- Task 6: virus propagation
- Task 7: scalability, parallelism and hadoop
- Conclusions

Detailed outline

- Motivation
- ➔ • Hard clustering – k pieces
- Hard co-clustering – (k, l) pieces
- Hard clustering – optimal # pieces
- Observations

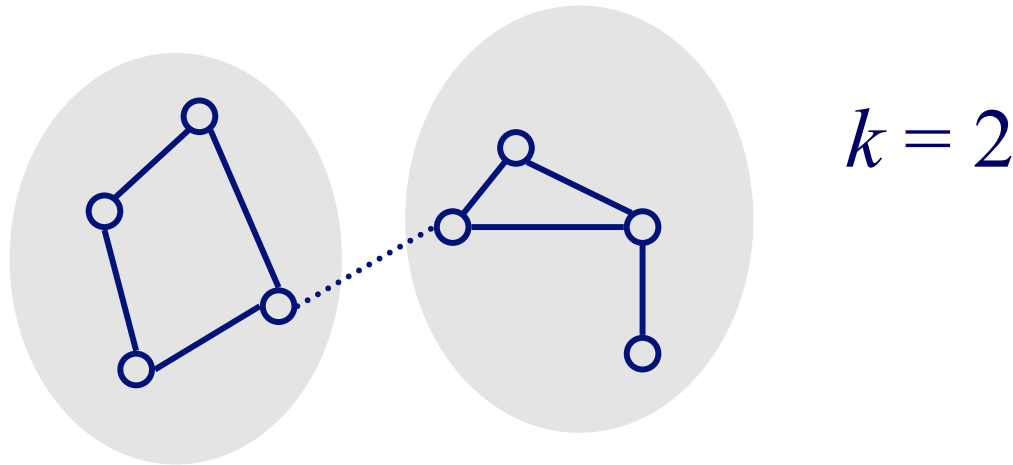
Problem

- Given a graph, and k
- Break it into k (disjoint) communities



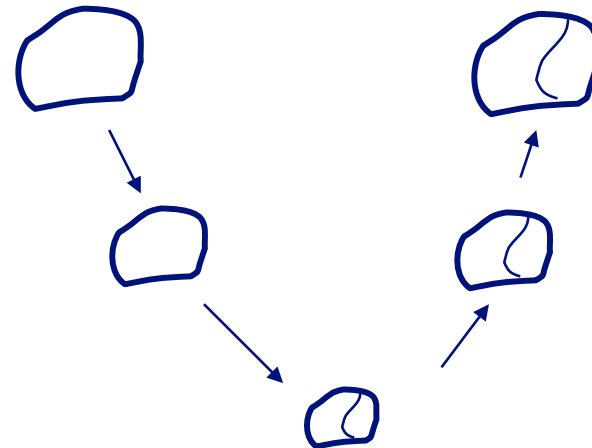
Problem

- Given a graph, and k
- Break it into k (disjoint) communities



Solution #1: METIS

- Arguably, the best algorithm
- Open source, at
 - <http://www.cs.umn.edu/~metis>
- and *many* related papers, at same url
- Main idea:
 - coarsen the graph;
 - partition;
 - un-coarsen



Solution #1: METIS

- G. Karypis and V. Kumar. *METIS 4.0: Unstructured graph partitioning and sparse matrix ordering system*. TR, Dept. of CS, Univ. of Minnesota, 1998.
- <and many extensions>



Solution #2

(problem: hard clustering, k pieces)

Spectral partitioning:

- Consider the 2nd smallest eigenvector of the (normalized) Laplacian

Solutions #3, ...

Many more ideas:

- Clustering on the A^2 (square of adjacency matrix) [Zhou, Woodruff, PODS'04]
- Minimum cut / maximum flow [Flake+, KDD'00]
- ...

Detailed outline

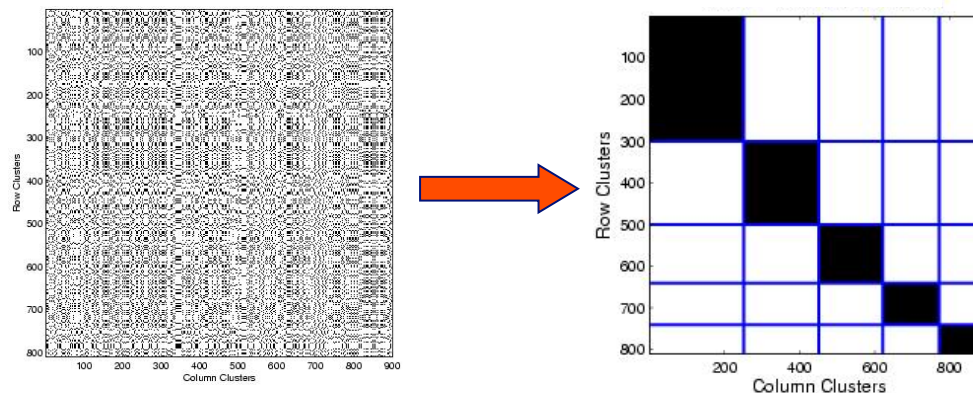
- Motivation
- Hard clustering – k pieces
- ➔ • Hard co-clustering – (k, l) pieces
- Hard clustering – optimal # pieces
- Soft clustering – matrix decompositions
- Observations

Problem definition

- Given a bi-partite graph, and k, l
- Divide it into k row groups and l row groups
- (Also applicable to uni-partite graph)

Co-clustering

- Given data matrix and the number of row and column groups k and l
- Simultaneously
 - Cluster rows into k disjoint groups
 - Cluster columns into l disjoint groups



Co-clustering

- Let X and Y be discrete random variables
 - X and Y take values in $\{1, 2, \dots, m\}$ and $\{1, 2, \dots, n\}$
 - $p(X, Y)$ denotes the joint probability distribution—if not known, it is often estimated based on co-occurrence data
 - Application areas: text mining, market-basket analysis, analysis of browsing behavior, etc.
- Key Obstacles in Clustering Contingency Tables
 - High Dimensionality, Sparsity, Noise
 - Need for robust and scalable algorithms

Reference:

1. Dhillon et al. Information-Theoretic Co-clustering, KDD'03

med. doc cs doc

term group x
doc. group

$$\begin{bmatrix} .05 & .05 & .05 & 0 & 0 & 0 \\ .05 & .05 & .05 & 0 & 0 & 0 \\ 0 & 0 & 0 & .05 & .05 & .05 \\ 0 & 0 & 0 & .05 & .05 & .05 \\ .04 & .04 & 0 & .04 & .04 & .04 \\ .04 & .04 & .04 & 0 & .04 & .04 \end{bmatrix}$$

| med. terms

| cs terms

| common terms

$$\begin{bmatrix} .5 & 0 & 0 \\ .5 & 0 & 0 \\ 0 & .5 & 0 \\ 0 & .5 & 0 \\ 0 & 0 & .5 \\ 0 & 0 & .5 \end{bmatrix}$$

$$\begin{bmatrix} .3 & 0 \\ 0 & .3 \\ .2 & .2 \end{bmatrix}$$

$$\begin{bmatrix} .36 & .36 & .28 & 0 & 0 & 0 \\ 0 & 0 & 0 & .28 & .36 & .36 \end{bmatrix} =$$

doc x
doc group

$$\begin{bmatrix} .054 & .054 & .042 & 0 & 0 & 0 \\ .054 & .054 & .042 & 0 & 0 & 0 \\ 0 & 0 & 0 & .042 & .054 & .054 \\ 0 & 0 & 0 & .042 & .054 & .054 \\ .036 & .036 & .028 & .028 & .036 & .036 \\ .036 & .036 & .028 & .028 & .036 & .036 \end{bmatrix}$$

term x
term-group

Co-clustering

Observations

- uses KL divergence, instead of L2
- the middle matrix is **not** diagonal
 - we saw that earlier in the Tucker tensor decomposition
- s/w at:

www.cs.utexas.edu/users/dml/Software/cocluster.html

Detailed outline

- Motivation
- Hard clustering – k pieces
- Hard co-clustering – (k,l) pieces
- ➔ • Hard clustering – optimal # pieces
- Soft clustering – matrix decompositions
- Observations

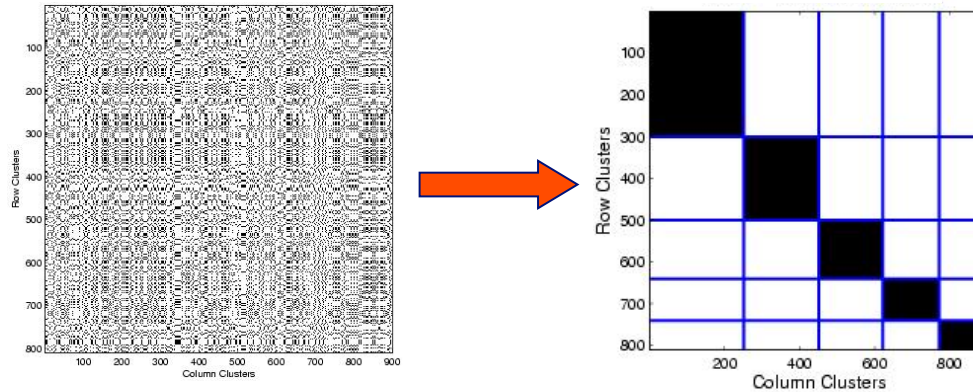
Problem with Information Theoretic Co-clustering

- Number of row and column groups must be specified

Desiderata:

- ✓ **Simultaneously discover** row and column groups
- ✗ **Fully Automatic:** No “magic numbers”
- ✓ **Scalable** to large graphs

Cross-association



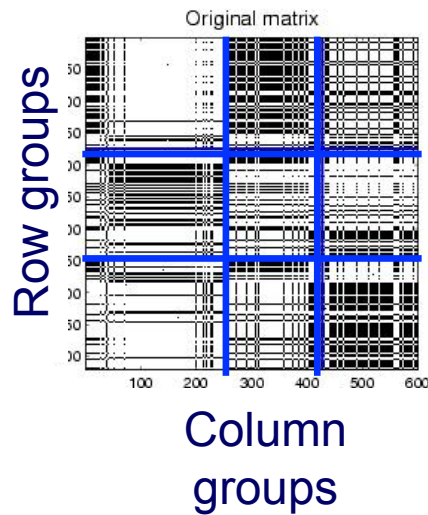
Desiderata:

- ✓ **Simultaneously discover** row and column groups
- ✓ **Fully Automatic:** No “magic numbers”
- ✓ **Scalable** to large matrices

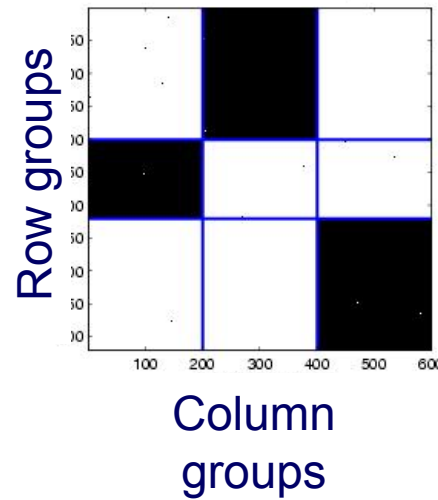
Reference:

1. Chakrabarti et al. Fully Automatic Cross-Associations, KDD'04

What makes a cross-association “good”?

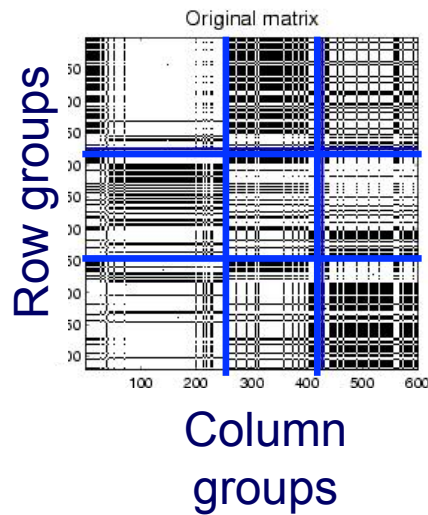


versus

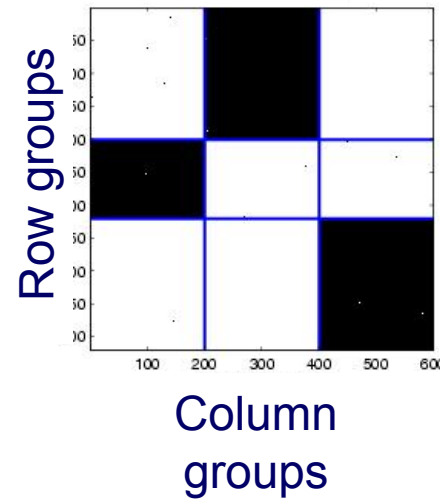


Why is this better?

What makes a cross-association “good”?



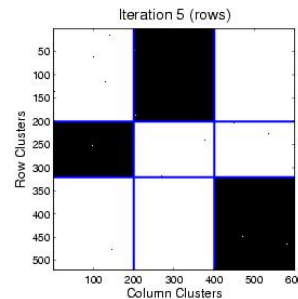
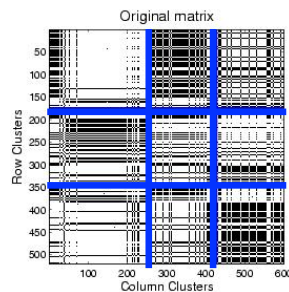
versus



Why is this
better?

simpler; easier to describe
easier to compress!

What makes a cross-association “good”?



Problem definition: given an encoding scheme

- decide on the # of col. and row groups k and l
- and reorder rows and columns,
- to achieve best compression

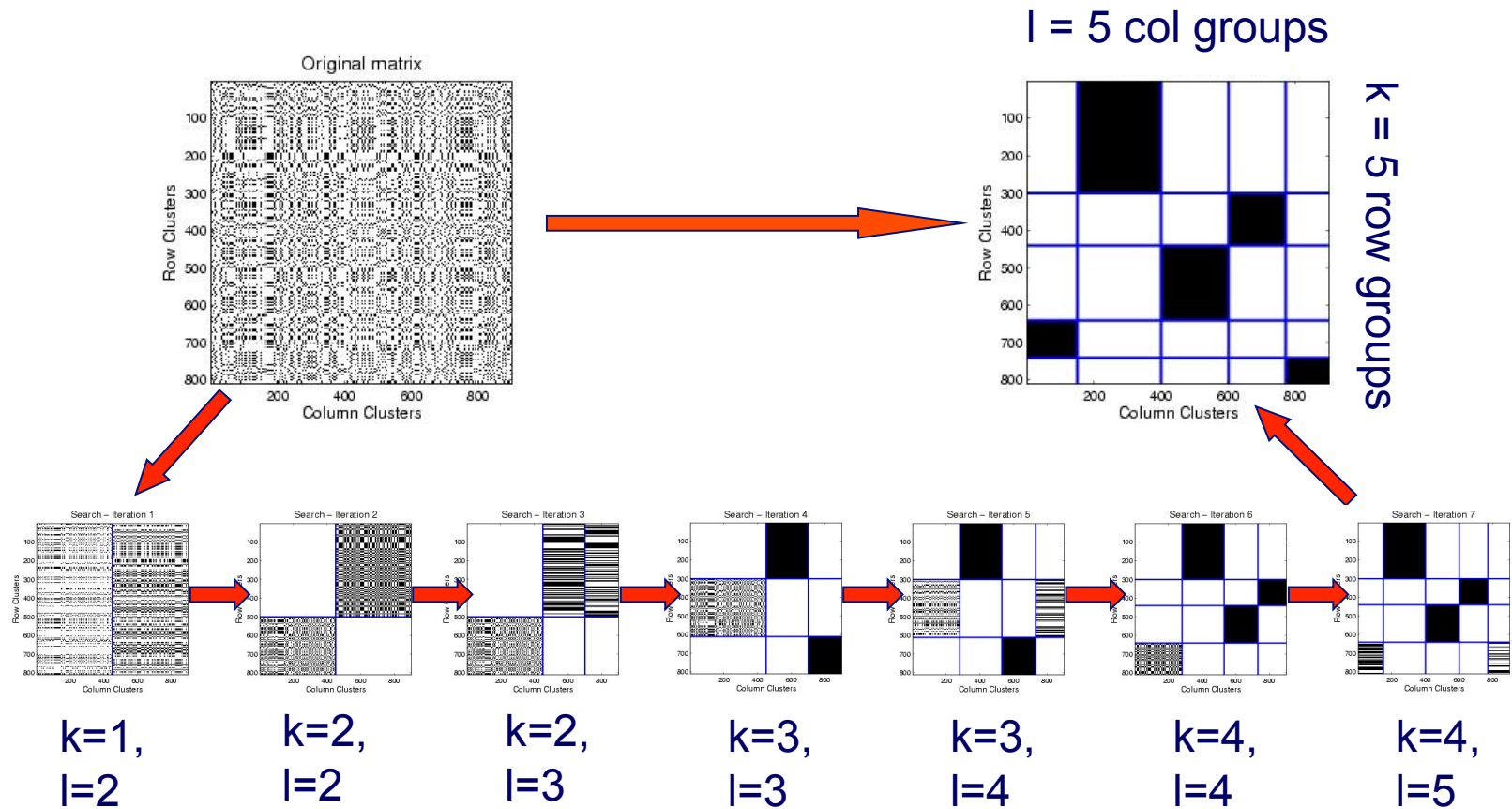
Main Idea



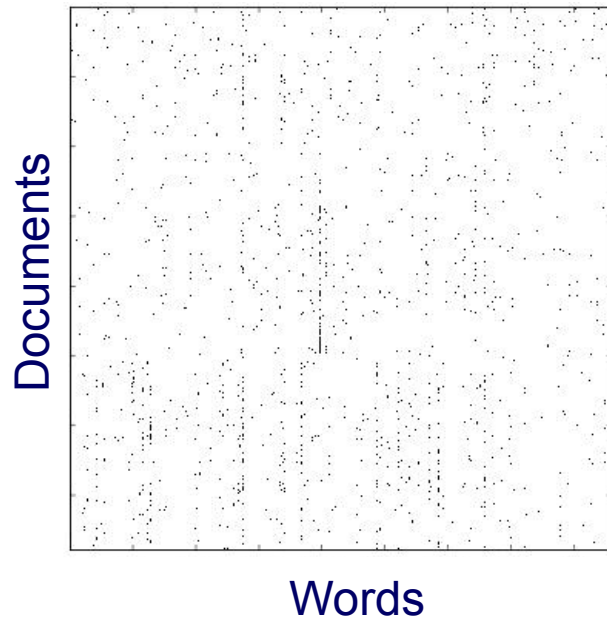
$$\text{Total Encoding Cost} = \underbrace{\sum_i \text{size}_i * H(x_i)}_{\text{Code Cost}} + \underbrace{\text{Cost of describing cross-associations}}_{\text{Description Cost}}$$

Minimize the total cost (# bits)
for lossless compression

Algorithm



Experiments



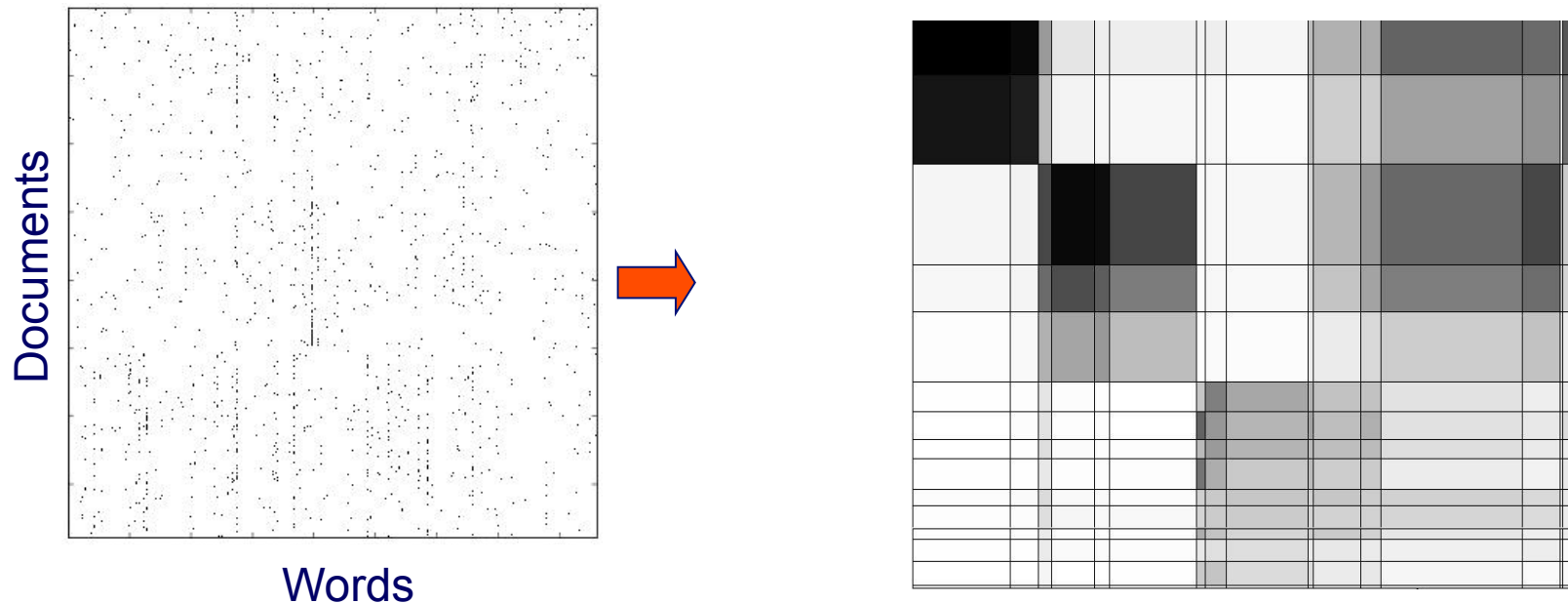
“CLASSIC”

- 3,893 documents
- 4,303 words
- 176,347 “dots”

Combination of 3 sources:

- MEDLINE (medical)
- CISI (info. retrieval)
- CRANFIELD (aerodynamics)

Experiments



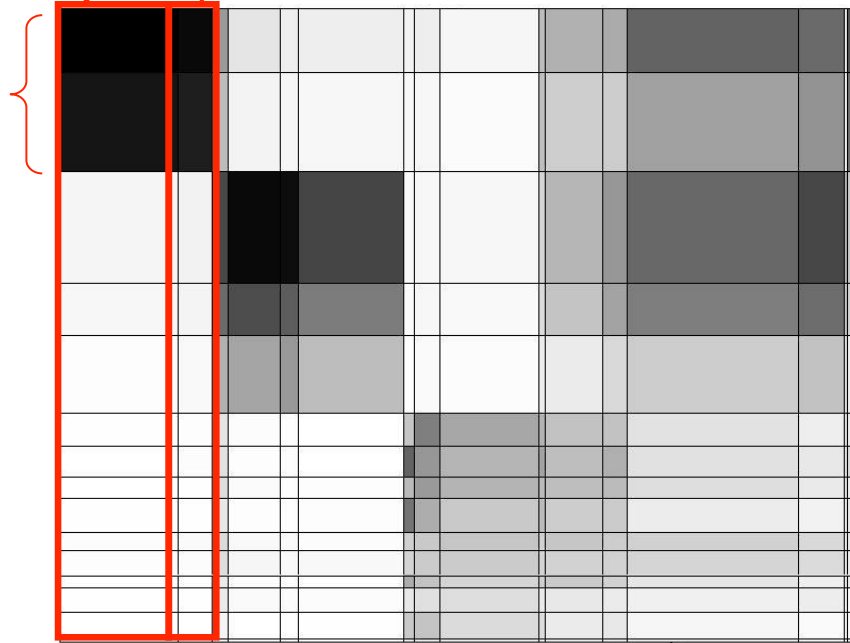
“CLASSIC” graph of documents & words: $k=15$, $l=19$

Experiments

insipidus, alveolar, aortic,
death, prognosis, intravenous

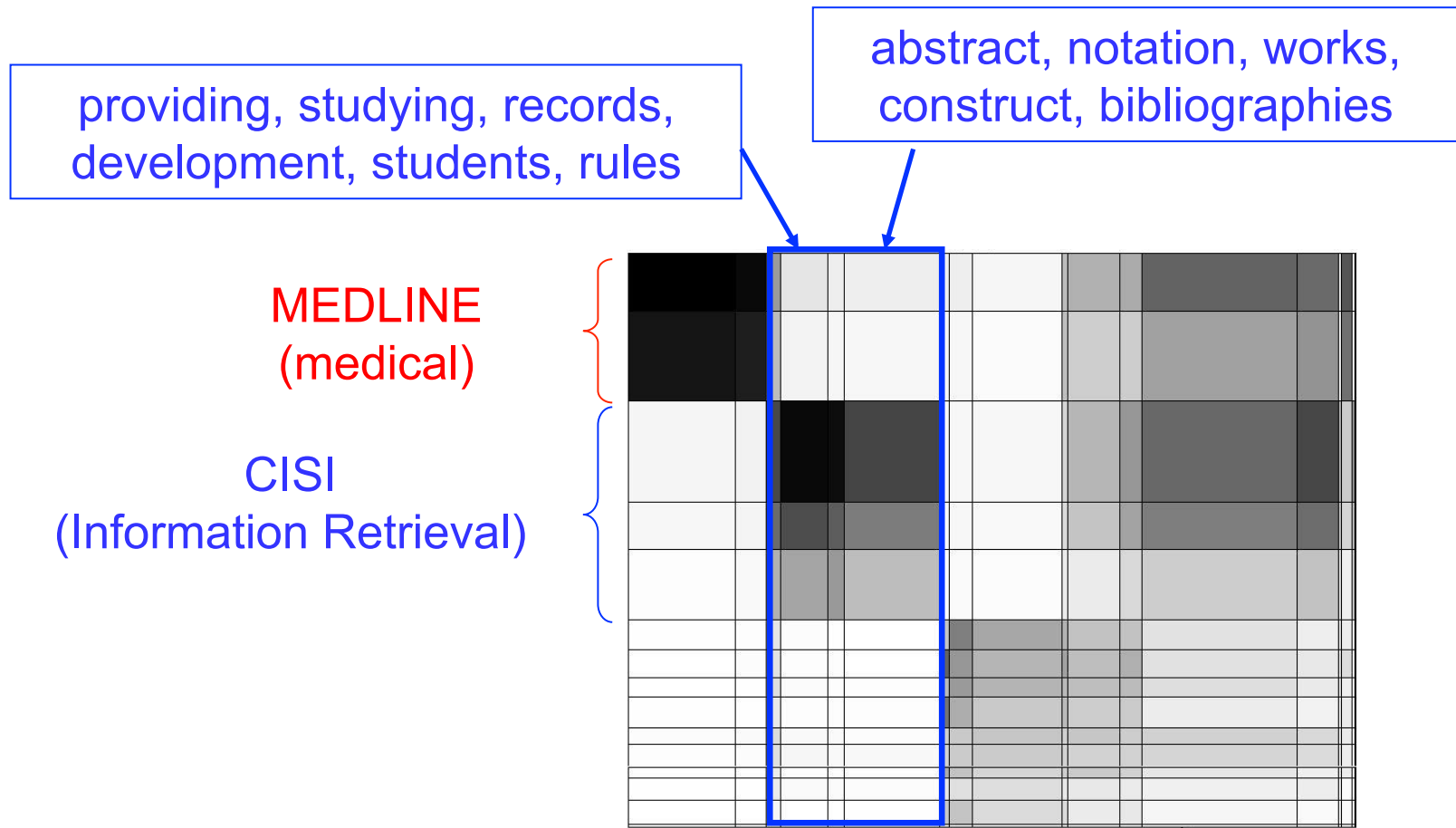
blood, disease, clinical, cell,
tissue, patient

MEDLINE
(medical)



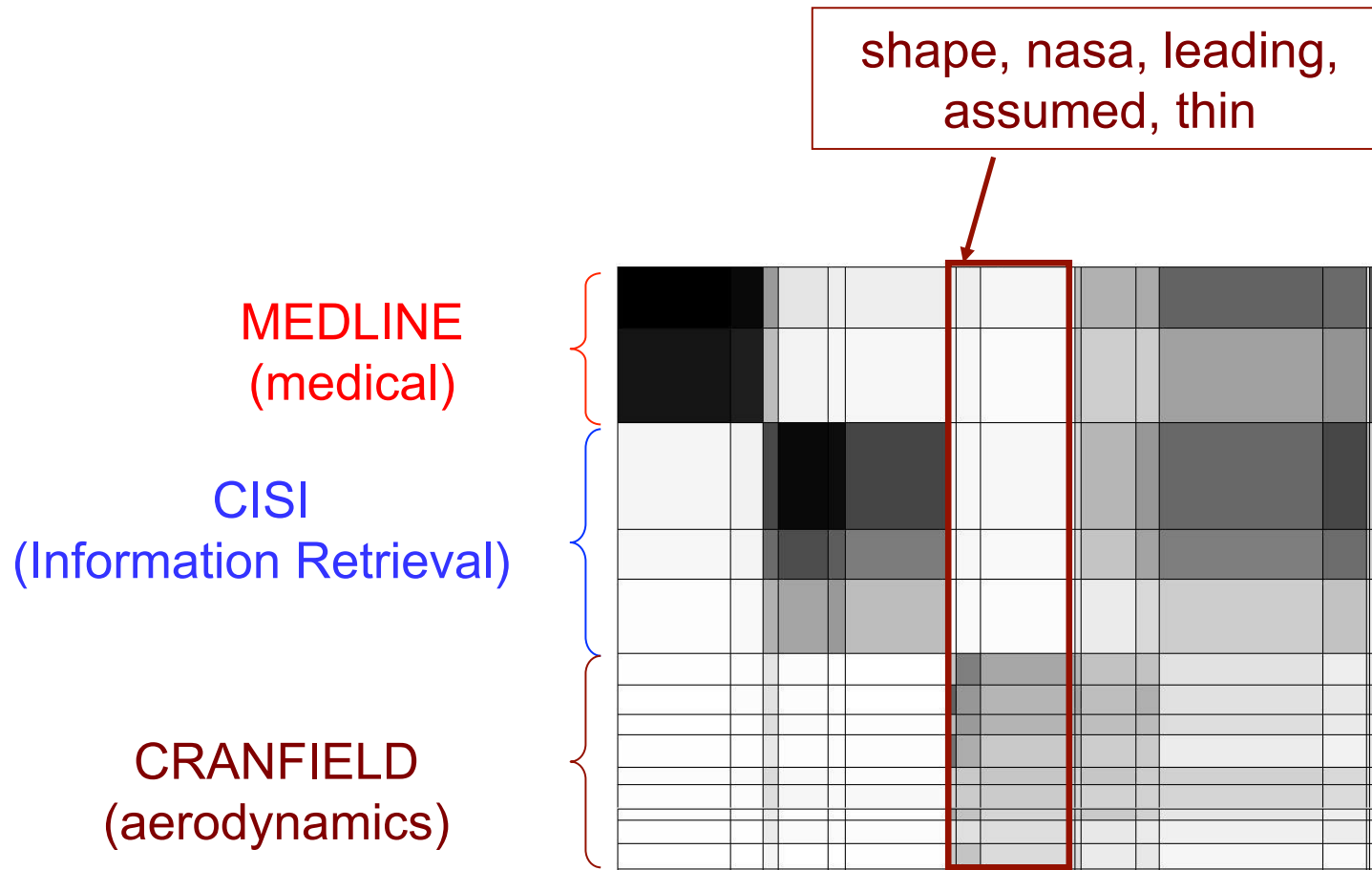
“CLASSIC” graph of documents &
words: $k=15$, $l=19$

Experiments



“CLASSIC” graph of documents & words: $k=15$, $l=19$

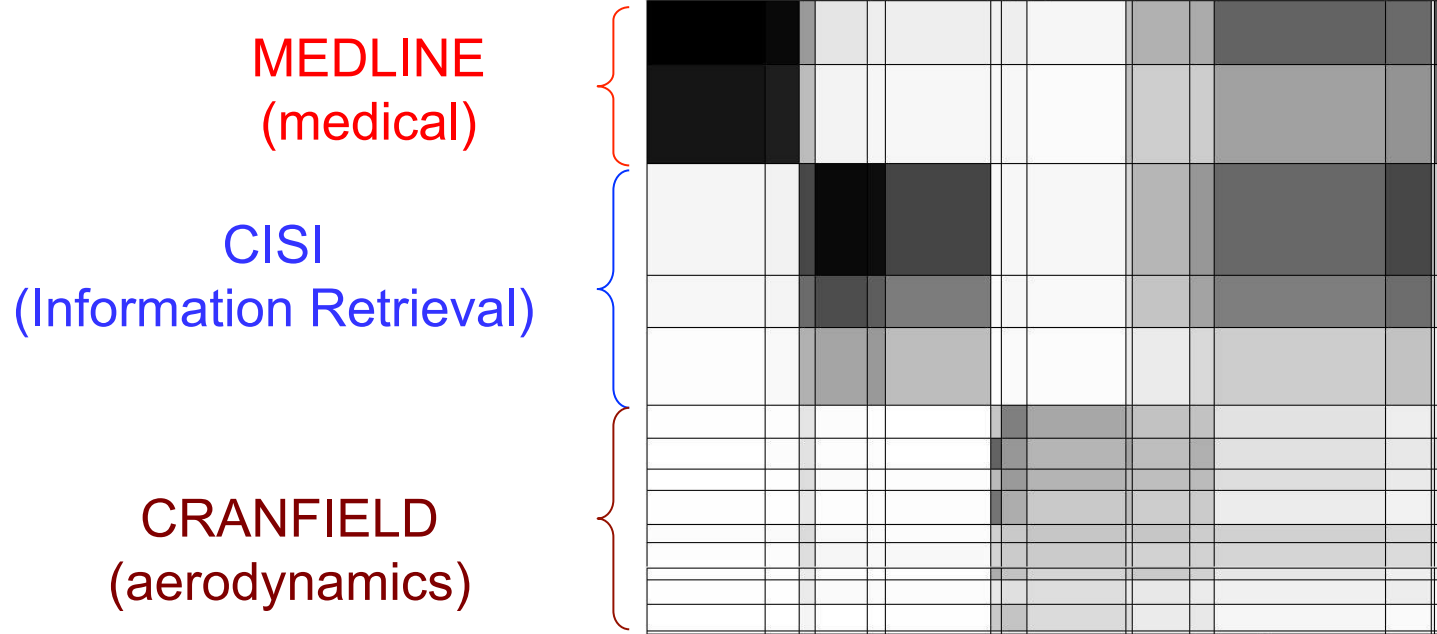
Experiments



“CLASSIC” graph of documents &
words: $k=15$, $l=19$

Experiments

paint, examination, fall,
raise, leave, based



“CLASSIC” graph of documents &
words: $k=15$, $l=19$

Algorithm

Code for cross-associations (matlab):

[www.cs.cmu.edu/~deepay/mywww/software/
CrossAssociations-01-27-2005.tgz](http://www.cs.cmu.edu/~deepay/mywww/software/CrossAssociations-01-27-2005.tgz)

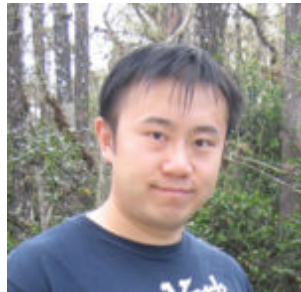
Variations and extensions:

- ‘Autopart’ [Chakrabarti, PKDD’04]
- www.cs.cmu.edu/~deepay



Algorithm

- Hadoop implementation [ICDM'08]



Spiros Papadimitriou, Jimeng Sun: DisCo: Distributed Co-clustering with Map-Reduce: A Case Study towards Petabyte-Scale End-to-End Mining. ICDM 2008: 512-521

Detailed outline

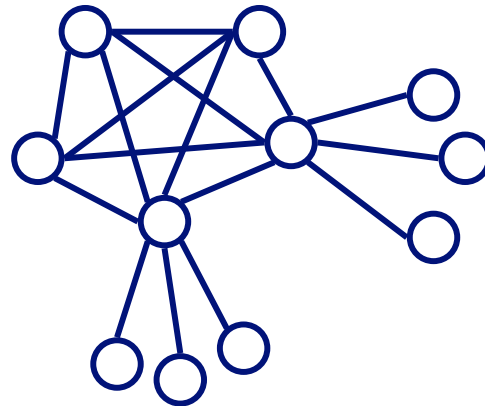
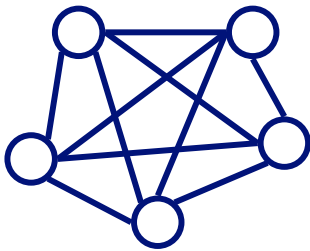
- Motivation
- Hard clustering – k pieces
- Hard co-clustering – (k, l) pieces
- Hard clustering – optimal # pieces
- ➔ • Observations

Observation #1

- Skewed degree distributions – there are nodes with huge degree ($>O(10^4)$, in facebook/linkedin popularity contests!)

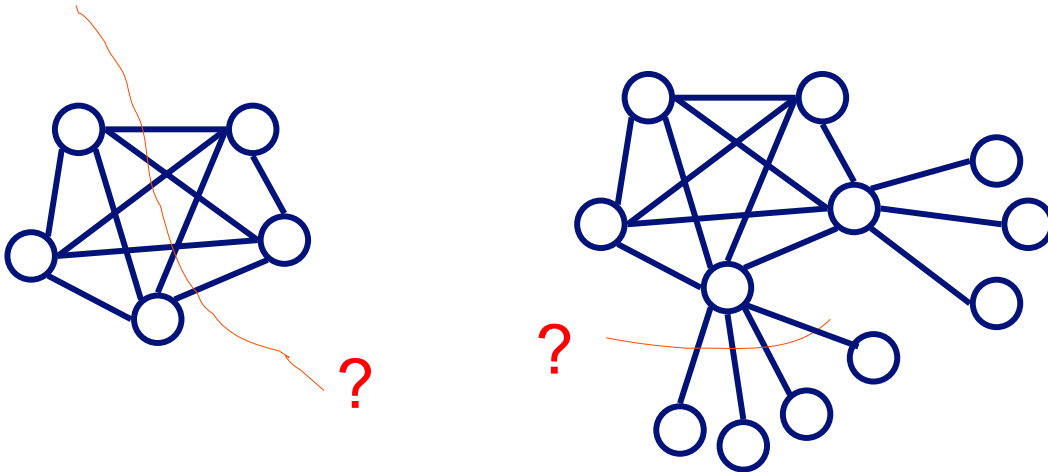
Observation #2

- Maybe there are no good cuts: “jellyfish” shape [Tauro+’01], [Siganos+,’06], strange behavior of cuts [Chakrabarti+’04], [Leskovec+,’08]

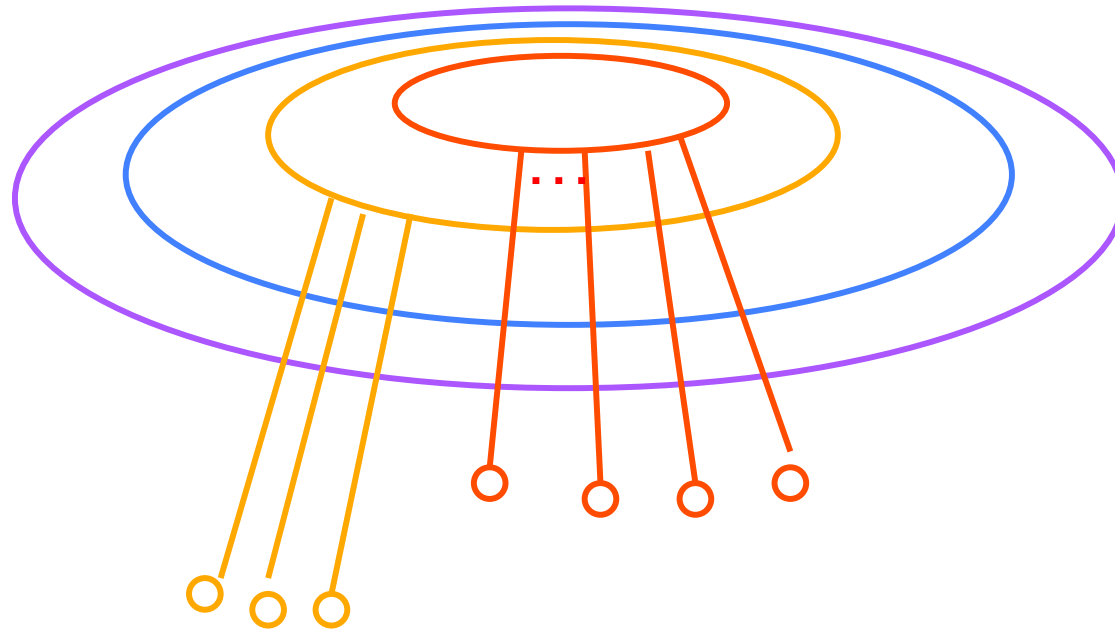


Observation #2

- Maybe there are no good cuts: “jellyfish” shape [Tauro+’01], [Siganos+,’06], strange behavior of cuts [Chakrabarti+,’04], [Leskovec+,’08]



Jellyfish model [Tauro+]



A Simple Conceptual Model for the Internet Topology, L. Tauro, C. Palmer, G. Siganos, M. Faloutsos, Global Internet, November 25-29, 2001

Jellyfish: A Conceptual Model for the AS Internet Topology G. Siganos, Sudhir L. Tauro, M. Faloutsos, J. of Communications and Networks, Vol. 8, No. 3, pp 339-350, Sept. 2006.

Strange behavior of min cuts

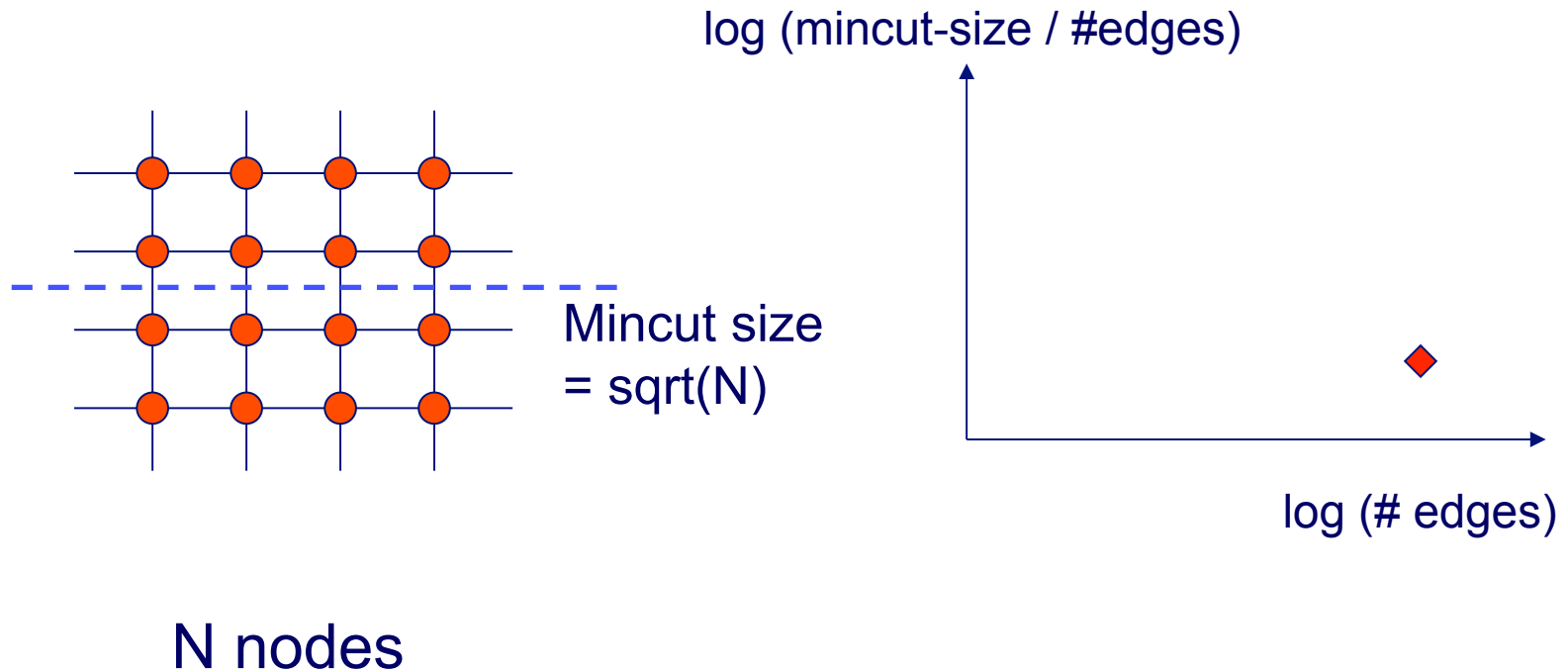
- ‘negative dimensionality’ (!)

NetMine: New Mining Tools for Large Graphs, by D. Chakrabarti, Y. Zhan, D. Blandford, C. Faloutsos and G. Blelloch, in the SDM 2004 Workshop on Link Analysis, Counter-terrorism and Privacy

Statistical Properties of Community Structure in Large Social and Information Networks, J. Leskovec, K. Lang, A. Dasgupta, M. Mahoney. WWW 2008.

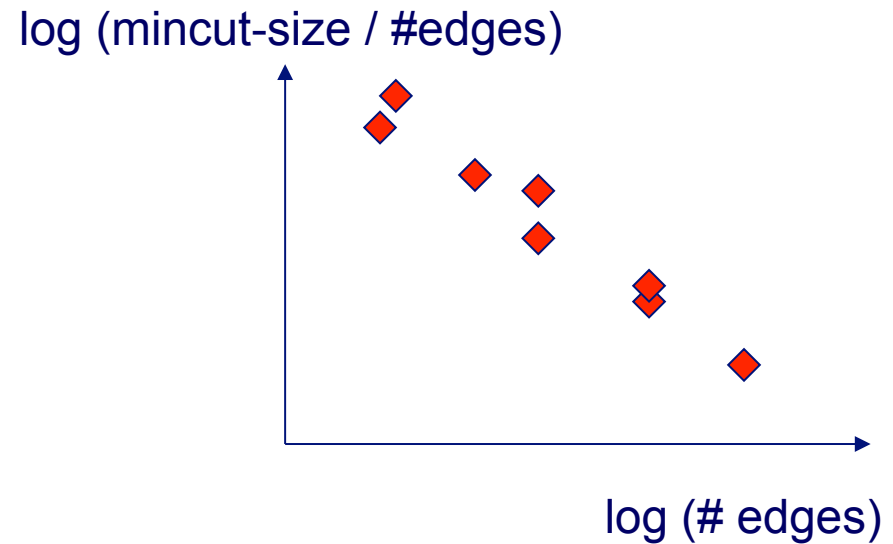
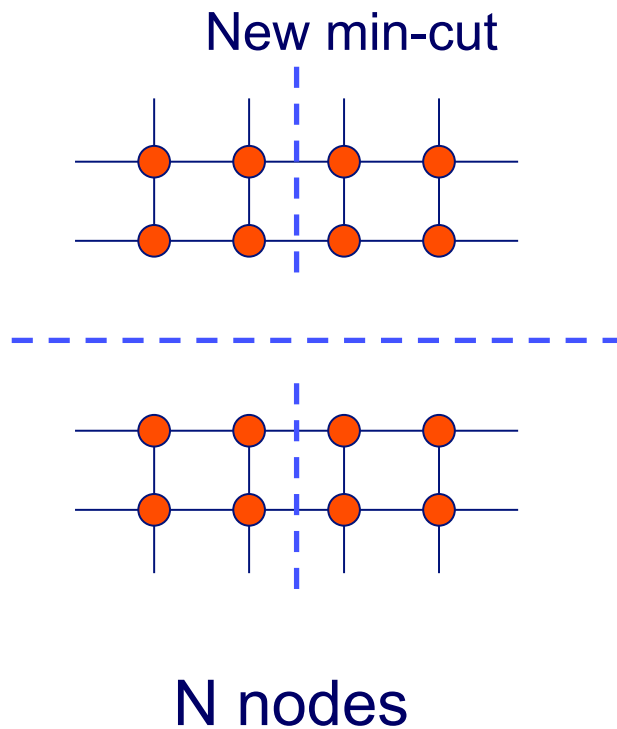
“Min-cut” plot

- Do min-cuts recursively.



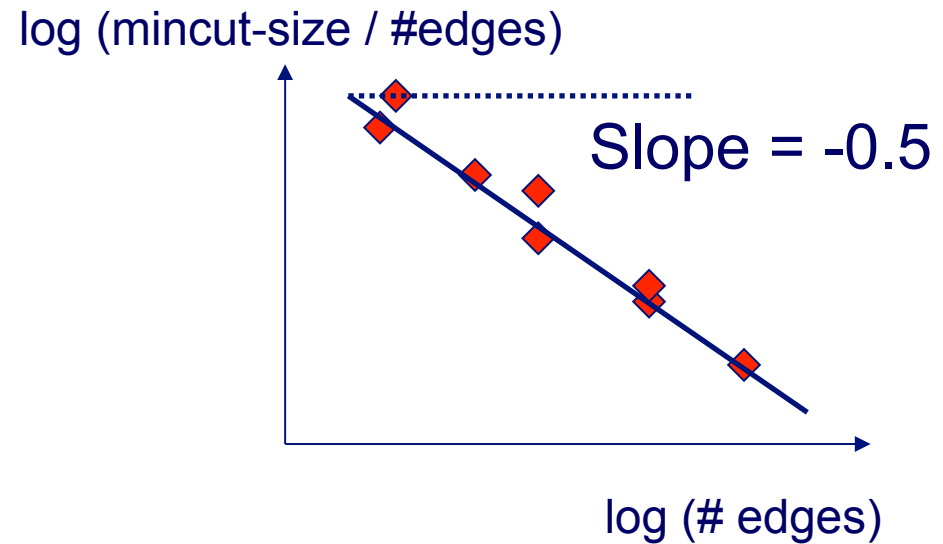
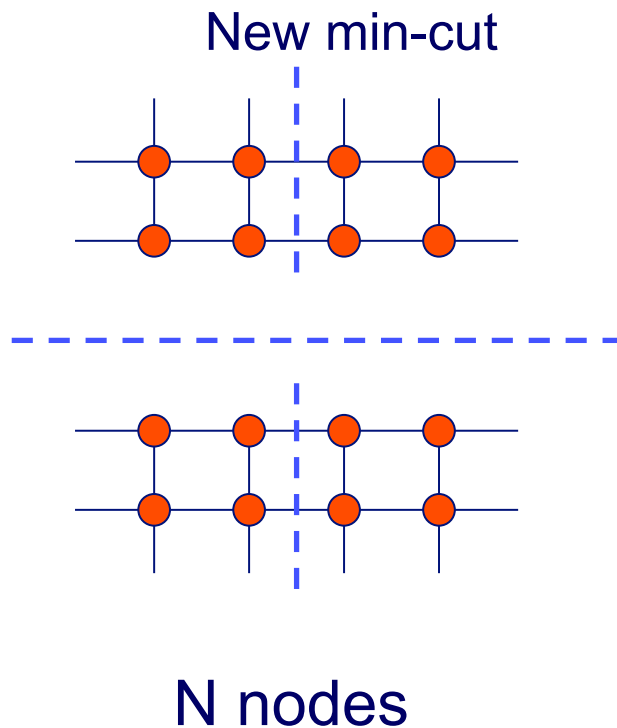
“Min-cut” plot

- Do min-cuts recursively.



“Min-cut” plot

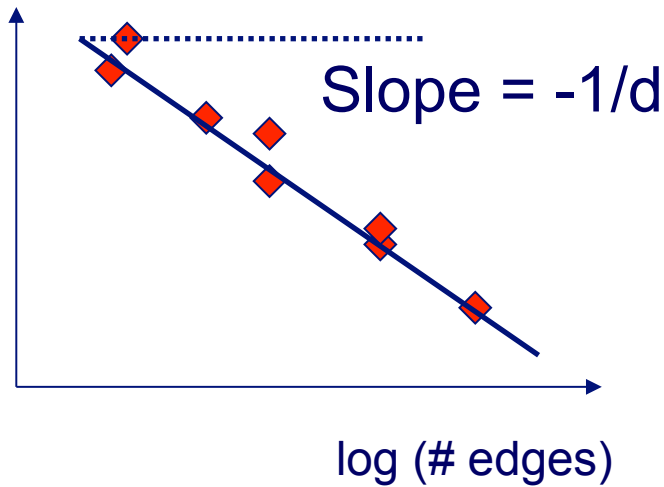
- Do min-cuts recursively.



For a d -dimensional grid, the slope is $-1/d$

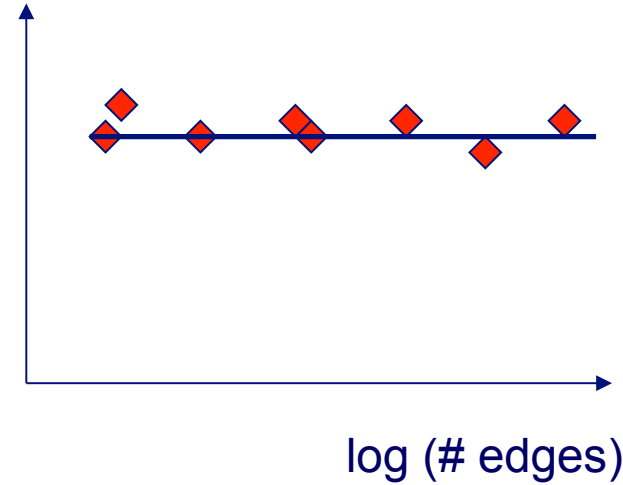
“Min-cut” plot

log (mincut-size / #edges)



For a d -dimensional grid, the slope is $-1/d$

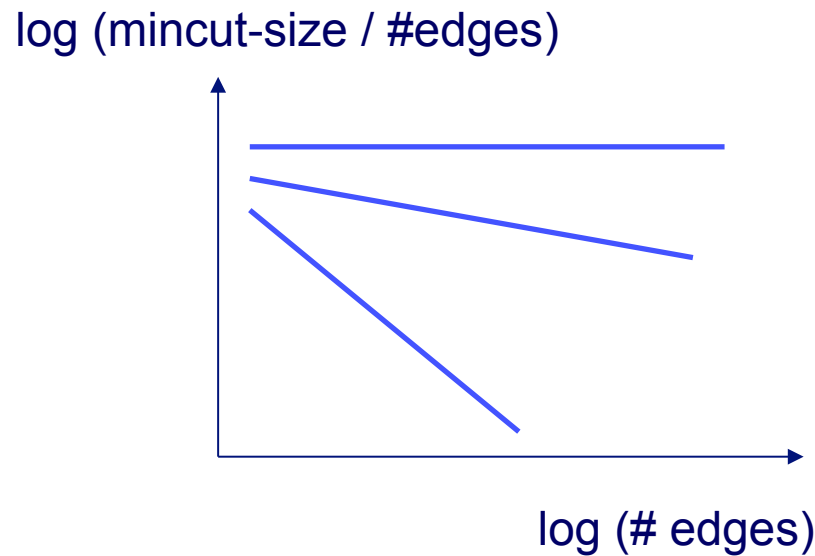
log (mincut-size / #edges)



For a random graph, the slope is 0

“Min-cut” plot

- What does it look like for a real-world graph?



?

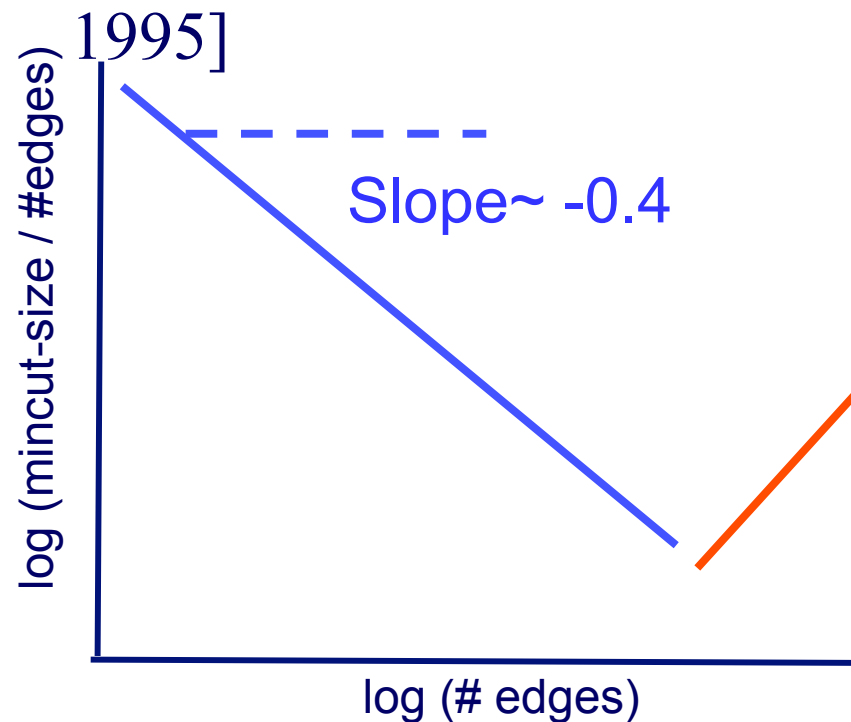
Experiments

- Datasets:
 - **Google Web Graph**: 916,428 nodes and 5,105,039 edges
 - **Lucent Router Graph**: Undirected graph of network routers from www.isi.edu/scan/mercator/maps.html; 112,969 nodes and 181,639 edges
 - **User → Website Clickstream Graph**: 222,704 nodes and 952,580 edges

NetMine: New Mining Tools for Large Graphs, by D. Chakrabarti, Y. Zhan, D. Blandford, C. Faloutsos and G. Blelloch, in the SDM 2004 Workshop on Link Analysis, Counter-terrorism and Privacy

Experiments

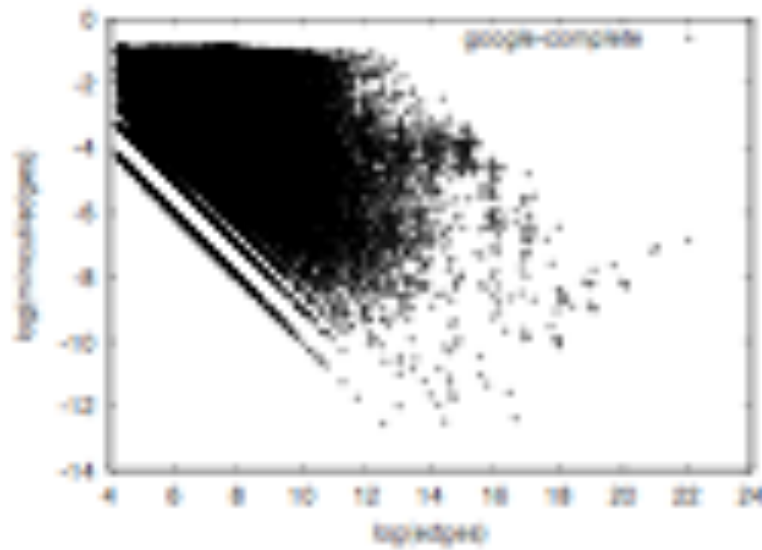
- Used the METIS algorithm [Karypis, Kumar, 1995]



- Google Web graph
- Values along the y-axis are averaged
- We observe a “lip” for large edges
- Slope of -0.4, corresponds to a 2.5-dimensional grid!

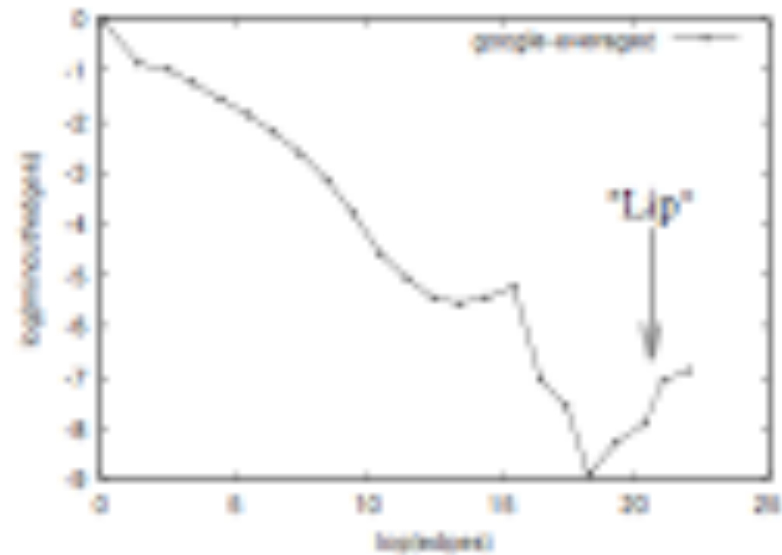
Google graph

$\log(\text{mincut-size} / \#\text{edges})$



Log(#edges)

All min-cuts



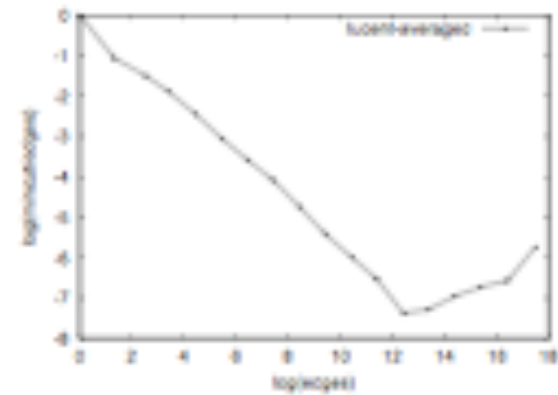
Log(#edges)

averaged

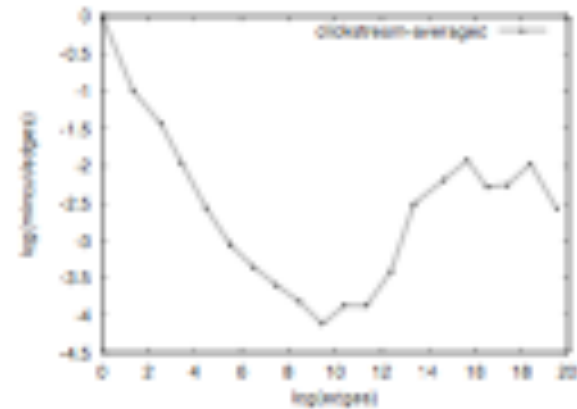
Experiments

- Same results for other graphs too...

Lucent Router graph



Clickstream graph



Conclusions – Practitioner’s guide

- Hard clustering – k pieces **METIS**
- Hard co-clustering – (k, l) pieces **Co-clustering**
- Hard clustering – optimal # pieces **Cross-associations**
- Observations **‘jellyfish’:
Maybe, there are
no good cuts**

Outline

- Task 4: time-evolving graphs – tensors
- Task 5: community detection
- ➔ • Task 6: virus propagation
- Task 7: scalability, parallelism and hadoop
- Conclusions

Detailed outline

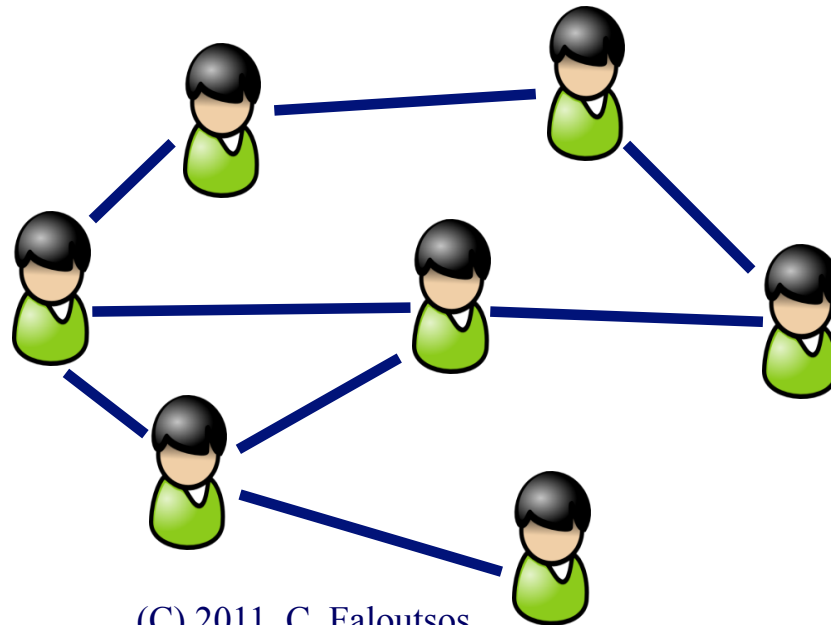
- Problem definition
- Analysis
- Experiments

Immunization and epidemic thresholds

- Q1: which nodes to immunize?
- Q2: will a virus vanish, or will it create an epidemic?

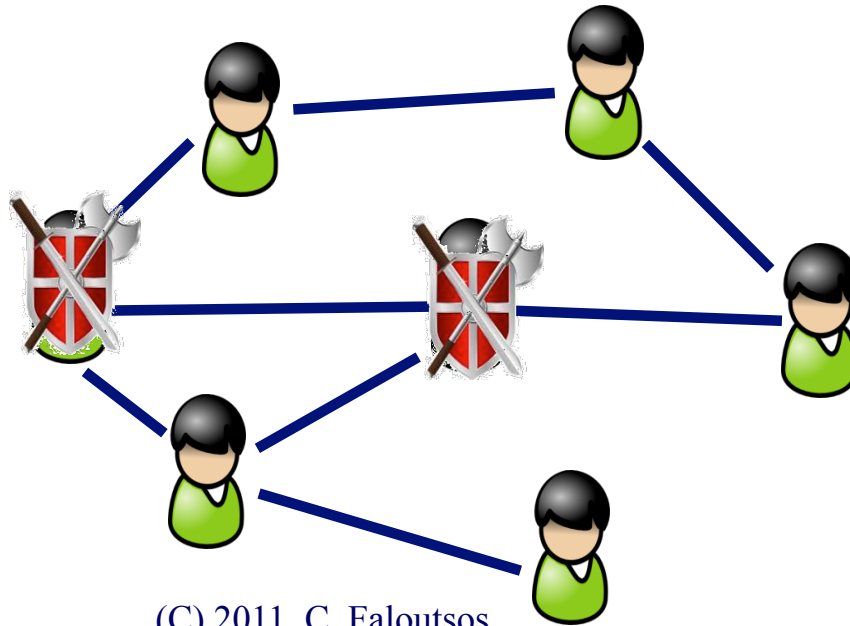
Q1: Immunization:

- Given
 - a network,
 - k vaccines, and
 - the virus details
- Which nodes to immunize?



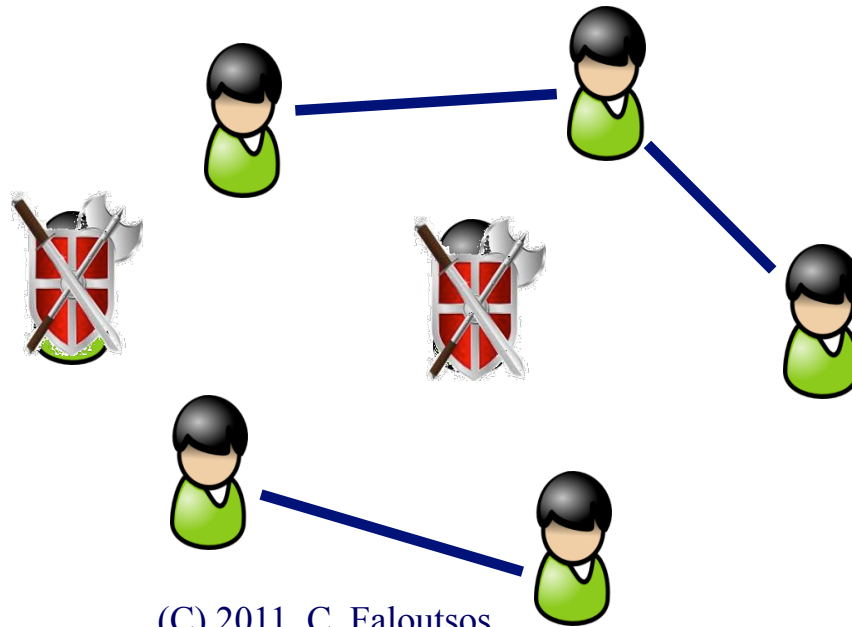
Q1: Immunization:

- Given
 - a network,
 - k vaccines, and
 - the virus details
- Which nodes to immunize?



Q1: Immunization:

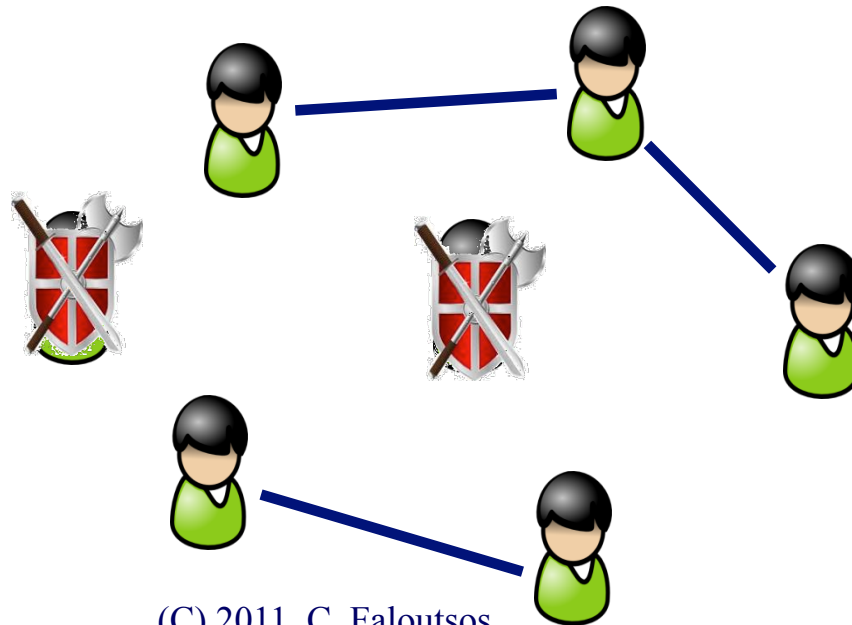
- Given
 - a network,
 - k vaccines, and
 - the virus details
- Which nodes to immunize?



Q1: Immunization:

- Given
 - a network,
 - k vaccines, and
 - the virus details
- Which nodes to immunize?

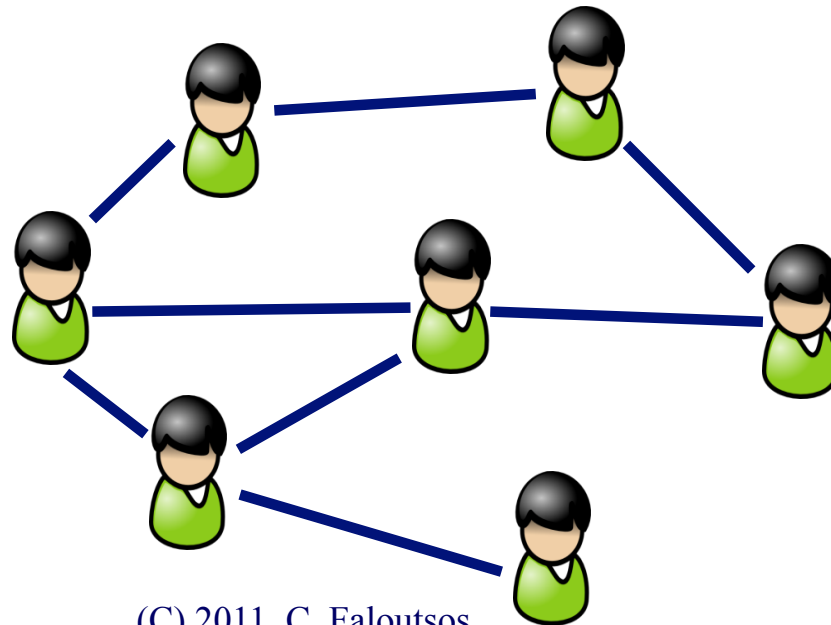
A: immunize the ones that maximally raise the `epidemic threshold' [Tong+, ICDM'10]



Q2: will a virus take over?

- Flu-like virus (no immunity, ‘SIS’)
- Mumps (life-time immunity, ‘SIR’)
- Pertussis (finite-length immunity, ‘SIRS’)

β : attack prob
 δ : heal prob



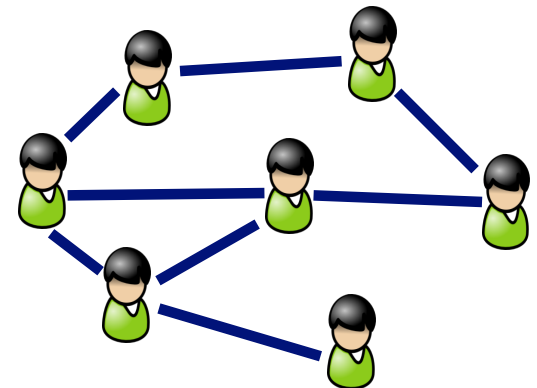
Q2: will a virus take over?

- Flu-like virus (no immunity, ‘SIS’)
- Mumps (life-time immunity, ‘SIR’)
- Pertussis (finite-length immunity, ‘SIRS’)

β : attack prob

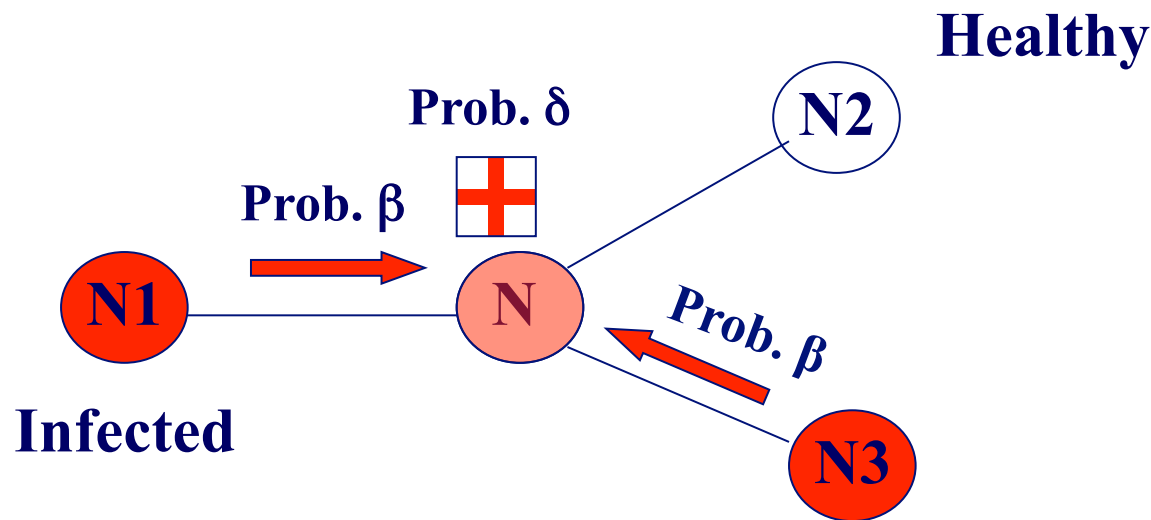
δ : heal prob

A: depends on connectivity
(avg degree? Max degree?
variance? Something else?)



The model: SIS

- ‘Flu’ like: Susceptible-Infected-Susceptible
- Virus ‘strength’ $s = \beta / \delta$



Epidemic threshold τ

of a graph: the value of τ , such that

if strength $s = \beta / \delta < \tau$

an epidemic can not happen

Thus,

- given a graph
- compute its epidemic threshold

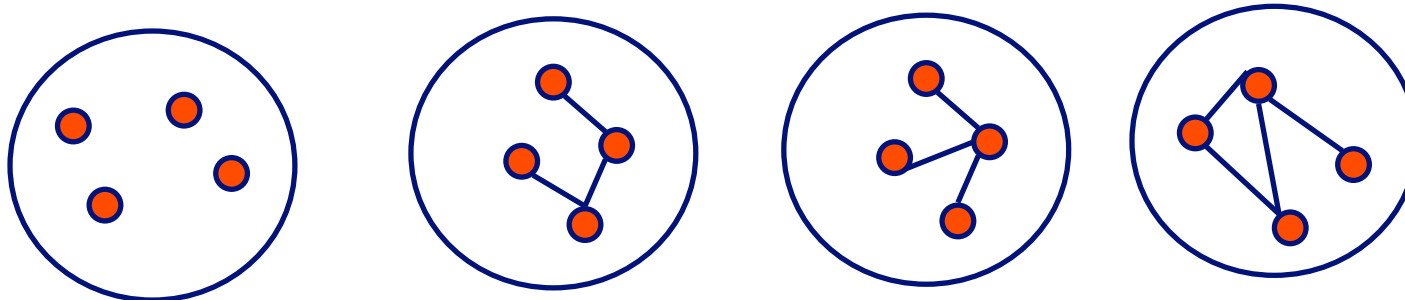
Detailed outline

- Problem definition
- ➔ • Analysis
- Experiments

Epidemic threshold τ

What should τ depend on?

- avg. degree? and/or highest degree?
- and/or variance of degree?
- and/or third moment of degree?
- and/or diameter?



Epidemic threshold

- [Theorem] We have no epidemic, if

$$\beta/\delta < \tau = 1/\lambda_{1,A}$$

Epidemic threshold

- [Theorem] We have no epidemic, if

recovery prob. $\beta/\delta < \tau = 1/\lambda_{1,A}$ epidemic threshold

attack prob. β/δ

largest eigenvalue of adj. matrix A

Proof: [Wang+03] (proof: for **SIS=flu only**)

Beginning of proof

Healthy @ t+1:

- (healthy or healed)
- and not attacked @ t

Let: $p(i, t) = \text{Prob node } i \text{ is sick @ } t+1$

$$1 - p(i, t+1) = (1 - p(i, t) + p(i, t) * \delta) * \prod_j (1 - \beta a_{ji} * p(j, t))$$

Below threshold, if the above *non-linear dynamical system* above is 'stable' (eigenvalue of Hessian < 1)

Epidemic threshold for various networks

Formula includes older results as special cases:

- Homogeneous networks [Kephart+White]
 - $\lambda_{I,A} = \langle k \rangle$; $\tau = 1/\langle k \rangle$ ($\langle k \rangle$: avg degree)
- Star networks (d = degree of center)
 - $\lambda_{I,A} = \text{sqrt}(d)$; $\tau = 1/\text{sqrt}(d)$
- Infinite power-law networks
 - $\lambda_{I,A} = \infty$; $\tau = 0$; [Barabasi]

Epidemic threshold

- [Theorem 2] Below the epidemic threshold, the epidemic dies out exponentially

Recent generalization

- [Prakash+, arxiv '10]: similar threshold, for almost **all** virus propagation models (VPM)
 - SIS -> flu
 - SIR -> mumps
 - SIRS -> whooping cough (temporary immunity)
 - SIIR (-> HIV)
 - ...

A2: will a virus take over?

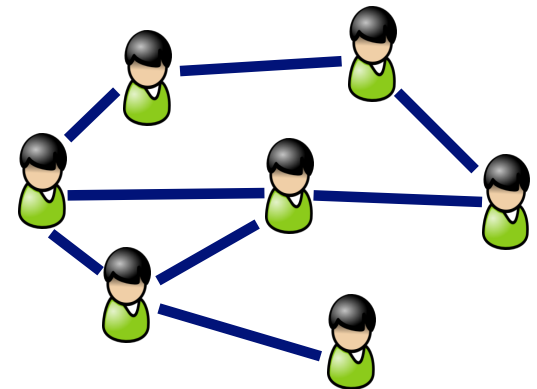
- For **all** typical virus propagation models (flu, mumps, pertussis, HIV, etc)
- The **only** connectivity measure that matters, is

$$1/\lambda_1$$


the first eigenvalue of the
adj. matrix

Proof for **all** VPM:

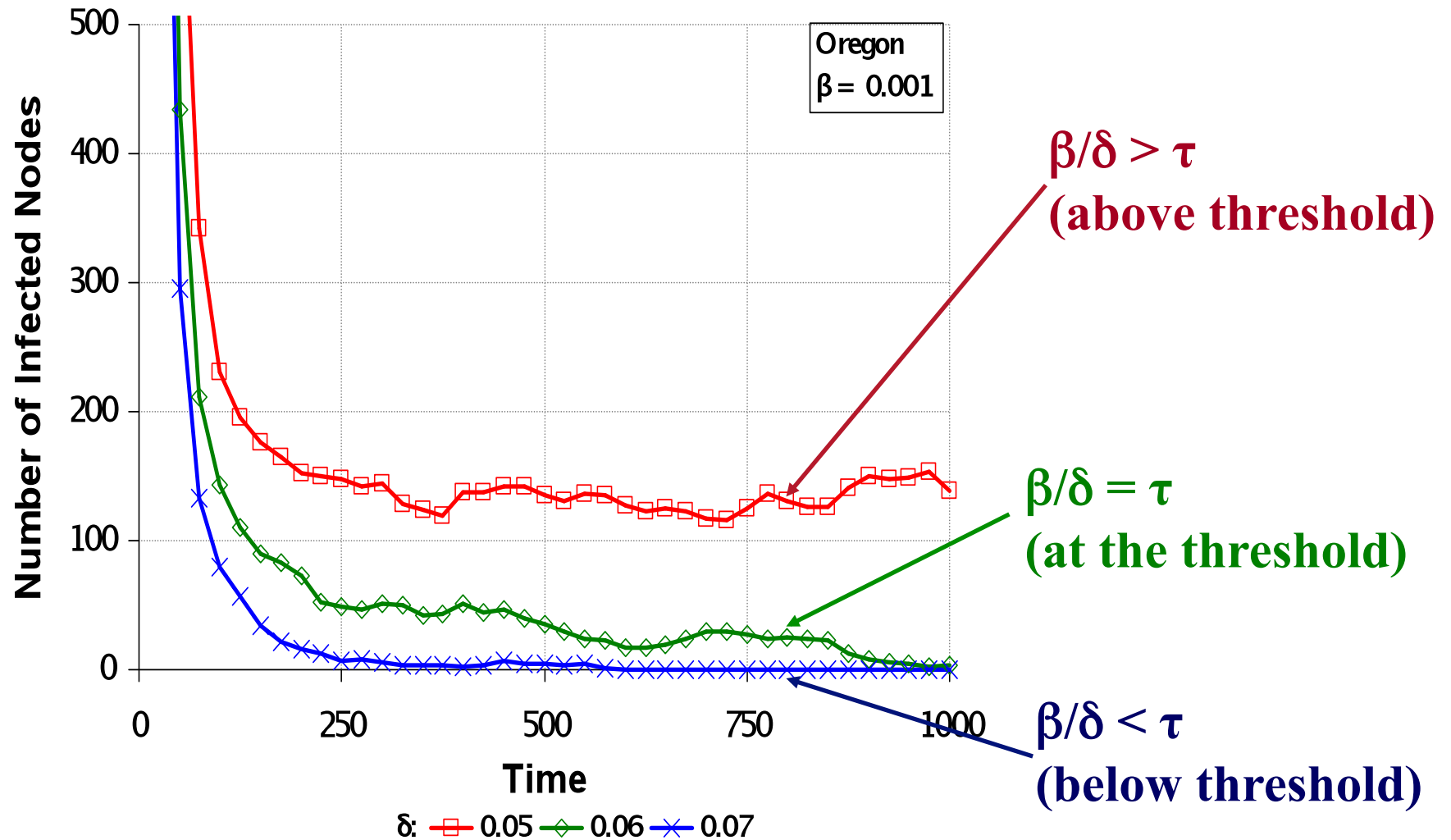
[Prakash+, '10, arxiv]



Detailed outline

- Epidemic threshold
 - Problem definition
 - Analysis
 -  – Experiments

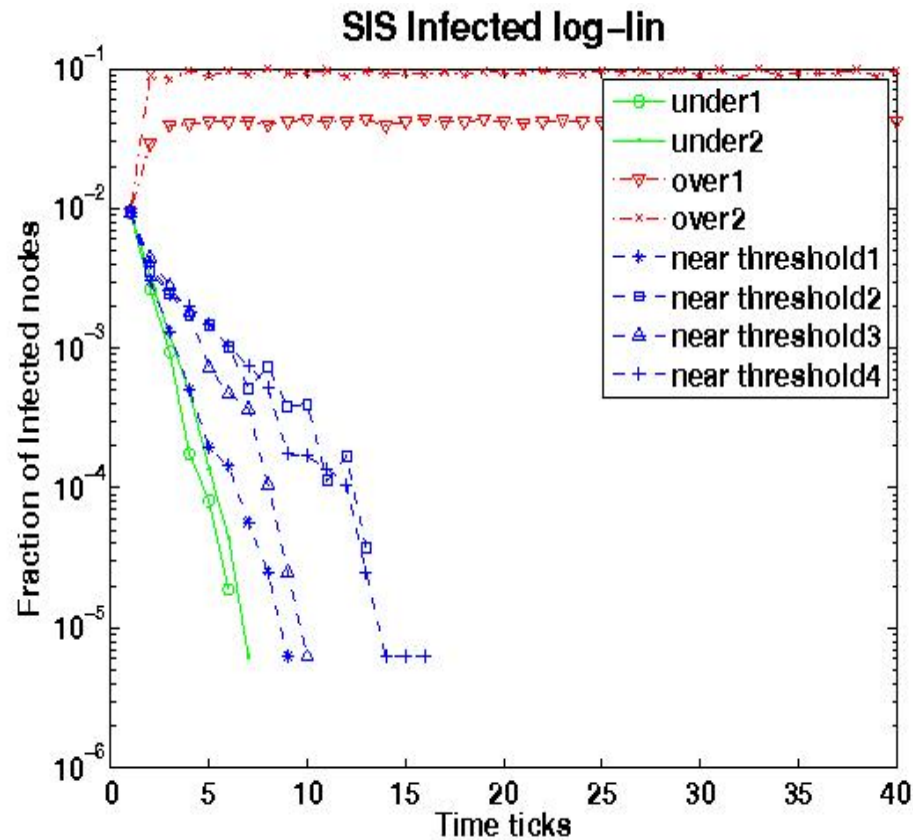
Experiments (Oregon)



SIS simulation - # infected nodes vs time

Log - Lin

#inf.
(log scale)



— above

— at

— below

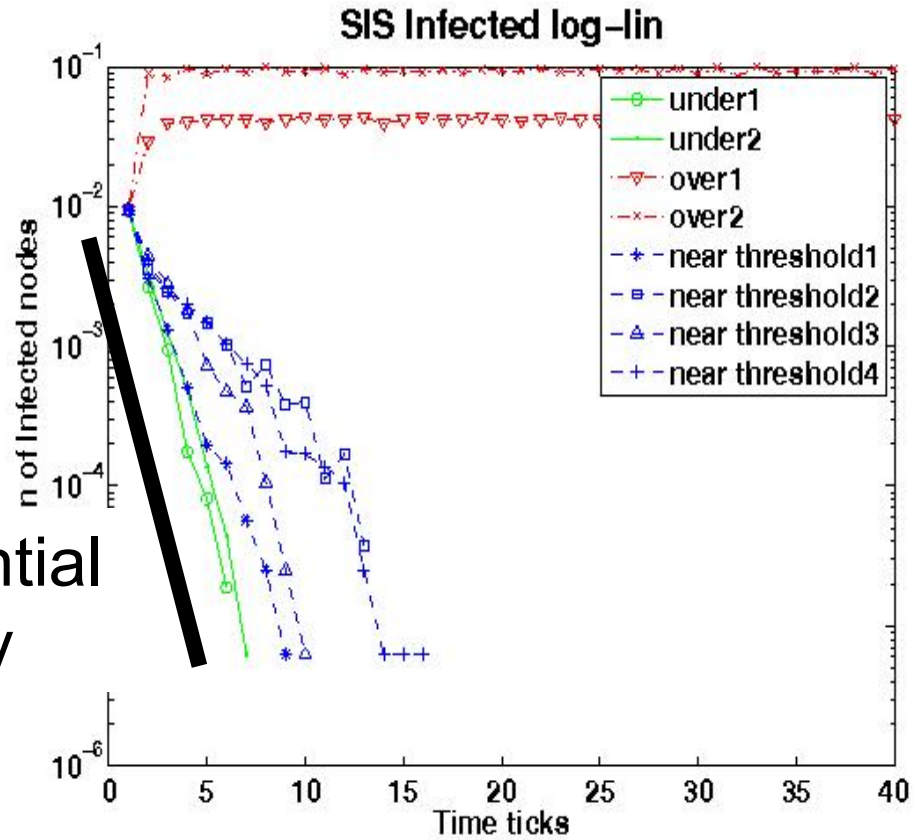
Time (linear scale)

SIS simulation - # infected nodes vs time

Log - Lin

#inf.
(log scale)

Exponential
decay



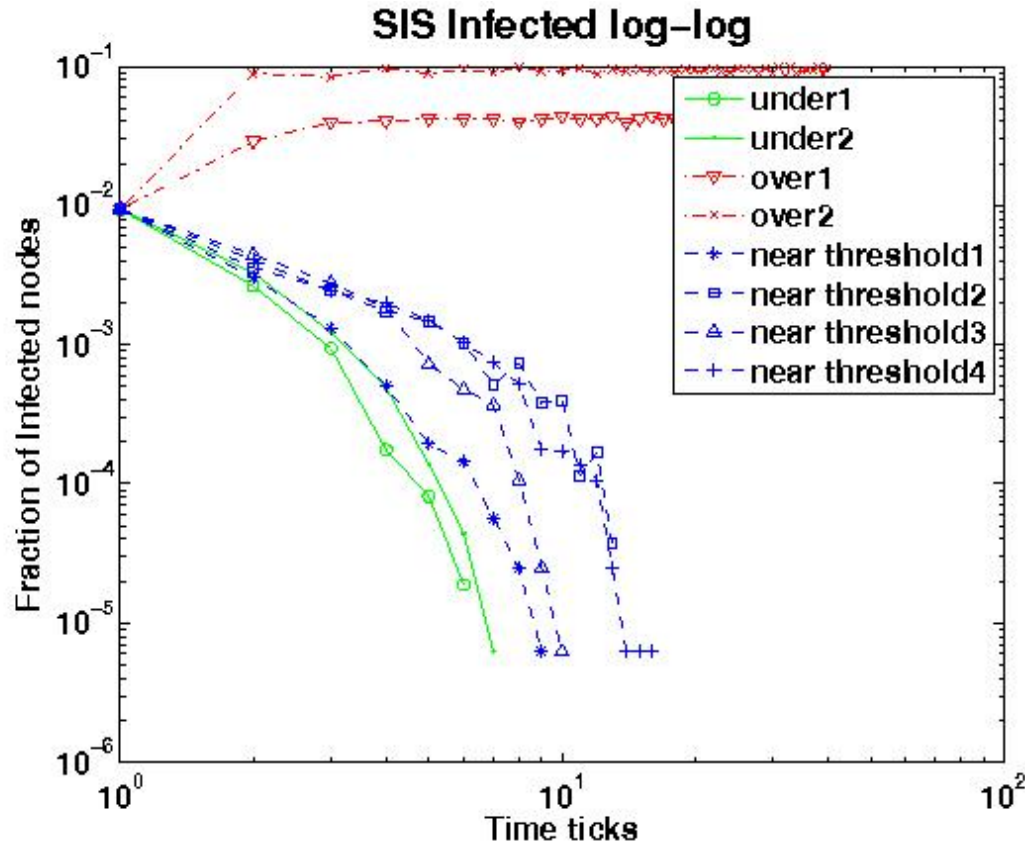
— above
— at
— below

Time (linear scale)

SIS simulation - # infected nodes vs time

Log - Log

#inf.
(log scale)



— above
— at
— below

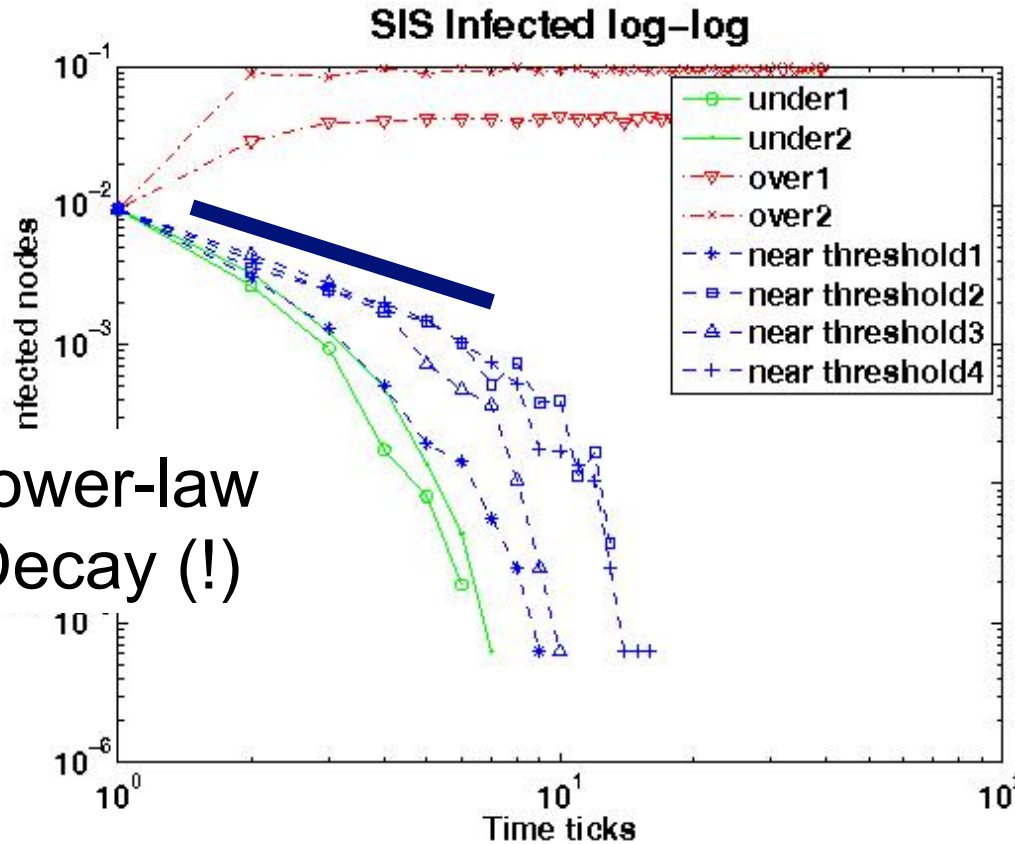
Time (log scale)

SIS simulation - # infected nodes vs time

Log - Log

#inf.
(log scale)

Power-law
Decay (!)

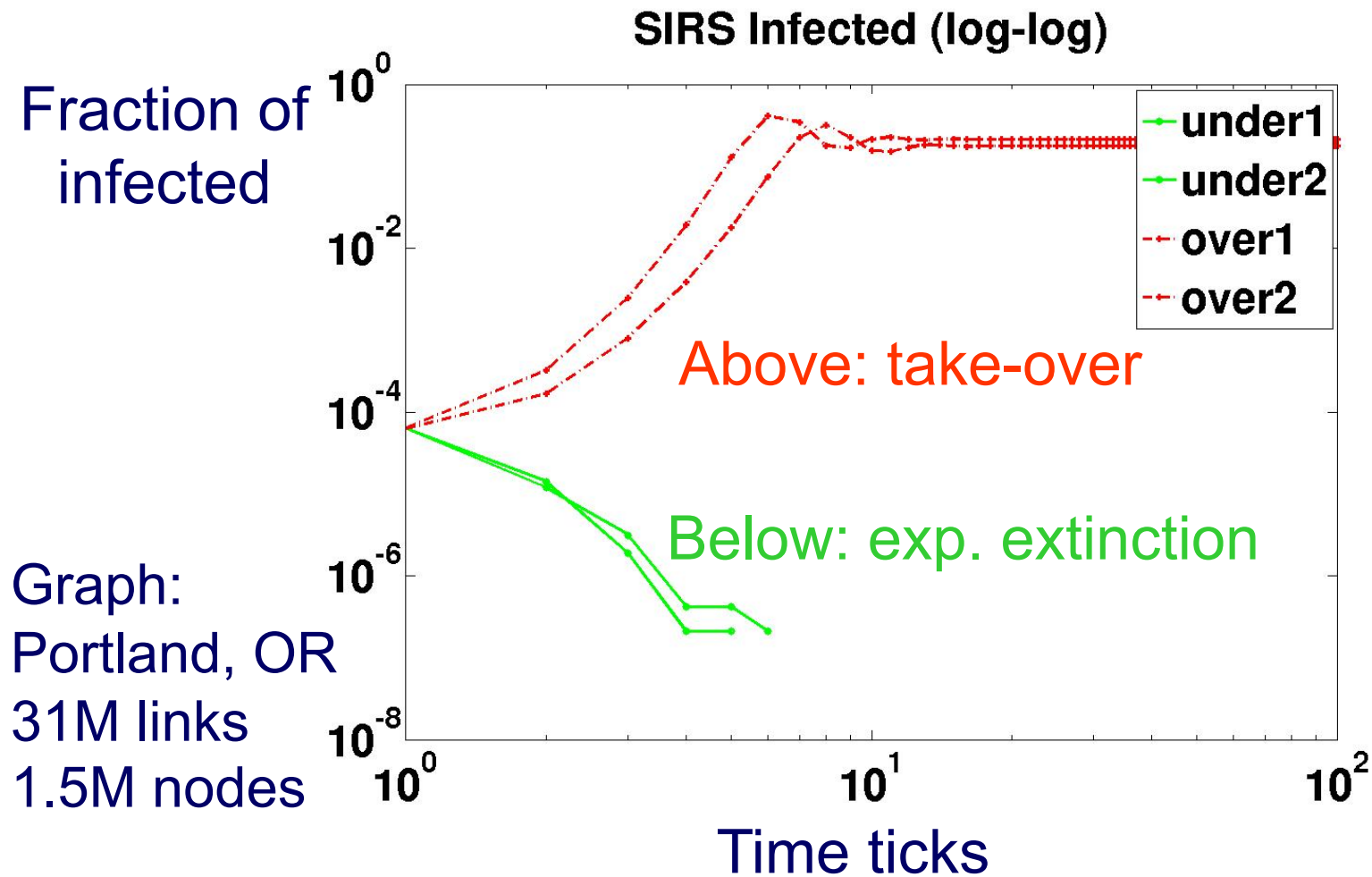


- above
- at
- below

Time (log scale)

How about other VPMs?

A2: will a virus take over? (SIRS case)



Conclusions

- $\lambda_{1,A}$: Eigenvalue of adjacency matrix
determines the survival of (almost) **any**
virus
- measure of connectivity (\sim # paths)
 - Can answer ‘what-if’ scenarios
 - May guide immunization policies
 - Can help us avoid expensive simulations

References

- D. Chakrabarti, Y. Wang, C. Wang, J. Leskovec, and C. Faloutsos, *Epidemic Thresholds in Real Networks*, in ACM TISSEC, 10(4), 2008
- Ganesh, A., Massoulié, L., and Towsley, D., 2005. The effect of network topology on the spread of epidemics. In *INFOCOM*.

References (cont'd)

- Hethcote, H. W. 2000. The mathematics of infectious diseases. *SIAM Review* 42, 599–653.
- Hethcote, H. W. AND Yorke, J. A. 1984. *Gonorrhea Transmission Dynamics and Control*. Vol. 56. Springer. Lecture Notes in Biomathematics.

References (cont'd)

- Y. Wang, D. Chakrabarti, C. Wang and C. Faloutsos, *Epidemic Spreading in Real Networks: An Eigenvalue Viewpoint*, in SRDS 2003 (pages 25-34), Florence, Italy

Outline

- Task 4: time-evolving graphs – tensors
- Task 5: community detection
- Task 6: virus propagation
- ➔ • Task 7: scalability, parallelism and hadoop
- Conclusions

Scalability

- How about if graph/tensor does not fit in core?
- How about handling huge graphs?

Scalability

- How about if graph/tensor does not fit in core?
- [‘MET’: Kolda, Sun, ICMD’08, best paper award]
- How about handling huge graphs?

Scalability

- Google: > 450,000 processors in clusters of ~2000 processors each [Barroso, Dean, Hölzle, “*Web Search for a Planet: The Google Cluster Architecture*” IEEE Micro 2003]
- Yahoo: 5Pb of data [Fayyad, KDD’07]
- Problem: machine failures, on a daily basis
- How to parallelize data mining tasks, then?

Scalability

- Google: > 450,000 processors in clusters of ~2000 processors each [Barroso, Dean, Hölzle, “*Web Search for a Planet: The Google Cluster Architecture*” IEEE Micro 2003]
- Yahoo: 5Pb of data [Fayyad, KDD’07]
- Problem: machine failures, on a daily basis
- How to parallelize data mining tasks, then?
- A: map/reduce – hadoop (open-source clone)
<http://hadoop.apache.org/>



2' intro to hadoop

- master-slave architecture; n-way replication (default n=3)
- ‘group by’ of SQL (in parallel, fault-tolerant way)
- e.g, find histogram of word frequency
 - compute local histograms
 - then merge into global histogram

```
select course-id, count(*)  
from ENROLLMENT  
group by course-id
```

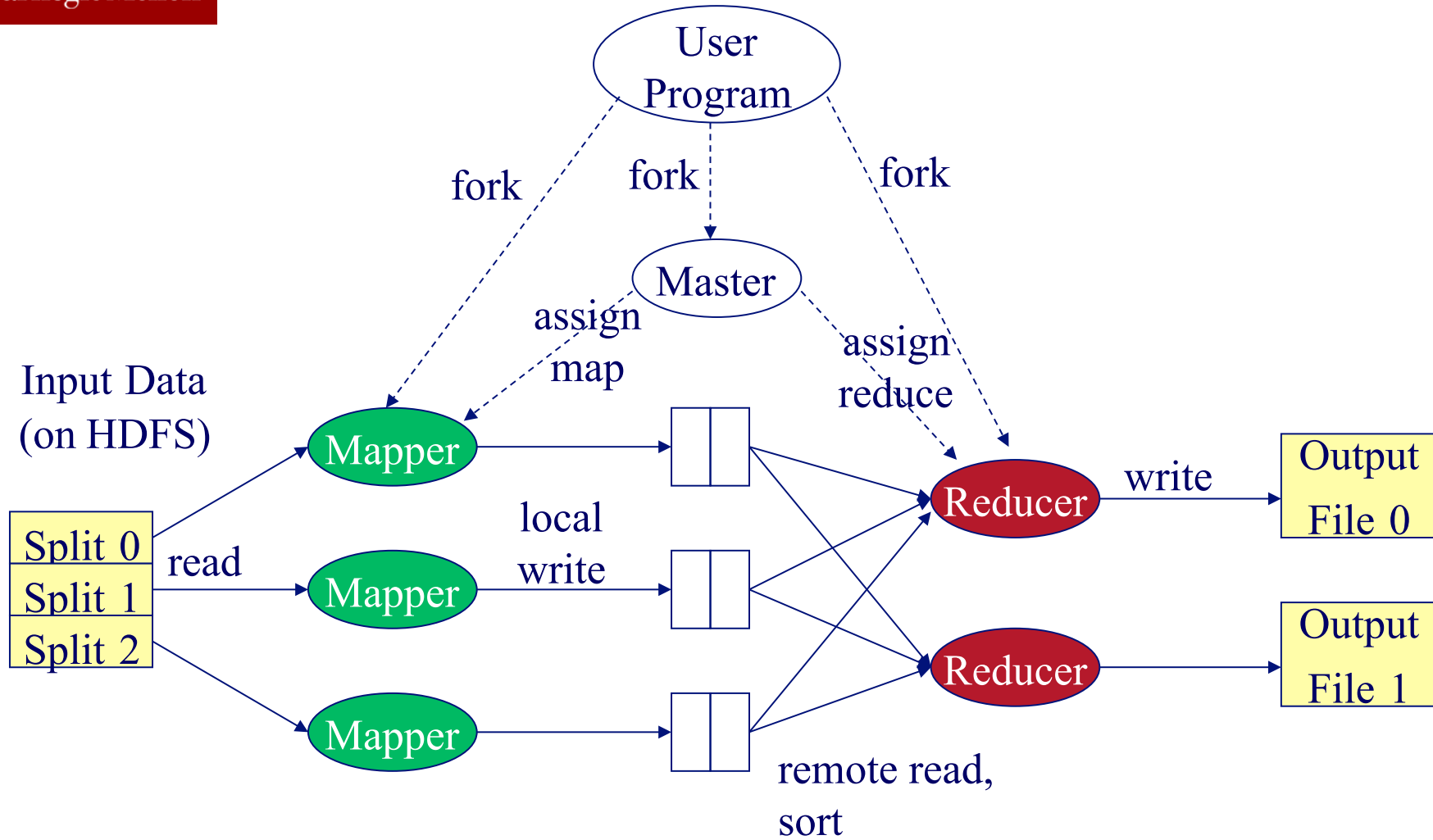
2' intro to hadoop

- master-slave architecture; n-way replication (default n=3)
- ‘group by’ of SQL (in parallel, fault-tolerant way)
- e.g, find histogram of word frequency
 - compute local histograms
 - then merge into global histogram

```
select course-id, count(*)  
from ENROLLMENT  
group by course-id
```

reduce

map



By default: 3-way replication;
 Late/dead machines: ignored, **transparently** (!)

D.I.S.C.

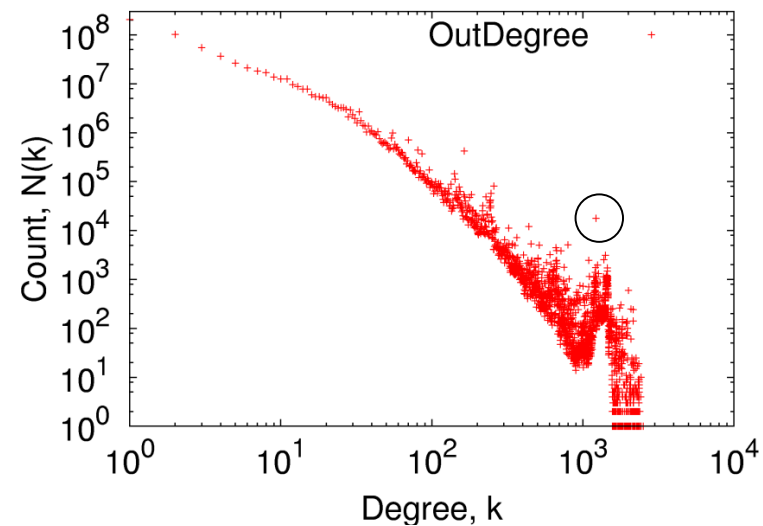
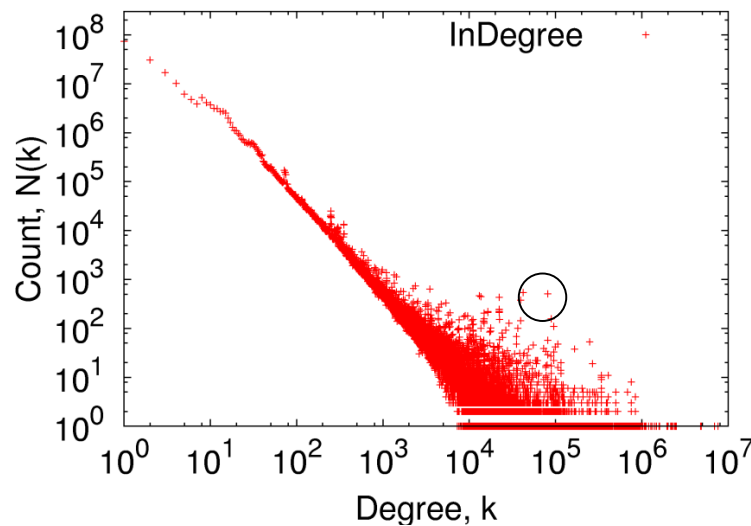


- ‘Data Intensive Scientific Computing’ [R. Bryant, CMU]
 - ‘big data’
 - www.cs.cmu.edu/~bryant/pubdir/cmu-cs-07-128.pdf

Analysis of a large graph

~200Gb (Yahoo crawl) - Degree Distribution:

- in 12 minutes with 50 machines
- Many (link spams ?) at out-degree 1200



Outline – Algorithms & results

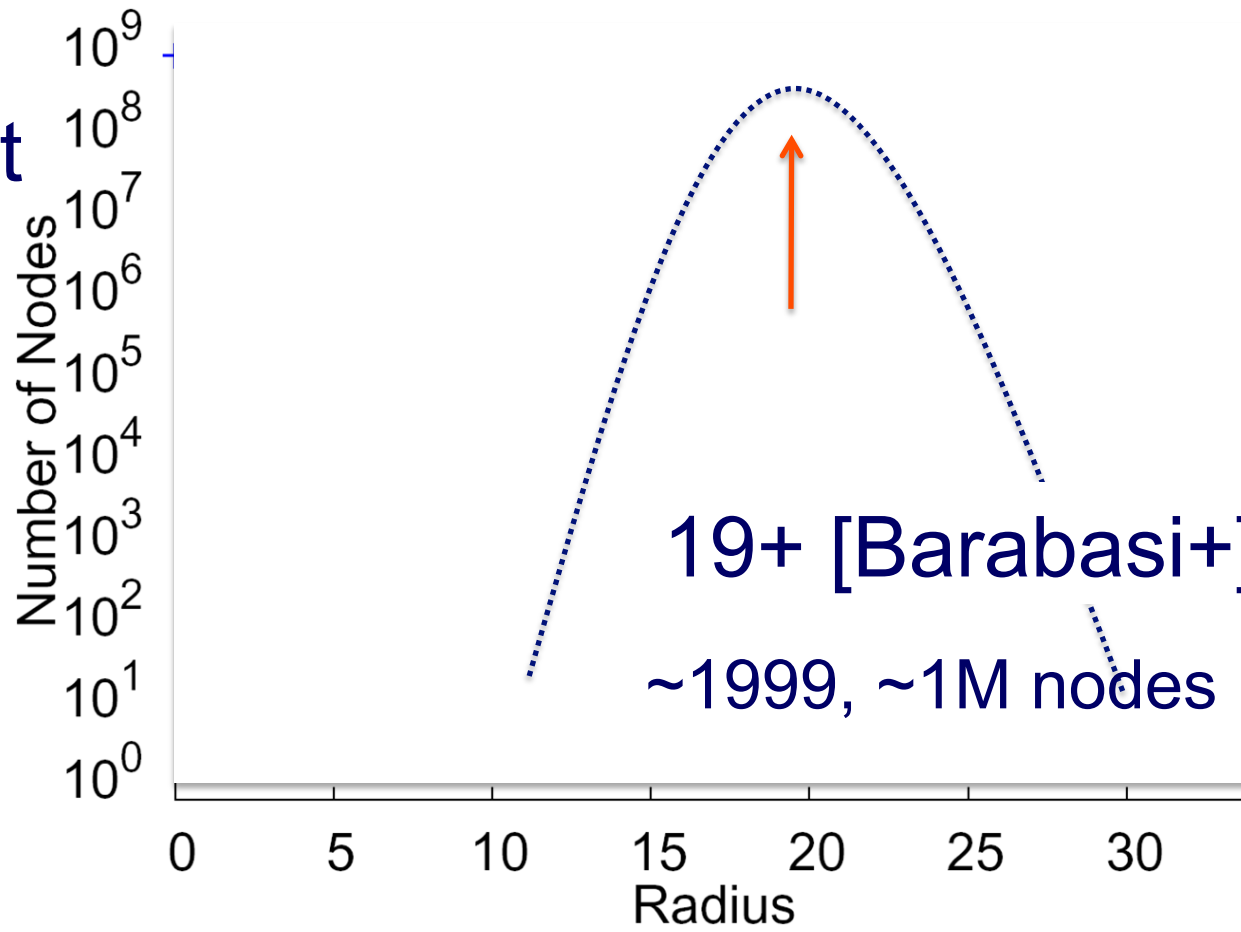
	Centralized	Hadoop/ PEGASUS
Degree Distr.	old	old
Pagerank	old	old
→ Diameter/ANF	old	DONE
Conn. Comp	old	DONE
Triangles	DONE	
Visualization	STARTED	



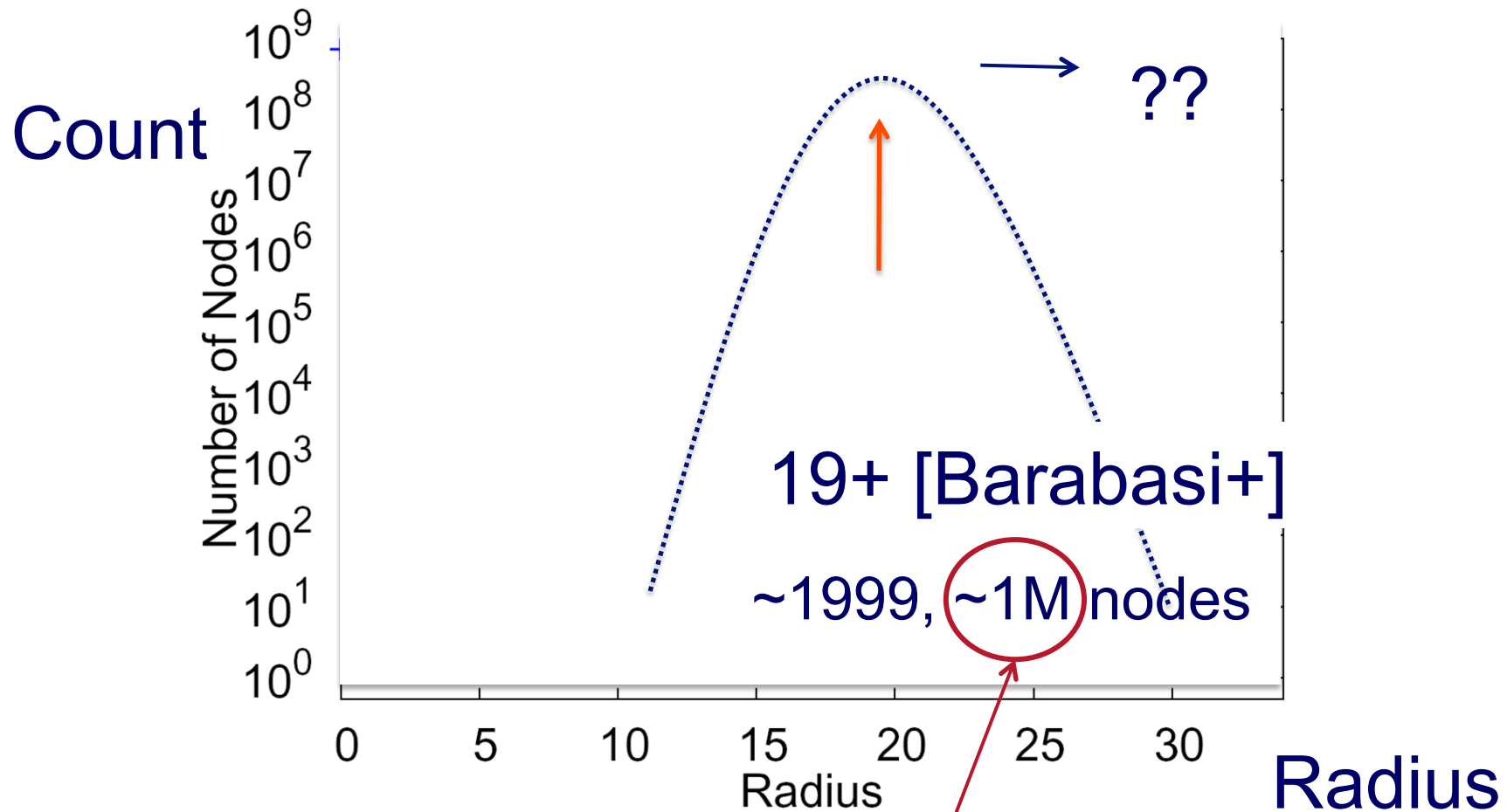
HADI for diameter estimation

- *Radius Plots for Mining Tera-byte Scale Graphs* **U Kang**, Charalampos Tsourakakis, Ana Paula Appel, Christos Faloutsos, Jure Leskovec, SDM'10
- Naively: diameter needs $O(N^2)$ space and up to $O(N^3)$ time – **prohibitive** ($N \sim 1B$)
- Our HADI: linear on E ($\sim 10B$)
 - Near-linear scalability wrt # machines
 - Several optimizations \rightarrow 5x faster

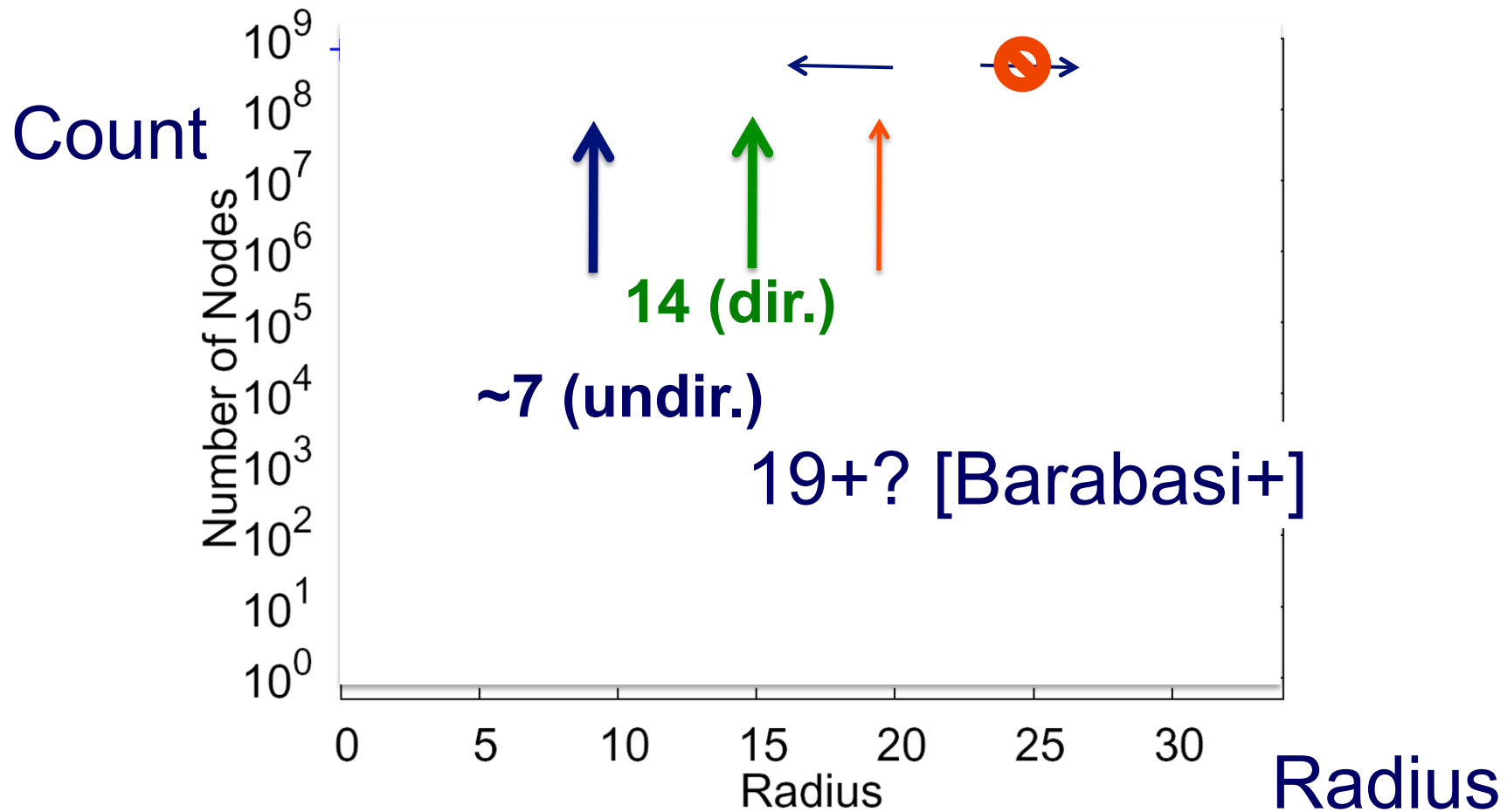
Count



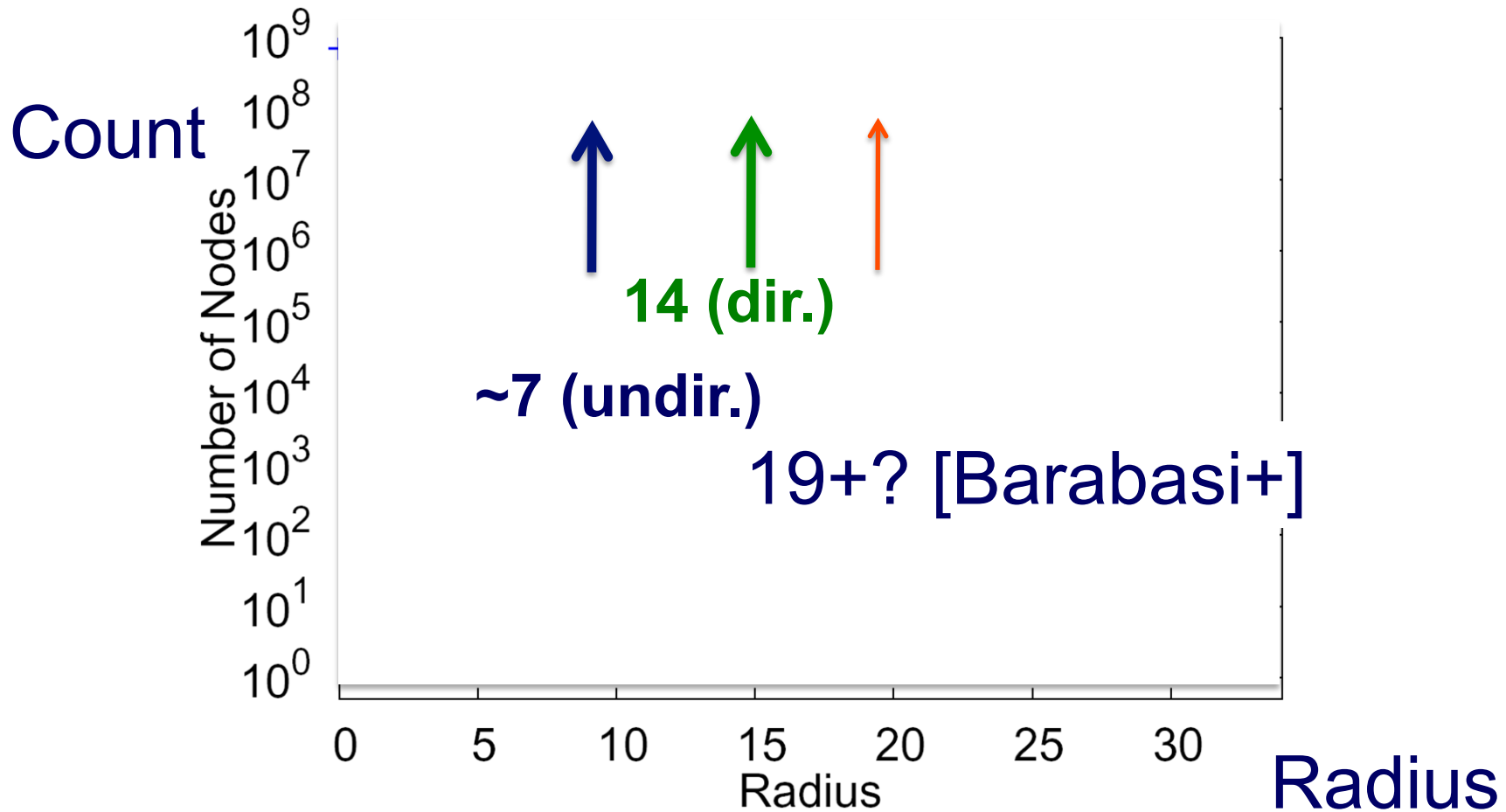
Radius



- YahooWeb graph (120Gb, 1.4B nodes, 6.6 B edges)
- Largest publicly available graph ever studied.

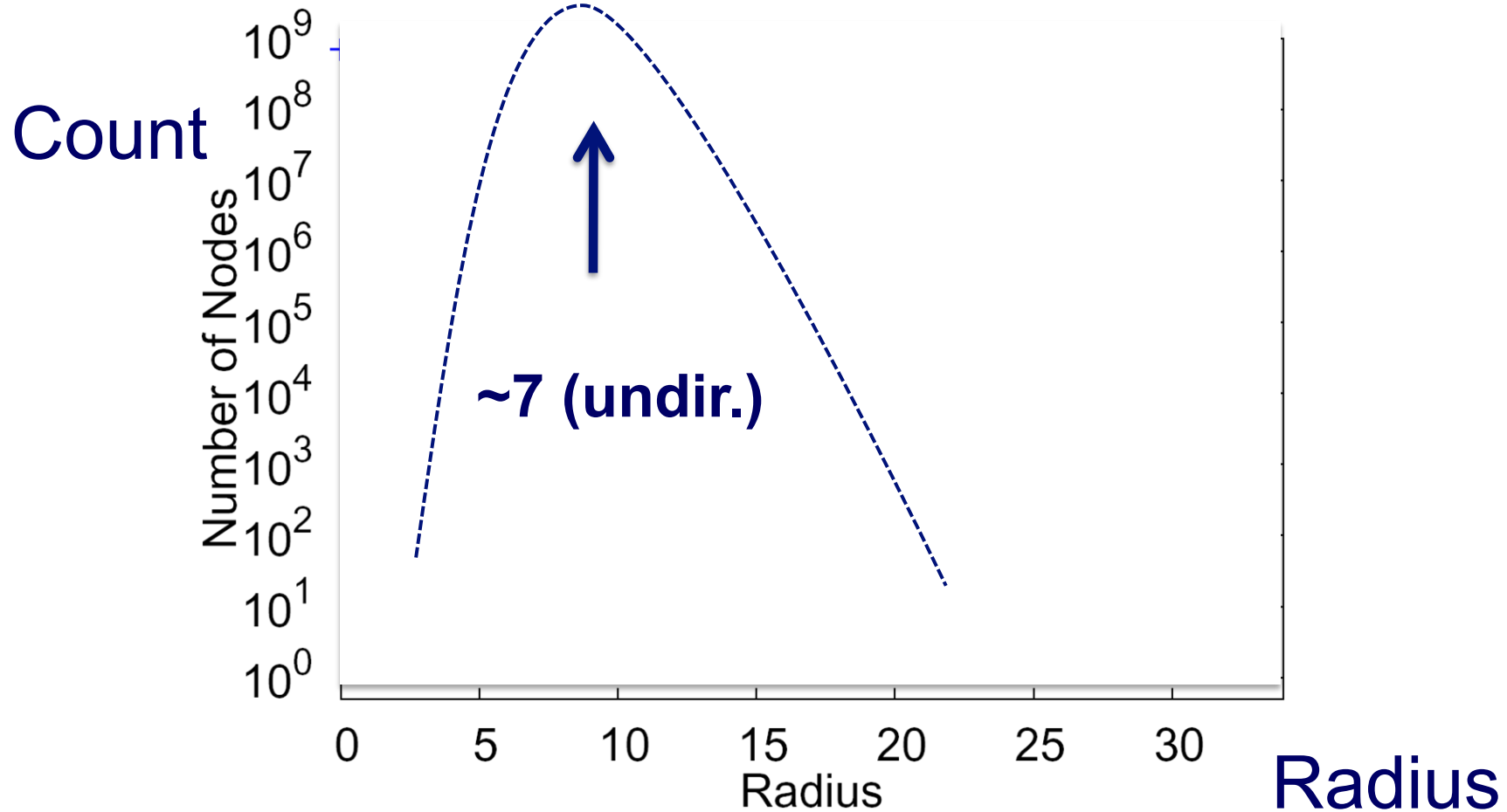


- YahooWeb graph (120Gb, 1.4B nodes, 6.6 B edges)
- Largest publicly available graph ever studied.

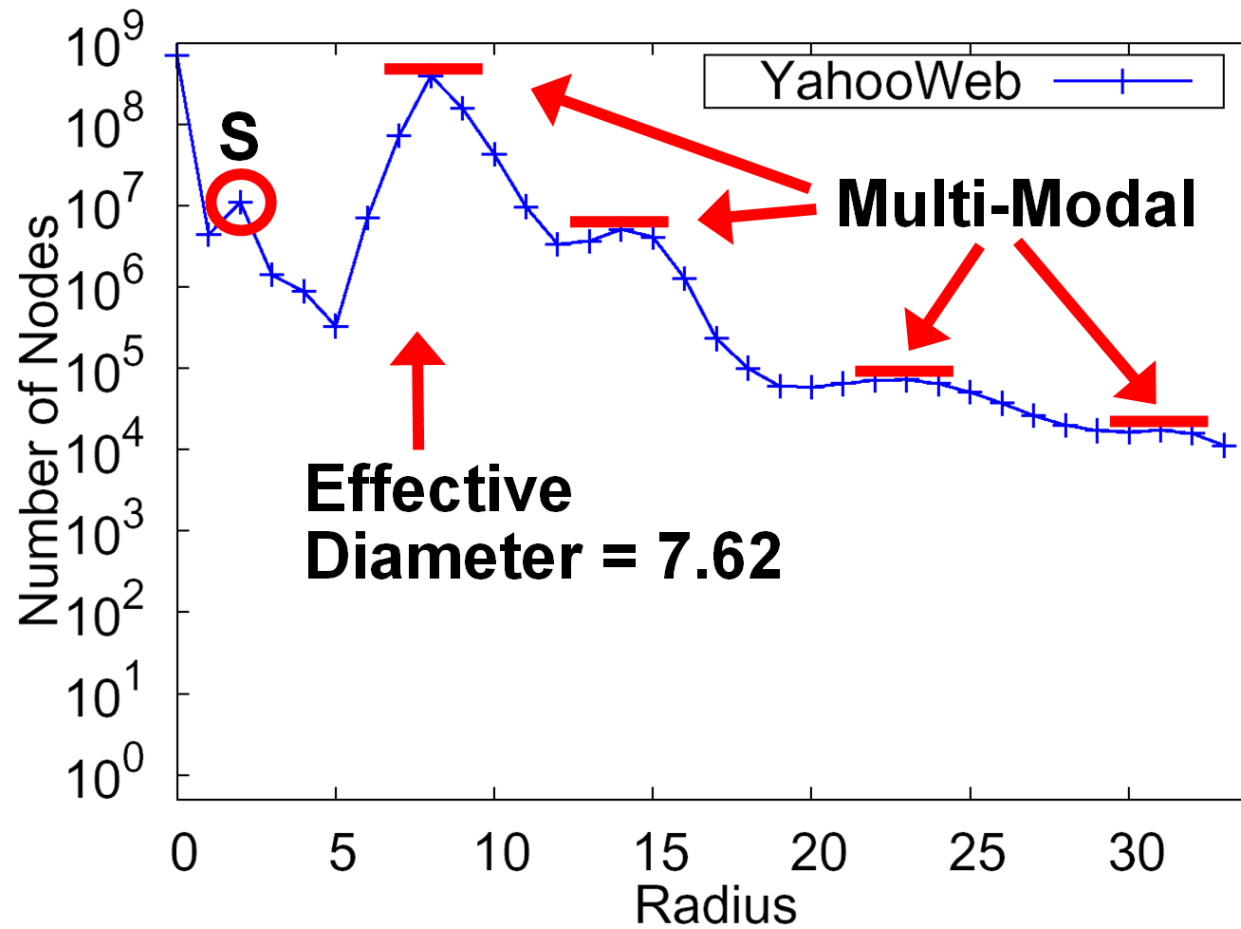


YahooWeb graph (120Gb, 1.4B nodes, 6.6 B edges)

- 7 degrees of separation (!)
- Diameter: shrunk

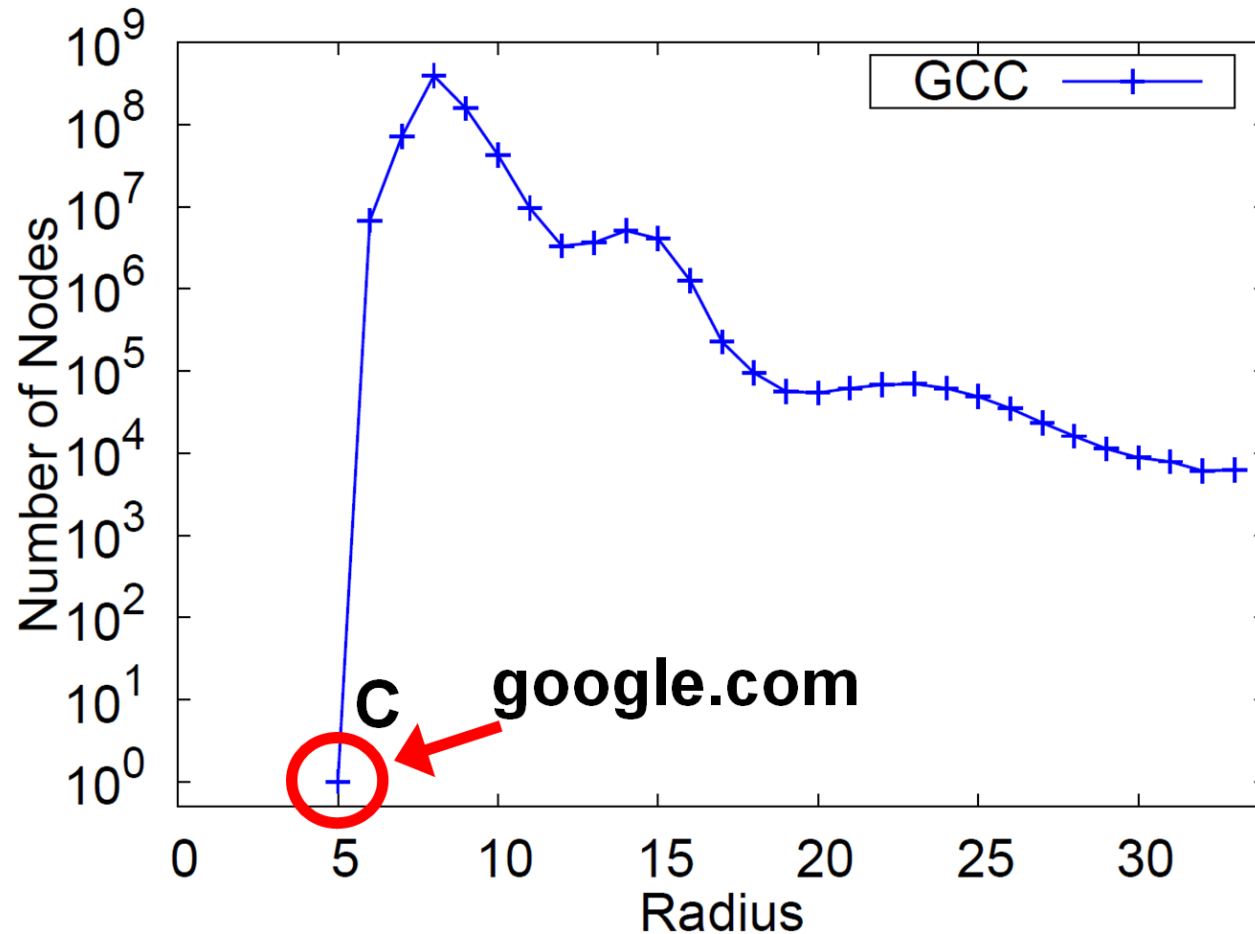


YahooWeb graph (120Gb, 1.4B nodes, 6.6 B edges)
Q: Shape?

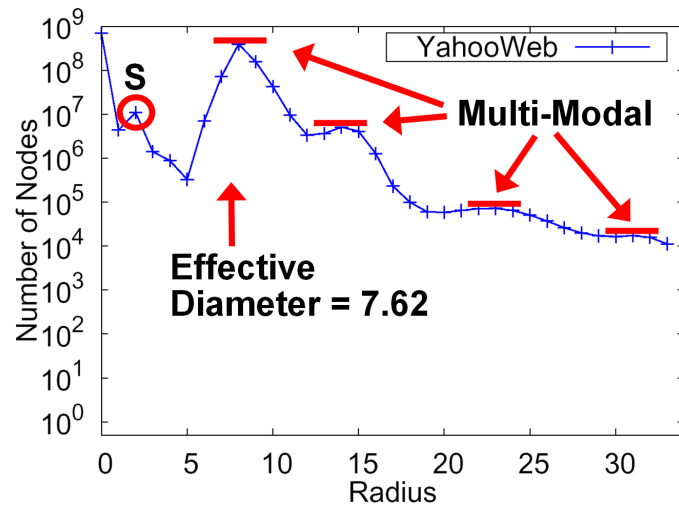


YahooWeb graph (120Gb, 1.4B nodes, 6.6 B edges)

- effective diameter: surprisingly small.
- Multi-modality (?!)

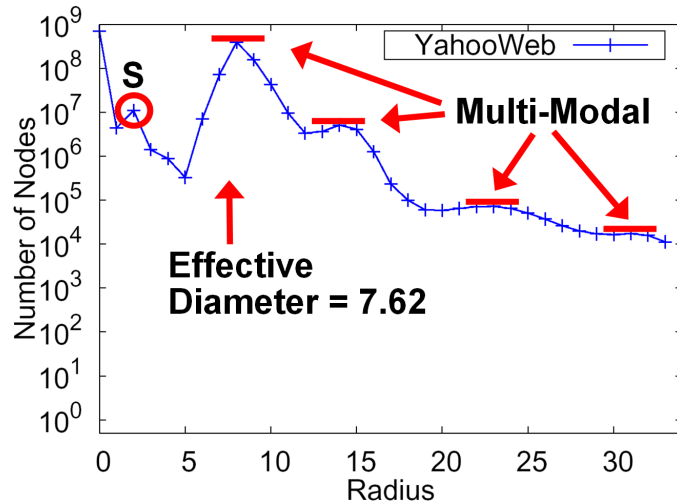


Radius Plot of **GCC** of YahooWeb.

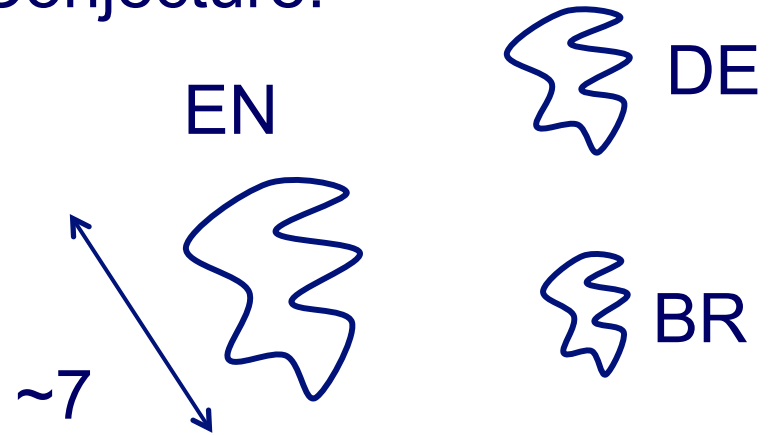


YahooWeb graph (120Gb, 1.4B nodes, 6.6 B edges)

- effective diameter: surprisingly small.
- Multi-modality: probably mixture of cores .

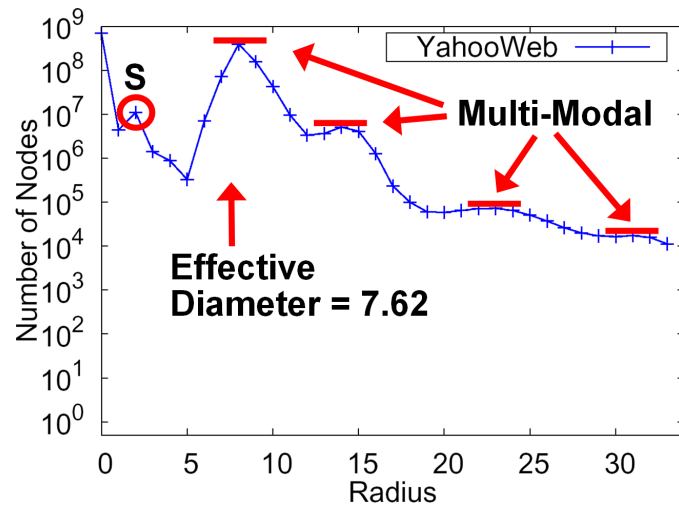


Conjecture:

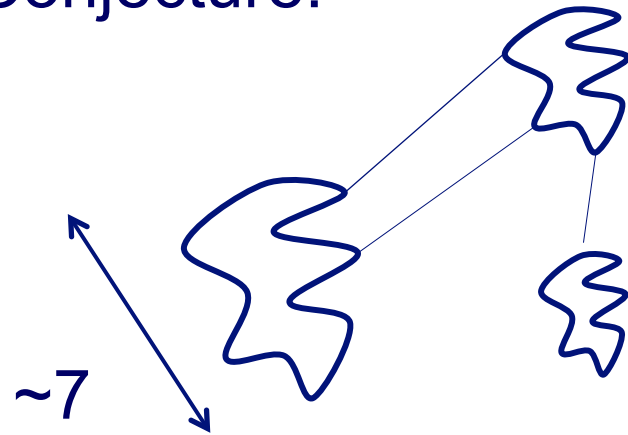


YahooWeb graph (120Gb, 1.4B nodes, 6.6 B edges)

- effective diameter: surprisingly small.
- Multi-modality: probably mixture of cores .

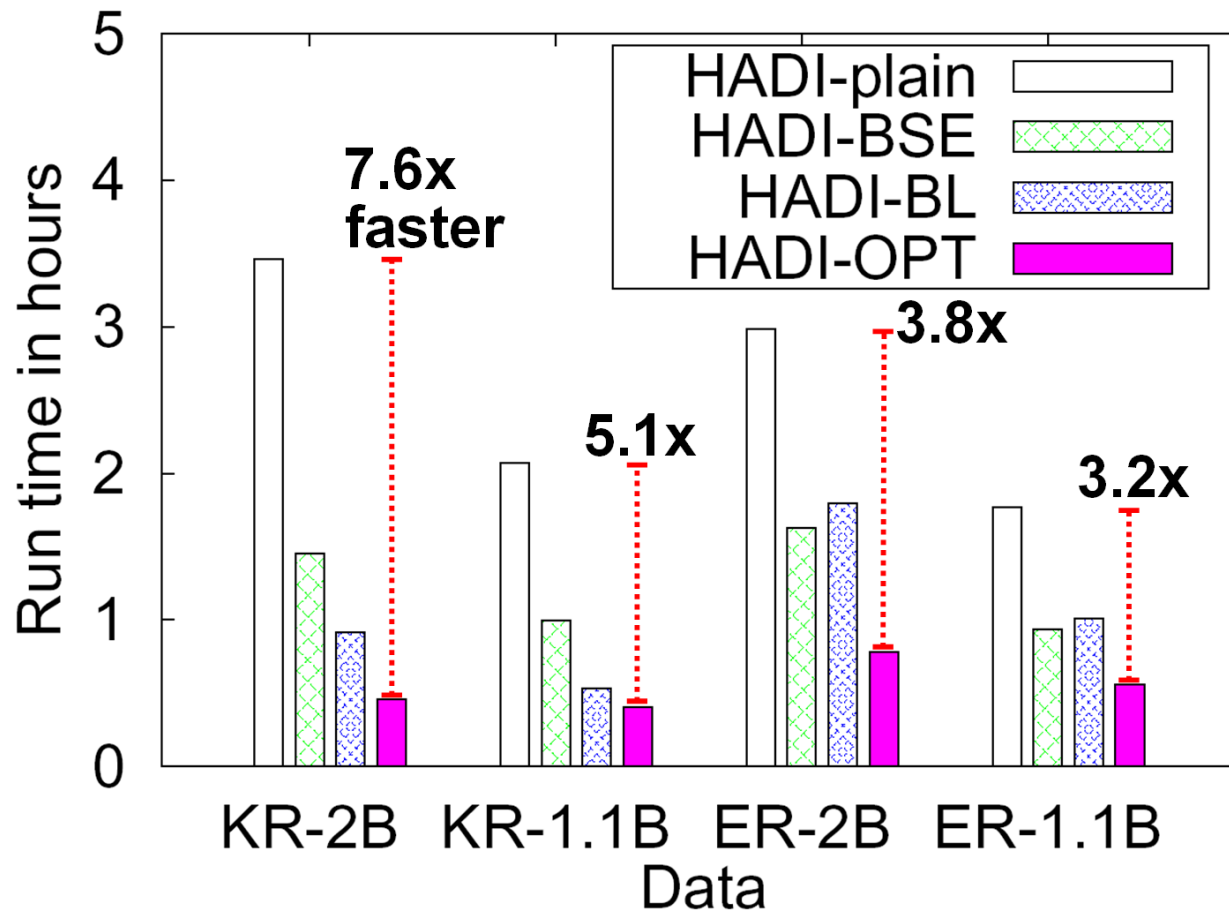


Conjecture:



YahooWeb graph (120Gb, 1.4B nodes, 6.6 B edges)

- effective diameter: surprisingly small.
- Multi-modality: probably mixture of cores .



Running time - Kronecker and Erdos-Renyi
Graphs with billions edges.

Outline – Algorithms & results

	Centralized	Hadoop/ PEGASUS
Degree Distr.	old	old
Pagerank	old	old
Diameter/ANF	old	DONE
→ Conn. Comp	old	DONE
Triangles	DONE	
Visualization	STARTED	

Generalized Iterated Matrix Vector Multiplication (GIMV)

*PEGASUS: A Peta-Scale Graph Mining
System - Implementation and Observations.*

U Kang, Charalampos E. Tsourakakis,
and Christos Faloutsos.

(ICDM) 2009, Miami, Florida, USA.

Best Application Paper (runner-up).

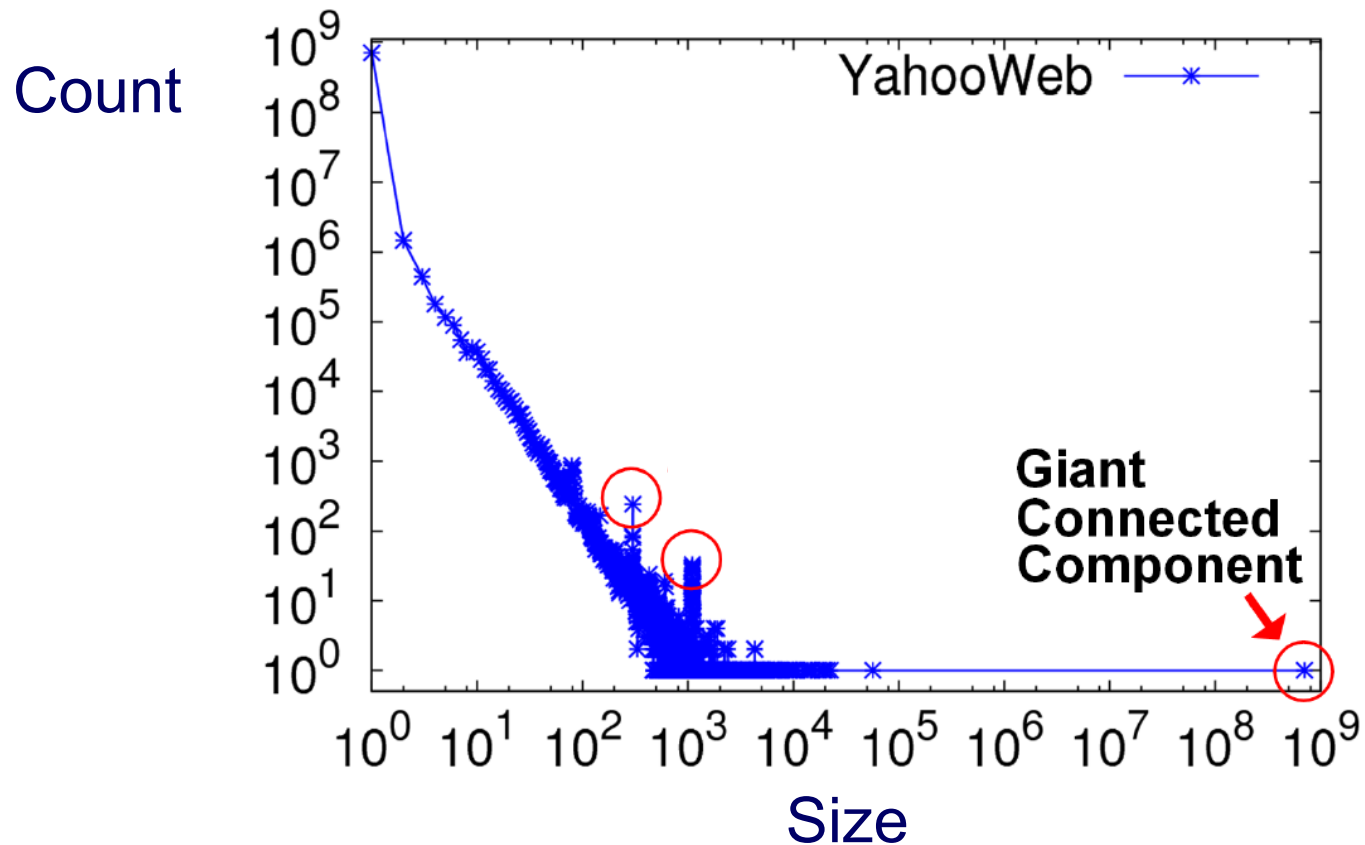
Generalized Iterated Matrix Vector Multiplication (GIMV)

- PageRank
- proximity (RWR)
- Diameter
- Connected components
- (eigenvectors,
- Belief Prop.
- ...)

Matrix – vector
Multiplication
(iterated)

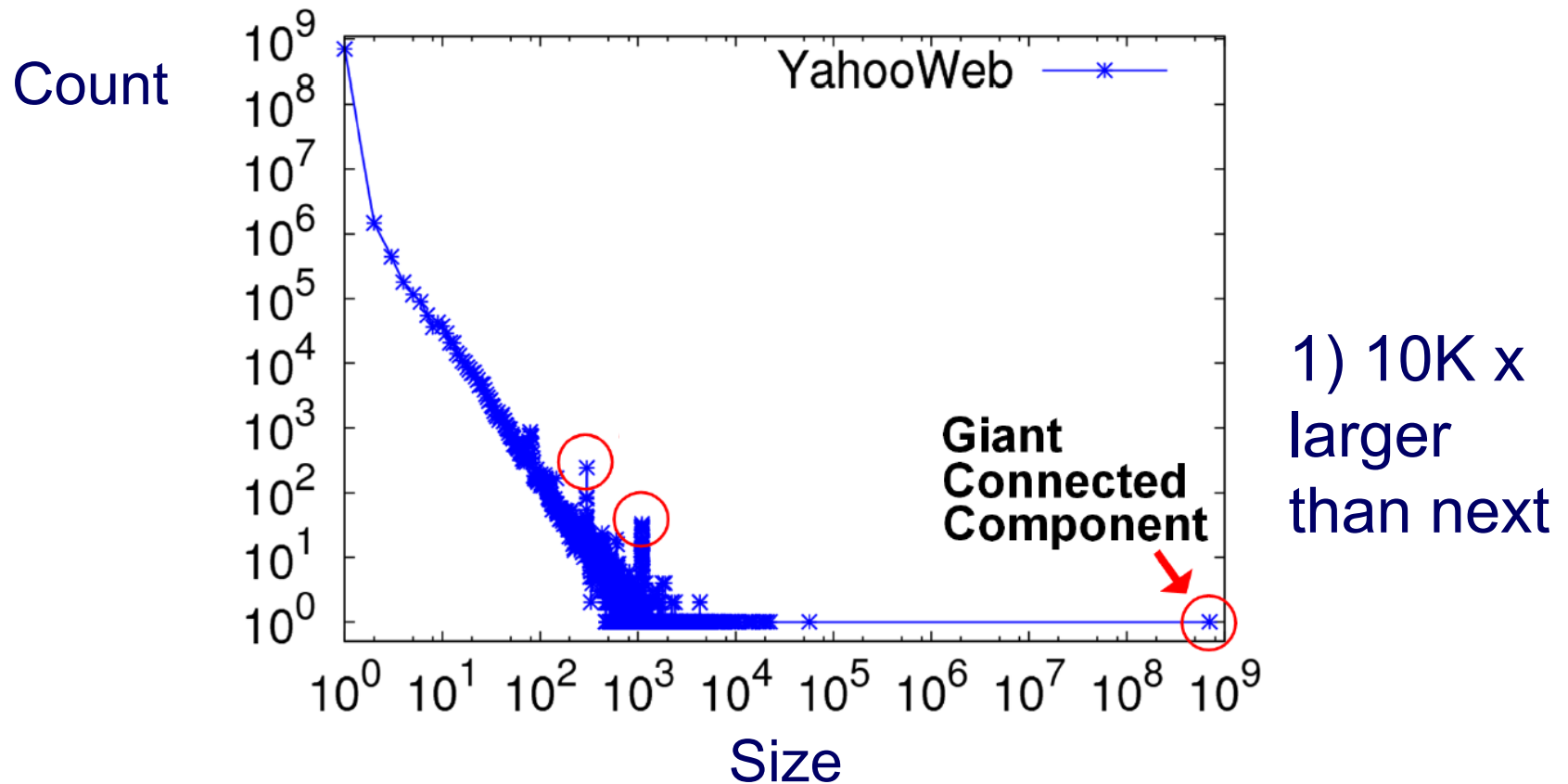
Example: GIM-V At Work

- Connected Components – 4 observations:



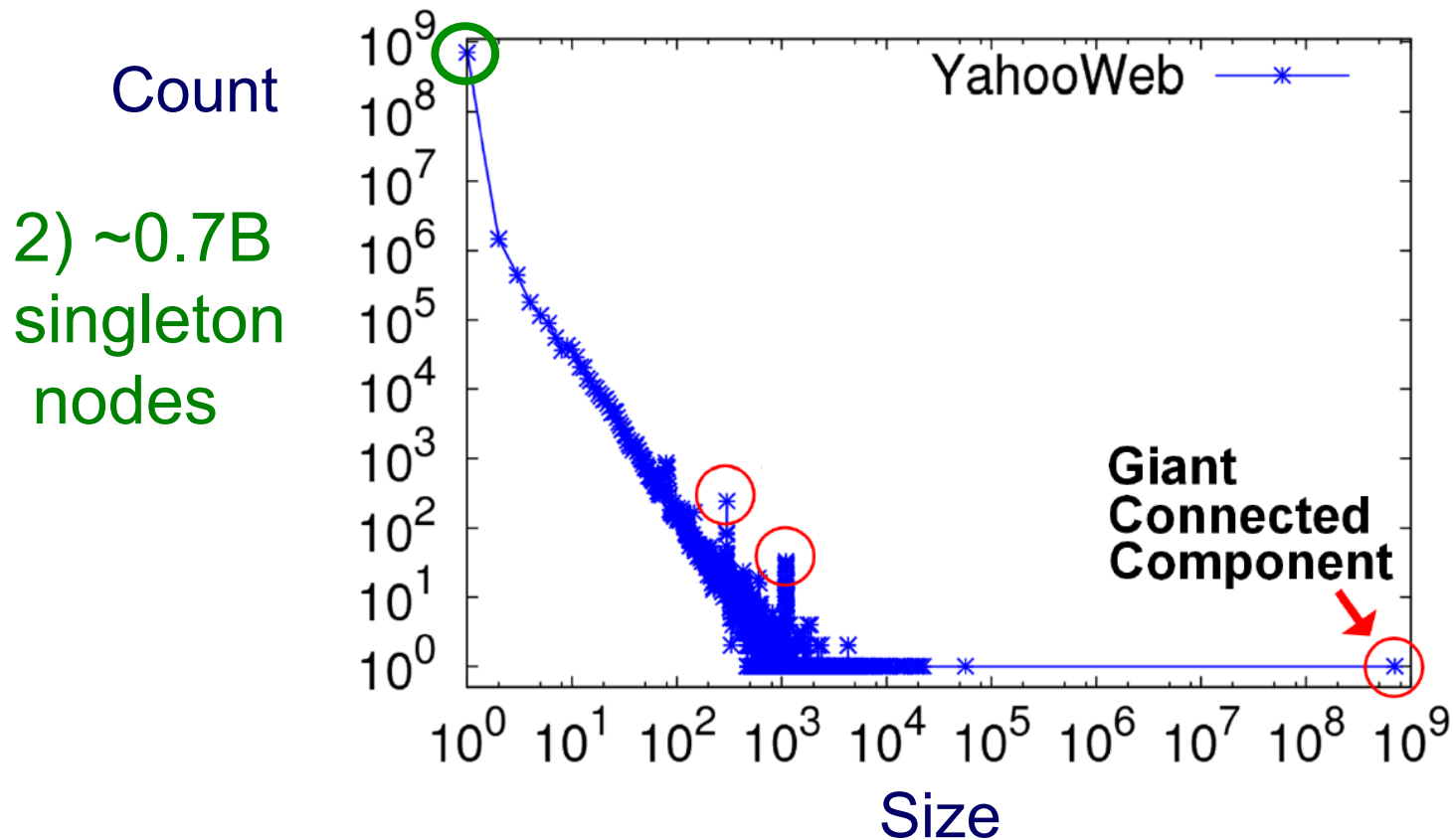
Example: GIM-V At Work

- Connected Components



Example: GIM-V At Work

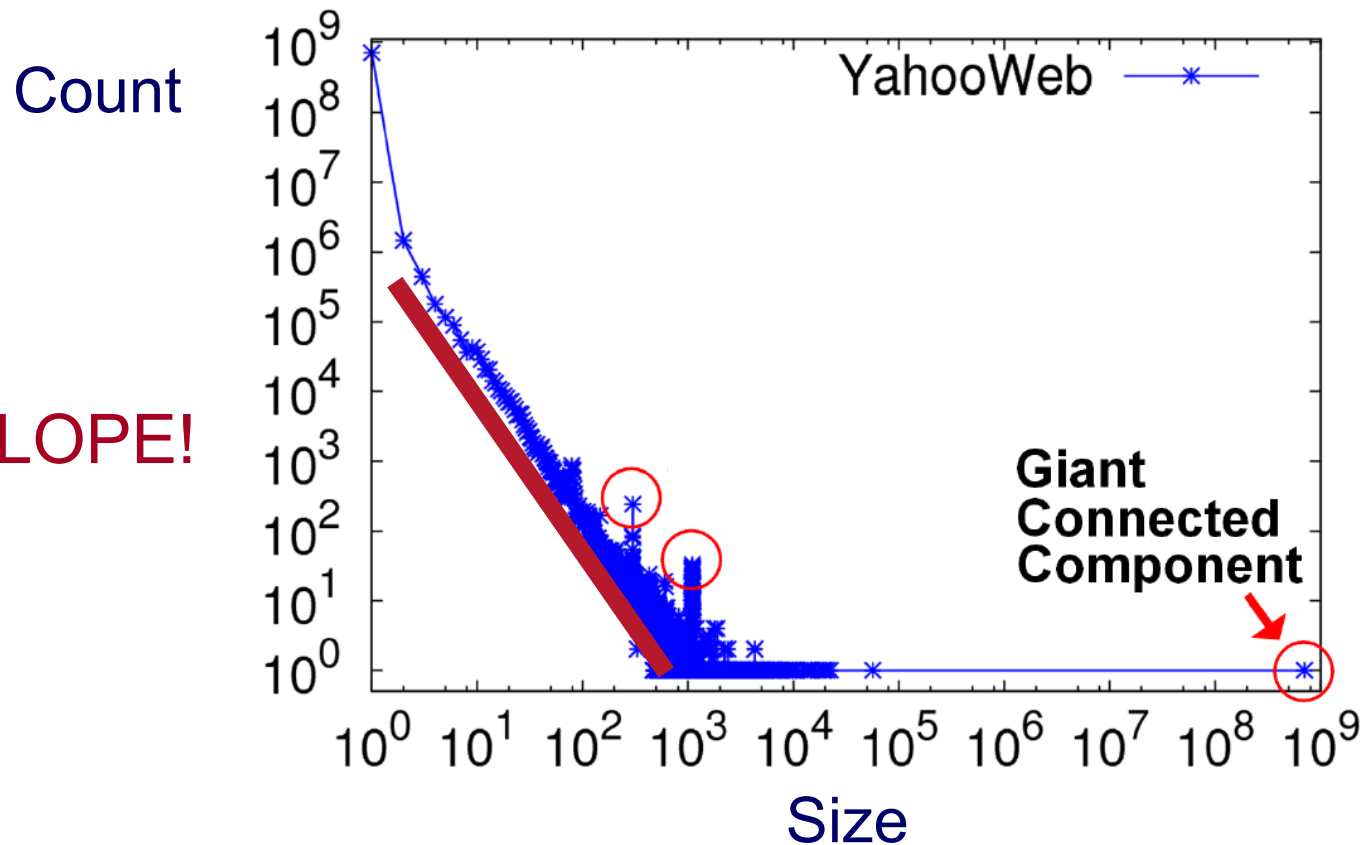
- Connected Components



Example: GIM-V At Work

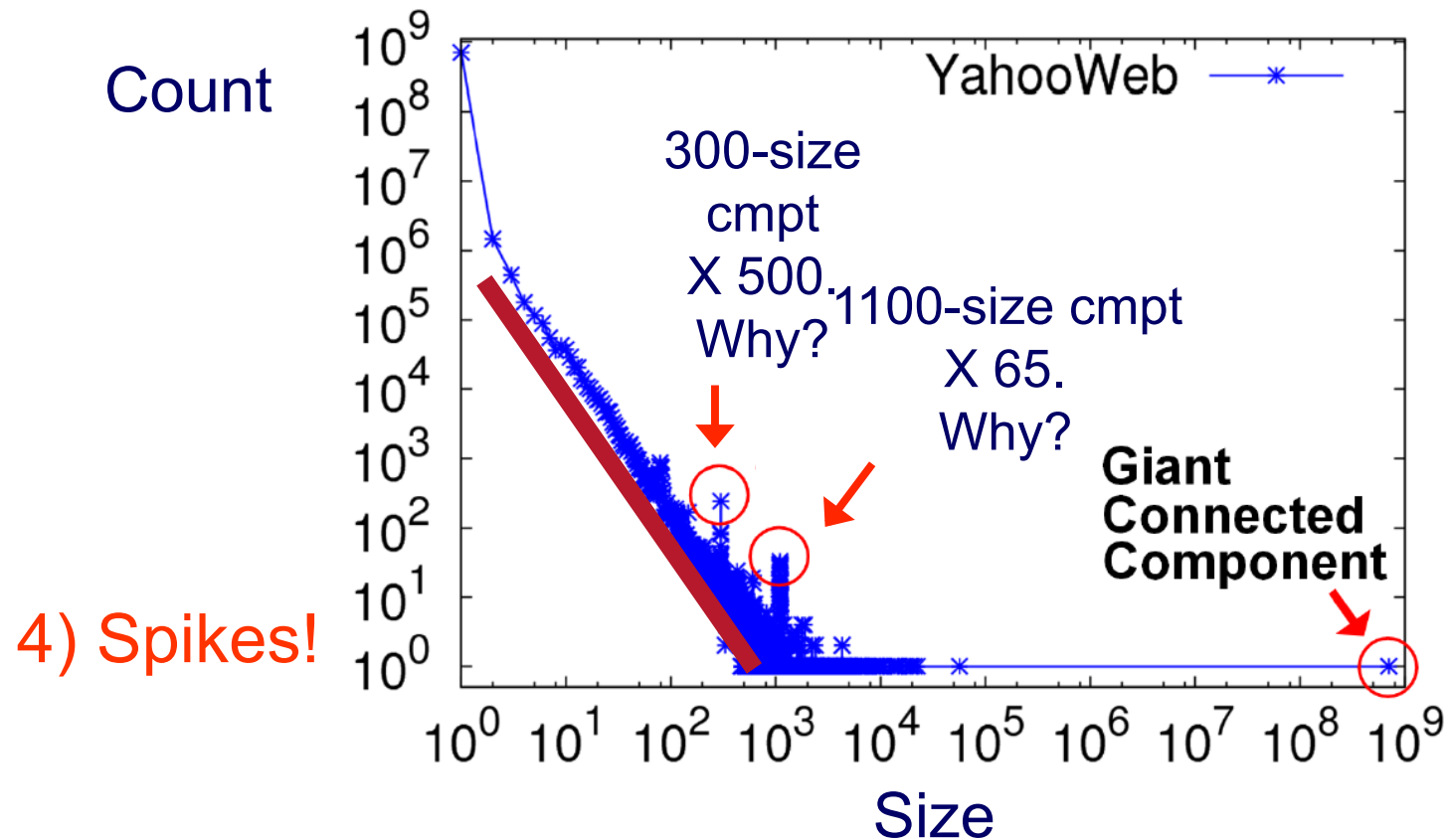
- Connected Components

3) SLOPE!



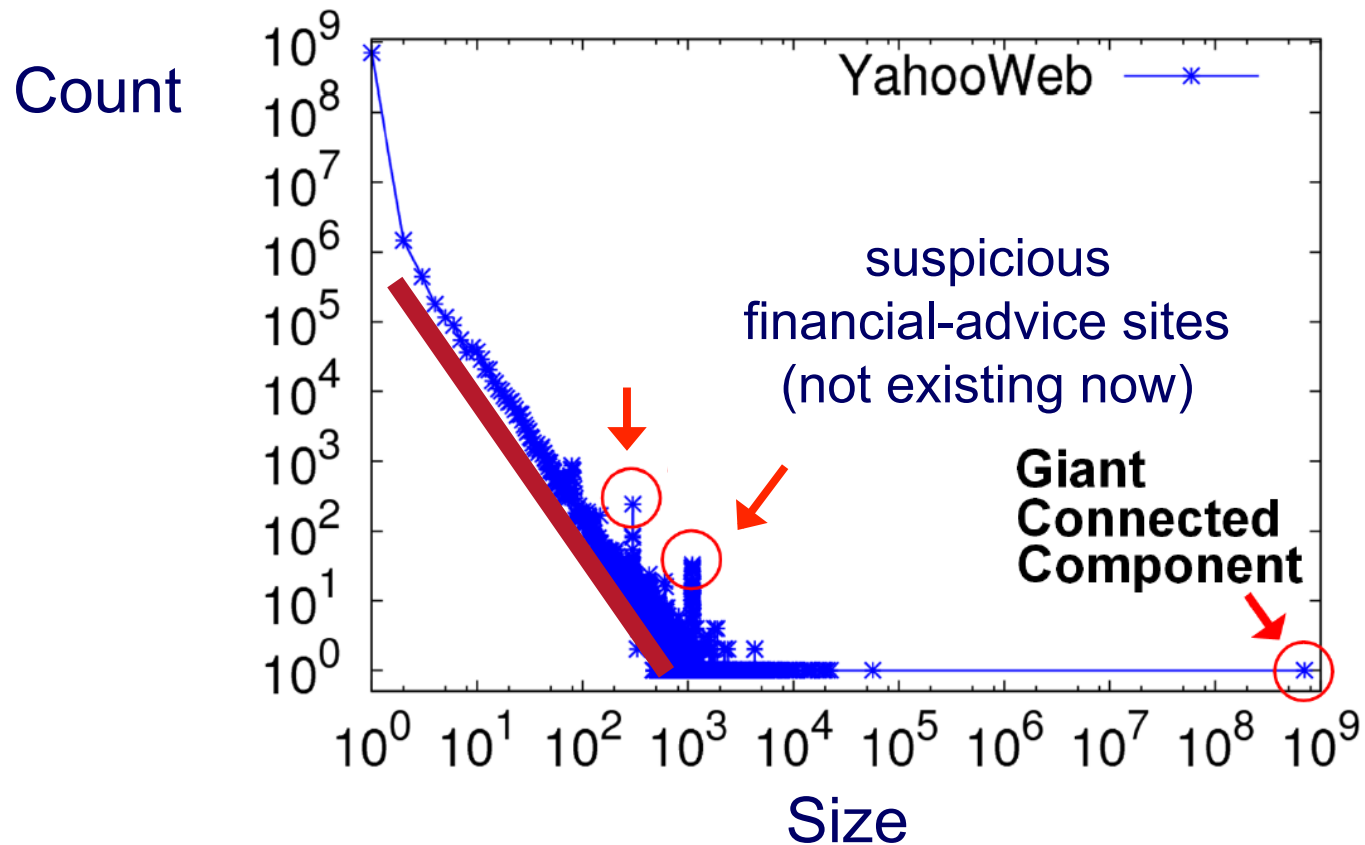
Example: GIM-V At Work

- Connected Components



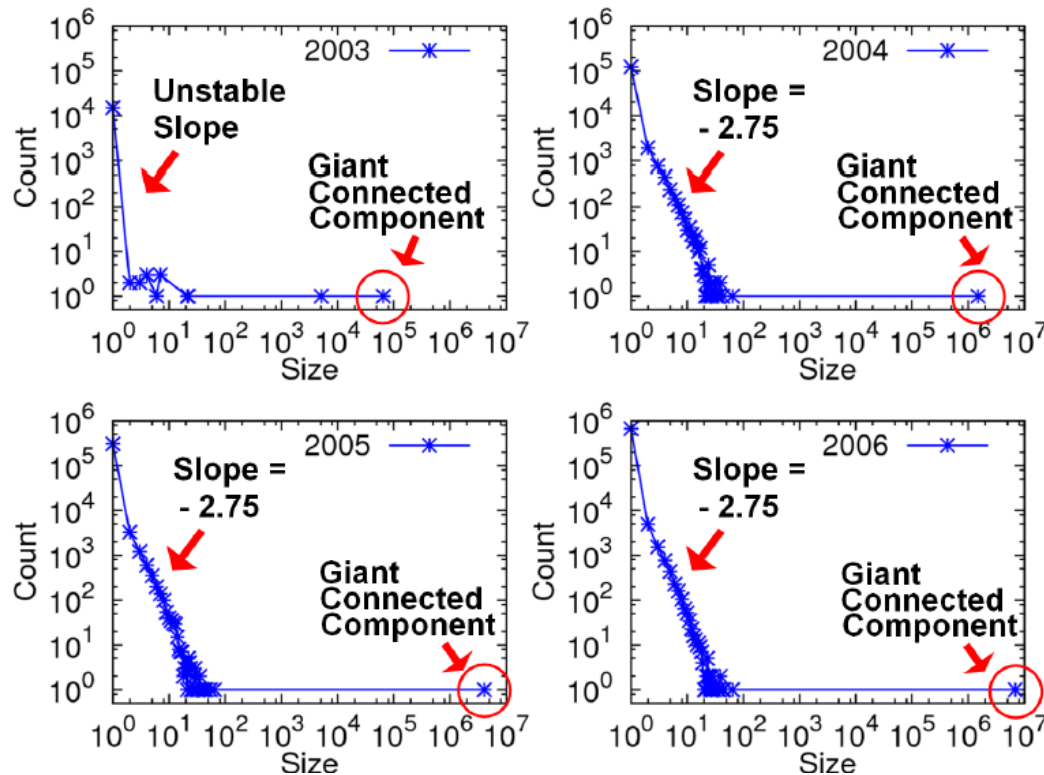
Example: GIM-V At Work

- Connected Components



GIM-V At Work

- Connected Components over Time
- **LinkedIn: 7.5M nodes and 58M edges**



Stable tail slope
after the gelling point

Conclusions

- Hadoop: promising architecture for Tera/Peta scale graph mining

Resources:

- <http://hadoop.apache.org/core/>
- <http://hadoop.apache.org/pig/>

Higher-level language for data processing

References

- [Jeffrey Dean](#) and [Sanjay Ghemawat](#), *MapReduce: Simplified Data Processing on Large Clusters*, OSDI'04
- Christopher Olston, [Benjamin Reed](#), [Utkarsh Srivastava](#), [Ravi Kumar](#), [Andrew Tomkins](#): *Pig latin: a not-so-foreign language for data processing*. [SIGMOD 2008](#): 1099-1110

Overall Conclusions

- Real graphs exhibit surprising **patterns** (power laws, shrinking diameter, super-linearity on edge weights, triangles etc)
- **SVD**: a powerful tool (HITS, PageRank)
- Several other tools: **tensors**, METIS, ...
 - But: good communities might **not** exist...
- Immunization: **first** eigenvalue
- Scalability: **hadoop**/parallelism

Our goal:

Open source system for mining huge graphs:

PEGASUS project (PEta GrAph mining System)

- www.cs.cmu.edu/~pegasus
- code and papers



Project info

www.cs.cmu.edu/~pegasus



Chau,
Polo



Koutra,
Danae



Prakash,
Aditya



Akoglu,
Leman

Kang, U



McGlohon,
Mary



Tong,
Hanghang



Thanks to: NSF IIS-0705359, IIS-0534205,
CTA-INARC; Yahoo (M45), LLNL, IBM, SPRINT,
Google, INTEL, HP, iLab

Extra material

- E-bay fraud detection
- Outlier detection

Detailed outline

extra

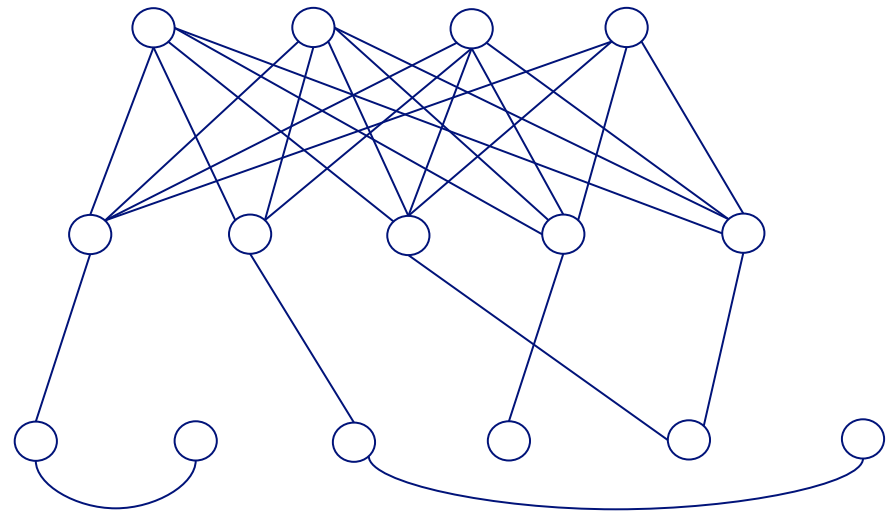
- ➔ • Fraud detection in e-bay
 - Anomaly detection

E-bay Fraud detection

extra



w/ Polo Chau &
Shashank Pandit, CMU

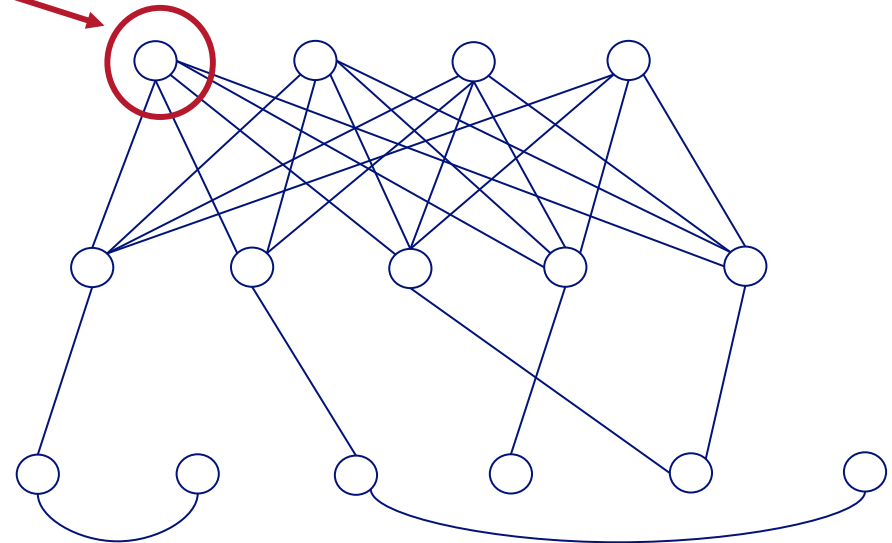


NetProbe: A Fast and Scalable System for Fraud Detection in Online Auction Networks, S. Pandit, D. H. Chau, S. Wang, and C. Faloutsos (*WWW'07*), pp. 201-210

E-bay Fraud detection

extra

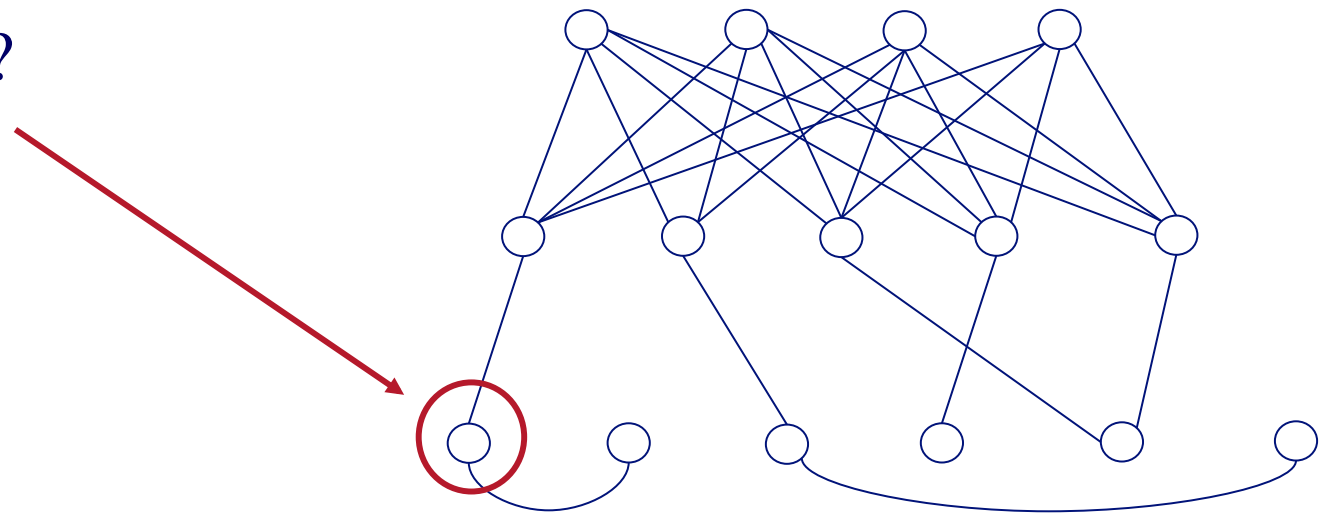
- lines: positive feedbacks
- would you buy from him/her?



E-bay Fraud detection

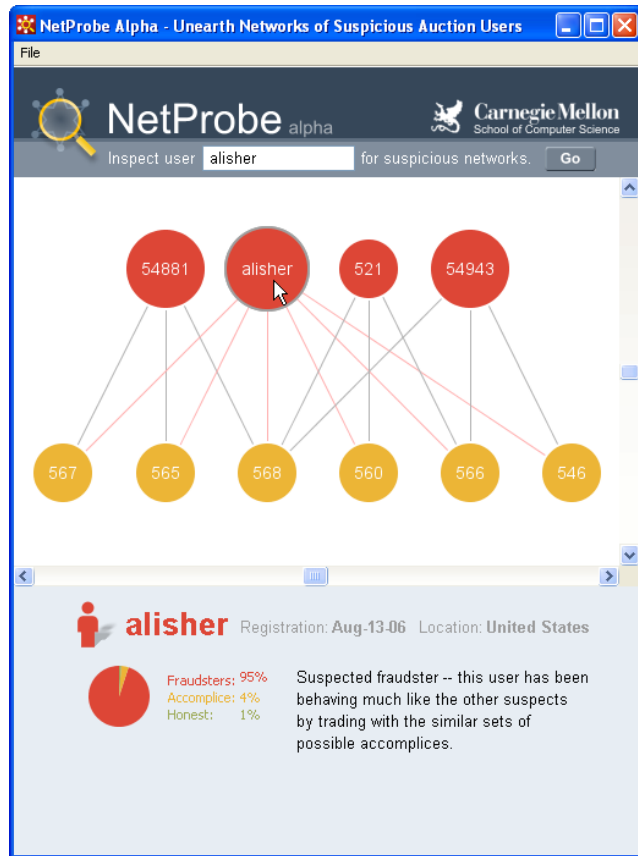
extra

- lines: positive feedbacks
- would you buy from him/her?
- or him/her?

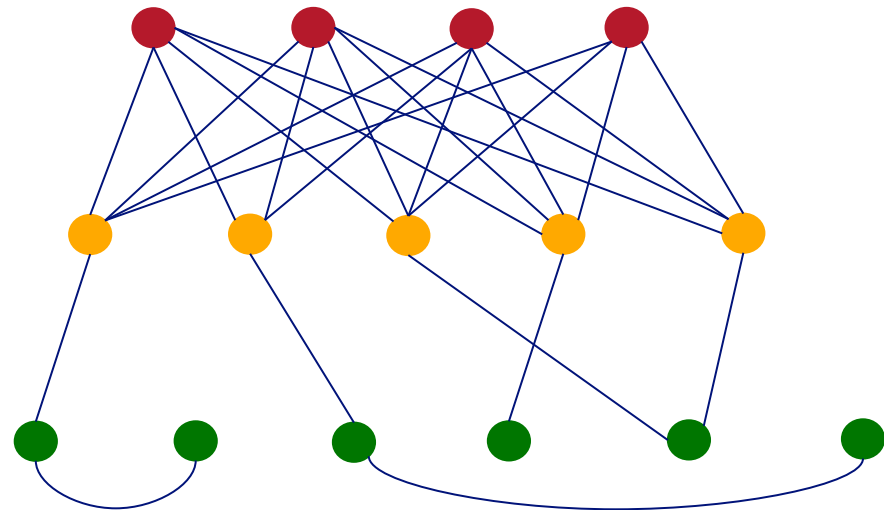


E-bay Fraud detection - NetProbe

extra



Belief Propagation gives:



Popular press



The Washington Post

Los Angeles Times

And less desirable attention:

- E-mail from ‘Belgium police’ (‘copy of your code?’)

Extra material

- E-bay fraud detection
- Outlier detection

OddBall: Spotting Anomalies in Weighted Graphs



Leman Akoglu, Mary McGlohon, Christos
Faloutsos

*Carnegie Mellon University
School of Computer Science*

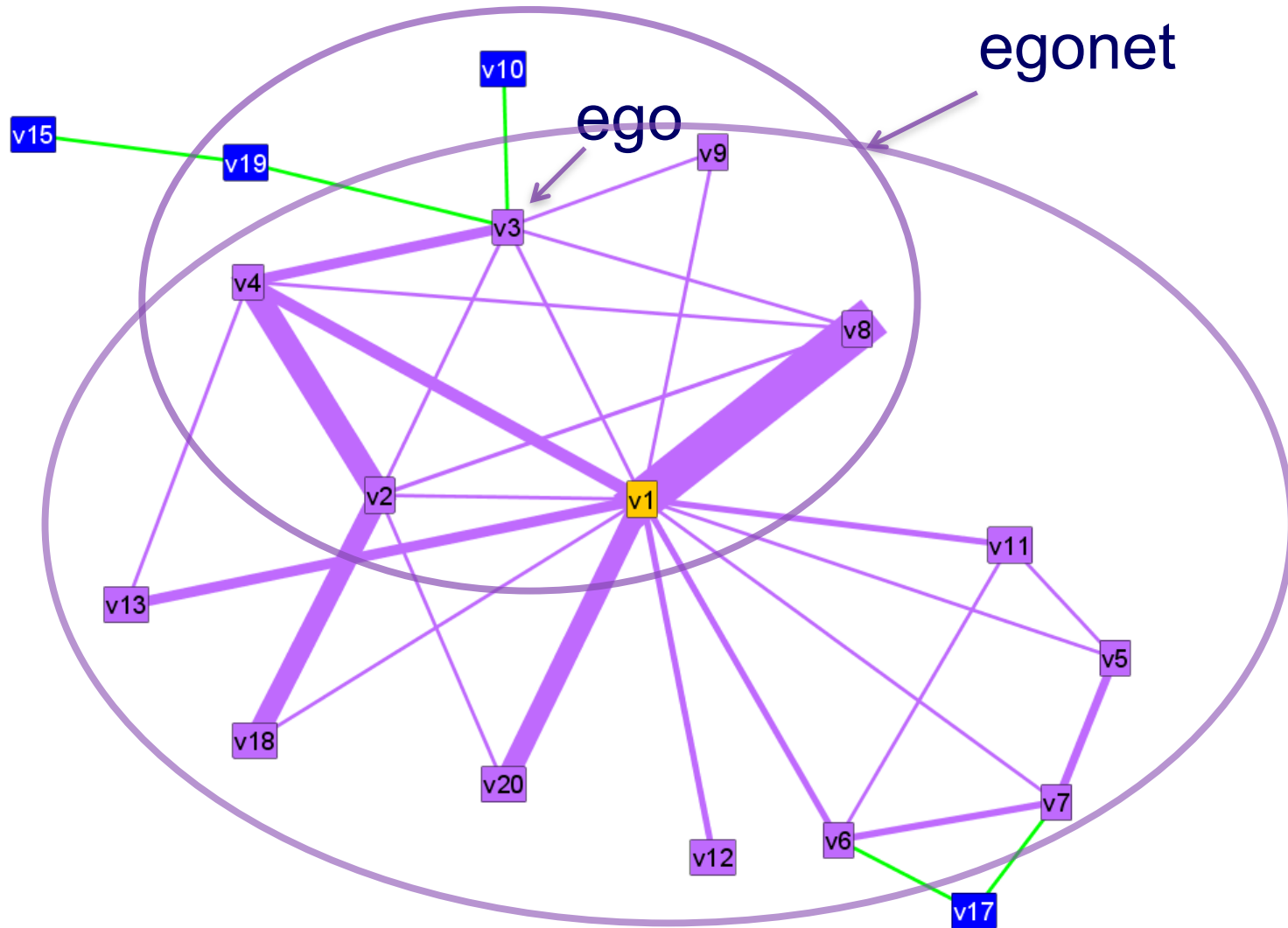
PAKDD 2010, Hyderabad, India

Main idea

For each node,

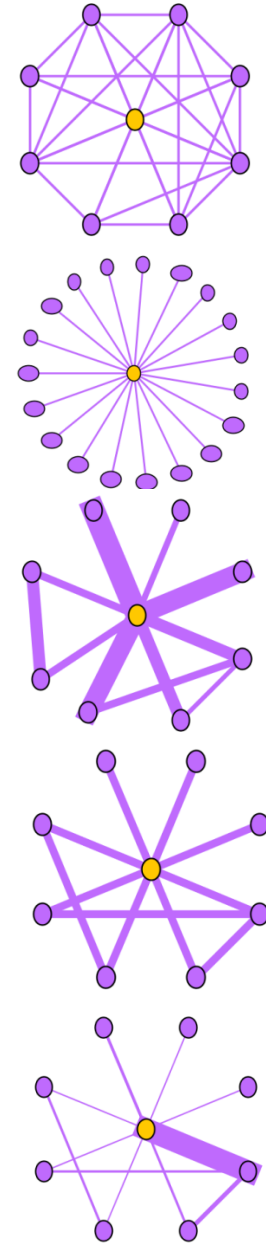
- extract ‘ego-net’ (=1-step-away neighbors)
- Extract features (#edges, total weight, etc etc)
- Compare with the rest of the population

What is an egonet?

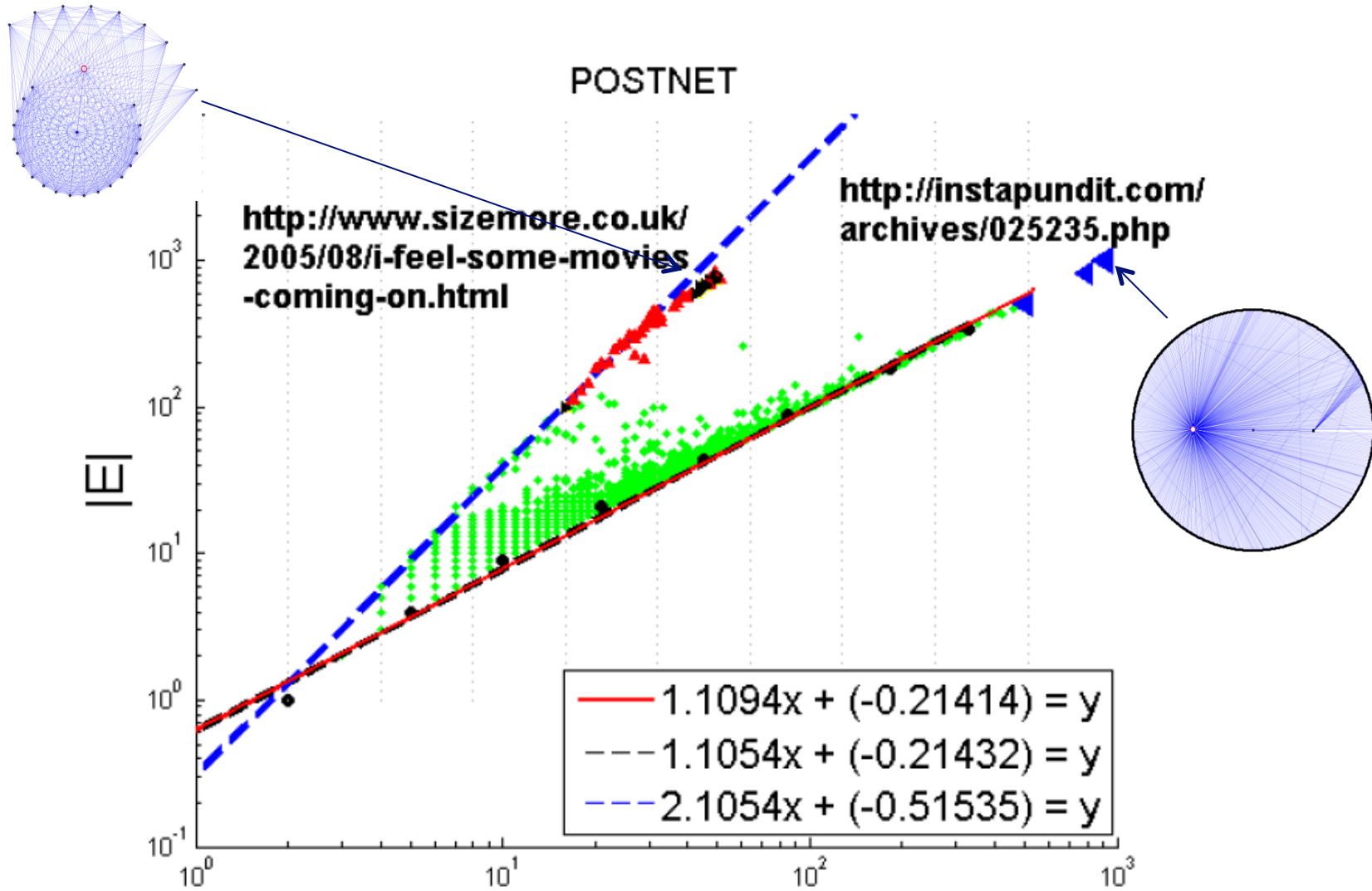


Selected Features

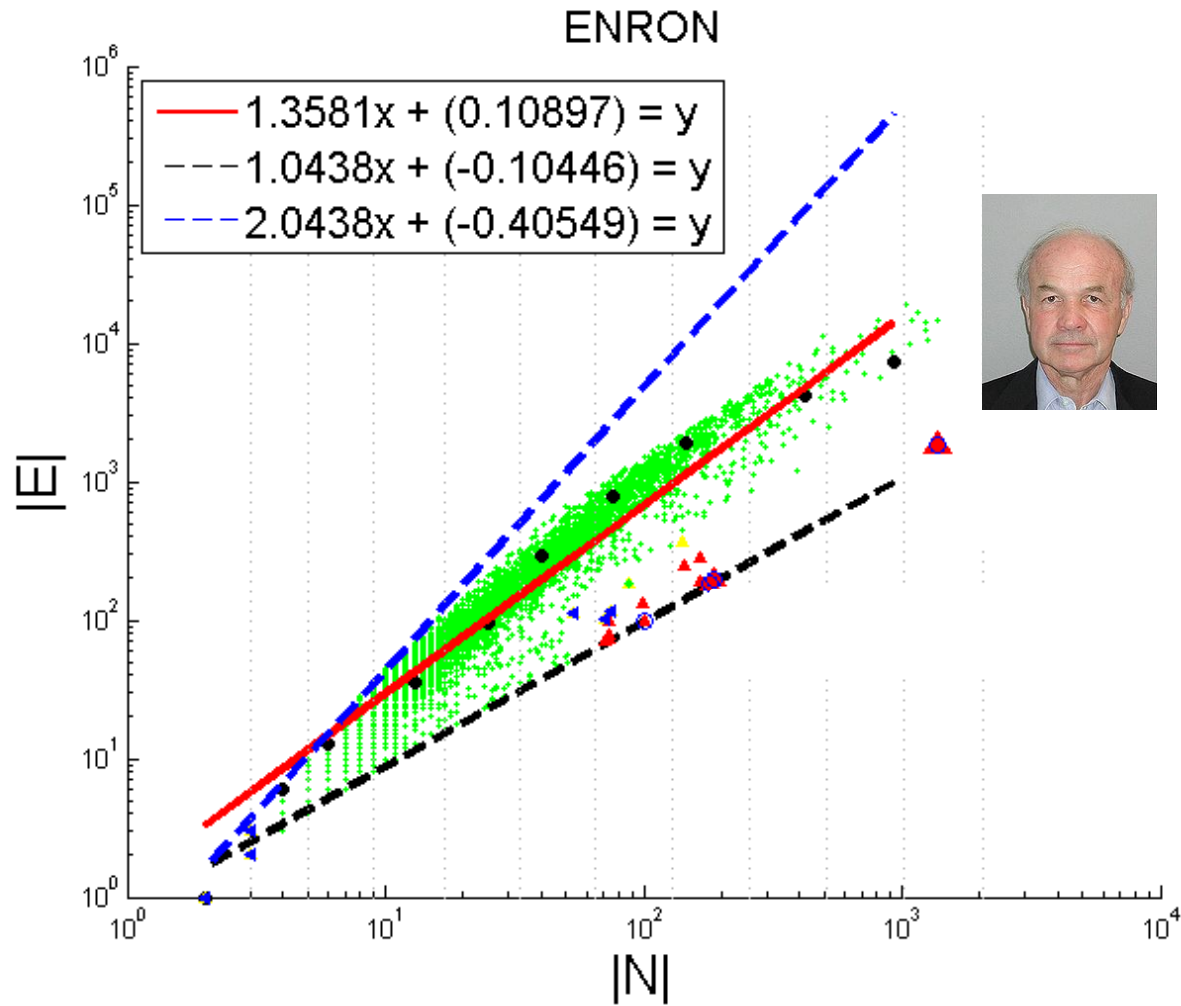
- N_i : number of neighbors (degree) of ego i
- E_i : number of edges in egonet i
- W_i : total weight of egonet i
- $\lambda_{w,i}$: principal eigenvalue of the **weighted** adjacency matrix of egonet I



Near-Clique/Star



Near-Clique/Star



END