

15-826: Multimedia Databases and Data Mining

Guest Lecture: Graph mining
Parts 2&3

Christos Faloutsos

Thanks

- Deepayan Chakrabarti (CMU)
- Soumen Chakrabarti (IIT-Bombay)
- Michalis Faloutsos (UCR)
- George Siganos (UCR)



PART 2: PageRank, HITS, and eigenvalues

Outline

Part 1: Topology, 'laws' and generators

➡ Part 2: PageRank, HITS and eigenvalues

Part 3: Influence, communities

Part 2: PageRank, HITS and eigenvalues

- How important is a node?
- Who is the best customer to advertise to?

Part 2: PageRank, HITS and eigenvalues

- How important is a node?
- Who is the best customer to advertise to?

Answers:

- PageRank (random surfer model)
- HITS (hubs and authorities)



(Published) PageRank

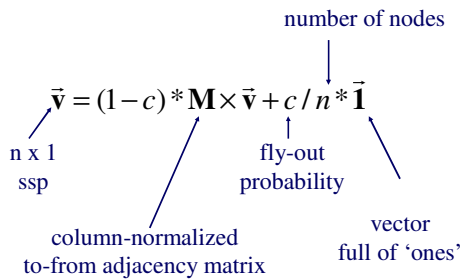
- Do a random walk, but
- with probability c , fly-out to a random node
- Then, the ssp vector \vec{v} obeys:



(Published) PageRank

$$\vec{v} = (1 - c) * \mathbf{M} \times \vec{v} + c / n * \vec{\mathbf{1}}$$

(Published) PageRank



CMU SCS

A clever variation

- Personalized PageRank:
- Who is the most important node in the vicinity of node i ?

15-826 Copyright: C. Faloutsos (2009) #10

CMU SCS

Personalized PageRank

- [Haveliwala+]

$$\vec{v}_i = (1-c) * \mathbf{M} \times \vec{v}_i + c * \vec{e}_i$$

ssp, when we restart from node ' i '

15-826 Copyright: C. Faloutsos (2009) #11

CMU SCS

Personalized PageRank

- [Haveliwala+]

$$\vec{v}_i = (1-c) * \mathbf{M} \times \vec{v}_i + c * \vec{e}_i$$

ssp, when we restart from node ' i '

i -th row \rightarrow

0
...
1
0
...

15-826 Copyright: C. Faloutsos (2009) #12

CMU SCS

Personalized PageRank

- [Haveliwala+]

$$\vec{v}_i = (1-c) * \mathbf{M} \times \vec{v}_i + c * \vec{e}_i \quad \text{new}$$

$$\vec{v} = (1-c) * \mathbf{M} \times \vec{v} + c/n * \mathbf{1} \quad \text{original}$$

15-826 Copyright: C. Faloutsos (2009) #13

CMU SCS

Personalized PageRank

- [Haveliwala+]

$$\vec{v}_i = (1-c) * \mathbf{M} \times \vec{v}_i + c * \vec{e}_i$$

- then $s_{i,j}$ = prob(a random walker with restarts from node i , will find itself at node j)

$$[s_{i,j}] = \mathbf{S} = c * [\mathbf{I} - (1-c) * \mathbf{M}]^{-1}$$

15-826 Copyright: C. Faloutsos (2009) #14

CMU SCS

Personalized PageRank

$s_{i,j}$ = prob(a random walker with restarts from node i , will find itself at node j)

$s_{i,j}$: good measure of how close is node j to node i

$$[s_{i,j}] = \mathbf{S} = c * [\mathbf{I} - (1-c) * \mathbf{M}]^{-1}$$

15-826 Copyright: C. Faloutsos (2009) #15

Our wish list:

- ✓ How important is a node?
- ✓ Who is the best customer to advertise to?

ssp values answer these questions

Outline

Part 1: Topology, 'laws' and generators

Part 2: PageRank, HITS and eigenvalues

- Eigenvalues and PageRank

➡ • SVD and HITS

Part 3: influence, virus prop., communities

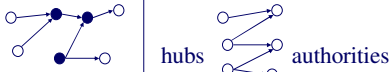
Kleinberg's algorithm ('HITS')

- Problem defn: given the web and a query
- find the most 'authoritative' web pages for this query

Kleinberg's algorithm



- give high score (= 'authorities') to nodes that many important nodes point to
- give high importance score ('hubs') to nodes that point to good 'authorities'



15-826

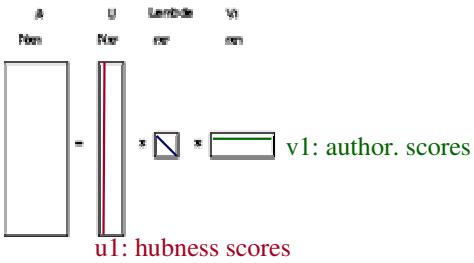
Copyright: C. Faloutsos (2009)

#19

SVD & HITS



- Adjacency matrix $A = U \Lambda V^T$



15-826

Copyright: C. Faloutsos (2009)

#20

Conclusions

eigenvalues/eigenvectors: vital for

- PageRank,
- (virus propagation - coming up next!)
- (graph partitioning - not mentioned here)

15-826

Copyright: C. Faloutsos (2009)

#21

Conclusions, cont'd

SVD

- closely related: HITS/Kleinberg
- (and also LSI, KLT, PCA, Least squares, ...)

Both eigen- and singular decompositions are **extremely useful, well understood** tools for graphs / matrices.

PART 3: Influence, virus propagation, communities

Outline

Part 1: Topology, 'laws' and generators

Part 2: PageRank, HITS and eigenvalues

- Eigenvalues and PageRank
- SVD and HITS

➡ Part 3: influence, virus prop., communities

CMU SCS

Problem definition

- Q1: How does a virus spread across an arbitrary network?
- Q2: will it create an epidemic?

15-826 Copyright: C. Faloutsos (2009) #25

CMU SCS

Framework

- Susceptible-Infected-Susceptible (SIS) model
 - Cured nodes immediately become susceptible

15-826 Copyright: C. Faloutsos (2009) #26

CMU SCS

The model

- (virus) Birth rate β : probability than an infected neighbor attacks
- (virus) Death rate δ : probability that an infected node heals

15-826 Copyright: C. Faloutsos (2009) #27

CMU SCS

The model

- Virus 'strength' $s = \beta/\delta$

15-826 Copyright: C. Faloutsos (2009) #28

CMU SCS

Other models:

- SIR: Susceptible - infected & infectious - recovered/removed
 - eg., mumps, chickenpox; black plague

15-826 Copyright: C. Faloutsos (2009) #29

CMU SCS

Other models:

- and many more:
- SEIR: Susceptible; Exposed (= infected, but not infectious yet); I; R
- variations:
 - M: passively immune, like infants
 - with births/newcomers
 - ...

15-826 Copyright: C. Faloutsos (2009) #30

Epidemic threshold τ

of a graph, defined as 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

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?



Epidemic threshold

- [Theorem] We have no epidemic, if

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

CMU SCS

Epidemic threshold

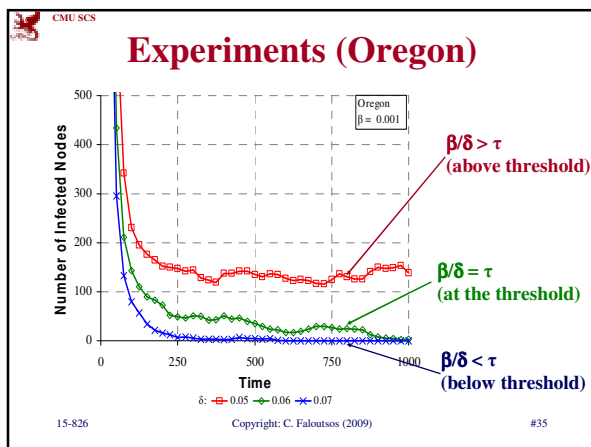
- [Theorem] We have no epidemic, if recovery prob. $\beta/\delta < \tau = 1/\lambda_{1,A}$

epidemic threshold
 attack prob.

largest eigenvalue of adj. matrix A

Proof: [Wang+03]

15-826 Copyright: C. Faloutsos (2009) #34



CMU SCS

Our wish list:

- Who is the best person/computer to immunize against a virus?

15-826 Copyright: C. Faloutsos (2009) #36

CMU SCS

Our wish list:

✓ Who is the best person/computer to immunize against a virus? Highest diff in λ_1

15-826 Copyright: C. Faloutsos (2009) #37

CMU SCS

Outline

Part 1: Topology, 'laws' and generators
Part 2: PageRank, HITS and eigenvalues

- Eigenvalues and PageRank
- SVD and HITS

➡ Part 3: influence, virus prop., communities

15-826 Copyright: C. Faloutsos (2009) #38

CMU SCS

Graph clustering & mining

- Q1: which edges/nodes are 'abnormal'?
- Q2: split a graph in k 'natural' communities - but how to determine k ?

15-826 Copyright: C. Faloutsos (2009) #39

CMU SCS

Graph partitioning

- Documents x terms
- Customers x products
- Users x web-sites

15-826 Copyright: C. Faloutsos (2009) #40

CMU SCS

Graph partitioning

- Documents x terms
- Customers x products
- Users x web-sites

- Q: HOW MANY PIECES?

15-826 Copyright: C. Faloutsos (2009) #41

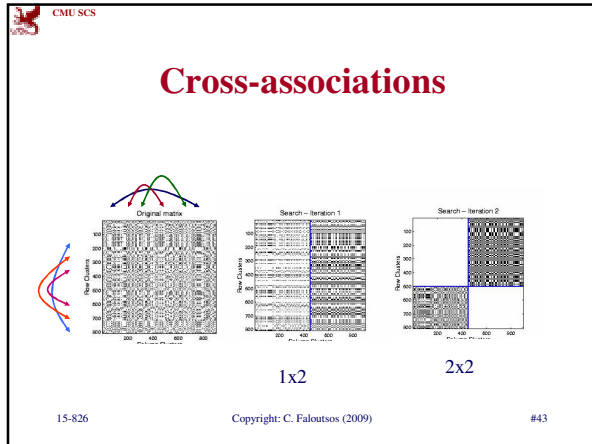
CMU SCS

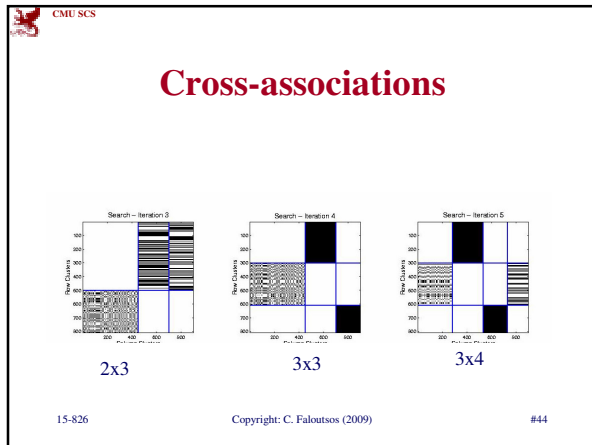
Graph partitioning

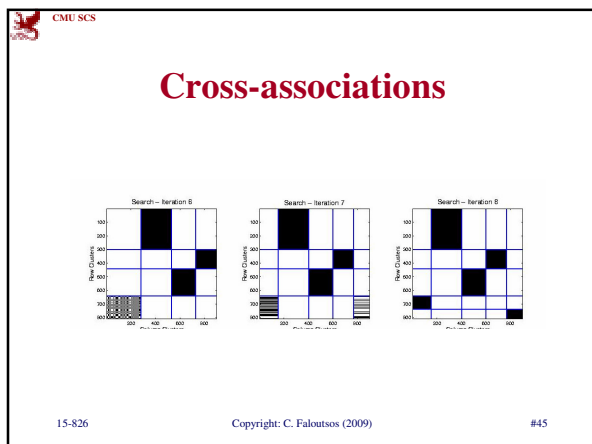
- Documents x terms
- Customers x products
- Users x web-sites

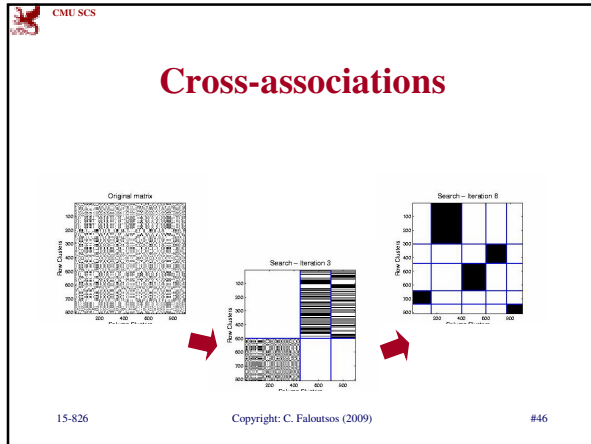
- Q: HOW MANY PIECES?
- A: MDL/ compression

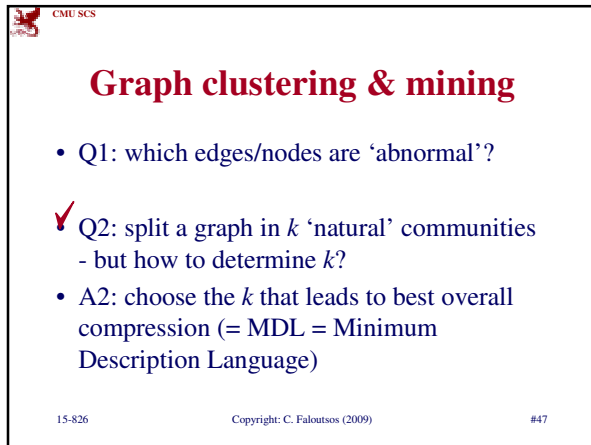
15-826 Copyright: C. Faloutsos (2009) #42

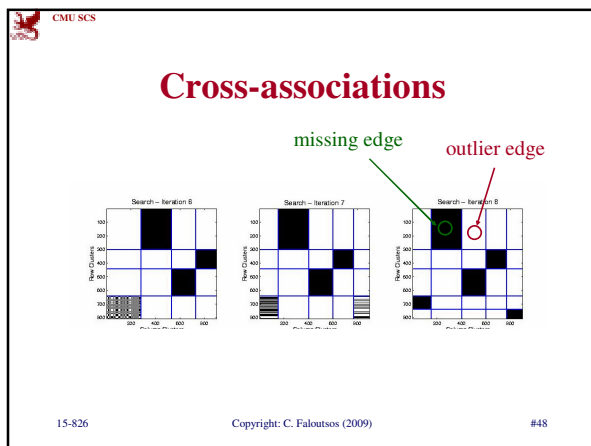












Conclusions

- virus propagation: eigenvalue determines the epidemic threshold (SIS model)
- communities/graph partitioning: MDL

Resources: Software and urls

- SVD packages: in **many** systems (matlab, mathematica, LINPACK, LAPACK)
- stand-alone, free code: SVDPACK from Michael Berry
<http://www.cs.utk.edu/~berry/projects.html>

Books

- Faloutsos, C. (1996). *Searching Multimedia Databases by Content*, Kluwer Academic Inc.
- Jolliffe, I. T. (1986). *Principal Component Analysis*, Springer Verlag.

Books

- [Press+92] William H. Press, Saul A. Teukolsky, William T. Vetterling and Brian P. Flannery: *Numerical Recipes in C*, Cambridge University Press, 1992, 2nd Edition. (Great description, intuition and code for SVD)

References

- Berry, Michael: <http://www.cs.utk.edu/~lsi/>
- [Brin+98] Brin, S. and L. Page (1998). *Anatomy of a Large-Scale Hypertextual Web Search Engine*. 7th Intl World Wide Web Conf.
- [Chakrabarti'04] D. Chakrabarti, *AutoPart: Parameter-Free Graph Partitioning and Outlier Detection*, PKDD 2004 (pages 112-124), Pisa, Italy

References (cont'd)

- [Chakrabarti+,04a] D. Chakrabarti, S. Papadimitriou, D. Modha and C. Faloutsos, *Fully Automatic Cross-Associations*, KDD 2004 (pp. 79-88), Washington, USA
- [Haveliwala02] Taher H. Haveliwala, *Topic-Sensitive PageRank* World Wide Web Conference, 2002

References (cont'd)

- Fukunaga, K. (1990). *Introduction to Statistical Pattern Recognition*, Academic Press.
- Kleinberg, J. (1998). *Authoritative sources in a hyperlinked environment*. Proc. 9th ACM-SIAM Symposium on Discrete Algorithms.

References (cont'd)

- [Wang+03] Yang Wang, Deepayan Chakrabarti, Chenxi Wang and Christos Faloutsos: *Epidemic Spreading in Real Networks: an Eigenvalue Viewpoint*, SRDS 2003, Florence, Italy.

Discussion

A lot of recent interest - topics we didn't cover:

- Relational learning, e.g., [David Jensen; Daphne Koller; Saso Dzeroski]
- Frequent sub-graphs, e.g., [Jiawei Han, Jian Pei; George Karypis, Vipin Kumar; Mohammed Zaki]

Discussion cont'd

- Graph partitioning, e.g., [METIS (Karypis)]
- Social networks, e.g., [Kathleen Carley; Wasserman+Faust]
- Web mining, e.g., [Soumen Chakrabarti]

Overall conclusions

- Surprising patterns in graphs
- Powerful tools exist:
 - Self-similarity, fractals, Kronecker
 - SVD, eigenvalues
 - MDL for partitioning
