# 15-826: Multimedia Databases and Data Mining

### Graph mining - Part 2&3

*Christos Faloutsos*

15-826          Copyright: C. Faloutsos (2006)          #1

---

# Thanks

- Deepayan Chakrabarti (CMU)

- Soumen Chakrabarti    (*IIT-Bombay*)

- Michalis Faloutsos (UCR)

- George Siganos (UCR)

15-826          Copyright: C. Faloutsos (2006)          #2

---

# PART 2:
# PageRank, HITS, and eigenvalues

15-826          Copyright: C. Faloutsos (2006)          #3

---

# Outline

Part 1: Topology, 'laws' and generators
➡ Part 2: PageRank, HITS and eigenvalues
Part 3: Influence, communities

15-826          Copyright: C. Faloutsos (2006)          #4

---

### Part 2: PageRank, HITS and eigenvalues

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

15-826          Copyright: C. Faloutsos (2006)          #5

---

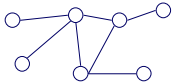# Outline

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 (2006)          #6

# Motivating problem
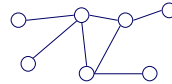
Given a graph, find its most interesting/central node

# Motivating problem

Given a graph, find its most interesting/central node

A node is important,
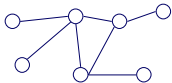if it is connected
with important nodes
(recursive, but OK!)

# Motivating problem – pageRank solution

Given a graph, find its most interesting/central node

Proposed solution: Random walk; spot most 'popular' node (-> steady state prob. (ssp))

A node has high ssp,
if it is connected
with high ssp nodes
(recursive, but OK!)

# Notational conventions

- bold capitals -> matrix (eg. $\mathbf{A}$, $\mathbf{U}$, $\mathbf{\Lambda}$, $\mathbf{V}$)
- bold lower-case -> <u>column</u> vector (eg., $\mathbf{x}$, $\mathbf{v}_1$, $\mathbf{u}_3$)
- regular lower-case -> scalars (eg., $\lambda_1$, $\lambda_r$ )

# (Simplified) PageRank algorithm

- Let $\mathbf{A}$ be the transition matrix (= adjacency matrix); let $\mathbf{M} = \mathbf{A}^\mathbf{T}$ and column-normalized - then

# (Simplified) PageRank algorithm

- $\mathbf{M} \, \mathbf{p} = \mathbf{p}$

2

# (Simplified) PageRank algorithm

- $\mathbf{M} \, \mathbf{p} = 1 * \mathbf{p}$
- thus, $\mathbf{p}$ is the eigenvector that corresponds to the highest eigenvalue (=1, since the matrix is column-normalized)

# (Simplified) PageRank algorithm

- In short: imagine a particle randomly moving along the edges
- compute its steady-state probabilities (ssp)

Full version of algo:  with occasional random jumps – see later

# Formal definition

If $\mathbf{M}$ is a (n x n) square matrix
$(\lambda \, , \mathbf{x})$ is an **eigenvalue/eigenvector** pair of $\mathbf{M}$ if

$$\mathbf{M} \, \mathbf{x} = \lambda \, \mathbf{x}$$

# (Published) PageRank

- Do a random walk, but
- with probability $c$, fly-out to a random node

- Then, the ssp vector $\mathbf{v}$ obeys:

# (Published) PageRank

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

# (Published) PageRank

number of nodes

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

n x 1
ssp

fly-out
probability

column-normalized
to-from adjacency matrix

vector
full of 'ones'

3

# Personalized PageRank

- [Haveliwala+]

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

ssp, when
we restart from node '*i*'

---

# Personalized PageRank

- [Haveliwala+]

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

ssp, when
we restart from node '*i*'

*i*-th row →

| 0 |
| --- |
| … |
| 1 |
| 0 |
| … |

---

# Personalized PageRank

- [Haveliwala+]

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

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

---

# Personalized PageRank

- [Haveliwala+]

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

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

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

---

# 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

4

# Kleinberg's algorithm ('HITS')

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

---

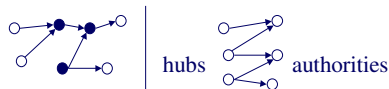# Kleinberg's algorithm

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

hubs          authorities

---

# Kleinberg's algorithm

Observations
- recursive definition!
- each node (say, '$i$'-th node) has both an authoritativeness score $a_i$ and a hubness score $h_i$

---

# Kleinberg's algorithm

Let $A$ be the adjacency matrix:
   the $(i,j)$ entry is 1 if the edge from $i$ to $j$ exists
Let $h$ and $a$ be [n x 1] vectors with the 'hubness' and 'authoritativiness' scores.
Then:

---

# Kleinberg's algorithm

In conclusion, we want vectors $h$ and $a$ such that:

$$h = A\ a$$
$$a = A^T\ h$$

That is:

$$a = A^T A\ a$$

---

# Kleinberg's algorithm

$a$ is a right- singular vector of the adjacency matrix $A$ (by dfn!)
== eigenvector of $A^T A$

# SVD & HITS

- $A = U \Lambda V^T$ - example:

A      U     Lambda    V1

Nxm     Nxr     rxr     rxm

v1: author. scores

u1: hubness scores

---

# Conclusions

eigenvalues/eigenvectors: vital for
- PageRank,
- (virus propagation - coming up next!)
- (graph partitioning - not mentioned here)

---

# Conclusions, cont'd

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

Both 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

---

# Problem definition

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

# Framework

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



Copyright: C. Faloutsos (2006) #37

---

# The model

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



Copyright: C. Faloutsos (2006) #38

---

# The model

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



Copyright: C. Faloutsos (2006) #39

---

# Other models:

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

Copyright: C. Faloutsos (2006) #40

---

# Other models:

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

Copyright: C. Faloutsos (2006) #41

---

# 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

Copyright: C. Faloutsos (2006) #42

# Epidemic threshold $\tau$

What should $\tau$ depend on?
- avg. degree? and/or highest degree?
- and/or variance of degree?
- and/or third moment of degree?

Copyright: C. Faloutsos (2006) #43

---

# Epidemic threshold

- [Theorem] We have no epidemic, if

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

Copyright: C. Faloutsos (2006) #44

---

# Epidemic threshold

- [Theorem] We have no epidemic, if

recovery prob.     epidemic threshold

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

attack prob.     largest eigenvalue of adj. matrix $A$

Proof: [Wang+03]

Copyright: C. Faloutsos (2006) #45

---

# Experiments (Oregon)



$\beta/\delta > t$ (above threshold)

$\beta/\delta = t$ (at the threshold)

$\beta/\delta < t$ (below threshold)

$\delta$: 0.05  0.06  0.07

Copyright: C. Faloutsos (2006) #46

---

# Our wish list:

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

Copyright: C. Faloutsos (2006) #47

---

# Our wish list:

✓Who is the best person/computer to immunize against a virus?   Highest diff in $\lambda1$

Copyright: C. Faloutsos (2006) #48

## Slide 1

# 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

## Slide 2

# Graph clustering & mining

• Q1: which edges/nodes are 'abnormal'?

• Q2: split a graph in $k$ 'natural' communities - but how to determine $k$?

## Slide 3

# Graph partitioning

• Documents x terms

• Customers x products

• Users x web-sites

## Slide 4

# Graph partitioning

• Documents x terms

• Customers x products

• Users x web-sites

• Q: HOW MANY PIECES?

## Slide 5

# Graph partitioning

• Documents x terms

• Customers x products

• Users x web-sites

• Q: HOW MANY PIECES?

• A: MDL/ compression

## Slide 6

# Cross-associations

1x2      2x2

# Cross-associations



2x3          3x3          3x4

# Cross-associations

# Cross-associations

# Graph clustering & mining

- Q1: which edges/nodes are 'abnormal'?

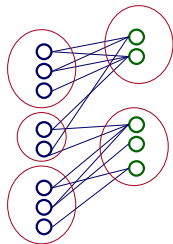- Q2: split a graph in $k$ 'natural' communities - but how to determine $k$?
- A2: choose the $k$ that leads to best overall compression (= MDL = Minimum Description Language)

# Cross-associations

missing edge          outlier edge

# Conclusions

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

**CMU SCS**

# 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**

15-826        Copyright: C. Faloutsos (2006)        #61

---

**CMU SCS**

# Books

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

15-826        Copyright: C. Faloutsos (2006)        #62

---

**CMU SCS**

# 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)

15-826        Copyright: C. Faloutsos (2006)        #63

---

**CMU SCS**

# 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

15-826        Copyright: C. Faloutsos (2006)        #64

---

**CMU SCS**

# 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

15-826        Copyright: C. Faloutsos (2006)        #65

---

**CMU SCS**

# 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.

15-826        Copyright: C. Faloutsos (2006)        #66

11

**CMU SCS**

# 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.

15-826                    Copyright: C. Faloutsos (2006)                    #67

---

**CMU SCS**

# 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]

15-826                    Copyright: C. Faloutsos (2006)                    #68

---

**CMU SCS**

# Discussion cont'd

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

15-826                    Copyright: C. Faloutsos (2006)                    #69

---

**CMU SCS**

# Overall conclusions

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

15-826                    Copyright: C. Faloutsos (2006)                    #70