## Einführung in Web- und Data-Science

**Community Analysis** 

Dr. Marcel Gehrke Universität zu Lübeck Institut für Informationssysteme



**IM FOCUS DAS LEBEN** 

#### Today's lecture

- Social Network Analysis
- Anchor text
- Link analysis for ranking
  - PageRank and variants
  - Hyperlink-Induced Topic Search (HITS)



#### Acknowledgements

- Slides are based on material provided for *CS276A, Stanford Univ., Text Information Retrieval, Mining, and Exploitation Chr. Manning, P. Raghavan, H. Schütze*
- Thanks also to other lecturers who provided their teaching material on the web



**Community Analysis** 

## SOCIAL NETWORK ANALYSIS



#### Social Network Analysis (SNA)

- Mapping and measuring of relationships and flows between people, groups, organizations, computers or other information/knowledge processing entities.
- The nodes in the network are the people and groups while the links show relationships or flows between the nodes.





#### Kite Network



- Who are connecters or hubs in the network?
- Who has control over what flows in the network?
- Who has best visibility of what is happening in the network?
- Who are peripheral players? Are they important?



#### IM FOCUS DAS LEBEN 7

#### **Measures**

#### 1. Degree Centrality:

The number of direct connections a node has. What really matters is where those connections lead to and how they connect the otherwise unconnected.

$$C_D(n_i) = d(n_i)$$

A node with high betweenness has great influence over what flows in the network indicating important links and single points of failure.

$$C_B(n_i) = \sum_{j < k} g_{jk}(n_i) / g_{jk}$$

#### 3. Closeness Centrality:

ERSITÄT ZU LÜBECK

The measure of closeness of a node to everyone else.

Determined by the sum of the length of the <u>shortest paths</u> between the node and all other nodes in the graph.

$$C_{C}(n_{i}) = \left[\sum_{j=1}^{g} d(n_{i}, n_{j})\right]^{-1}$$

$$C_{C}'(n_{i}) = \frac{g-1}{\sum_{j=1}^{g} d(n_{i}, n_{j})} = (g-1)C_{C}(n_{i})$$

 $C'_D(n_i) = \frac{d(n_i)}{\sigma-1}$ 

 $C'_{B}(n_{i}) = \frac{C_{B}(n_{i})}{(g-1)(g-2)/2}$ 

Undirected graph

Rescaling by dividing through by the number of pairs of nodes not including n<sub>i</sub>

#### Legend

- g = size of graph (number of nodes)
- d(.) = (in)degree
- $g_{ik}$  = number of minimal paths between nodes j and k
- g<sub>jk</sub>(n) = number of minimal paths between nodes j and k that contain n
- (g-1)(g-2)/2 = number of potential paths  $\sum_{x=1}^{u} x = (u+1)u/2 \text{ für } u = (g-2)$
- d(.,.)= distance between two nodes
- Scaling with (g 1)(g 2): For every node n except n<sub>i</sub> pair the node with all other nodes except n and n<sub>i</sub>



### Example: Kite-Network





### Example

	Α	В	C	D	Е
А	0	1	1	0	0
В	1	0	0	1	1
С	1	0	0	1	0
D	0	1	1	0	1
Е	0	1	0	1	0

#### Adjacency

	C <sub>B</sub>	C <sub>C</sub>	C <sub>D</sub>
А	1	1/6	2
В	3	1/5	3
С	1	1/6	2
D	3	1/5	3
Е	0	1/6	2

	А	В	С	D	E
А	0	1	1	2	2
В	1	0	2	1	1
С	1	2	0	1	2
D	2	1	1	0	1
Е	2	1	2	1	0

Distance

	А	В	С	D	Е
А	0	А	А	BC	В
В	В	0	AD	В	В
С	С	AD	0	С	D
D	BC	D	D	0	D
E	В	Е	D	Е	0

Paths



IM FOCUS DAS LEBEN 10

Community Analysis
ANCHOR TEXT



IM FOCUS DAS LEBEN 11

#### The Web as a Directed Graph



**Assumption 1:** A hyperlink between pages denotes author perceived relevance (quality signal)

Assumption 2: The anchor of the hyperlink describes the target page (textual context)



#### **Anchor Text**

RSITÄT ZU LÜBECK

RMATIONSSYSTEM

- For IBM how to distinguish between:
  - IBM's home page (mostly graphical)
  - IBM's copyright page (high term freq. for 'ibm')
  - Rival's spam page (arbitrarily high term freq.)



Oliver A. McBryan. GENVL and WWWW: Tools for Taming the Web. Research explained at First International Conference on the World Wide Web. CERN, Geneva (Switzerland), May 25-26-27 **1994** (WWWW=World Wide Web Worm, first serach engine for the web)

#### Indexing anchor text

• When indexing a document *D*, include anchor text from links pointing to *D*.





#### The Web as a Resource for NLP





Community Analysis
PAGE RANK



#### The Web as a Resource for Ranking

- First generation: using link counts as simple measures of popularity.
- Two basic suggestions:
  - <u>Undirected popularity:</u>
    - Each page gets a score = the number of in-links plus the number of out-links (3+2=5).
  - Directed popularity:
    - Score of a page = number of its in-links (3).





#### Query processing

- First retrieve all pages matching the text query (say *venture capital*).
- Order these by their link popularity (either variant on the previous page).

How to organize for "Search Engine Optimization"?



#### PageRank scoring

- Imagine a browser doing a random walk on web pages:
  - Start at a random page



- At each step, go out of the current page along one of the links on that page, equiprobably
- Each page has a long-term visit rate use this as the page's score



#### Not quite enough

- The web is full of dead-ends.
  - Random walk can get stuck in dead-ends.
  - Makes no sense to talk about long-term visit rates.





## Teleporting / damping

- At a dead end, jump to a random web page.
- At any non-dead end, with probability 10%, jump to a random web page.
  - With remaining probability (90%), go out on a random link.
  - 10% a parameter.
- There is a long-term rate at which any page is visited.
  - How do we compute this visit rate?



#### Markov chains

- A Markov chain consists of *n* states, plus an *n*×*n* <u>transition matrix</u> P.
- At each step, we are in exactly one of the states.
- For 1 ≤ i,j ≤ n, the matrix entry P<sub>ij</sub> tells us the relative frequency of j being the next state, given we are currently in state i.







 $P_{ii} > 0$ 

#### **Ergodic Markov chains**

- A Markov chain is ergodic if
  - you have a path from any state to any other (reducibility)
  - returns to states occur at irregular times (aperiodicity)
  - For any start state, after a finite transient time  $T_0$ , <u>the</u> probability of being in any state at a fixed time  $T>T_0$  is <u>nonzero.</u> (positive recurrence)





#### Ergodic Markov chains

 For any ergodic Markov chain, there is a unique long-term visit rate for each state.

- "Steady-state" distribution.

- Over a long time-period, we visit each state in proportion to this rate.
- It doesn't matter where we start.



#### State vectors

- A (row) vector (state vector) x = (x<sub>1</sub>, ... x<sub>n</sub>) tells us where the walk is at any point.
- E.g., (000...1...000) means we're in state i.

1 i n

More generally, the vector  $\mathbf{x} = (x_1, \dots, x_n)$  means the walk is in state i with relative frequency  $x_i$ .

$$\sum_{i=1}^n x_i = 1.$$



- If the state vector is  $\mathbf{x} = (x_1, \dots, x_n)$  at this step, what is it at the next step?
- Recall that row *i* of the transition matrix P tells us where we go next from state *i*
- So from **x**, our next state is distributed as **xP**.



The steady state looks like a vector of probabilities a = (a<sub>1</sub>, ... a<sub>n</sub>):

 $-a_i$  is the relative frequency that we are in state i.



For this example,  $a_1 = 1/4$  and  $a_2 = 3/4$ .



## How do we compute this vector?

- Let  $\mathbf{a} = (a_{1}, \dots, a_{n})$  denote the row vector of steadystate rates.
- If we our current position is described by **a**, then the next step is described as **aP**.
- But **a** is the steady state, so **a**=**aP**.
- Solving this matrix equation gives us **a**.
  - So **a** is the (left) eigenvector for **P**.
  - (Corresponds to the "principal" eigenvector of P with the largest eigenvalue)
  - Transition matrices always have largest eigenvalue 1.



#### Eigenvectors and Eigenvalues $Mx = \lambda x$







[Wikipedia]

#### One way of computing a

- Recall, regardless of where we start, we eventually reach the steady state a.
- Start with any distribution (say **x**=(*10...0*)).
- After one step, we're at **xP**;
- after two steps at **xP**<sup>2</sup>, then **xP**<sup>3</sup> and so on.
- "Eventually" means for "large" k,  $xP^k = a$ .
- Algorithm: multiply x by increasing powers of P until the product looks stable.



#### PageRank Summary

- Preprocessing:
  - Given graph of links, build matrix P
  - From it compute **a**
  - The entry a<sub>i</sub> is a number between 0 and 1: the pagerank of page *i*.
- Query processing:
  - Retrieve pages meeting query
  - Rank them by their pagerank
  - Order is query-*independent*
- A variant of PageRank is used in Google, but also many other clever heuristics



#### PageRank: Issues and Variants

- How realistic is the random surfer model?
  - What if we modeled the back button?
  - Surfer behavior sharply skewed towards short paths
  - Search engines, bookmarks & directories make jumps non-random
- Biased Surfer Models
  - Weight edge traversal probabilities based on match with topic/query (non-uniform edge selection)
  - Bias jumps to pages on topic (e.g., based on personal bookmarks & categories of interest)



#### Google PageRank

- Links are also weighted according to the importance of the source node
  - Page C has a higher
     PageRank than Page E,
     even though there
     are fewer links to C;
     the one link to C
     comes from an
     important page
     and hence is
     of high value.





[Wikipedia]

**Community Analysis** 

## HYPERLINK-INDUCED TOPIC SEARCH



### Hyperlink-Induced Topic Search (HITS)

- In response to a query, instead of an ordered list of pages each matching the query, find <u>two</u> sets of inter-related pages:
  - Hub pages are good lists of links on a subject
    - e.g., "Bob's list of cancer-related links."
  - Authority pages occur recurrently on good hubs for the subject
- Best suited for "broad topic" queries rather than for page-finding queries



#### Hubs and Authorities

- Thus, a good hub page for a topic *points* to many authoritative pages for that topic
- A good authority page for a topic is *pointed* to by many good hubs for that topic
- Circular definition will turn this into an iterative computation



#### The hope





- Extract from the web a <u>base set</u> of pages that could be good hubs or authorities
- From these, identify a small set of top hub and authority pages;

 $\rightarrow$ iterative algorithm



#### Base set

- Given text query (say "*browser*"), use a text index to get all pages containing "*browser*"
  - Call this the root set of pages
- Add in any page that either
  - points to a page in the root set, or
  - is pointed to by a page in the root set
- Call this the base set



#### Visualization





#### Assembling the base set

- Root set typically 200-1000 nodes
- Base set may have up to 5000 nodes
- How do you find the base set nodes?
  - Follow out-links by parsing root set pages
  - Get in-links (and out-links) from a connectivity server
  - Actually, suffices to text-index strings of the form href="<u>URL</u>" to get in-links to <u>URL</u>



### Distilling hubs and authorities

Compute, for each page x in the base set, a hub score h(x) and an authority score a(x)

- Initialize: for all x,  $h(x) \leftarrow 1$ ;  $a(x) \leftarrow 1$ ;
- Iteratively update all h(x), a(x);

$$h(x) \leftarrow \sum_{x \mapsto y} a(y)$$
$$a(x) \leftarrow \sum_{y \mapsto x} h(y)$$



- After iterations output pages with
  - highest h() scores as top hubs
  - highest a() scores as top authorities

### Scaling

- To prevent the *h()* and *a()* values from getting too big, can scale down after each iteration
- Scaling factor doesn't really matter:
  - we only care about the *relative* values of the scores



#### How many iterations?

- Claim: relative values of scores will converge after a few iterations:
  - In fact, suitably scaled, *h()* and *a()* scores settle into a steady state!
- We only require the relative orders of the h() and a() scores - not their absolute values
- In practice, ~5 iterations get you close to stability



#### Deeper understanding

- Determination of conditional dependencies between attribute values of objects in a social network
- Identification of clusters
- . . .



# Data Models vs. Algorithmic Models

Data Modeling	VS.	Algorithmic Modeling
Y ← F( X, random noise, parameters	Y ← Black Box ← X ↑ Random Forests	
We understand the world How well 'my data model' works	Data Science	We don't understand the world The world produces data in a black-box
Linear Regression Logistic Regression Known Distributions Confidence Intervals Predictor Variables & Goodness of Fit	AI	Machine Learning, Al Random Forests, SVM Unknown Multivariate Distributions Iterative Predictive Accuracy

