Note to other teachers and users of these slides: We would be delighted if you found this our material useful in giving your own lectures. Feel free to use these slides verbatim, or to modify them to fit your own needs. If you make use of a significant portion of these slides in your own lecture, please include this message, or a link to our web site: <u>http://www.mmds.org</u>

Analysis of Large Graphs: TrustRank and WebSpam

Mining of Massive Datasets Jure Leskovec, Anand Rajaraman, Jeff Ullman Stanford University http://www.mmds.org



Example: PageRank Scores



J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org

Random Teleports ($\beta = 0.8$)



PageRank: The Complete Algorithm

Input: Graph G and parameter β

- Directed graph G with spider traps and dead ends
- Parameter β
 Output: PageRank vector r

• Set:
$$r_j^{(0)} = \frac{1}{N}, t = 1$$

do:

•
$$\forall j: \mathbf{r}'_{j}^{(t)} = \sum_{i \to j} \boldsymbol{\beta} \; \frac{r_{i}^{(t-1)}}{d_{i}}$$

 $\mathbf{r}'_{j}^{(t)} = \mathbf{0} \; \text{if in-degree of } \mathbf{j} \text{ is } \mathbf{0}$

Now re-insert the leaked PageRank:

$$\forall j: r_j^{(t)} = r'_j^{(t)} + \frac{1-S}{N}$$

$$t = t + 1$$

where:
$$S = \sum_{j} r'^{(t)}_{j}$$

• while
$$\sum_{j} \left| r_{j}^{(t)} - r_{j}^{(t-1)} \right| > \varepsilon$$

If the graph has no deadends then the amount of leaked PageRank is $1-\beta$. But since we have dead-ends the amount of leaked PageRank may be larger. We have to explicitly account for it by computing **S**.

N.B. S = (1-P) 181 G has no deadends

Some Problems with PageRank

- Measures generic popularity of a page
 - Will ignore/miss topic-specific authorities
 - Solution: Topic-Specific PageRank (next)
- Uses a single measure of importance
 - Other models of importance
 - Solution: Hubs-and-Authorities
- Susceptible to Link spam
 - Artificial link topographies created in order to boost page rank
 - Solution: TrustRank

Topic-Specific PageRank

Topic-Specific PageRank

- Instead of generic popularity, can we measure popularity within a topic?
- Goal: Evaluate Web pages not just according to their popularity, but by how close they are to a particular topic, e.g. "sports" or "history"
- Allows search queries to be answered based on interests of the user
 - Example: Query "Trojan" wants different pages depending on whether you are interested in sports, history and computer security

Topic-Specific PageRank

- Random walker has a small probability of teleporting at any step
- Teleport can go to:
 - Standard PageRank: Any page with equal probability
 - To avoid dead-end and spider-trap problems
- Topic Specific PageRank: A topic-specific set of "relevant" pages (teleport set) = \$ = \$ (TOPK)
 Idea: Bias the random walk
 - When walker teleports, she pick a page from a set S
 - S contains only pages that are relevant to the topic
 - E.g., Open Directory (DMOZ) pages for a given topic/query
 - For each teleport set S, we get a different vector r_s

Matrix Formulation

To make this work all we need is to update the teleportation part of the PageRank formulation:

$$A_{ij} = \begin{cases} \beta M_{ij} + (1 - \beta) / |S| & \text{if } i \in S = S(\text{ToPic } S) \\ \beta M_{ij} + 0 & \text{otherwise} \end{cases}$$

- A is stochastic!
- We weighted all pages in the teleport set S equally
 - Could also assign different weights to pages!
- Compute as for regular PageRank:
 - Multiply by *M*, then add a vector
 - Maintains sparseness

Example: Topic-Specific PageRank



J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org

Discovering the Topic Vector S

Create different PageRanks for different topics

- The 16 DMOZ top-level categories:
 - arts, business, sports,...

Which topic ranking to use?

- User can pick from a menu
- Classify query into a topic
- Can use the context of the query
 - E.g., query is launched from a web page talking about a known topic
 - History of queries e.g., "basketball" followed by "Jordan"
- User context, e.g., user's bookmarks, ...

Application to Measuring Proximity in Graphs

Random Walk with Restarts: S is a single element

[Tong-Faloutsos, 'o6]

Proximity on Graphs



Good proximity measure?



No effect of degree-1 nodes (E, F, G)!
 Multi-faceted relationships

Good proximity measure?

Network flow is not good:



Does not punish long paths

[Tong-Faloutsos, '06] What is good notion of proximity?



It must be seusihive to:

- Multiple connections
- Quality of connection
 - Direct & Indirect
 - connections
 - Length, Degree,
 - Weight...

SimRank: Idea (For K-PARTITE GRAPHS)

SimRank: Random walks from a fixed node on **k**-partite graphs Conferences Tags **Authors** H11 B1 Setting: k-partite graph * H12 B2 with **k** types of nodes B3 H21 E.g.: Authors, Conferences, Tags
Topic Specific PageRank A2 • B4 H22 **B**5 H23 · Apply T-Page-Rank. from node u: teleport set $S = \{u\}$ Resulting scores measures similarity to node u **Problem:** Must be done once for each node u Suitable for sub-Web-scale applications

SimRank: Example



Q: What is most related conference to ICDM? ^ Arply A: Topic-Specific PageRank with teleport set S={ICDM}

SimRank: Example



PageRank: Summary

"Normal" PageRank:

- Teleports uniformly at random to any node ([3, (1-p))
- Topic-Specific PageRank also known as Personalized PageRank:
 - Teleports to a topic specific set of pages
 - Nodes can have different probabilities of surfer landing there: S = [0.1, 0, 0, 0.2, 0, 0, 0.5, 0, 0, 0.2] Random Walk with Restarts:
 - Topic-Specific PageRank where teleport is always to the same node. S=[0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0]
 NODE CROXING

TrustRank: Combating the Web Spam

What is Web Spam?

Spamming:

 Any deliberate action to boost a web page's position in search engine results, incommensurate with page's real value

Spam:

- Web pages that are the result of spamming
- This is a very broad definition
 - **SEO** industry might disagree!
 - SEO = search engine optimization
- Approximately 10-15% of web pages are spam

Web Search

Early search engines:

- Crawl the Web
- Index pages by the words they contained
- Respond to search queries (lists of words) with the pages containing those words

Early page ranking:

- Attempt to order pages matching a search query by "importance" - Corulating Score
 - First search engines considered:
 - (1) Number of times query words appeared
 - (2) Prominence of word position, e.g. title, header

 As people began to use search engines to find things on the Web, those with commercial interests tried to exploit search engines to bring people to their own site – whether they wanted to be there or not

Example:

- Shirt-seller might pretend to be about "movies"
- Techniques for achieving high relevance/importance for a web page

First Spammers: Term Spam

- How do you make your page appear to be about movies?
 - (1) Add the word movie 1,000 times to your page
 - Set text color to the background color, so only search engines would see it
 - (2) Or, run the query "movie" on your target search engine
- See what page came first in the listings
 Copy it into your page, make it "invisible"
 These and similar techniques are term spam

Google's Solution to Term Spam

- Believe what people say about you, rather than what you say about yourself
 - Use words in the anchor text (words that appear underlined to represent the link) and its surrounding text
- PageRank as a tool to measure the "importance" of Web pages

Why It Works?

Our hypothetical shirt-seller looses

- Saying he is about movies doesn't help, because others don't say he is about movies
- His page isn't very important, so it won't be ranked high for shirts or movies

Example:

- Shirt-seller creates 1,000 pages, each links to his with "movie" in the anchor text
- These pages have no links in, so they get little PageRank
- So the shirt-seller can't beat truly important movie pages, like IMDB

Why it does not work?

Google	Web	Images	Groups	News	Froogle	Local	more	e »
	miserable failure						Search Advanced Search Preferences	

Web

Results 1 - 10 of about 969,000 for miserable failure. (0.06 seconds)

Biography of President George W. Bush

Biography of the president from the official White House web site. www.whitehouse.gov/president/gwbbio.html - 29k - <u>Cached</u> - <u>Similar pages</u> <u>Past Presidents</u> - <u>Kids Only</u> - <u>Current News</u> - <u>President</u> <u>More results from www.whitehouse.gov »</u>

Welcome to MichaelMoore.com!

Official site of the gadfly of corporations, creator of the film Roger and Me and the television show The Awful Truth. Includes mailing list, message board, ... www.michaelmoore.com/ - 35k - Sep 1, 2005 - Cached - Similar pages

BBC NEWS | Americas | 'Miserable failure' links to Bush

Web users manipulate a popular search engine so an unflattering description leads to the president's page. news.bbc.co.uk/2/hi/americas/3298443.stm - 31k - Cached - Similar pages

Google's (and Inktomi's) Miserable Failure

A search for **miserable failure** on Google brings up the official George W. Bush biography from the US White House web site. Dismissed by Google as not a ... searchenginewatch.com/sereport/article.php/3296101 - 45k - Sep 1, 2005 - Cached - Similar pages



Google vs. Spammers: Round 2!

- Once Google became the dominant search engine, spammers began to work out ways to fool Google
- Spam farms were developed to concentrate
 PageRank on a single page
- Link spam:
- SUBGRAPHS
- Creating link structures that boost PageRank of a particular page



Link Spamming

- Three kinds of web pages from a spammer's point of view :
 - Inaccessible pages
 - Accessible pages
 - e.g., blog comments pages
 - spammer can post links to his pages
 - Owned pages
 - Completely controlled by spammer
 - May span multiple domain names

Link Farms

Spammer's goal:

Maximize the PageRank of target page t

Technique:

 Get as many links from accessible pages as possible to target page *t*

 Construct "link farm" to get PageRank multiplier effect

Link Farms







J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org





Multiplier effect for acquired PageRank for the
By making *M* large, we can make *y* as presence large as we want

TrustRank: Combating the Web Spam

Combating Spam

Combating term spam

- Analyze text using statistical methods
- Similar to email spam filtering
- Also useful: Detecting approximate duplicate pages

Combating link spam



- Leads to another war hiding and detecting spam farms
- TrustRank = topic-specific PageRank with a teleport set of trusted pages

Example: .edu domains, similar domains for non-US schools. Teleport is not any more "UNIFORH" over all pages PP

Le Owned

SUBGRAPHS

TrustRank: Idea

- Basic principle: Approximate isolation
 - It is rare for a "good" page to point to a "bad" (spam) page
- Sample a set of seed pages from the web
- Have an oracle (human) to identify the good pages and the spam pages in the seed set
 - Expensive task, so we must make seed set as small as possible

Trust Propagation

- Call the subset of seed pages that are identified as good the trusted pages
- Perform a topic-sensitive PageRank with TOPIC Page-RANK teleport set = trusted pages Apply
 - **Propagate trust through links:**
 - Each page gets a trust value between 0 and 1

Solution 1: Use a threshold value and mark all pages below the trust threshold as spam

STRUST

SE

Simple Model: Trust Propagation

- Set trust of each trusted page to 1
- Suppose trust of page *p* is *t_p*
 - Page *p* has a set of out-links *o_p*
- For each *q* ∈ *o_p*, *p* confers the trust to *q*
 - $\beta t_p / |o_p|$ for $0 < \beta < 1$
- Trust is additive
 - Trust of *p* is the sum of the trust conferred on *p* by all its in-linked pages
- Note similarity to Topic-Specific PageRank
 - Within a scaling factor, TrustRank = PageRank with trusted pages as teleport set

Why is it a good idea?

Trust attenuation:

The degree of trust conferred by a trusted page decreases with the distance in the graph

Trust splitting:

- The larger the number of out-links from a page, the less scrutiny the page author gives each outlink
- Trust is split across out-links

Picking the Seed Set

Two conflicting considerations:

- Human has to inspect each seed page, so seed set must be as small as possible
- Must ensure every good page gets adequate trust rank, so need make all good pages reachable from seed set by short paths

Approaches to Picking Seed Set

- Suppose we want to pick a seed set of k pages
- How to do that?
- (1) PageRank:
 - Pick the top k pages by PageRank
 - Theory is that you can't get a bad page's rank really high
- (2) Use trusted domains whose membership
 is controlled, like .edu, .mil, .gov

- In the TrustRank model, we start with good pages and propagate trust
- Complementary view: What fraction of a page's PageRank comes from spam pages?
- In practice, we don't know all the spam pages, so we need to estimate

Trusted

set

Web

Spam Mass Estimation : TWO Scores

Solution 2:

- r_p = PageRank of page p
 r⁺_p = PageRank of p with teleport into trusted pages only
- Then: What fraction of a page's PageRank comes from spam pages?

•
$$r_p^- = r_p - r_p^+$$

- Spam mass of $p = \frac{r_p^-}{r_p}$
 - Pages with high spam mass are spam.

Trusted

set

Web

HITS: Hubs and Authorities

HITS (Hypertext-Induced Topic Selection)

- Is a measure of importance of pages or documents, similar to PageRank
- Proposed at around same time as PageRank ('98)
- Goal: Say we want to find good newspapers
 - Don't just find newspapers. Find "experts" people who link in a coordinated way to good newspapers
- Idea: Links as votes
 - Page is more important if it has more links
 - In-coming links? Out-going links?

Finding newspapers

Hubs and Authorities

Each page has 2 scores:

- Quality as an expert (hub):
 - Total sum of votes of authorities pointed to
- Quality as a content (authority):
 - Total sum of votes coming from experts





Interesting pages fall into two classes:

- 1. Authorities are pages containing useful information
 - Newspaper home pages
 - Course home pages
 - Home pages of auto manufacturers
- 2. Hubs are pages that link to authorities
 - List of newspapers
 - Course bulletin
 - List of US auto manufacturers



Counting in-links: Authority



(Note this is idealized example. In reality graph is not bipartite and each page has both the hub and authority score)

Counting in-links: Authority



(Note this is idealized example. In reality graph is not bipartite and each page has both the hub and authority score)

Expert Quality: Hub



(Note this is idealized example. In reality graph is not bipartite and each page has both the hub and authority score)

J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org

Reweighting



J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org

Mutually Recursive Definition

- A good hub links to many good authorities
- A good authority is linked from many good hubs
- Model using two scores for each node:
 - Hub score and Authority score
 - Represented as vectors h and a

Each page *i* has 2 scores: Authority score: a_i Hub score: h_i **HITS algorithm:** $a_i = \sum h_j$ • Initialize: $a_i^{(0)} = 1/\sqrt{N}$, $h_i^{(0)} = 1/\sqrt{N}$ Then keep iterating until convergence: • $\forall i$: Authority: $a_i^{(t+1)} = \sum_{i \to i} h_i^{(t)}$ • $\forall \mathbf{i}$: Hub: $h_i^{(t+1)} = \sum_{\mathbf{i} \to \mathbf{i}} a_i^{(t)}$ • $\forall i$: Normalize: $h_i = \sum a_j$ $\sum_{i} \left(a_{i}^{(t+1)} \right)^{2} = 1, \sum_{j} \left(h_{j}^{(t+1)} \right)^{2} = 1$ J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org 55

- HITS converges to a single stable point
- Notation:
 - Vector $a = (a_1 ..., a_n), \quad h = (h_1 ..., h_n)$
 - Adjacency matrix A (NxN): $A_{ij} = 1$ if $i \rightarrow j$, 0 otherwise
- Then $h_i = \sum_{i \to j} a_j$ can be rewritten as $h_i = \sum_j A_{ij} \cdot a_j$
 - So: $h = A \cdot a$
- Similarly, $a_i = \sum_{j \to i} h_j$

can be rewritten as $a_i = \sum_j A_{ji} \cdot h_j = A^T \cdot h$

HITS algorithm in vector notation:

• Set:
$$a_i = h_i = \frac{1}{\sqrt{n}}$$

Repeat until convergence:

•
$$h = A \cdot a$$

•
$$\boldsymbol{a} = \boldsymbol{A}^T \cdot \boldsymbol{h}$$

• Normalize a and h• Then: $a = A^T \cdot (A \cdot a)$ • Normalize a and h• Normalize a and h Convergence criterion: $\sum_{i} \left(h_{i}^{(t)} - h_{i}^{(t-1)} \right)^{2} < \varepsilon$ $\sum_{i} \left(a_{i}^{(t)} - a_{i}^{(t-1)} \right)^{2} < \varepsilon$

a is updated (in 2 steps): $a = A^T (A \ a) = (A^T A) a$ *h* is updated (in 2 steps): $h = A (A^T h) = (A \ A^T) h$

Repeated matrix powering

Existence and Uniqueness

- h = λ A a
- $\bullet a = \mu A^T h$
- $h = \lambda \mu A A^{T} h$
- $= a = \lambda \mu A^T A a$

$$\lambda = 1 / \sum h_i$$

$$\mu = 1 / \sum a_i$$

- Under reasonable assumptions about A, HITS converges to vectors h* and a*:
 - h^{*} is the principal eigenvector of matrix A A^T
 - *a*^{*} is the principal eigenvector of matrix A^TA

Example of HITS

$$\mathbf{A} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \qquad \mathbf{A}^{\mathrm{T}} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}$$



h(yahoo)	=	.58	.80	.80	.79	• • •	.788
h(amazon)	=	.58	.53	.53	.57	• • •	.577
h(m'soft)	=	.58	.27	.27	.23	•••	.211
a(yahoo)	=	.58	.58	.62	.62	•••	.628
a(amazon)	=	.58	.58	.49	.49	• • •	.459
a(m'soft)	=	.58	.58	.62	.62	• • •	.628

PageRank and HITS

- PageRank and HITS are two solutions to the same problem:
 - What is the value of an in-link from u to v?
 - In the PageRank model, the value of the link depends on the links into u
 - In the HITS model, it depends on the value of the other links out of u

The destinies of PageRank and HITS post-1998 were very different