# Data-Intensive Distributed Computing
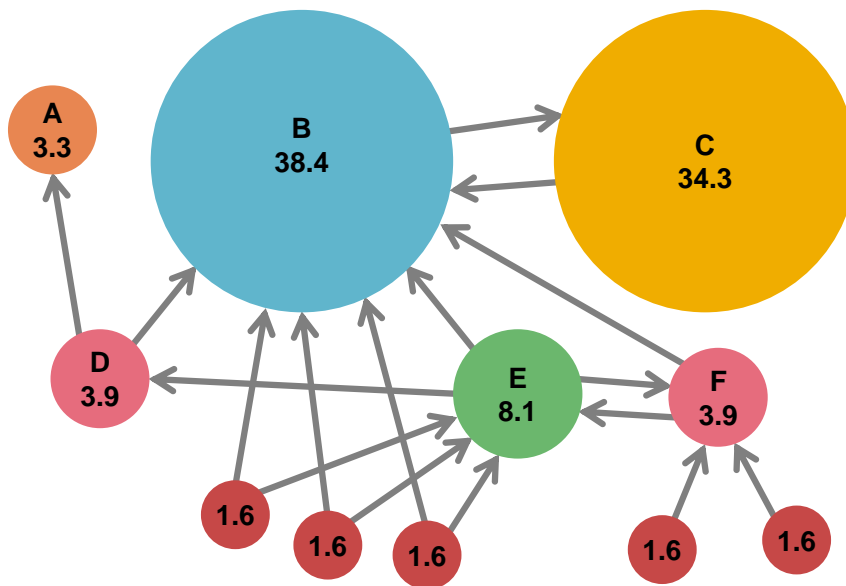
CS 431/631 451/651 (Fall 2021)

## Part 10b: Analyzing Graphs, Redux

Ali Abedi

Thanks to Jure Leskovec, Anand Rajaraman, Jeff Ullman (Stanford University)
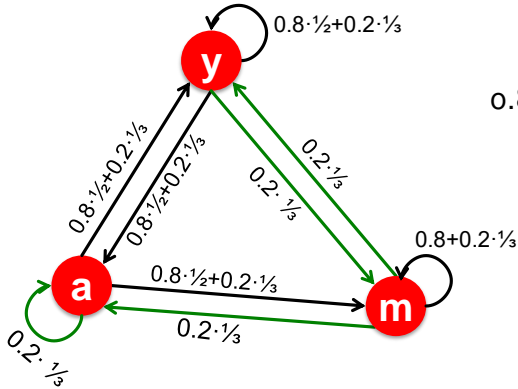
These slides are available at https://www.student.cs.uwaterloo.ca/~cs451

# Example: PageRank Scores

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

# Random Teleports (β = 0.8)



**M**

$$0.8 \begin{bmatrix} 1/2 & 1/2 & 0 \\ 1/2 & 0 & 0 \\ 0 & 1/2 & 1 \end{bmatrix} + 0.2 \begin{bmatrix} 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \end{bmatrix}$$

**[1/N]$_{NxN}$**

|   | | | |
|---|---|---|---|
| y | 7/15 | 7/15 | 1/15 |
| a | 7/15 | 1/15 | 1/15 |
| m | 1/15 | 7/15 | 13/15 |

**A**

| | | | | | | |
|---|---|---|---|---|---|---|
| y | | 1/3 | 0.33 | 0.24 | 0.26 | 7/33 |
| a | = | 1/3 | 0.20 | 0.20 | 0.18 | . . . 5/33 |
| m | | 1/3 | 0.46 | 0.52 | 0.56 | 21/33 |

**r   =        A r**

**Equivalently:** $r = \beta \, M \cdot r + \left[ \frac{1-\beta}{N} \right]_N$

3

# Some Problems with PageRank

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

# 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

Photo credit to Huffington Post

# 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"
  - **First search engines considered:**
    - **(1)** Number of times query words appeared
    - **(2)** Prominence of word position, e.g. title, header

# First Spammers

- 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**

Do not forget to tell the joke "I used to do this when I was young" when talking about point (1) :D

# 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

SPAM FARMING

# 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:**
  - 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



**One of the most common and effective organizations for a link farm**

# Analysis



Accessible     Owned

Inaccessible

t

1
2

M

N…# pages on the web
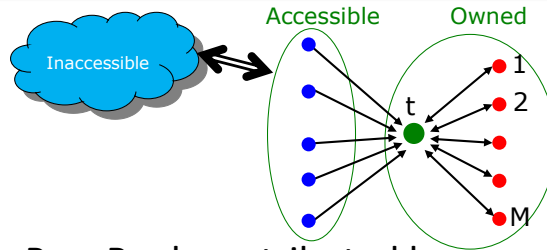M…# of pages spammer owns

- **x**: PageRank contributed by accessible pages
- **y**: PageRank of target page **t**
- Rank of each "farm" page $= \frac{\beta y}{M} + \frac{1-\beta}{N}$
- $y = x + \beta M \left[ \frac{\beta y}{M} + \frac{1-\beta}{N} \right] + \frac{1-\beta}{N}$

  $= x + \beta^2 y + \frac{\beta(1-\beta)M}{N} + \boxed{\frac{1-\beta}{N}}$     Very small; ignore
  
  Now we solve for **y**
- $y = \frac{x}{1-\beta^2} + c\frac{M}{N}$    where $c = \frac{\beta}{1+\beta}$

# Analysis



Accessible  Owned

Inaccessible

t

1
2

M

N…# pages on the web
M…# of pages spammer owns

- $y = \dfrac{x}{1-\beta^2} + c\,\dfrac{M}{N}$   where $c = \dfrac{\beta}{1+\beta}$
- For $\beta = 0.85$, $1/(1-\beta^2)= 3.6$

- Multiplier effect for acquired PageRank
- By making **M** large, we can make **y** as **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**
  - **Detection and blacklisting of structures that look like spam farms**
    - 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

# 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 **teleport set = trusted pages**
  - **Propagate trust through links:**
    - Each page gets a trust value between **0** and **1**

- <u>**Solution 1:**</u> **Use a threshold value and mark all pages below the trust threshold as spam**

# 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 $q \in o_p$, $p$ **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 Personalized 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 out-link
  - Trust is **split** across out-links
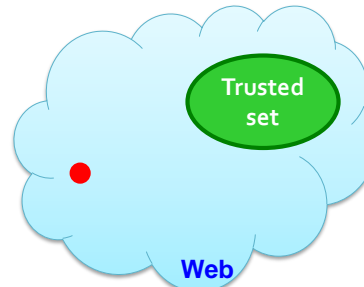
# Picking the Seed Set

- **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

# Spam Mass

- 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
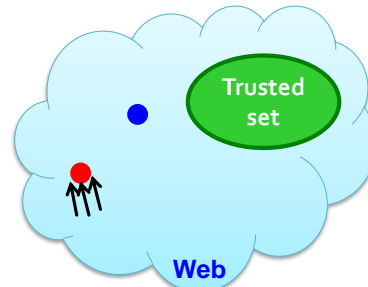


Trusted set

Web

# Spam Mass Estimation

**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 $p$ = $\dfrac{r_p^-}{r_p}$**

  - Pages with high spam mass are spam.



Trusted set

Web

The End!

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)**
- **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
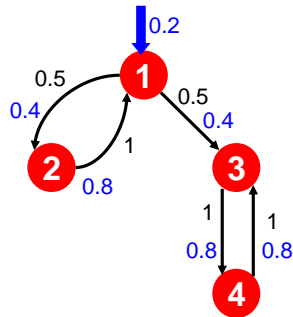  - For each teleport set $S$, we get a different vector $r_S$

# Matrix Formulation

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

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

  - *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



Suppose **S = {1}**, **β = 0.8**

| Node | Iteration | | | |
|------|-----------|-----|------|--------|
|      | 0         | 1   | 2 …  | stable |
| 1    | 0.25      | 0.4 | 0.28 | 0.294  |
| 2    | 0.25      | 0.1 | 0.16 | 0.118  |
| 3    | 0.25      | 0.3 | 0.32 | 0.327  |
| 4    | 0.25      | 0.2 | 0.24 | 0.261  |

**S={1}, β=0.90:**
r=[0.17, 0.07, 0.40, 0.36]
**S={1} , β=0.8:**
r=[0.29, 0.11, 0.32, 0.26]
**S={1}, β=0.70:**
r=[0.39, 0.14, 0.27, 0.19]

**S={1,2,3,4}, β=0.8:**
r=[0.13, 0.10, 0.39, 0.36]
**S={1,2,3} , β=0.8:**
r=[0.17, 0.13, 0.38, 0.30]
**S={1,2} , β=0.8:**
r=[0.26, 0.20, 0.29, 0.23]
**S={1} , β=0.8:**
r=[0.29, 0.11, 0.32, 0.26]