

CMSC5741 Big Data Tech. & Apps.

Lecture Assignment Project Exam Help

k Analysis

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Prof. Michael R. Lyu

Computer Science & Engineering Dept.
The Chinese University of Hong Kong

What's the Mechanism Distinguished Google?

Assignment Project Exam Help

How Google return such kind of rankings (e.g. Charles Kao Wikipedia first then his NobelPrize. Even his passing away)?

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

One important factor is PageRank score.

How to make PageRank computation scalable

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

- There are **billions** of web pages and hyperlinks between them, how to compute their ranking score (e.g., PageRank) **efficiently**?

Outline

- Web as a Graph
- PageRank [Assignment](#) [Project](#) [Exam](#) [Help](#)
- Topic-Specific <https://eduassistpro.github.io/>
- Appendix: Trust-Rank [Add WeChat edu_assist_pro](#)

Outline

- Web as a Graph
- PageRank Assignment Project Exam Help
- Topic-Specific <https://eduassistpro.github.io/>
- Appendix: Trust-Rank Add WeChat edu_assist_pro

Web as a Graph

- Web as a directed graph:
 - Nodes: Web pages
Assignment Project Exam Help
 - Edges: Hype
<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro

Web as a Directed Graph

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Broad Question

- How to **organize** the Web?
- First try: Human curated **Web directories**
 - Yahoo, DM <https://eduassistpro.github.io/>
- Second try: **Web Search**
 - Information Retrieval inv : Find relevant docs in a small and trusted set
 - Newspaper articles, Patents, etc.
 - But: Web is **huge**, full of untrusted documents, random things, web spam, etc.

Web Search: 2 Challenges

- Web contains many sources of information:
Who to “trust”?
Assignment Project Exam Help
– Trick: Trust points to each other
- What is the “<https://eduassistpro.github.io/> query”?
“newspaper”? Add WeChat edu_assist_pro
– No single right answer
– Trick: Pages that actually know about newspapers might all be pointing to many newspapers

Ranking Nodes on the Graph

- All web pages are **not equally** “important”
 - www.cuhk. Assignment Project Exam Help
www.joe-sc <https://eduassistpro.github.io/>
- There is large ~~Adversity~~ **iversity** in web-graph node connectivity.
 - Rank the pages by the link structure!

Link Analysis Algorithms

- We will cover the following **Link Analysis approaches** for computing importance of nodes in a graph
 - PageRank <https://eduassistpro.github.io/>
 - Topic-Specific (Personalized) Rank [Add WeChat edu_assist_pro](https://eduassistpro.github.io/)
 - Web Spam Detection Algorithms, e.g. TrustRank

Outline

- Web as a Graph
- PageRank
Assignment Project Exam Help
- Topic-Specific
<https://eduassistpro.github.io/>
- Appendix: Trust-Rank
Add WeChat edu_assist_pro

Links as Votes

- Idea: Links as votes
 - Page is more important if it has more links
Assignment Project Exam Help
 - In-coming
- Think of in-links
 - <https://eduassistpro.github.io/>
 - www.cuhk.edu.hk has 1 in-links
 - www.joe-schmoe.com has 1 in-link
- Are all in-links equal?
 - Link from important pages count more
 - Recursive question!

Example: PageRank Scores

Assignment Project Exam Help

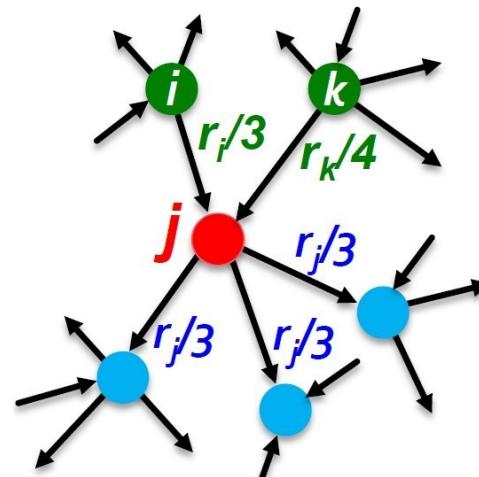
<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Simple Recursive Formulation

- Each link's vote is proportional to the **importance** of its source page
[Assignment](#) [Project](#) [Exam](#) [Help](#)
- If page j with n out-links, each link get
- Page j 's own importance is the sum of the votes on its in-links

$$r_j = r_i/3 + r_k/4$$

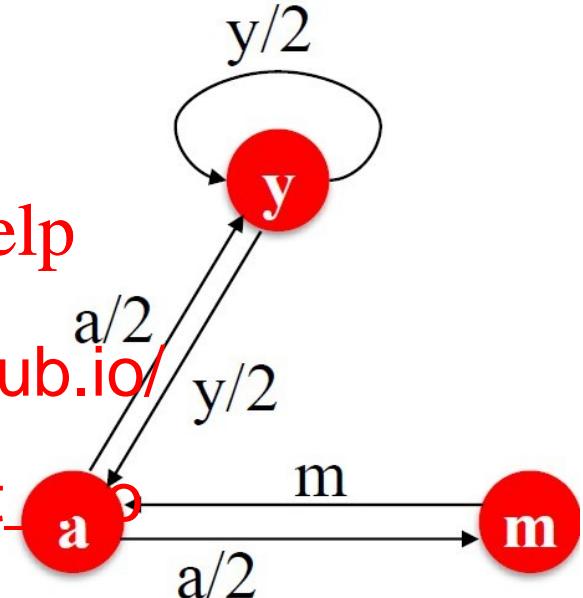


PageRank: The “Flow” Model

- A “vote” from an important page is worth more
- A page is important if it is pointed to by important pages
- Define a “rank” for page j

$$r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$$

d_i . . . out-degree of node i



“Flow” equations:

$$\mathbf{r}_y = \mathbf{r}_y/2 + \mathbf{r}_a/2$$

$$\mathbf{r}_a = \mathbf{r}_y/2 + \mathbf{r}_m$$

$$\mathbf{r}_m = \mathbf{r}_a/2$$

Solving the Flow Equations

- 3 equations, 3 unknowns,
no constants
 - No unique solution
 - All solutions <https://eduassistpro.github.io/factor>
- Additional constraint for ~~Add WeChat edu_assistness:~~
- $r_y + r_a + r_m = 1$
- Solution: $r_y = \frac{2}{5}, r_a = \frac{2}{5}, r_m = \frac{1}{5}$
- Gaussian Elimination method works for small examples, but we need a better method for ~~large web-size graphs~~

PageRank: Matrix Formulation

- Stochastic adjacency matrix M
 - Let page i have d_i out-links
[Assignment](#) [Project](#) [Exam](#) [Help](#)
 - If $i \rightarrow j$, then
 - M is a column vector with sum to 1
- Rank vector r : vector with one entry per page
 - r_i is the importance score of page i
 - $\sum_i r_i = 1$
- The flow equations can be written
$$r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$$
$$r = M \cdot r$$

Example

- Remember the flow equation: $r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$
- Flow equation in the matrix form:
Assignment Project Exam Help

<https://eduassistpro.github.io/>

- Suppose page i links to 3 pages including j
Add WeChat edu_assist_pro

Eigenvector Formulation

- The flow equations can be written $r = M \cdot r$
- So the rank vector r is an eigenvector of the stochastic w
 - In fact, its first component is the largest eigenvalue of M since M is column stochastic
- We can now efficiently solve for r !
 - The method is called Power iteration.

NOTE: x is an eigenvector with the corresponding eigenvalue λ if:
 $Ax = \lambda x$

Example: Flow Equation & M

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Power Iteration Method

- Given a web graph with n nodes, where the nodes are **pages** and edges are **hyperlinks**
Assignment Project Exam Help
- Power Iterative scheme
<https://eduassistpro.github.io/>
 - Suppose th
 - Initialize:
$$r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$$
 Add WeChat edu_assist_pro
 - Iterate: $\mathbf{r}^{(t+1)} = \mathbf{M} \cdot \mathbf{r}^{(t)}$ $d_i \dots$ out-degress of node i
 - Stop when $|\mathbf{r}^{(t+1)} - \mathbf{r}^{(t)}|_1 < \epsilon$
 - $|\mathbf{x}|_1 = \sum_{1 \leq i \leq N} |x_i|$ is the L1 norm

Why Power Iteration Works?

- Power iteration:
 - A method for finding dominant eigenvector (the largest eigenvalue)
Assignment Project Exam Help
 - <https://eduassistpro.github.io/>
 - Add WeChat edu_assist_pro
- Claim:
 - Sequence $M \cdot r^{(0)}, M^2 \cdot r^{(0)}, \dots, M^k \cdot r^{(0)}, \dots$ approaches the dominant eigenvector of M

Why Power Iteration Works?

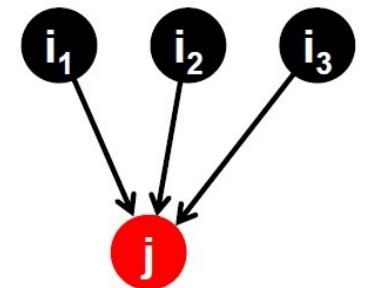
- Proof:
 - Assume M has n linearly independent eigenvectors x_1, x_2, \dots, x_n with corresponding eigenvalues <https://eduassistpro.github.io/>
 - Vectors x_1, x_2, \dots, x_n form a basis thus we can write:
 - $$\begin{aligned} Mr^{(0)} &= M(c_1x_1 + c_2x_2 + \dots + c_nx_n) \\ &= c_1(Mx_1) + c_2(Mx_2) + \dots + c_n(Mx_n) \\ &= c_1(\lambda_1x_1) + c_2(\lambda_2x_2) + \dots + c_n(\lambda_nx_n) \end{aligned}$$
 - Repeated multiplication on both sides:
$$M^k r^{(0)} = c_1(\lambda_1^k x_1) + c_2(\lambda_2^k x_2) + \dots + c_n(\lambda_n^k x_n)$$

Why Power Iteration Works?

- Proof: (cont.)
 - Repeated multiplication on both sides produces
[Assignment Project Exam Help](https://eduassistpro.github.io/)
 - Since $\lambda_1 > \lambda_2$ then
[Add WeChat edu_assist_pro](#)
 - Thus $M^k r^{(0)} \approx c_1(\lambda_1^k x_1)$
 - Note if $c_1 = 0$, then the method won't converge

Random Walk Interpretation

- Imagine a random web surfer:
 - At any time t , surfer is on some page i
 - At time $t+1$, the surfer follows an out-link from i uniformly at <https://eduassistpro.github.io/>
 - Ends up on some page j
 - Process repeats indefinitely



- Let:
 - $p(t)$... vector whose i^{th} coordinate is the prob. that the surfer is at page i at time t
 - So $p(t)$ is a probability distribution over pages

$$r_j = \sum_{i \rightarrow j} \frac{r_i}{d_{\text{out}}(i)}$$

The Stationary Distribution

- Where is the surfer at time $t+1$
 - Follows a link uniformly at random
[Assignment Project Exam Help](https://eduassistpro.github.io/)
<https://eduassistpro.github.io/>
- Suppose the random walk $p(t)$ is a state
 $p(t + 1) = M \cdot p(t) = p(t)$
Then $p(t)$ is **stationary distribution** of a random walk
- Our original rank vector r satisfies $r = M \cdot r$
 - So r is a stationary distribution for the random walk

PageRank

- Three questions:
 - Does this converge?
Assignment Project Exam Help
 - Does it converge?
<https://eduassistpro.github.io/>
 - Are results reasonable
Add WeChat edu_assist_pro

Does This Converge?

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Does it Converge to What We Want?

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

PageRank: Problems

- Two problems:
 - Spider traps: all out-links are within the group
 - Eventually all importance <https://eduassistpro.github.io/>
 - Some pages cause importance leak out
 - Such pages cause importance leak out”

Problem: Spider Traps

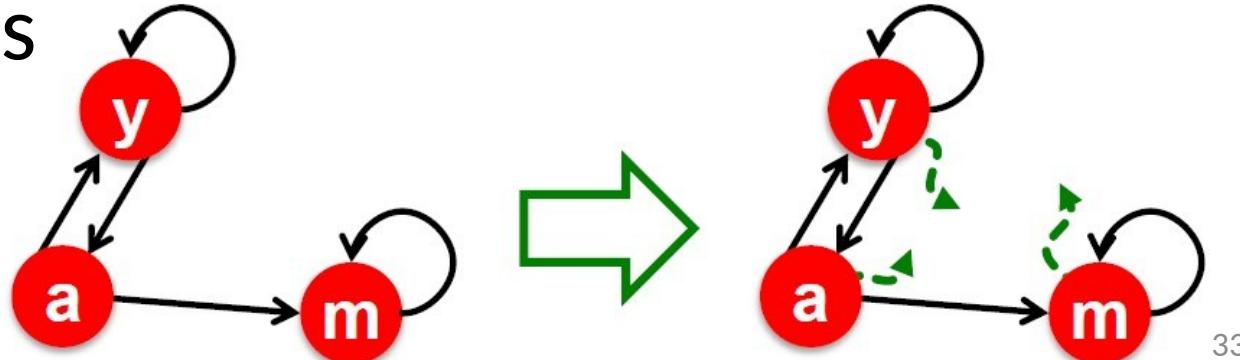
- Power Iteration:
 - Set $r_j = 1$
 - Assignment Project Exam Help
 - - And iterate
- Example

$$\begin{bmatrix} r_y \\ r_a \\ r_m \end{bmatrix} = \begin{matrix} 1/3 & 2/6 & 3/12 & 5/24 & 0 \\ 1/3 & 1/6 & 2/12 & 3/24 & \dots & 0 \\ 1/3 & 3/6 & 7/12 & 16/24 & & 1 \end{matrix}$$

Iteration 0, 1, 2, ...

Solution: Random Teleports

- The Google solution for spider traps: At each time step, the random surfer has two options
 - With prob. $\frac{1}{3}$, follow a link at random
 - With prob. $\frac{2}{3}$, teleport to a random page
 - Common values: $\alpha = 0.8$ to 0.9
- Surfer will teleport out of spider trap within a few time steps



Problem: Dead Ends

- Power Iteration:

- Set $r_j = 1$

Assignment Project Exam Help

-

- And iterate

<https://eduassistpro.github.io/>

- Example

$$\begin{bmatrix} r_y \\ r_a \\ r_m \end{bmatrix} = \begin{bmatrix} 1/3 & 2/6 & 3/12 & 5/24 & \dots & 0 \\ 1/3 & 1/6 & 2/12 & 3/24 & \dots & 0 \\ 1/3 & 1/6 & 1/12 & 2/24 & \dots & 0 \end{bmatrix}$$

Iteration 0, 1, 2, ...

Solution: Always Teleport

- Teleports: Follow random teleport links with probability 1.0 from dead-ends
 - Adjust matr

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Why Teleports Solve the Problem?

$$\mathbf{r}^{(t+1)} = M\mathbf{r}^{(t)}$$

Assignment Project Exam Help

- Markov chain <https://eduassistpro.github.io/>
 - Set of states X [Add WeChat edu_assist_pro](#)
 - Transition matrix P where $P_{ij} = P(X_t = i | X_{t-1} = j)$
 - π specifying the stationary probability of being at each state $x \in X$
 - Goal is to find π such that $\pi = P\pi$

Why Is This Analogy Useful?

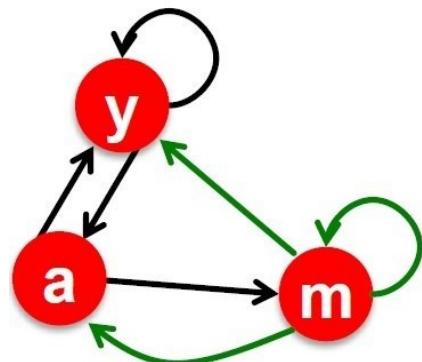
- Theory of Markov chains
- Fact: For ~~any start vector, Assignment Project, Exam Help~~, the power method applied to a ~~matrix P will converge to a stationary vector as long as P is Strongly irreducible and aperiodic.~~ <https://eduassistpro.github.io/>

Make M Stochastic

- **Stochastic:** Every column sums to 1
- A possible solution: add green links
Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro



	y	a	m
y	$\frac{1}{2}$	$\frac{1}{2}$	$\textcolor{green}{1/3}$
a	$\frac{1}{2}$	0	$\textcolor{green}{1/3}$
m	0	$\frac{1}{2}$	$\textcolor{green}{1/3}$

$$\mathbf{r}_y = \mathbf{r}_y/2 + \mathbf{r}_a/2 + \mathbf{r}_m/3$$

$$\mathbf{r}_a = \mathbf{r}_y/2 + \mathbf{r}_m/3$$

$$\mathbf{r}_m = \mathbf{r}_a/2 + \mathbf{r}_m/3$$

Make M Aperiodic

- A chain is periodic if there exists $k > 1$ such that the interval between two visits to some state s is always k .
- A possible solution is $\frac{1}{k}$.
<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro

Make M Irreducible

- From any state, there is a non-zero probability of going from any one state to any another
- A possible so

[Assignment](#) [Project](#) [Exam](#) [Help](#)

[links
https://eduassistpro.github.io/](https://eduassistpro.github.io/)

[Add WeChat edu_assist_pro](#)

Solution: Random Jumps

- Google's solution that does it all:
 - Makes \mathbf{M} stochastic, aperiodic, irreducible
- At each step, s two options:
 - With probability β , follow random
 - With probability $1 - \beta$, jump to random page
- PageRank equation [Brin-Page, 98]

$$r_j = \sum_{i \rightarrow j} \beta \frac{r_i}{d_i} + (1 - \beta) \frac{1}{n}$$

This formulation assumes that \mathbf{M} has no dead ends. We can either preprocess matrix \mathbf{M} to remove all dead ends or explicitly follow random teleport links with probability 1.0 from dead-ends.

The Google Matrix

- PageRank equation [Brin-Page,98]

$$r_j = \sum_{i \rightarrow j} \beta \frac{r_i}{d_i} + (1 - \beta) \frac{1}{n}$$

- The Google <https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

- A is stochastic, aperiodic and irreducible, so

$$\mathbf{r}^{(t+1)} = \mathbf{A} \cdot \mathbf{r}^{(t)}$$

- In practice $\beta = 0.8, 0.9$ (make 5 steps and jump)

In-class Practice

- Go to Practice

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Computing PageRank

- Key step is matrix-vector multiplication
 - $r^{new} = A \cdot r^{old}$
- Easy if we have A, r^{old}, r^{new}
 - <https://eduassistpro.github.io/>
 - Add WeChat edu_assist_pro
- Say $N=1$ billion pages
 - We need 4 bytes for each entry (say)
 - 2 billion entries for vectors, approx. 8GB
 - Matrix A has N^2 entries: 10^{18} is a large number!

Matrix Formulation

- Suppose there are N pages
 - Consider page j , with d_j out-links
[Assignment](#) [Project](#) [Exam](#) [Help](#)
 - We have $M_{ij} = \frac{1}{d_j}$ and $M_{ij} = 0$ otherwise
- The random teleport is $\frac{1}{N}$ to:
 - Adding a teleport link from j to every other page and setting transition prob. to $(1 - \beta)/N$
 - Reducing the prob. of following each out-link from $1/|d_j|$ to $\beta/|d_j|$

Rearranging the Equation

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Note: Here we assumed **M** has no **dead-ends**

$[x]_N$... a vector of length N with all entries x

Spare Matrix Formulation

- We just rearranged the PageRank equation

$$r = \beta M \cdot r + \left[\frac{1-\beta}{N} \right]$$

Assignment Project Exam Help

- M is a sparse matrix (dead-ends)

– 10 links per node, approximately 10% dead-ends

- So in each iteration, we need to

– Compute $r^{new} = A \cdot r^{old}$

– Add a constant $(1 - \beta)/N$ to each entry in r^{new}

- Note: if M contains dead-ends then $\sum_i r_i^{new} < 1$ and we also have to renormalize r^{new} so that it sums to 1

PageRank: The Complete Algorithm

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Sparse Matrix Encoding

- Encode sparse matrix using only nonzero entries

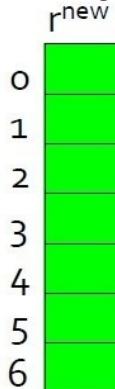
Assignment Project Exam Help

- Space prop umber of links
- Say $10N$, or <https://eduassistpro.github.io/>
- Still won't fit in memory ~~Add WeChat edu_assist_pro~~ ~~t on disk~~

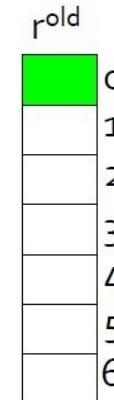
source node	degree	destination nodes
0	3	1, 5, 7
1	5	17, 64, 113, 117, 245
2	2	13, 23

Basic Algorithm: Update Step

- Assume enough RAM to fit r^{new} into memory
 - Store r^{old} and matrix M on disk
- Then 1 step ~~Assignment Project Exam Help~~
 - Initialize all <https://eduassistpro.github.io/>
 - For each page p (of out-degree ≥ 1)
 - Read into memory: $p, n, \text{st}_n, r^{old}(p)$
 - For $j=1\dots n$: $r^{new}(\text{dest}_j) += \beta r^{old}(p)/n$



src	degree	destination
0	3	1, 5, 6
1	4	17, 64, 113, 117
2	2	13, 23



Analysis

- Assume enough RAM to fit r^{new} into memory
 - Store r^{old} and matrix M on disk
[Assignment](#) [Project](#) [Exam](#) [Help](#)
- In each iteration
 - Read r^{old} and M
<https://eduassistpro.github.io/>
[Add WeChat edu_assist_pro](#)
 - Write r^{new} back to disk
 - IO cost = $2|r| + |M|$
- Question:
 - What if we could not even fit r^{new} in memory

Block-based Update Algorithm

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Analysis of Block Update

- Similar to nested-loop join in databases
 - Break r^{new} into k blocks that fit in memory
[Assignment](#) [Project](#) [Exam](#) [Help](#)
 - Scan M and r k times, each with a lock
<https://eduassistpro.github.io/>
- k scans of M
 - $k(|M| + |r|) + |r| = k|M|$
[Add WeChat](#) [edu_assist_pro](#)
- Can we do better?
 - Hint: M is much bigger than r (approx. 10-20x), so we must avoid reading it k times per iteration

Block-Strip Update Algorithm

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Block-Strip Analysis

- Break M into stripes
 - Each strip contains only destination nodes in the corresponding row
- Some additional overhead
 - But it is usually worth it
- Cost per iteration

$$|M|(1 + \epsilon) + (k + 1)|r|$$

Some Problems with PageRank

- Measures generic popularity of a page
 - Biased against topic-specific authorities
Assignment Project Exam Help
 - Solution: To k (next)
<https://eduassistpro.github.io/>
- Susceptible to
 - Artificial link topography in order to boost page rank
Add WeChat edu_assist_pro
 - Solution: TrustRank (next)
- Uses a single measure of importance
 - Solution: Hubs-and-Authorities (in further reading) □

Outline

- Web as a Graph
- PageRank [Assignment](#) [Project](#) [Exam](#) [Help](#)
- Topic-Specific <https://eduassistpro.github.io/>
- Appendix: Trust-Rank [Add WeChat edu_assist_pro](#)

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 according to whether they are related to a particular topic, e.g. “edu_assist_pro” 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
 - Assignment Project Exam Help
- Teleport can
 - Standard PageRank with equal probability
 - To avoid dead-end and spider-trap problems
 - **Topic Specific PageRank**: A topic-specific set of “relevant” pages (teleport set)

Topic-Specific PageRank

- Idea: Bias the random walk
 - When walker teleports, she picks a page from a set S
 - S contains only pages that are relevant to the topic
 - e.g., Open <https://eduassistpro.github.io/> on 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:
Assignment Project Exam Help

<https://eduassistpro.github.io/>

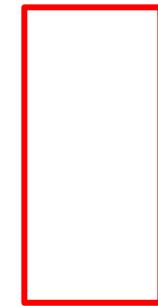
- A is stochastic!
- We have weighted all pages in the teleport set S equally
 - Could also assign different weights to pages!
- Compute as for standard PageRank

Example

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro



Discovering the Topic

- Create different PageRanks for different topics
 - The 16 DMOZ top-level categories: arts, business, sports, ... **Assignment Project Exam Help**
- Which topic <https://eduassistpro.github.io/>
 - User can pick ~~Add from WeChat~~ **edu_assist_pro**
 - 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
 - User context, e.g., user's bookmarks,...

SimiRank: An Application of Personalized PageRank

- SimRank: Random walks from a fixed node on k-partite graphs
[Assignment](#) [Project](#) [Exam](#) [Help](#)
- Setting: k-pa
 - Example: pi <https://eduassistpro.github.io/> types of nodes
- How to find nodes similar to u ?
[Add WeChat edu_assist_pro](#)
- Do a Random-Walk with Restarts from node u
 - i.e. teleport set $S = \{u\}$

SimiRank(cont.)

- Resulting scores measures similarity/proximity to node u
 - Generally applied (typically unweighted)
 - Problems:
 - Must be done once for each node u
 - Suitable for sub-Web-scale applications
- Assignment Project Exam Help
etworks
<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro

SimRank: Example

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

SimRank: Example

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Outline

- Web as a Graph
- PageRank [Assignment](#) [Project](#) [Exam](#) [Help](#)
- Topic-Specific <https://eduassistpro.github.io/>
- Appendix: Trust-Rank [Add WeChat edu_assist_pro](#)
- [Skip to conclusion](#)

What is Web Spam?

- Spaming:
 - Any deliberate action to boost a web page's position in search results, incommensurate with page's <https://eduassistpro.github.io/>
- Spam:
 - Web pages that are the result of spamming
- Approximately 10-15% of web pages are spam

Web Search

- Early Search Engines
 - Crawl the Web
 - Index page
 - Respond to search query
- Assignment Project Exam Help
<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro

Web Search

- Early page ranking
 - Attempt to order pages matching a search query by “i”
<https://eduassistpro.github.io/>:
Assignment Project Exam Help
 - First search
 - Number of times query “i” has been searched
 - Prominence of word position, e.g. title, header

First Spammers

- Those with commercial interests tried to exploit search engines to bring people to their own site – would be there or not
<https://eduassistpro.github.io/>
- Example Add WeChat edu_assist_pro
 - 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?
Assignment Project Exam Help
 - Add the word “movie” to your page, set the text color to <https://eduassistpro.github.io/>, or, so only search engines would see it.
 - 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 that appear und surrounding text
 - PageRank as a tool to measure the “importance” of Web pages

Why It Works?

- Our hypothetical shirt-seller loses
 - Saying he is about movies doesn't help, because others don't say he is about movies, his page isn't very important, it's not ranked high for shirts or movies <https://eduassistpro.github.io/>
- Example: Add WeChat edu_assist_pro
 - 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 Farms

- Once Google became the dominant search engine, spammers began to work out ways to fool Google
- Spam farms
- 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
 - Accessible p
 - E.g. blog comments page
 - Spammer can post links to his pages
 - Own pages
 - Completely controlled by spammer
 - May span multiple domain names

Link Farms

- Spammer's goal:
 - Maximize the PageRank of target page t
[Assignment](#) [Project](#) [Exam](#) [Help](#)
- Technique:
 - Get as many links from <https://eduassistpro.github.io/> pages as possible to target page t
[Add WeChat](#) [edu_assist_pro](#)
 - Construct “link farm” to get PageRank multiplier effect

Link Farms

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Analysis

Assignment Project Exam Help

- x : PageRank cont <https://eduassistpro.github.io/ges>
- y : PageRank of target page [Add WeChat edu_assist_pro](#)
- Rank of each “farm” page = $\frac{\beta y}{M} + \frac{1 - \beta}{N}$
- $$\begin{aligned}y &= x + \beta M \left[\frac{\beta y}{M} + \frac{1 - \beta}{N} \right] + \frac{1 - \beta}{N} \\&= x + \beta^2 y + \frac{\beta(1 - \beta)M}{N} + \frac{1 - \beta}{N}\end{aligned}$$

Very small: ignore now
we solve for y


$$= \frac{x}{1 - \beta^2} + c \frac{M}{N}, \text{ where } c = \frac{\beta}{1 + \beta}$$

Analysis

Assignment Project Exam Help

<https://eduassistpro.github.io/>

- Add WeChat edu_assist_pro
- Multiplier effect for “acquired” PageRank
- By making M large, we can make y as large as we want

TrustRank: Combating Spam

- Combating term spam
 - Analyze text using statistical methods
 - Similar to email spam filtering
 - Also useful: D <https://eduassistpro.github.io/>
- 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”
Assignment Project Exam Help
(spam) pag
- Sample a set <https://eduassistpro.github.io/> eb
- Have an oracle (human) *Add WeChat edu_assist_pro* fy 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**
Assignment Project Exam Help
- Perform a to **Rank with**
teleport set <https://eduassistpro.github.io/>
– **Propagate trust through**
 • Each page gets a trust value between 0 and 1
- Use a **threshold value** and mark all pages below the trust threshold as spam

Why is It a Good Idea?

- Trust attenuation:
 - The degree of trust conferred by a trusted page decreases w Assignment Project Exam Help e graph
- Trust splitting
 - The larger the number o Add WeChat edu_assist_pro s 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 [Assignment](#) [Project](#) [Exam](#) [Help](#)
 - Must ensure <https://eduassistpro.github.io/> is adequate rank, so need [Add](#) [WeChat](#) [edu_assist_pro](#) reachable from seed set by short paths

Approaches to Picking Seed Set

- Suppose we want to pick a good set of k pages
- How to do that?
[Assignment Project Exam Help](#)
- **PageRank** <https://eduassistpro.github.io/>
 - Pick the top k pages by P
[Add WeChat edu_assist_pro](#)
 - Theory is that you can't have a page's rank really high
- **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.
[Assignment](#) [Project](#) [Exam](#) [Help](#)
- Complement
 - What fraction of the rank comes from spam pages?
<https://eduassistpro.github.io/>
[Add WeChat edu_assist_pro](#)
- In practice, we don't know all the spam pages, so we need to estimate.

Spam Mass Estimation

- r_p = PageRank of page p
- r_p^+ = PageRank of p with ~~Assignment Project Exam Help~~ into trusted pages only <https://eduassistpro.github.io/>
Add WeChat edu_assist_pro
- Then: What fraction of a page's PageRank comes from spam pages?
$$p = \frac{r_p^-}{r_p}$$
- Spam mass of

One-slide Takeaway

- Web as a Graph
 - Denote the web structure as a graph
- PageRankAssignment Project Exam Help
 - PageRank s <https://eduassistpro.github.io>of web pages
- Topic-SpecificPageRank
 - Evaluate web pages by their popularity as well as particular topic
- Trust-Rank
 - Deal with link spams

Further Reading

- Original PageRank paper: [http://
ilpubs.stanford.edu:8090/422/1/1999-66.pdf](http://ilpubs.stanford.edu:8090/422/1/1999-66.pdf)
Assignment Project Exam Help
- An Analytical approach to Personalizing
<https://eduassistpro.github.io/>
[http://www-cs-student.d.edu/~taherh/
papers/comparison.pdf](http://www-cs-student.d.edu/~taherh/papers/comparison.pdf)
Add WeChat edu_assist_pro

Further Reading

- HITS algorithm: [http://
www.cs.cornell.edu/home/kleinber/auth.pdf](http://www.cs.cornell.edu/home/kleinber/auth.pdf)
Assignment Project Exam Help
- Parallel Page [https://eduassistpro.github.io/
http://link.springer.com/pdf/10.1007%2F11735106_22.pdf](https://eduassistpro.github.io/nt/pdf/10.1007%2F11735106_22.pdf)
Add WeChat edu_assist_pro

Reference

- <http://www.stanford.edu/class/cs246/slides/09-pagerank.pdf>
- <http://www.stanford.edu/class/cs246/slides/10-spam.pdf>
- [Assignment Project Exam Help
http://i.stanford.edu/~ullman/mmds/ch5.pdf](http://i.stanford.edu/~ullman/mmds/ch5.pdf)
- <http://en.wikipe> <https://eduassistpro.github.io/>
- <http://en.wikipedia.org/wiki/HI> [Add WeChat edu_assist_pro](#)

In-class Practice

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

- Compute the final PageRank Score of the given graph, with Google matrix A and assume $\beta = 0.8$
- Show the matrix A, and solve by both Gaussian Elimination and Power Iteration methods