

CMSC5741 Big Data Tech. & Apps.

Lecture 8: Massive Link Analysis

<https://powcoder.com>

Add WeChat powcoder

Prof. Michael R. Lyu

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

What's the Mechanism Distinguishing Google?

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

One important factor is PageRank score.

Charles K. Kao - Wikipedia
https://en.wikipedia.org/wiki/Charles_K._Kao
Sir Charles Kuen Kao GBM KBE FRS FREng (4 November 1933 – 23 September 2018) was a Hong Kong electrical engineer and physicist who pioneered the ...
Known for: Fibre optics; Fibre-optic communic... Doctoral advisor: Harold Barlow
Fields: Physics Citizenship: United States; United Kingdom
Early life and education · Academic career · Honours and awards · Awards

Charles K. Kao - Facts - NobelPrize.org
<https://www.nobelprize.org/prizes/physics/2009/kao/facts/>
Charles Kuen Kao, Born: November 1933, Shanghai, China, Died: 23 September 2018, Affiliation at the time of the award: Standard Telecommunications Laboratories, London, United Kingdom, Prize motivation: "for their ground-breaking achievements concerning the transmission of optical fiber in light and in particular for their contribution to the development of optical fiber for communication"

Sir Charles Kao: Fibre optics genius passes away - BBC News
<https://www.bbc.com/news/uk-england-essex-45647549>
Sep 26, 2018 - Tributes have been paid to Sir Charles Kao, the scientist whose work in Essex "transformed the world". Sir Charles, who won the 2009 Nobel

https://powcoder.com

Mourning Professor Sir Charles Kao, former Vice-Chancellor of CUHK ...
https://www.cpr.cuhk.edu.hk/en/press_detail.php?id=2857
Sep 23, 2018 - Professor Sir Charles K. Kao, the third Vice-Chancellor, Honorary Professor of Engineering, and Doctor of Science, honoris causa, of The ...

Hong Kong mourns passing of Nobel Prize winner and father of fibre optics
<https://www.scmp.com › News › Hong Kong › Education>
Sep 23, 2018 - Hong Kong on Sunday mourned the passing of the city's Nobel Prize winner in physics, Professor Charles Kao Kuen, whose seminal work on ...

Add WeChat powcoder

Map data ©2018 Google

Rating -

The Chinese University of Hong Kong Charles Kuen Kao Building Scie...
5.0 ★★★★★ (2)
Open now

Charles K. Kao Auditorium

Charles K. Kao

Charles K. Kao, Electrical engineer

Sir Charles Kuen Kao GBM KBE FRS FREng was a Hong Kong electrical engineer and physicist who pioneered the development and use of fibre optics in telecommunications. [Wikipedia](#)

Born: November 4, 1933, Shanghai, China
Died: September 23, 2018, Sha Tin, Hong Kong
Spouse: May-Wan Kao (m. 1959)
Education: University of Greenwich (1957), St Joseph's College (1952), UCL, University of London International Programmes
Awards: Nobel Prize in Physics, Grand Bauhinia Medal, MORE

Books

A Time And A Tide 2011
Optical Fiber Systems... 1982
A Choice Fulfilled: The Business Of Optical Fiber Technology 1991
Optical Fiber Technolo... 1981
Nonlinear Photonics: Nonlinear... 2002

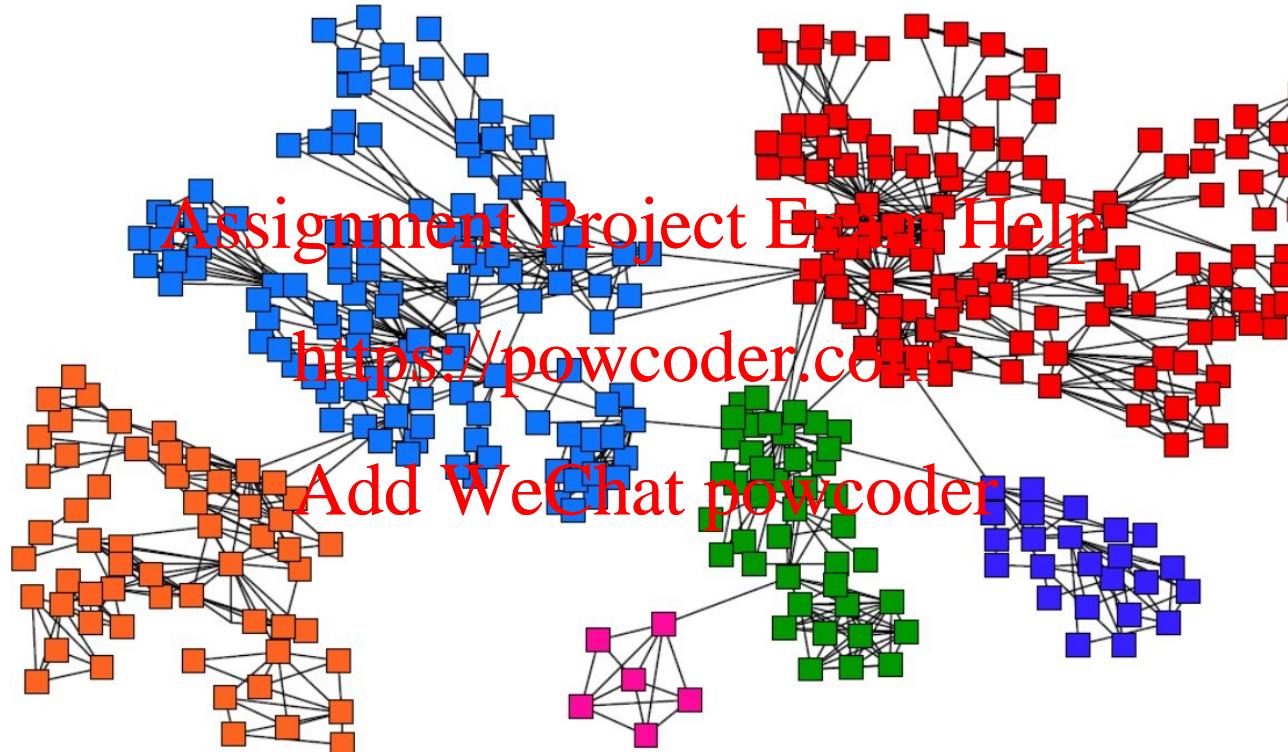
View 1+ more

People also search for

Willard Boyle Daniel C. Tsui Yoichiro Nambu Samuel C. C. Ting Subrahmanyan Chandrasekhar

Feedback

How to make PageRank computation scalable



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

Outline

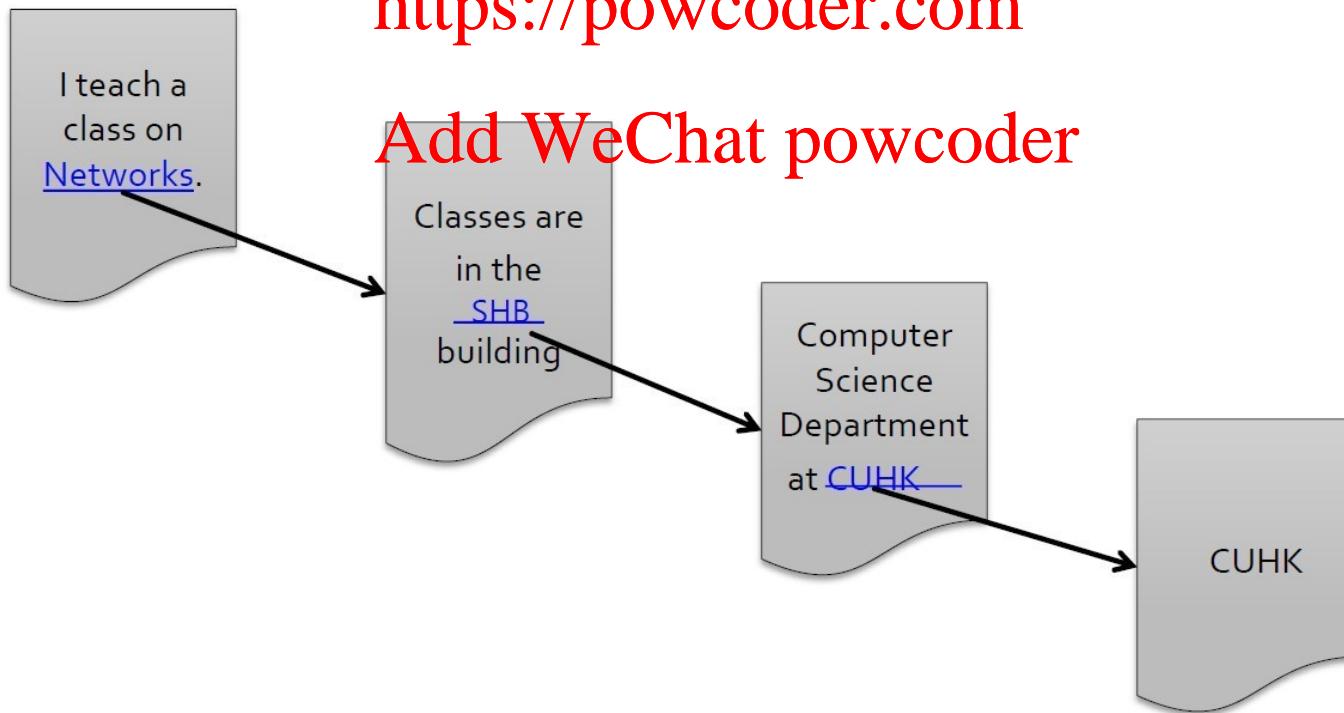
- Web as a Graph
- PageRank [Assignment](#) [Project](#) [Exam](#) [Help](#)
- Topic-Specific PageRank <https://powcoder.com>
- Appendix: Trust-Rank [Add WeChat](#) [powcoder](#)

Outline

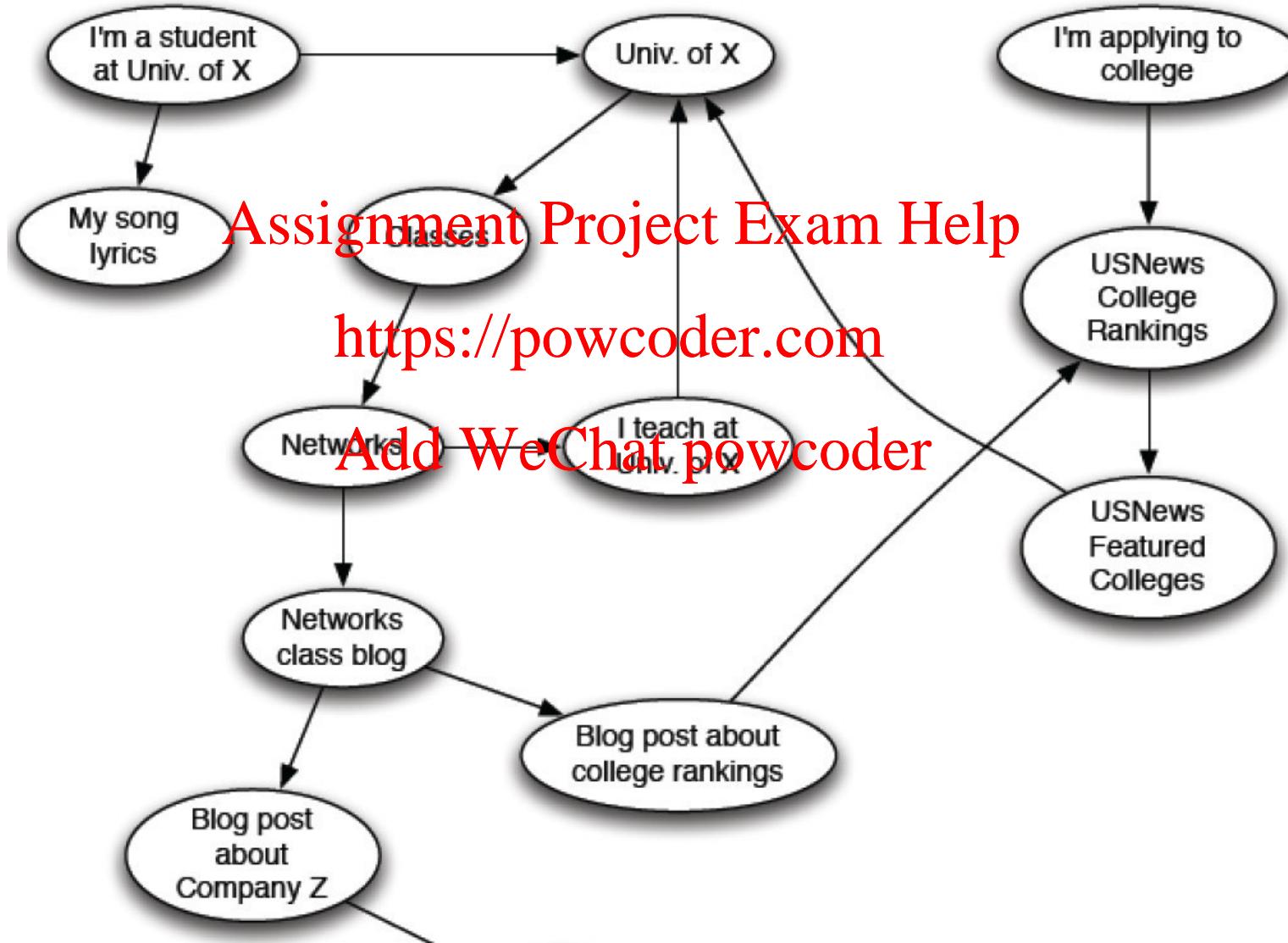
- Web as a Graph
- PageRank [Assignment](#) [Project](#) [Exam](#) [Help](#)
- Topic-Specific PageRank <https://powcoder.com>
- Appendix: Trust-Rank [Add WeChat](#) [powcoder](#)

Web as a Graph

- Web as a directed graph:
 - Nodes: Web pages
Assignment Project Exam Help
 - Edges: Hyperlinks
https://powcoder.com



Web as a Directed Graph



Broad Question

- How to **organize** the Web?
- First try: Human curated **Web directories**
 - Yahoo, DMOZ, LookSmart
<https://powcoder.com>
- Second try: **Web Search**
 - **Information Retrieval** investigates: 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: Trustworthy pages may point to each other
<https://powcoder.com>
- What is the “best” answer to query
“newspaper”?
Add WeChat powcoder
 - No single right answer
 - Trick: Pages that actually know about newspapers might all be pointing to many newspapers

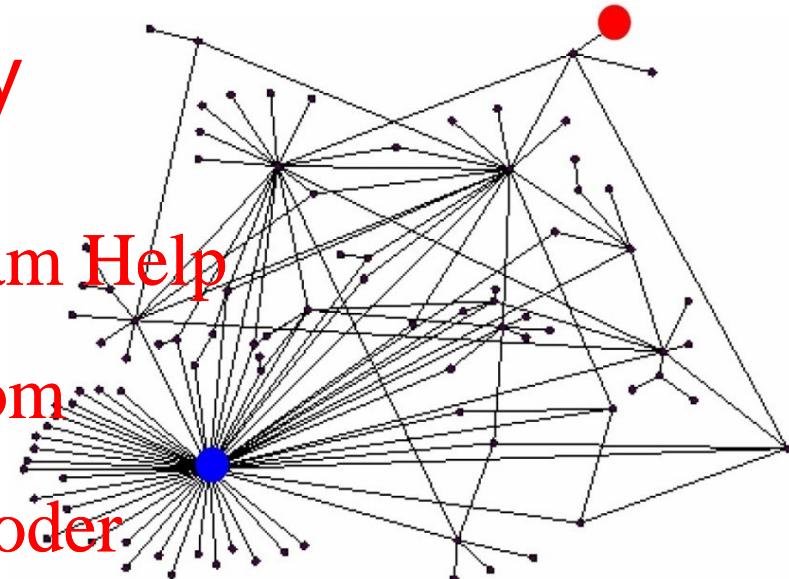
Ranking Nodes on the Graph

- All web pages are **not equally** “important”

– www.cuhk.edu.hk vs.
www.joe-schimpo.com

- There is large **Adversity** in the 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
 - Topic-Specific (Personalized) PageRank
 - Web Spam Detection Algorithms, e.g. TrustRank

<https://powcoder.com>

Add WeChat powcoder

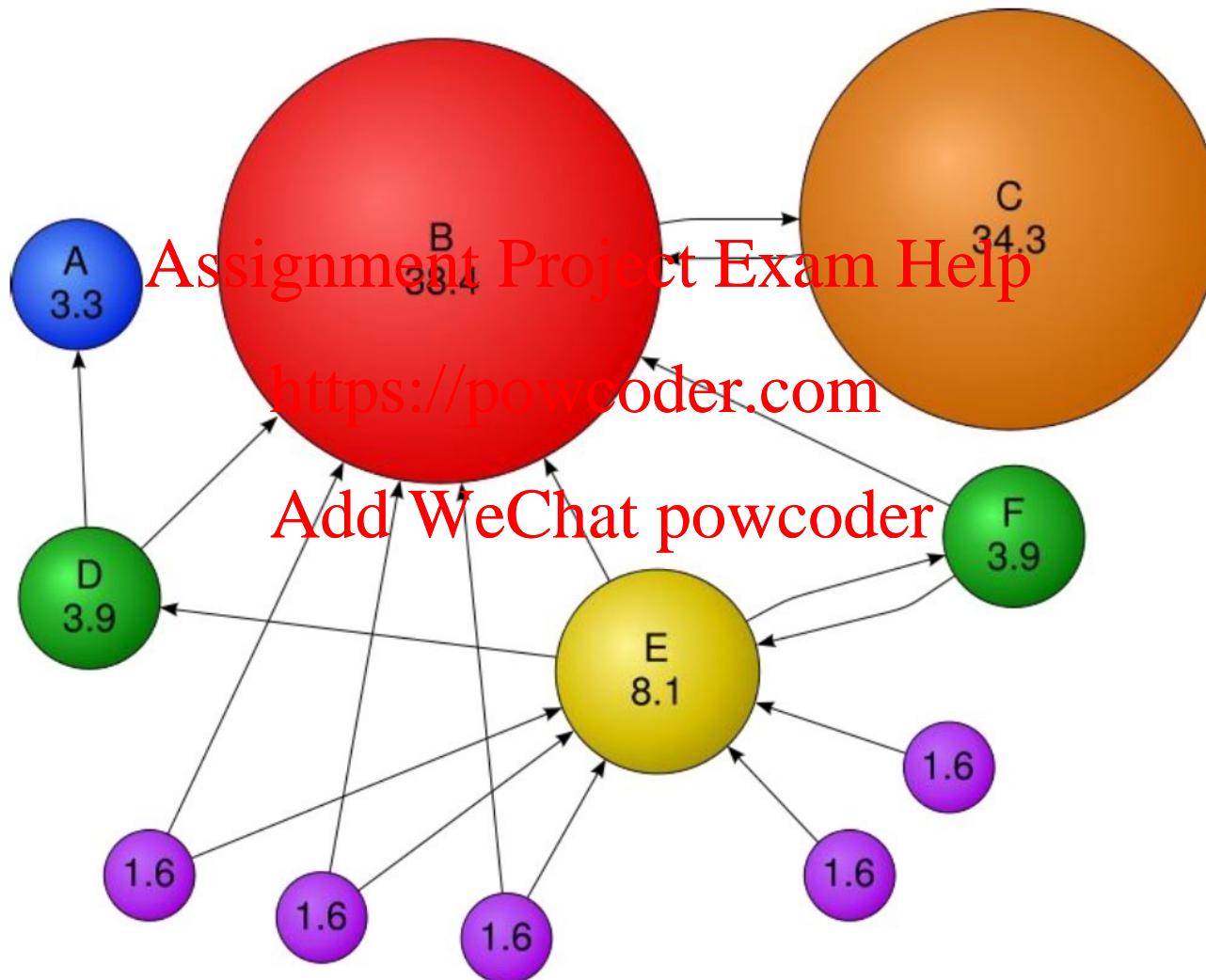
Outline

- Web as a Graph
- PageRank
Assignment Project Exam Help
- Topic-Specific PageRank
<https://powcoder.com>
- Appendix: Trust-Rank
Add WeChat powcoder

Links as Votes

- Idea: Links as votes
 - Page is more important if it has more links
Assignment Project Exam Help
 - In-coming links? Out-going links?
https://powcoder.com
- Think of in-links as votes:
 - www.cuhk.edu.hk has 15,432 in-links
Add WeChat powcoder
 - www.joe-schmoe.com has 1 in-link
- Are all in-links equal?
 - Link from important pages count more
 - Recursive question!

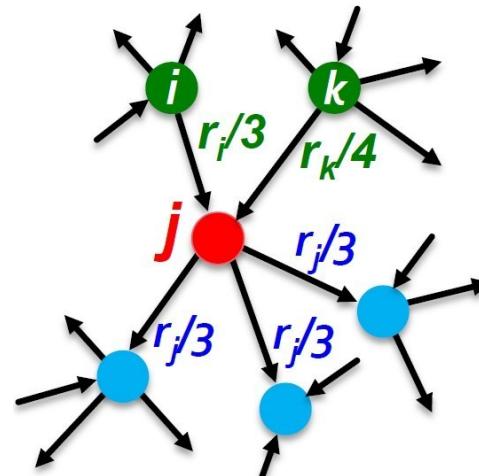
Example: PageRank Scores



Simple Recursive Formulation

- Each link's vote is proportional to the **importance** of its source page
[Assignment](#) [Project](#) [Exam](#) [Help](#)
- If page j with importance r_j has n out-links,
<https://powcoder.com>
each link gets r_j/n votes.
- Page j 's own importance is the sum of the
votes on its in-links

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

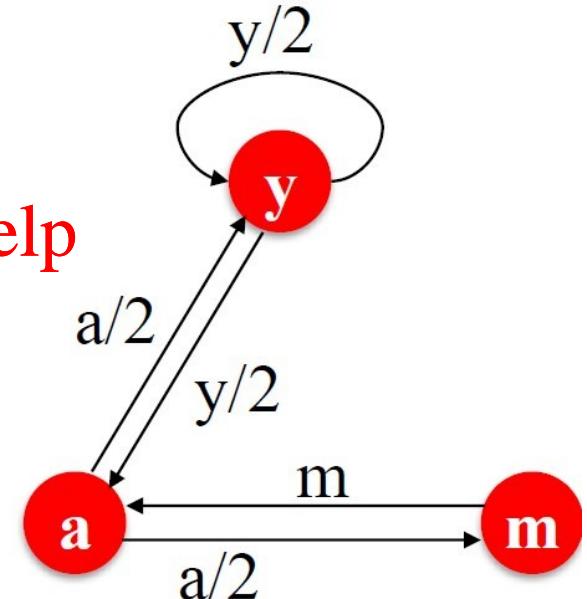


PageRank: The “Flow” Model

- A “vote” from an important page is worth more
- A page is important if it is pointed to by other 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 equivalent modulo the scale factor

- Additional constraint forces uniqueness:

- $r_y + r_a + r_m = 1$
 - 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

Flow equations:

$$r_y = r_y/2 + r_a/2$$

$$r_a = r_y/2 + r_m$$

$$r_m = r_a/2$$

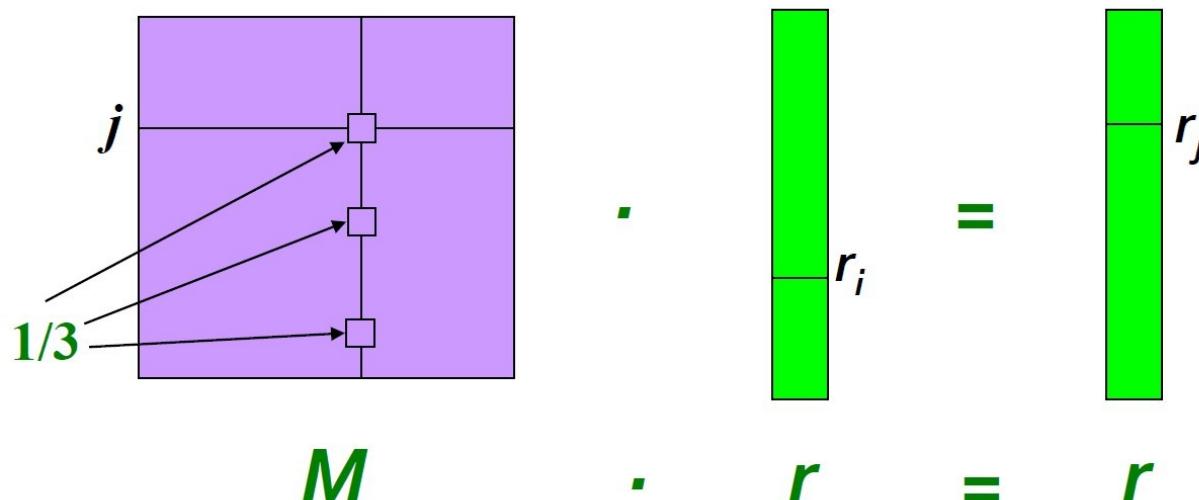
PageRank: Matrix Formulation

- Stochastic adjacency matrix M
 - Let page i have d_i out-links
 - If $i \rightarrow j$, then $M_{ji} = \frac{1}{d_i}$, else $M_{ji} = 0$
 - M is a column stochastic matrix, i.e., columns sum to 1
- Rank vector r : vector with an 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://powcoder.com

– Suppose page i links to 3 pages, including j

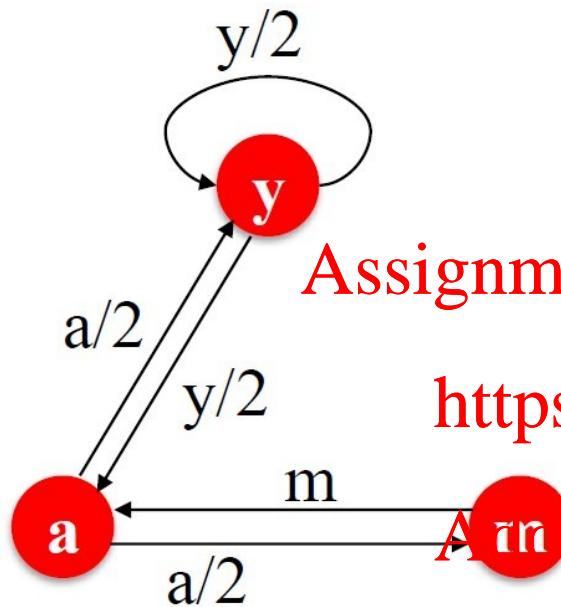


Eigenvector Formulation

- The flow equations can be written $r = M \cdot r$
- So the rank vector r is an eigenvector of the stochastic web matrix M
 - In fact, its first or principal eigenvector, with corresponding eigenvalue 1
 - Largest eigenvalue of M is 1 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://powcoder.com>

WeChat powcoder $r = M \cdot r$

y	a	m	
y	1/2	1/2	0
a	1/2	0	1
m	0	1/2	0

“Flow” equations:

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

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

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

$$\begin{bmatrix} y \\ a \\ m \end{bmatrix} = \begin{bmatrix} 1/2 & 1/2 & 0 \\ 1/2 & 0 & 1 \\ 0 & 1/2 & 0 \end{bmatrix} \begin{bmatrix} y \\ a \\ m \end{bmatrix}$$

Power Iteration Method

- Given a web graph with n nodes, where the nodes are pages and edges are hyperlinks
Assignment Project Exam Help
- Power Iteration: a simple iterative scheme
<https://powcoder.com>
 - Suppose there are N web pages
 - Initialize: $\mathbf{r}^{(0)} = [1/N, \dots, 1/N]^T$
 - 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

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

Why Power Iteration Works?

- Power iteration:
 - A method for finding **dominant eigenvector** (the vector corresponding to the largest eigenvalue)
 - $r^{(1)} = M \cdot r^{(0)}$
 - $r^{(2)} = M \cdot r^{(1)}$
 - $r^{(3)} = M \cdot r^{(2)} = M(M^2 r^{(0)}) = M^3 \cdot r^{(0)}$
- 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 $\lambda_1, \lambda_2, \dots, \lambda_n$, where $\lambda_1 > \lambda_2 > \dots > \lambda_n$
 - Vectors x_1, x_2, \dots, x_n form a basis and thus we can write: $r^{(0)} = c_1x_1 + c_2x_2 + \dots + c_nx_n$
 - $$\begin{aligned} Mr^{(0)} &= M(c_1x_1 + c_2x_2 + \dots + c_nx_n) \\ &= c_1(Mx_1) + c_2(Mx_2) + \dots + c_n(Mx_n) \\ &= c_1(\lambda_1x_1) + c_2(\lambda_2x_2) + \dots + c_n(\lambda_nx_n) \end{aligned}$$
 - 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

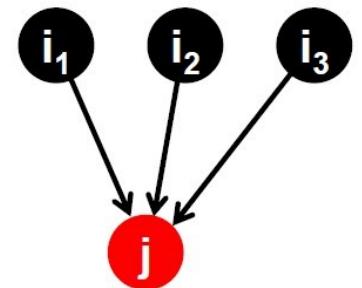
Assignment Project Exam Help
https://powcoder.com

$$M^k r^{(0)} = \lambda_1^k \left[c_1 x_1 + c_2 \left(\frac{\lambda_2}{\lambda_1} \right)^k x_2 + \dots + c_n \left(\frac{\lambda_2}{\lambda_1} \right)^k x_n \right]$$

- Since $\lambda_1 > \lambda_2$ then $\frac{\lambda_2}{\lambda_1}, \frac{\lambda_3}{\lambda_1}, \dots < 1$, so $\frac{\lambda_i}{\lambda_1} = 0$ as $k \rightarrow \infty$
 - 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 random
 - Ends up on some page j linked from i
 - 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

$$p(t + 1) = M \cdot p(t)$$

<https://powcoder.com>

- Suppose the random walk reaches a state

Add WeChat powcoder

$$p(t + 1) = M \cdot p(t) = p(t)$$

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 to what we want?
<https://powcoder.com>
 - Are results reasonable?
Add WeChat powcoder

Does This Converge?

Assignment Project Exam Help

$r_j^{(t+1)} = \sum_{i \rightarrow j} \frac{r_i^{(t)}}{d_i}$

- Example: Add WeChat powcoder

$$\begin{array}{lcl} r_a & = & 1 \quad 0 \quad 1 \quad 0 \\ r_b & = & 0 \quad 1 \quad 0 \quad 1 \end{array}$$

Iteration 0, 1, 2, ...

Does it Converge to What We Want?

Assignment Project Exam Help
https://powcoder.com

$$r_j^{(t+1)} = \sum_{i \rightarrow j} \frac{r_i^{(t)}}{d_i}$$

- Example:

$$\begin{array}{lcl} r_a & = & 1 & 0 & 0 & 0 \\ r_b & & 0 & 1 & 0 & 0 \end{array}$$

Iteration 0, 1, 2, ...

PageRank: Problems

- Two problems:
 - Spider traps: all out-links are within the group
 - Eventually spider traps absorb all importance
<https://powcoder.com>
 - Some pages are dead ends (have no out-links)
 - Such pages cause importance to “leak out”

Problem: Spider Traps

- Power Iteration:

- Set $r_j = 1$
- $r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$

- And iterate

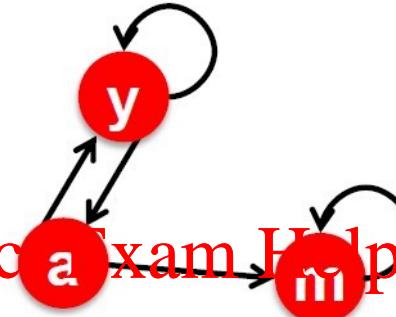
- Example

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Iteration 0, 1, 2, ...



	y	a	m
y	$\frac{1}{2}$	$\frac{1}{2}$	0
a	$\frac{1}{2}$	0	0
m	0	$\frac{1}{2}$	1

$$r_y = r_y/2 + r_a/2$$

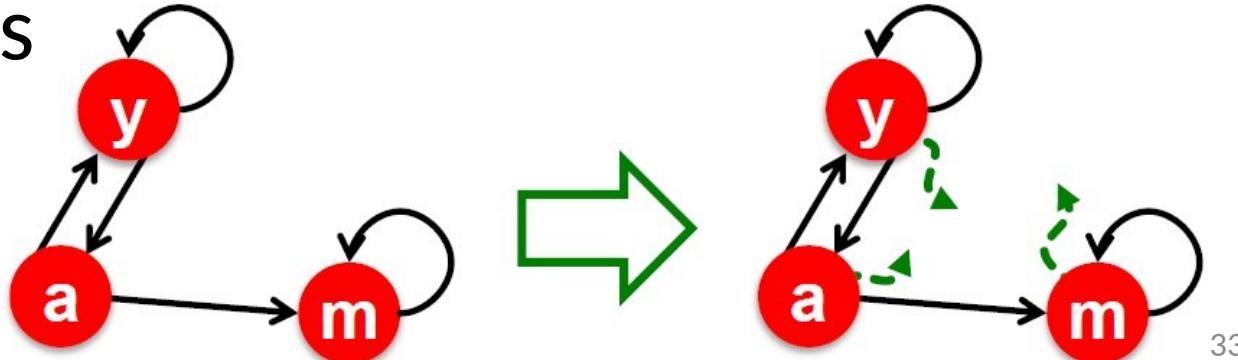
$$r_a = r_y/2$$

$$r_m = r_a/2 + r_m$$

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

Solution: Random Teleports

- The Google solution for spider traps: At each time step, the random surfer has two options
 - With prob. β , follow a link at random
 - With prob. $1 - \beta$, jump to some random page
 - Common values for β range 0.8 to 0.9
- Surfer will teleport out of spider trap within a few time steps



Problem: Dead Ends

- Power Iteration:

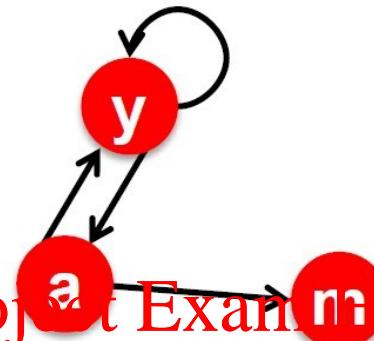
- Set $r_j = 1$
- $r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$

- And iterate

- Example

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

Iteration 0, 1, 2, ...



	y	a	m
y	1/2	1/2	0
a	1/2	0	0
m	0	1/2	0

$$r_y = r_y/2 + r_a/2$$

$$r_a = r_y/2$$

$$r_m = r_a/2$$

Solution: Always Teleport

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

Assignment Project Exam Help

<https://powcoder.com>



	y	a	m
y	$\frac{1}{2}$	$\frac{1}{2}$	0
a	$\frac{1}{2}$	0	0
m	0	$\frac{1}{2}$	0

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

Why Teleports Solve the Problem?

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

Assignment Project Exam Help

- Markov chains <https://powcoder.com>
 - Set of states X [Add WeChat powcoder](#)
 - 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 Markov transition matrix P will ~~https://powcoder.com~~ converge to a ~~unique positive stationary~~ vector as long as P is ~~Add WeChat powcoder~~ ~~stochastic~~, irreducible and ~~aperiodic~~.

Make M Stochastic

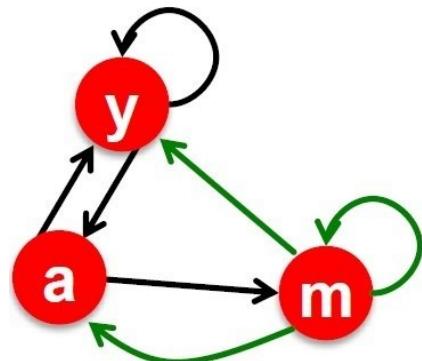
- **Stochastic:** Every column sums to 1
- A possible solution: add green links
Assignment Project Exam Help

$$A = M + a^T \left(\frac{1}{n} e \right)$$

Add WeChat powcoder

- $a_i = 1$ if node i has out degree 0, otherwise 0.

- e ... vector of all 1s



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

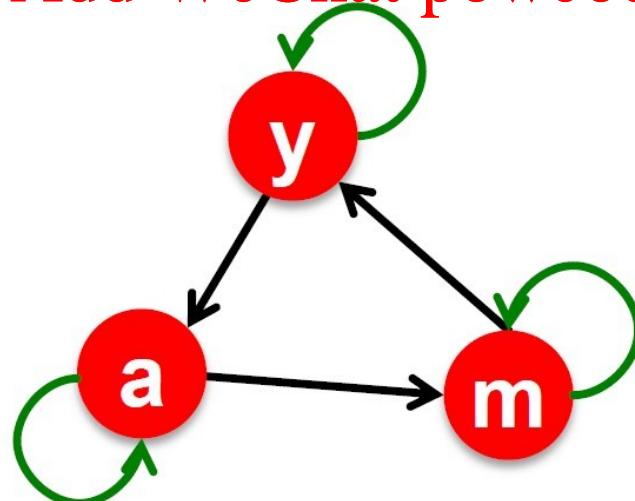
$$\mathbf{r}_y = \mathbf{r}_y/2 + \mathbf{r}_a/2 + \mathbf{r}_m/3$$

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

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

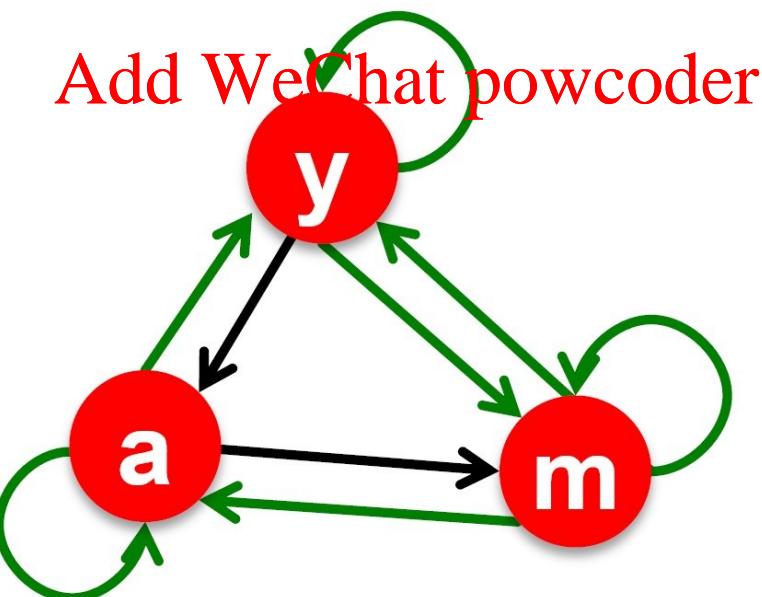
Make M Aperiodic

- A chain is periodic if there exists $k > 1$ such that the interval between two visits to some state s is always a multiple of k .
[Assignment Project Exam Help](https://powcoder.com)
<https://powcoder.com>
[Add WeChat powcoder](#)
- A possible solution: add green links



Make M Irreducible

- From any state, there is a non-zero probability of going from any one state to any another
Assignment Project Exam Help
- A possible solution: add green links
https://powcoder.com



Solution: Random Jumps

- Google's solution that does it all:
 - Makes \mathbf{M} stochastic, aperiodic, irreducible
- At each step, random surfer has two options:
 - With probability β , follow a link at random
 - With probability $1 - \beta$, jump to some 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 Matrix A :
<https://powcoder.com>

$$A = \beta M + (1 - \beta) \frac{1}{n} e \cdot e^T$$

- A is stochastic, aperiodic and irreducible, so

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

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

In-class Practice

- Go to Practice

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Computing PageRank

- Key step is matrix-vector multiplication
 - $r^{new} = A \cdot r^{old}$
- Easy if we have enough main memory to hold
 - A, r^{old}, r^{new}
- 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](#)
<https://powcoder.com>
 - We have $M_{ij} = 1/|d_j|$ when $j \rightarrow i$ and $M_{ij} = 0$ otherwise
- The random teleport is equivalent 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

- $r = A \cdot r$, where $A_{ij} = \beta M_{ij} + \frac{1-\beta}{N}$
- $r_i = \sum_{j=1}^N A_{ij} \cdot r_j$
- $$\begin{aligned} r_i &= \sum_{j=1}^N \left[\beta M_{ij} + \frac{1-\beta}{N} \right] \cdot r_j \\ &= \sum_{j=1}^N \beta M_{ij} \cdot r_j + \frac{1-\beta}{N} \sum_{j=1}^N r_j \\ &= \sum_{j=1}^N \beta M_{ij} \cdot r_j + \frac{1-\beta}{N}, \text{ since } \sum r_j = 1 \end{aligned}$$
- So we get: $r = \beta M \cdot r + \left[\frac{1-\beta}{N} \right]_N$

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! (with no dead-ends)

– 10 links per node, approx. $10N$ entries

Add WeChat powcoder

- 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

- **Input:** Graph \mathbf{G} and parameter β
 - Directed graph \mathbf{G} with **spider traps** and **dead ends**
 - Parameter β
- **Assignment Project Exam Help**
- **Output:** PageRank vector r
 - Set: $r_j^{(0)} = \frac{1}{N}, t = 1$
 - do:
 - * **Add WeChat powcoder**
 - * $\forall j: r'_j^{(t)} = \sum_{i \rightarrow j} \beta \frac{r_i^{(t-1)}}{d_i}$
 - * $r'_j^{(t)} = 0$ if in-deg. of j is 0
 - * Now re-insert the leaked PageRank:
 $\forall j : r_j^{(t)} = r'_j^{(t)} + \frac{1-S}{N}$ where $S = \sum_j r'_j^{(t)}$
 - * $t = t + 1$
 - while $\sum_j |r_j^{(t)} - r_j^{(t-1)}| > \epsilon$

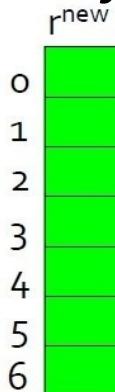
Sparse Matrix Encoding

- Encode sparse matrix using only nonzero entries
 - Space proportional roughly to number of links
<https://powcoder.com>
 - Say $10N$, or $4 * 10 * 1 \text{ billion} = 40\text{GB}$
 - Still won't fit in memory, but will fit 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 entries of r^{new} to $(1 - \beta)/N$
 - For each page p (of out-degree n):
 - Read into memory: $p, n, \text{dest}_1, \dots, \text{dest}_n, r^{old}(p)$
 - 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, we have to:
<https://powcoder.com>
 - Read r^{old} and M
[Add WeChat](#) [powcoder](#)
 - 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

	src	degree	destination	
0	0	4	0, 1, 3, 5	
1	1	2	0, 5	
2	2	2	3, 4	

r^{new}

r^{old}

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

0
1
2
3
4
5

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](#)
<https://powcoder.com>
 - Scan M and r^{old} once for each block
- k scans of M and r^{old}
 - $k(|M| + |r|) + |r| = k|M| + (k+1)|r|$
[Add WeChat](#) [powcoder](#)
- 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

	src	degree	destination
r^{new} 0 1	0	4	0, 1
	1	3	0
	2	2	1
	Assignment Project Exam Help https://powcoder.com		
2 3	0	4	3 Add WeChat powcoder
	2	2	3
4 5	0	4	5
	1	3	5
	2	2	4

Block-Strip Analysis

- Break M into stripes
 - Each strip contains only destination nodes in the corresponding block of r^{new}
- Some additional overhead per stripe
 - 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: Topic-Specific PageRank (next)
https://powcoder.com
- Susceptible to Link spam
 - Artificial link topographies created in order to boost page rank
Add WeChat powcoder
 - Solution: TrustRank (next)
- Uses a single measure of importance
 - Solution: Hubs-and-Authorities (in further reading) □
56

Outline

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

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
 - Assignment Project Exam Help
- Teleport can go to:
 - Standard PageRank: Any page with equal probability
 - Add WeChat powcoder
 - 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://powcoder.com> for a given topic/query
 - For each teleport set S , we get a different vector r_s

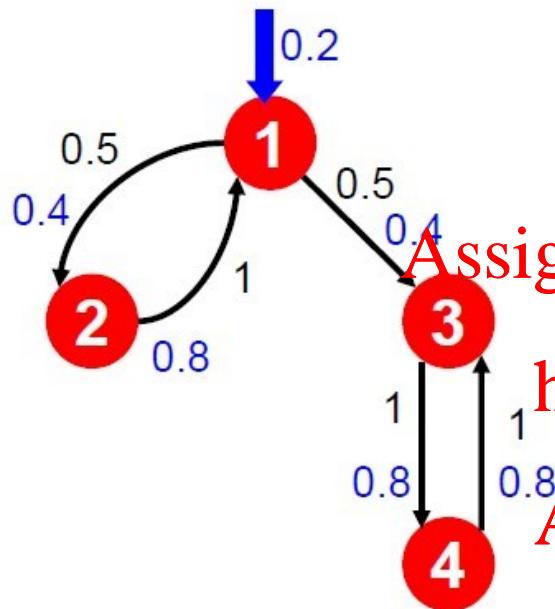
Matrix Formulation

- To make this work all we need is to update the teleportation part of the PageRank formulation:
Assignment Project Exam Help

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

- 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



Suppose $S = \{1\}$, $\beta = 0.8$

Node	Iteration 0	Iteration 1	Iteration 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

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

$S=\{1\}$, $\beta=0.90$:

$r=[0.17, 0.07, 0.40, 0.36]$

$S=\{1\}$, $\beta=0.8$:

$r=[0.29, 0.11, 0.32, 0.26]$

$S=\{1\}$, $\beta=0.70$:

$r=[0.39, 0.14, 0.27, 0.19]$

$S=\{1,2,3,4\}$, $\beta=0.8$:

$r=[0.13, 0.10, 0.39, 0.36]$

$S=\{1,2,3\}$, $\beta=0.8$:

$r=[0.17, 0.13, 0.38, 0.30]$

$S=\{1,2\}$, $\beta=0.8$:

$r=[0.26, 0.20, 0.29, 0.23]$

$S=\{1\}$, $\beta=0.8$:

$r=[0.29, 0.11, 0.32, 0.26]$

Discovering the Topic

- Create different PageRanks for different topics
 - The 16 DMOZ top-level categories: arts, business, sports, ... *Assignment Project Exam Help*
- Which topic ranking to use?
 - User can pick *Add from WeChat* [powcoder](https://powcoder.com)
 - 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-partite graph with k types of nodes
<https://powcoder.com>
 - Example: picture nodes and tag nodes
[Add WeChat powcoder](#)
- How to find nodes similar to node u ?
- Do a Random-Walk with Restarts from node u
 - i.e. teleport set $S = \{u\}$

SimiRank(cont.)

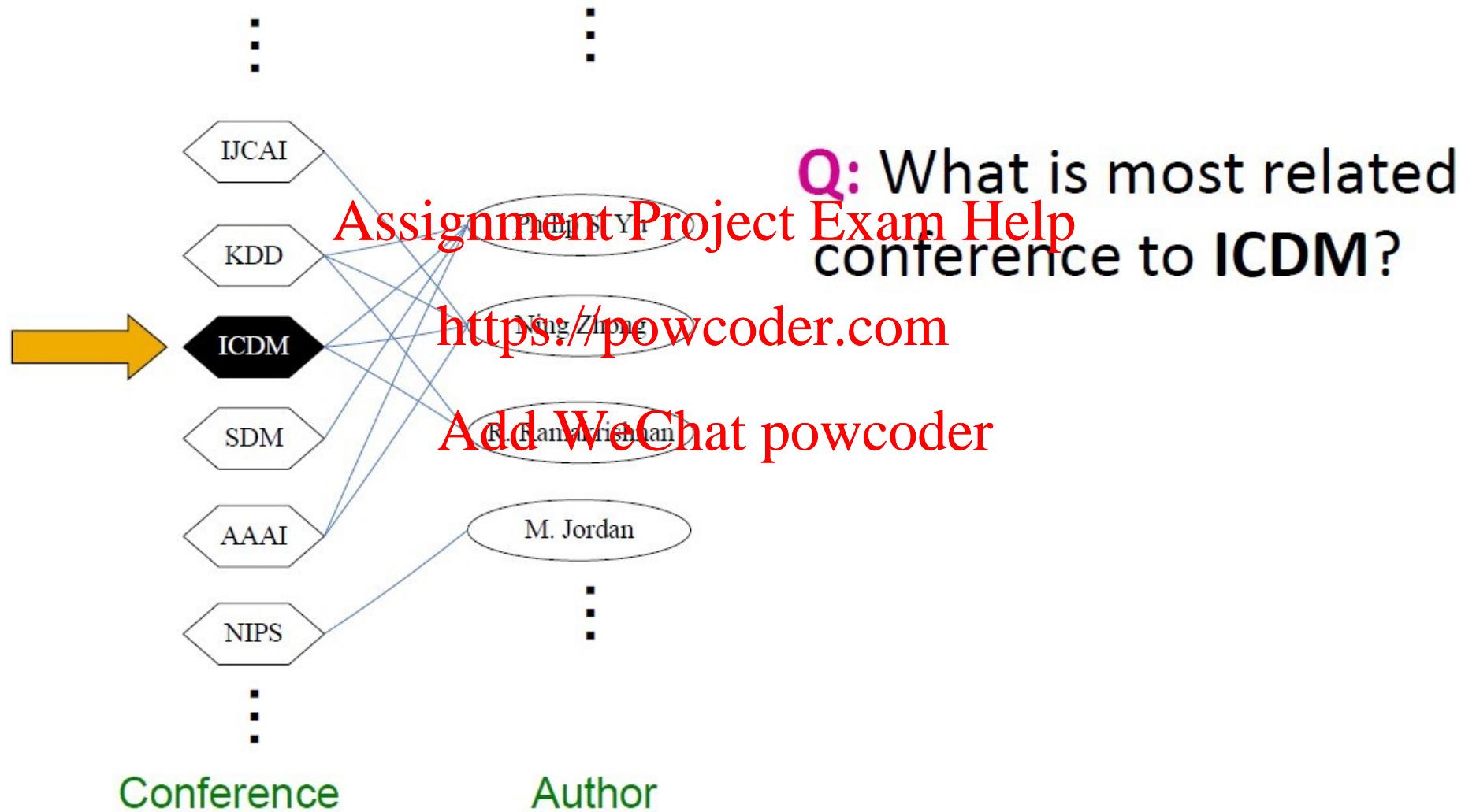
- Resulting scores measures similarity/proximity to node u
- Generally applicable to social networks (typically undirected graphs)
- Problems:
 - Must be done once for each node u
 - Suitable for sub-Web-scale applications

Assignment Project Exam Help

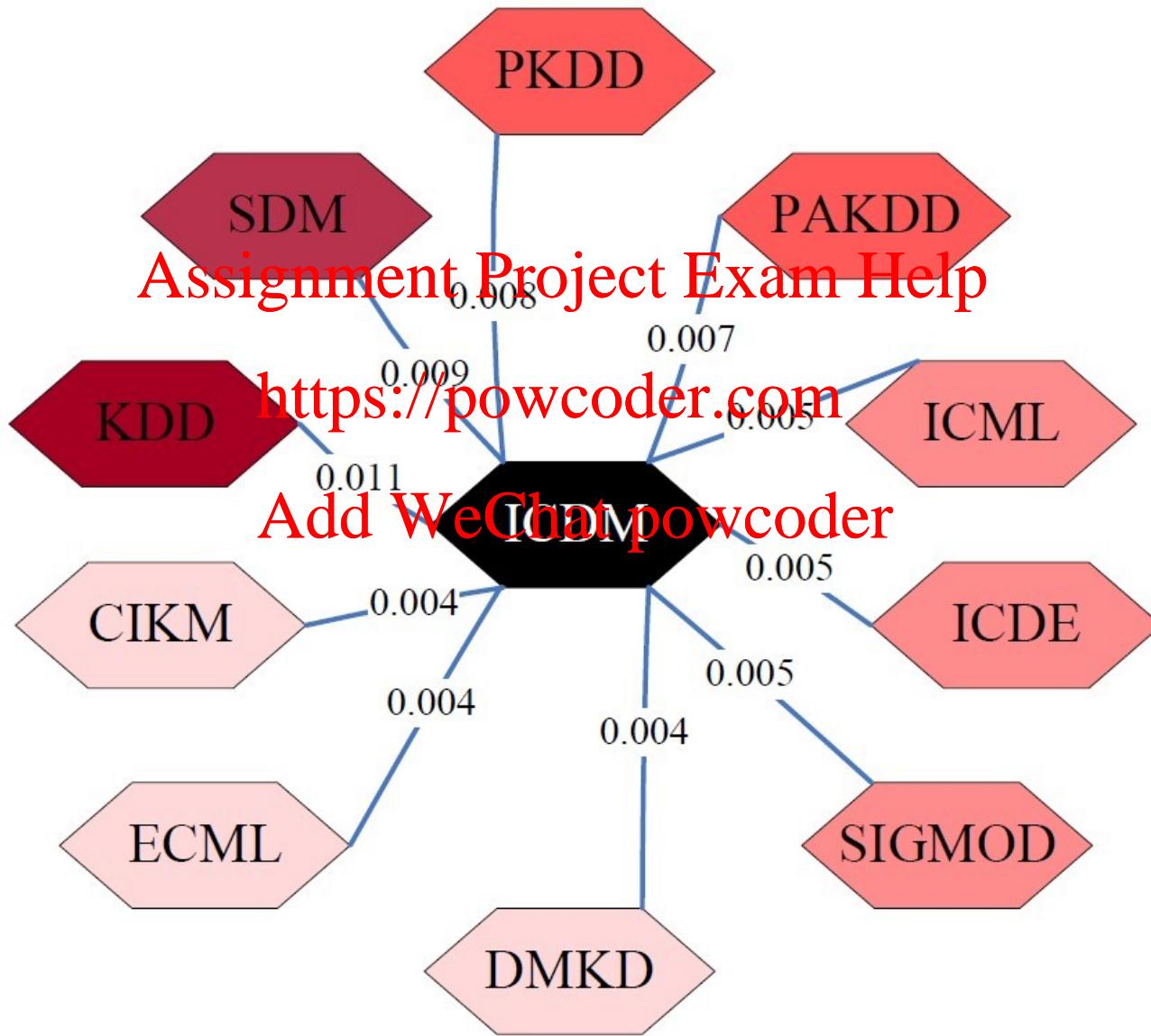
<https://powcoder.com>

Add WeChat powcoder

SimRank: Example



SimRank: Example



Outline

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

What is Web Spam?

- Spammering:
 - 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
- 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 with pages containing those words

Web Search

- Early page ranking
 - Attempt to order pages matching a search query by “importance”
<https://powcoder.com>
 - First search engines considered:
 - Number of times query words appeared
 - 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 – whether they wanted to be there or not
<https://powcoder.com>
- 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?
Assignment Project Exam Help
 - Add the word movie 1,000 times to your page, set text color to the background color, 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 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 loses
 - 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: [Add WeChat powcoder](#)
 - 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
 - [Assignment Project Exam Help](https://powcoder.com)
- Spam farms were developed to concentrate PageRank on a single page
 - Add WeChat powcoder
- Link spam:
 - Creating link structures that boost PageRank of a particular page

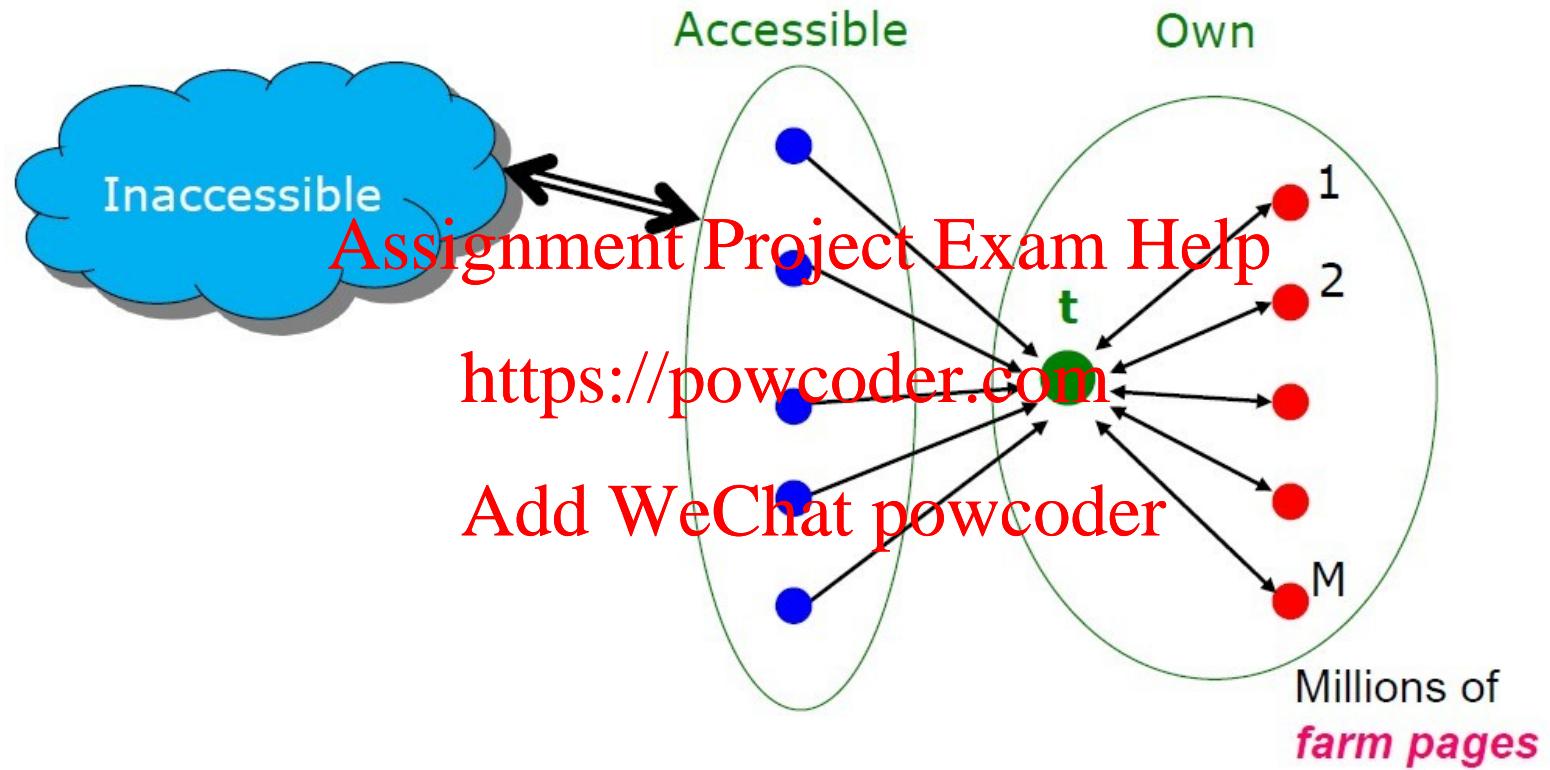
Link Spamming

- Three kinds of web pages from a spammer's point of view
 - Inaccessible pages
<https://powcoder.com>
 - Accessible pages
 - E.g. blog comments pages
 - Spammer can post links to his pages
 - Own pages
 - Completely controlled by spammer
 - May span multiple domain names

Link Farms

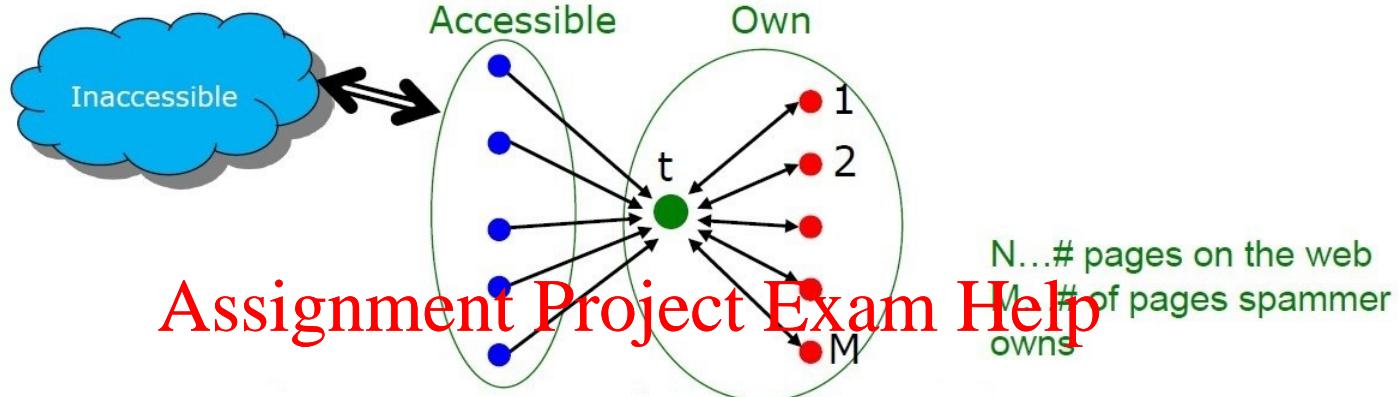
- Spammer's goal:
 - Maximize the PageRank of target page t
[Assignment](#) [Project](#) [Exam](#) [Help](#)
- Technique:
 - Get as many links from accessible pages as possible to target page t
<https://powcoder.com>
[Add WeChat](#) [powcoder](#)
 - Construct “link farm” to get PageRank multiplier effect

Link Farms



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

Analysis

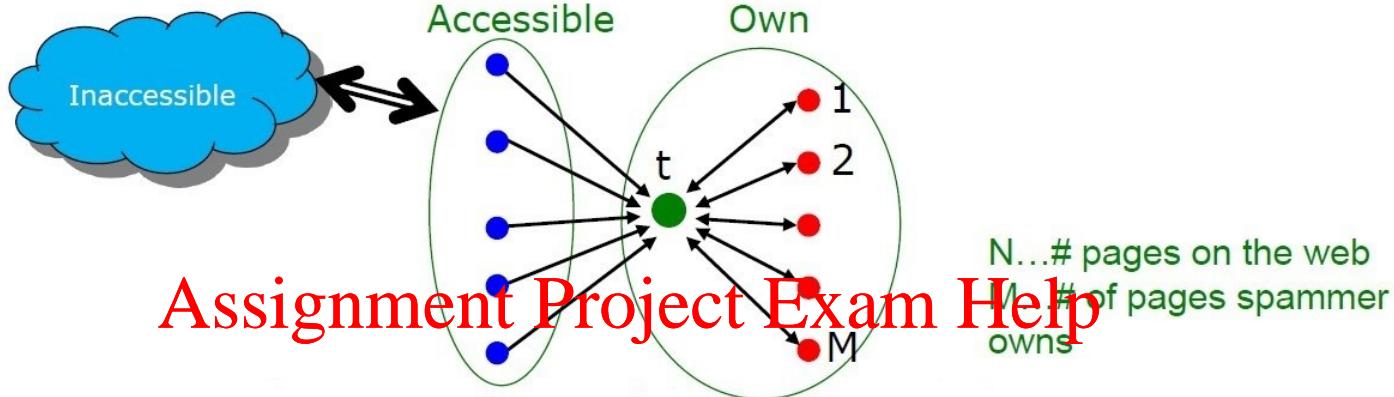


<https://powcoder.com>

- x : PageRank contributed by accessible pages
- y : PageRank of target page t
- Rank of each “farm” page = $\frac{\beta y}{M} + \frac{1 - \beta}{N}$
- $$\begin{aligned} y &= x + \beta M \left[\frac{\beta y}{M} + \frac{1 - \beta}{N} \right] + \frac{1 - \beta}{N} \\ &= x + \beta^2 y + \frac{\beta(1 - \beta)M}{N} + \frac{1 - \beta}{N} \end{aligned}$$

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



<https://powcoder.com>

- $y = \frac{x}{1 - \beta^2} + c \frac{M}{N}$, where $c = \frac{\beta}{1 + \beta}$ Add WeChat powcoder
- 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: 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”
Assignment Project Exam Help
(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**
Assignment Project Exam Help
- Perform a topic-sensitive PageRank with
https://powcoder.com
teleport set = trusted pages
Add WeChat powcoder
 - Propagate trust through **links**:
 - 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 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](https://powcoder.com) page gets adequate trust rank, so need to make all good pages 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://powcoder.com>
 - Pick the top k pages by PageRank
Add WeChat powcoder
 - Theory is that you can't get a bad 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](#)
- Complementary view:
<https://powcoder.com>
 - What fraction of a page's PageRank comes from spam pages?
[Add WeChat](#) [powcoder](#)
- 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 teleport into trusted pages only <https://powcoder.com>

Assignment Project Exam Help
Add WeChat powcoder
- Then: What fraction of a page's PageRank comes from spam pages?

$r_p - r_p^+$

$$p = \frac{r_p^-}{r_p}$$
- Spam mass of

One-slide Takeaway

- Web as a Graph
 - Denote the web structure as a graph
- PageRankAssignment Project Exam Help
 - PageRank score <https://powcoder.com> measures the importance 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 Comparison of Approaches to
Personalizing PageRank:
[https://powcoder.com
Add WeChat powcoder
http://www-cs-students.stanford.edu/~taherh/
/papers/comparison.pdf](http://www-cs-students.stanford.edu/~taherh/papers/comparison.pdf)

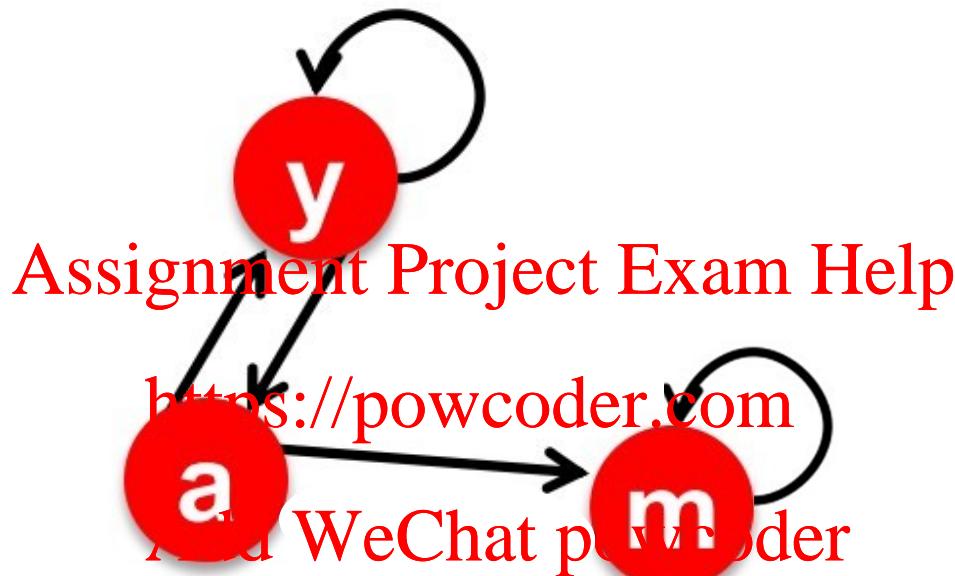
Further Reading

- 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 PageRank:
<https://powcoder.com>
http://link.springer.com/content/pdf/10.1007%2F11735106_22.pdf
~~Add WeChat powcoder~~

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.wikipedia.org/wiki/PageRank>
- http://en.wikipedia.org/wiki/HITS_algorithm
Add WeChat powcoder

In-class Practice



- 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