Massive Link Analysis

1

What's the Mechanism Behind Google?

How google return such kind of rankings (e.g. cuhk homepage first then Wikipedia page)? One important factor is PageRank score.

The Chinese University of Hong Kong (CUHK) is a top Hong Kong university with strong research emphasis. The university aims to bring together China and the ...

MyCUHK MyCUHK not only is your personal portal to a wide range of ...

University Library System

CUHK Library Catalogue - Opening Hours - Title - ...

Job Vacancies

Job Vacancies. ... About Personnel Office · Job Vacancies ... Studying at CUHK If you would like to study at a Hong Kong university, The Chinese ...

Contact Us CUHK address and contact telephone and fax number as ...

Campus Map Users should beware that this map may not be drawn to scale and ...

Search cuhk.edu.hk

Chinese University of Hong Kong - Wikipedia, the free encyclopedia en.wikipedia.org/wiki/Chinese_University_of_Hong_Kong *

The Chinese University of Hong Kong (Abbreviation: CUHK or Chinese University) is the second oldest university in Hong Kong. The CUHK is internationally Tradition and history - Funding - Libraries and museum - Faculties

The Chinese University of Hong Kong | Top Universities

www.topuniversities.com/universities/chinese-university-hong-kong ▼ Log in · Top Universities- Worldwide University rankings, guides & events ... The Chinese University of Hong Kong. Primary tabs... Hong Kong S.A.R., China ...

Chinese University of Hong Kong

maps.google.com/.../place?...chinese+university...hong+kong...chinese+uni... Score: 21 / 30 · 19 Google reviews · Write a review

◇ 大學道 Hong Kong

The Chinese University of Hong Kong | Coursera https://www.coursera.org/cuhk *

The Chinese University of Hong Kong. Founded in 1963, The Chinese University of Hong Kong (CUHK) is a forward looking comprehensive research university ...

The Chinese University of Hong Kong | Worldwide Universities ...



Chinese University of Hong Kong

Directions

The Chinese University of Hong Kong is the second oldest university in Hong Kong. The CUHK is internationally known for its achievements in physics and mathematics. Wikipedia

Address: 大學道, Hong Kong

Enrollment: 14,315 (2010)

Founded: 1963

Colors: Gold, Purple

Motto: 博文約禮, To broaden one's intellectual horizon and keep within the bounds of propriety.

People also search for



Feedback / More info

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
- Topic-Specific PageRank
- Trust-Rank
- Further Reading

Outline

- Web as a Graph
- PageRank
- Topic-Specific PageRank
- Trust-Rank
- Further Reading

Web as a Graph

- Web as a directed graph:
 - Nodes: Web pages
 - Edges: Hyperlinks



Web as a Directed Graph



Broad Question

- How to organize the Web?
- First try: Human curated Web directories

 Yahoo, DMOZ, LookSmart
- 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"?

- Trick: Trustworthy pages may point to each other

- What is the "best" answer to query "newspaper"?
 - 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"
 - <u>www.cuhk.edu.hk</u> vs. <u>www.joe-schmoe.com</u>
- There is large diversity 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

Outline

- Web as a Graph
- PageRank
- Topic-Specific PageRank
- Trust-Rank
- Further Reading

Links as Votes

- Idea: Links as votes
 - Page is more important if it has more links
 - In-coming links? Out-going links?
- Think of in-links as votes:
 - www.cuhk.edu.hk has 15,432 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



Simple Recursive Formulation

- Each link's vote is proportional to the importance of its source page
- If page j with importance r_j has n out-links, 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 r_j
- Define a "rank" for page j

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

 $d_i \dots$ out-degress of node i





Solving the Flow Equations

- 3 equations, 3 unknowns, no constants
 - No unique solution
 - All solutions equivalent modulo the scale factor
- Additional constraint forces uniqueness:
 - $r_y + r_a + r_m = 1$ $\text{Solution: } r_y = \frac{2}{5}, r_a = \frac{2}{5}, r_m = \frac{1}{5}$ Flow equations: $r_y = r_y/2 + r_a/2$ $r_a = r_y/2 + r_m$ $r_m = r_a/2$
- 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
 - 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 $r_j = \sum_{i \to j} \frac{r_i}{d_i}$ $r = M \cdot r$

Example

- Remember the flow equation: $r_j = \sum_{i \to j} \frac{r_i}{d_i}$
- Flow equation in the matrix form:

 $r = M \cdot r$

- 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!

NOTE: \boldsymbol{x} is an eigenvector with the corresponding eigenvalue λ if:

 $Ax = \lambda x$

– The method is called **Power iteration**.

Example: Flow Equation & M



	У	a	m
у	1/2	1/2	0
a	1/2	0	1
m	0	1/2	0

 $r = M \cdot r$



 $r_{y} = r_{y}/2 + r_{a}/2$ $r_{a} = r_{y}/2 + r_{m}$ $r_{m} = r_{a}/2$

Power Iteration Method

- Given a web graph with *n* nodes, where the nodes are pages and edges are hyperlinks
- Power Iteration: a simple iterative scheme
 - 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

Why Power Iteration Works?

- Power iteration:
 - A method for finding dominant eigenvector (the vector corresponding to the largest eigenvalue)

•
$$r^{(1)} = M \cdot r^{(0)}$$

•
$$\mathbf{r}^{(2)} = \mathbf{M} \cdot \mathbf{r}^{(1)} = \mathbf{M}(\mathbf{Mr}^{(1)}) = \mathbf{M}^2 \cdot \mathbf{r}^{(0)}$$

 $\mathbf{r}^{(3)} = \mathbf{M} \cdot \mathbf{r}^{(2)} = \mathbf{M}(\mathbf{M}^2 \mathbf{r}^{(0)}) = \mathbf{M}^3 \cdot \mathbf{r}^{(0)}$

- Claim:
 - Sequence $\mathbf{M} \cdot \mathbf{r^{(0)}}, \mathbf{M^2} \cdot \mathbf{r^{(0)}}, \dots, \mathbf{M^k} \cdot \mathbf{r^{(0)}}, \dots$ approaches the dominant eigenvector of \boldsymbol{M}

Why Power Iteration Works?

• Proof:

- Assume *M* has *n* linearly independent eigenvectors $x_1, x_2, ..., x_n$ with corresponding eigenvalues $\lambda_1, \lambda_2, ..., \lambda_n$, where $\lambda_1 > \lambda_2 > ... > \lambda_n$

- Vectors \mathbf{x}_1 , \mathbf{x}_2 , ..., \mathbf{x}_n form a basis and thus we can write: $r^{(0)} = c_1 x_1 + c_2 x_2 + \ldots + c_n x_n$

$$Mr^{(0)} = M(c_1x_1 + c_2x_2 + \dots + c_nx_n)$$

= $c_1(Mx_1) + c_2(Mx_2) + \dots + c_n(Mx_n)$
= $c_1(\lambda_1x_1) + c_2(\lambda_2x_2) + \dots + c_n(\lambda_nx_n)$

Repeated multiplication on both sides:

$$M^{k}r^{(0)} = c_{1}(\lambda_{1}^{k}x_{1}) + c_{2}(\lambda_{2}^{k}x_{2}) + \ldots + c_{n}(\lambda_{n}^{k}x_{n})$$

Why Power Iteration Works?

• Proof: (cont.)

- Repeated multiplication on both sides produces

$$M^{k}r^{(0)} = \lambda_{1}^{k} \left[c_{1}x_{1} + c_{2} \left(\frac{\lambda_{2}}{\lambda_{1}}\right)^{k} x_{2} + \ldots + c_{n} \left(\frac{\lambda_{2}}{\lambda_{1}}\right)^{k} x_{n} \right]$$

- 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 $k \to \infty$ - 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 ith coordinate is the prob. that the surfer is at page i at time t
 - So p(t) is a probability distribution over pages

The Stationary Distribution

• Where is the surfer at time t+1

- Follows a link uniformly at random

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

• Suppose the random walk reaches 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?
 - Does it converge to what we want?
 - Are results reasonable?

Does This Converge?



Does it Converge to What We Want?



PageRank: Problems

- Two problems:
 - Spider traps: all out-links are within the group
 - Eventually spider traps absorb all importance
 - Some pages are dead ends (have no out-links)
 - Such pages cause importance to "leak out"

Problem: Spider Traps

- Power Iteration:
 - Set $r_j = 1$ - $r_j = \sum_{i \to j} \frac{r_i}{d_i}$
 - And iterate
- Example



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

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 are in the range 0.8 to 0.9
- Surfer will teleport out of spider trap within a few time steps



33

Problem: Dead Ends

• Power Iteration:

- Set
$$r_j = 1$$

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

- And iterate
- Example



	у	a	m
у	1/2	1/2	0
a	1/2	0	0
m	0	1/2	0

$$r_{y} = r_{y}/2 + r_{a}/2$$
$$r_{a} = r_{y}/2$$
$$r_{m} = r_{a}/2$$

r_v		1/3	2/6	3/12	5/24		0
r _a	=	1/3	1/6	2/12	3/24	•••	0
r _m		1/3	1/6	1/12	2/24		0
		Iteratio	n 0, 1, 2,				

Solution: Always Teleport

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

Adjust matrix accordingly



Why Teleports Solve the Problem?

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

- Markov chains
 - Set of states X
 - 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, the power method applied to a Markov transition matrix P will converge to a unique positive stationary vector as long as P is stochastic, irreducible and aperiodic.

Make M Stochastic

- Stochastic: Every column sums to 1
- A possible solution: add green links

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

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

• $e \dots$ vector of all 1s



_	У	a	m	
y	1/2	1/2	1/3	
a	1/2	0	1/3	
m	0	1/2	1/3	

$$r_y = r_y/2 + r_a/2 + r_m/3$$

 $r_a = r_y/2 + r_m/3$
 $r_m = r_a/2 + r_m/3$

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.
- 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
- A possible solution: add green links



Solution: Random Jumps

- Google's solution that does it all:
 Makes M stochastic, aperiodic, irreducible
- At each step, random surfer has two options:
 With probability β, follow a link at random
 - With probability 1β , jump to some random page
- PageRank equation [Brin-Page,98]

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

This formulation assumes that **M** has no dead ends. We can either preprocess matrix **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]

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

• The Google Matrix A:

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

- 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)

Random Teleports ($\beta=0.8$)



In-class Practice

• Go to Practice

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_i out-links
 - 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 $A_{ij} = \beta M_{ij} + \frac{1-\beta}{N}$

•
$$r_i = \sum_{j=1}^N A_{ij} \cdot r_j$$

•
$$r_i = \sum_{j=1}^{N} \left[\beta M_{ij} + \frac{1-\beta}{N} \right] \cdot r_j$$

 $= \sum_{j=1}^{N} \beta M_{ij} \cdot r_j + \frac{1-\beta}{N} \sum_{j=1}^{N} r_j$
 $= \sum_{j=1}^{N} \beta M_{ij} \cdot r_j + \frac{1-\beta}{N}$, since $\sum r_j = 1$

• So we get:
$$r = \beta M \cdot r + \left[\frac{1-\beta}{N}\right]_N$$

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

 $[x]_N \dots$ a vector of length *N* with all entries x

Spare Matrix Formulation

- We just rearranged the PageRank equation $r = \beta M \cdot r + \left[\frac{1-\beta}{N}\right]_N$
- M is a sparse matrix! (with no dead-ends)
 10 links per node, approx. 10N entries
- 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 **G** and parameter β
 - Directed graph ${\bf G}$ with ${\bf spider \ traps}$ and ${\bf dead \ ends}$
 - Parameter β
- **Output**: PageRank vector r

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

$$- \text{ do:}$$

$$* \forall j: r'_j^{(t)} = \sum_{i \to j} \beta \frac{r_i^{(t-1)}}{d_i}$$

$$r'_j^{(t)} = 0 \text{ if in-deg. of } j \text{ is } 0$$

$$* \text{ Now re-insert the leaked PageRank:}$$

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

$$* t = t + 1$$

$$- \text{ while } \sum_j |r_j^{(t)} - r_j^{(t-1)}| > \epsilon$$

Sparse Matrix Encoding

- Encode sparse matrix using only nonzero entries
 - Space proportional roughly to number of links
 - Say 10N, or 4*10*1 billion = 40GB
 - 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 of power-iteration is:
 - Initialize all entries of r^{new} to $(1 \beta)/N$
 - For each page p (of out-degree n):

0

1

2

3

4

5

6

• Read into memory: *p*, *n*, dest₁, ..., dest_n, *r*^{old} (*p*)



• For j=1...n: r^{new} (dest_i) += $\beta r^{old}(p)/n$

Analysis

- Assume enough RAM to fit r^{new} into memory - Store r^{old} and matrix M on disk
- In each iteration, we have to:
 - Read r^{old} and M
 - Write r^{new} back to disk
 - $|O \cos t = 2|r| + |M|$
- Question:

– What if we could not even fit r^{new} in memory

Block-based Update Algorithm



Analysis of Block Update

- Similar to nested-loop join in databases
 - Break r^{new} into k blocks that fit in memory
 - 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|
- 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
0	4	0, 1
1	3	0
2	2	1





0	4	3
2	2	3



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

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

In-class Practice



• Compute the final PageRank Score of the given graph.

