

CS345 Data Mining

Link Analysis Algorithms Page Rank

Anand Rajaraman, Jeffrey D. Ullman

Link Analysis Algorithms

- Page Rank
- Hubs and Authorities
- Topic-Specific Page Rank
- Spam Detection Algorithms
- Other interesting topics we won't cover
 - Detecting duplicates and mirrors
 - Mining for communities
 - Classification
 - Spectral clustering

Ranking web pages

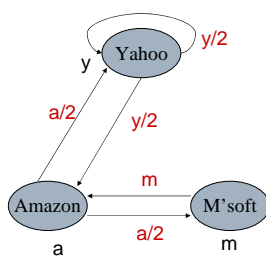
- Web pages are not equally "important"
 - www.joe-schmoe.com v www.stanford.edu
- Inlinks as votes
 - www.stanford.edu has 23,400 inlinks
 - www.joe-schmoe.com has 1 inlink
- Are all inlinks equal?
 - Recursive question!

Simple recursive formulation

- Each link's vote is proportional to the importance of its source page
- If page P with importance x has n outlinks, each link gets x/n votes

Simple "flow" model

The web in 1839



$$\begin{aligned}y &= y/2 + a/2 \\ a &= y/2 + m \\ m &= a/2\end{aligned}$$

Solving the flow equations

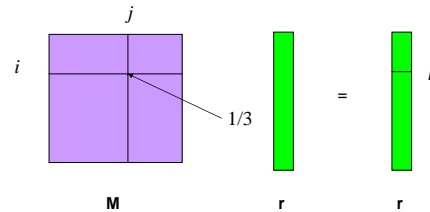
- 3 equations, 3 unknowns, no constants
 - No unique solution
 - All solutions equivalent modulo scale factor
- Additional constraint forces uniqueness
 - $y+a+m = 1$
 - $y = 2/5, a = 2/5, m = 1/5$
- Gaussian elimination method works for small examples, but we need a better method for large graphs

Matrix formulation

- Matrix \mathbf{M} has one row and one column for each web page
- Suppose page j has n outlinks
 - If $j \rightarrow i$, then $M_{ij} = 1/n$
 - Else $M_{ij} = 0$
- \mathbf{M} is a **column stochastic matrix**
 - Columns sum to 1
- Suppose \mathbf{r} is a vector with one entry per web page
 - r_i is the importance score of page i
 - Call it the **rank vector**

Example

Suppose page j links to 3 pages, including i

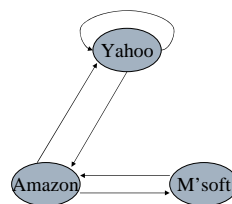


Eigenvector formulation

- The flow equations can be written

$$\mathbf{r} = \mathbf{M}\mathbf{r}$$
- So the rank vector is an eigenvector of the stochastic web matrix
 - In fact, its first or principal eigenvector, with corresponding eigenvalue 1

Example



$$\begin{aligned} y &= y/2 + a/2 \\ a &= y/2 + m \\ m &= a/2 \end{aligned}$$

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

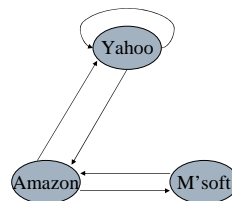
$$\mathbf{r} = \mathbf{M}\mathbf{r}$$

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

Power Iteration method

- Simple iterative scheme (aka **relaxation**)
- Suppose there are N web pages
- Initialize: $\mathbf{r}^0 = [1/N, \dots, 1/N]^T$
- Iterate: $\mathbf{r}^{k+1} = \mathbf{M}\mathbf{r}^k$
- Stop when $|\mathbf{r}^{k+1} - \mathbf{r}^k|_1 < \epsilon$
 - $|\mathbf{x}|_1 = \sum_{i=1}^N |x_i|$ is the L_1 norm
 - Can use any other vector norm e.g., Euclidean

Power Iteration Example



$$\begin{matrix} y & 1/3 & 1/3 & 5/12 & 3/8 & 2/5 \\ a & 1/3 & 1/2 & 1/3 & 11/24 & \dots & 2/5 \\ m & 1/3 & 1/6 & 1/4 & 1/6 & & 1/5 \end{matrix}$$

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

Random Walk Interpretation

- Imagine a **random web surfer**
 - At any time t , surfer is on some page P
 - At time $t+1$, the surfer follows an outlink from P uniformly at random
 - Ends up on some page Q linked from P
 - Process repeats indefinitely
- Let $\mathbf{p}(t)$ be a vector whose i^{th} component is the probability that the surfer is at page i at time t
 - $\mathbf{p}(t)$ is a probability distribution on pages

The stationary distribution

- Where is the surfer at time $t+1$?
 - Follows a link uniformly at random
 - $\mathbf{p}(t+1) = \mathbf{M}\mathbf{p}(t)$
- Suppose the random walk reaches a state such that $\mathbf{p}(t+1) = \mathbf{M}\mathbf{p}(t) = \mathbf{p}(t)$
 - Then $\mathbf{p}(t)$ is called a **stationary distribution** for the random walk
- Our rank vector \mathbf{r} satisfies $\mathbf{r} = \mathbf{M}\mathbf{r}$
 - So it is a stationary distribution for the random surfer

Existence and Uniqueness

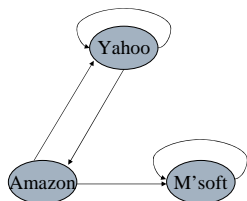
A central result from the theory of random walks (aka Markov processes):

For graphs that satisfy certain conditions, the stationary distribution is unique and eventually will be reached no matter what the initial probability distribution at time $t = 0$.

Spider traps

- A group of pages is a **spider trap** if there are no links from within the group to outside the group
 - Random surfer gets trapped
- Spider traps violate the conditions needed for the random walk theorem

Microsoft becomes a spider trap



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

y	=	1	1	3/4	5/8	0
a		1	1/2	1/2	3/8	...
m		1	3/2	7/4	2	3

Random teleports

- The Google solution for spider traps
- At each time step, the random surfer has two options:
 - With probability β , follow a link at random
 - With probability $1-\beta$, jump to some page uniformly at random
 - 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

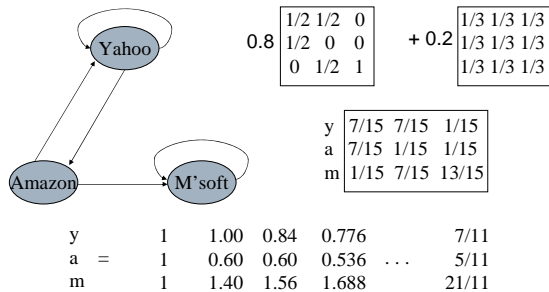
Matrix formulation

- Suppose there are N pages
 - Consider a page j , with set of outlinks $O(j)$
 - We have $M_{ij} = 1/|O(j)|$ 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 with probability $(1-\beta)/N$
 - reducing the probability of following each outlink from $1/|O(j)|$ to $\beta/|O(j)|$
 - Equivalent: tax each page a fraction $(1-\beta)$ of its score and redistribute evenly

Page Rank

- Construct the $N \times N$ matrix \mathbf{A} as follows
 - $A_{ij} = \beta M_{ij} + (1-\beta)/N$
- Verify that \mathbf{A} is a stochastic matrix
- The **page rank vector** \mathbf{r} is the principal eigenvector of this matrix
 - satisfying $\mathbf{r} = \mathbf{A}\mathbf{r}$
- Equivalently, \mathbf{r} is the stationary distribution of the random walk with teleports

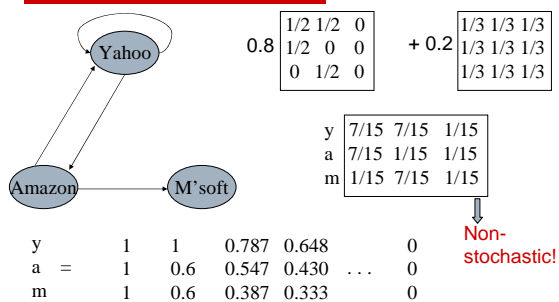
Previous example with $\beta=0.8$



Dead ends

- Pages with no outlinks are "**dead ends**" for the random surfer
 - Nowhere to go on next step

Microsoft becomes a dead end



Dealing with dead-ends

- Teleport
 - Follow random teleport links with probability 1.0 from dead-ends
 - Adjust matrix accordingly
- Prune and propagate
 - Preprocess the graph to eliminate dead-ends
 - Might require multiple passes
 - Compute page rank on reduced graph
 - Approximate values for deadends by propagating values from reduced graph

Computing page rank

- Key step is matrix-vector multiply
 - $r^{new} = Ar^{old}$
- 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!

Sparse matrix formulation

- Although A is a dense matrix, it is obtained from a sparse matrix M
 - 10 links per node, approx 10N entries
- We can restate the page rank equation
 - $r = \beta M r + [(1-\beta)/N]_N$
 - $[(1-\beta)/N]_N$ is an N -vector with all entries $(1-\beta)/N$
- So in each iteration, we need to:
 - Compute $r^{new} = \beta M r^{old}$
 - Add a constant value $(1-\beta)/N$ to each entry in r^{new}

Sparse matrix encoding

- Encode sparse matrix using only nonzero entries
 - Space proportional roughly to number of links
 - say 10N, or $4 \times 10^9 \times 10 = 40GB$
 - 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

- Assume we have enough RAM to fit r^{new} , plus some working memory
 - Store r^{old} and matrix M on disk

Basic Algorithm:

- Initialize: $r^{old} = [1/N]_N$
- Iterate:
 - **Update:** Perform a sequential scan of M and r^{old} and update r^{new}
 - Write out r^{new} to disk as r^{old} for next iteration
 - Every few iterations, compute $|r^{new} - r^{old}|$ and stop if it is below threshold
 - Need to read in both vectors into memory

Update step

Initialize all entries of r^{new} to $(1-\beta)/N$
 For each page p (out-degree n):
 Read into memory: $p, n, dest_1, \dots, dest_n, r^{old}(p)$
 for $j = 1..n$:
 $r^{new}(dest_j) += \beta * r^{old}(p)/n$

r^{new}	src	degree	destination	r^{old}
0	0	3	1, 5, 6	0
1	1	4	17, 64, 113, 117	1
2	2	2	13, 23	2
3				3
4				4
5				5
6				6

Analysis

- In each iteration, we have to:
 - Read r^{old} and M
 - Write r^{new} back to disk
 - IO Cost = $2|r| + |M|$
- What if we had enough memory to fit both r^{new} and r^{old} ?
- What if we could not even fit r^{new} in memory?
 - 10 billion pages

Block-based update algorithm

r^{new}	src	degree	destination	r^{old}
0	0	4	0, 1, 3, 5	0
1	1	2	0, 5	1
2	2	2	3, 4	2
3				3
4				4
5				5

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-Stripe Update algorithm

r^{new}	src	degree	destination	r^{old}
0	0	4	0, 1	0
1	1	3	0	1
2	2	2	1	2
3	0	4	3	3
4	2	2	3	4
5	0	4	5	5
	1	3	5	
	2	2	4	

Block-Stripe Analysis

- Break M into stripes
 - Each stripe contains only destination nodes in the corresponding block of r^{new}
- Some additional overhead per stripe
 - But usually worth it
- Cost per iteration
 - $|M|(1+\epsilon) + (k+1)|r|$

Next

- Topic-Specific Page Rank
- Hubs and Authorities
- Spam Detection