

Link Analysis

CSE 454 Advanced Internet Systems
University of Washington

10/21/2010 5:30 PM

Copyright © 2000-2009 D.S.Weld

1

Ranking Search Results

- **TF / IDF or BM25**
- **Tag Information**
 - Title, headers
- **Font Size / Capitalization**
- **Anchor Text on Other Pages**
- **Classifier Predictions**
 - Spam, Adult, Review, Celebrity, ...
- **Link Analysis**
 - HITS – (Hubs and Authorities)
 - PageRank

10/21/2010 5:30 PM

Copyright © 2000-2009 D.S.Weld

2

Pagerank Intuition

Think of Web as a big graph.

Suppose surfer keeps **randomly** clicking on the links.
Importance of a page = probability of being on the page

Derive **transition matrix** from adjacency matrix

Suppose $\exists N$ forward links from page P
Then the probability that surfer clicks on any one is $1/N$

10/21/2010 5:30 PM

Copyright © 2000-2009 D.S.Weld

3

Matrix Representation

Let M be an $N \times N$ matrix

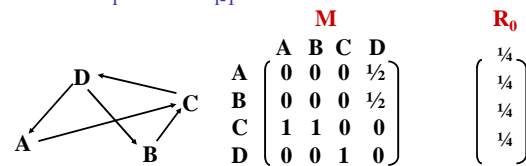
$$m_{uv} = 1/N_v \text{ if page } v \text{ has a link to page } u$$

$$m_{uv} = 0 \text{ if there is no link from } v \text{ to } u$$

Let R_0 be the initial rank vector

Let R_i be the $N \times 1$ rank vector for i^{th} iteration

Then $R_i = M \times R_{i-1}$



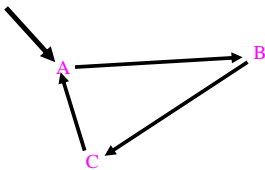
10/21/2010 5:30 PM

Copyright © 2000-2009 D.S.Weld

4

Problem: Page Sinks.

- **Sink = node (or set of nodes) with no out-edges.**
- **Why is this a problem?**



10/21/2010 5:30 PM

Copyright © 2000-2009 D.S.Weld

5

Solution to Sink Nodes

Let:

$(1-c)$ = chance of random transition from a sink.

N = the number of pages

$$K = \begin{pmatrix} \dots \\ \dots 1/N \dots \\ \dots \end{pmatrix}$$

$$M^* = cM + (1-c)K$$

$$R_i = M^* \times R_{i-1}$$

10/21/2010 5:30 PM

Copyright © 2000-2009 D.S.Weld

6

Computing PageRank - Example

$$M = \begin{matrix} & \begin{matrix} A & B & C & D \end{matrix} \\ \begin{matrix} A \\ B \\ C \\ D \end{matrix} & \begin{pmatrix} 0 & 0 & 0 & \frac{1}{2} \\ 0 & 0 & 0 & \frac{1}{2} \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix} \end{matrix}$$

$$M^* = \begin{pmatrix} 0.05 & 0.05 & 0.05 & 0.45 \\ 0.05 & 0.05 & 0.05 & 0.45 \\ 0.85 & 0.85 & 0.05 & 0.05 \\ 0.05 & 0.05 & 0.85 & 0.05 \end{pmatrix}$$

$R_0 = \begin{pmatrix} \frac{1}{4} \\ \frac{1}{4} \\ \frac{1}{4} \\ \frac{1}{4} \end{pmatrix}$

$R_{30} = \begin{pmatrix} 0.176 \\ 0.176 \\ 0.332 \\ 0.316 \end{pmatrix}$

10/21/2010 5:30 PM Copyright © 2000-2009 D.S.Weld 7

Ooops

- What About Sparsity?

$$M^* = \begin{pmatrix} 0.05 & 0.05 & 0.05 & 0.45 \\ 0.05 & 0.05 & 0.05 & 0.45 \\ 0.85 & 0.85 & 0.05 & 0.05 \\ 0.05 & 0.05 & 0.85 & 0.05 \end{pmatrix}$$

$M^* = cM + (1-c)K$

$$K = \begin{pmatrix} \dots & & & \\ \dots & \frac{1}{N} & \dots & \\ & & \dots & \\ \dots & & & \dots \end{pmatrix}$$

Copyright © 2000-2009 D.S.Weld

Authority and Hub Pages (1)

- A page is a good **authority** (with respect to a given query) if it is pointed to by many good hubs (with respect to the query).
- A page is a good **hub page** (with respect to a given query) if it points to many good authorities (for the query).
- Good authorities & hubs reinforce

10/21/2010 5:30 PM Copyright © 2000-2009 D.S.Weld 9

Authority and Hub Pages (2)

Authorities and hubs for a query *tend* to form a bipartite subgraph of the web graph.

(A page can be a good authority *and* a good hub)

10/21/2010 5:30 PM Copyright © 2000-2009 D.S.Weld 10

Linear Algebraic Interpretation

- PageRank = principle eigenvector of M^*
 - in limit
- HITS = principle eigenvector of $M^* \times (M^*)^T$
 - Where $[]^T$ denotes transpose $\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}^T = \begin{pmatrix} 1 & 3 \\ 2 & 4 \end{pmatrix}$
- **Stability**
 - Small changes to graph \rightarrow small changes to weights.
 - Can prove PageRank is stable
 - And HITS isn't

10/21/2010 5:30 PM Copyright © 2000-2009 D.S.Weld 11

Stability Analysis (Empirical)

- Make 5 subsets by deleting 30% randomly

1	1	3	1	1	1
2	2	5	3	3	2
3	3	12	6	6	3
4	4	52	20	23	4
5	5	171	119	99	5
6	6	135	56	40	8
7	10	179	159	100	7
8	8	316	141	170	6
9	9	257	107	72	9
10	13	170	80	69	18

- PageRank much more stable

10/21/2010 5:30 PM Copyright © 2000-2009 D.S.Weld 12

Practicality

- **Challenges**
 - M no longer sparse (don't represent explicitly!)
 - Data too big for memory (be sneaky about disk usage)
- **Stanford Version of Google :**
 - 24 million documents in crawl
 - 147GB documents
 - 259 million links
 - Computing pagerank “few hours” on single 1997 workstation
- **But How?**
 - Next discussion from Haveliwala paper...

10/21/2010 5:30 PM

Copyright © 2000-2009 D.S.Weld

13

Efficient Computation: Preprocess

- **Remove ‘dangling’ nodes**
 - Pages w/ no children
- **Then repeat process**
 - Since now more danglers
- **Stanford WebBase**
 - 25 M pages
 - 81 M URLs in the link graph
 - After two prune iterations: 19 M nodes

10/21/2010 5:30 PM

Copyright © 2000-2009 D.S.Weld

14

Representing ‘Links’ Table

- **Stored on disk in binary format**

Source node (32 bit integer)	Outdegree (16 bit int)	Destination nodes (32 bit integers)
0	4	12, 26, 58, 94
1	3	5, 56, 69
2	5	1, 9, 10, 36, 78

- **Size for Stanford WebBase: 1.01 GB**

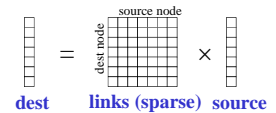
- Assumed to exceed main memory
- (But source & dest assumed to fit)

10/21/2010 5:30 PM

Copyright © 2000-2009 D.S.Weld

15

Algorithm 1



```

∀s Source[s] = 1/N
while residual > τ {
  ∀d Dest[d] = 0
  while not Links.eof() {
    Links.read(source, n, dest1, ... destn)
    for j = 1 ... n
      Dest[destj] = Dest[destj] + Source[source]/n
  }
  ∀d Dest[d] = (1-c) * Dest[d] + c/N /* damping c= 1/N */
  residual = || Source - Dest || /* recompute every few iterations */
  Source = Dest
}

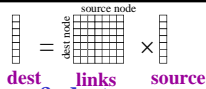
```

10/21/2010 5:30 PM

Copyright © 2000-2009 D.S.Weld

16

Analysis



- **If memory can hold both source & dest**
 - IO cost per iteration is | Links |
 - Fine for a crawl of 24 M pages
 - But web > 8 B pages in 2005 [Google]
 - Increase from 320 M pages in 1997 [NEC study]
- **If memory only big enough to hold just dest...?**
 - Sort Links on source field
 - Read Source sequentially during rank propagation step
 - Write Dest to disk to serve as Source for next iteration
 - IO cost per iteration is | Source | + | Dest | + | Links |
- **But What if memory can't even hold dest?**
 - Random access pattern will make working set = | Dest |
 - Thrash!!!

.....???

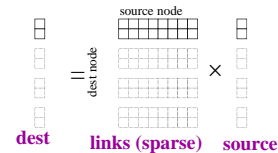
10/21/2010 5:30 PM

Copyright © 2000-2009 D.S.Weld

17

Block-Based Algorithm

- **Partition Dest into B blocks of D pages each**
 - If memory = P physical pages
 - $D < P/2$ since need input buffers for Source & Links
- **Partition (sorted) Links into B files**
 - Links_i only has some of the dest nodes for each source
 - Specifically, Links_i only has dest nodes such that
 - $DD * i <= dest < DD * (i+1)$
 - Where DD = number of 32 bit integers that fit in D pages



10/21/2010 5:30 PM

Copyright © 2000-2009 D.S.Weld

18

Partitioned Link File

Source node (32 bit int)	Outdegr (16 bit)	Num out (16 bit)	Destination nodes (32 bit integer)
0	4	2	12, 26
1	3	1	5
2	5	3	1, 9, 10

Buckets 0-31

0	4	1	58
1	3	1	56
2	5	1	36

Buckets 32-63

0	4	1	94
1	3	1	69
2	5	1	78

Buckets 64-95

10/21/2010 5:30 PM

Copyright © 2000-2009 D.S.Weld

19

Analysis of Block Algorithm

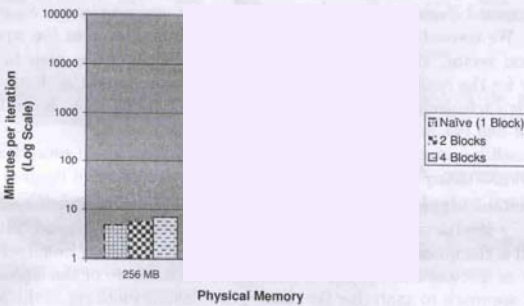
- **IO Cost per iteration =**
 - $B * |Source| + |Dest| + |Links| * (1+e)$
 - e is factor by which Links increased in size
 - Typically 0.1-0.3
 - Depends on number of blocks
- **Algorithm ~ nested-loops join**

10/21/2010 5:30 PM

Copyright © 2000-2009 D.S.Weld

20

Comparing the Algorithms

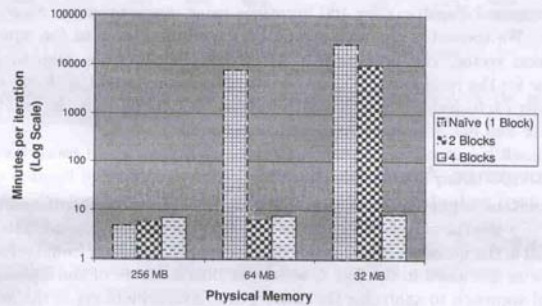


10/21/2010 5:30 PM

Copyright © 2000-2009 D.S.Weld

21

Comparing the Algorithms



10/21/2010 5:30 PM

Copyright © 2000-2009 D.S.Weld

22

Adding PageRank to a SearchEngine

- Weighted sum of importance+similarity with query
- $Score(q, d) = w * sim(q, p) + (1-w) * R(p)$, if $sim(q, p) > 0$
= 0, otherwise
- Where
 - $0 < w < 1$
 - $sim(q, p), R(p)$ must be normalized to $[0, 1]$.

10/21/2010 5:30 PM

Copyright © 2000-2009 D.S.Weld

23

Summary of Key Points

- **PageRank Iterative Algorithm**
- **Sink Pages**
- **Efficiency of computation – Memory!**
 - Don't represent M^* explicitly.
 - Minimize IO Cost.
 - Break arrays into Blocks.
 - Single precision numbers ok.
- **Number of iterations of PageRank.**
- **Weighting of PageRank vs. doc similarity.**

10/21/2010 5:30 PM

Copyright © 2000-2009 D.S.Weld

24