

COMP4434 Big Data Analytics

Lecture 12 PageRank

HUANG Xiao xiaohuang@comp.polyu.edu.hk

PageRank Motivation: How to organize the Web?

- First try: Human curated Web directories
 - Yahoo, DMOZ, LookSmart
- Second try: Web search
 - Information Retrieval investigates: Finding relevant docs in a small and trusted set
 - Newspaper articles, Patents, etc.



<u>But:</u> Web is huge, full of untrusted documents, random things, web spam, etc.

Challenges in Web Search

- (1) Web contains many sources of information Who to "trust"?
 - **Trick:** Trustworthy pages may point to each other!
- (2) 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

Hint: Web as a Directed Graph

- Nodes: Webpages
- Edges: Hyperlinks



Web as a Directed Graph



Ranking Nodes on the Graph

 All web pages are not equally "important" <u>https://xhuang31.github.io</u> vs. <u>https://www.polyu.edu.hk</u>

 There is large diversity in the web-graph node connectivity.
 Let's rank the pages by the link structure!



Example of Node Ranking

- Page Ranking
- Social Ranking
- Paper Ranking

.....

Scholar Ranking

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.stanford.edu has 23,400 in-links
 - <u>https://xhuang31.github.io</u> has 0 in-link
- Are all in-links are equal?
 - Links from important pages count more
 - Recursive question!

Google PageRank

- In-coming links! Out-going links?
- A page with high PageRank value
 - Many pages pointing to it, or
 - There are some pages that point to it and have high PageRank values
- Example:
 - Page C has a higher PageRank than Page E, even though it has fewer links to it
 - The link it has is of a much higher value



Is Page == "Webpage"?

- Born in March 26, 1973
- Found Google at September 4, 1998
- As of Nov 2024, own an estimated net worth of \$163 billion (No.15 Richest)
- Begins from "Larry Page and Sergey Brin developed PageRank at Stanford University in 1996" as part of a research project about a new kind of search engine.



Larry Page Co-founder of Google

Simple Recursive Formulation

- Each link's vote is proportional to the importance of its source page
- If page j with importance PR(j) has n out-links, each link gets
 PR(j) / n votes
- Page j's own importance is the sum of the votes on its in-links

PR(j) = PR(i)/3 + PR(k)/4



How to Represent a Graph

• Graph model G = (V, E)

- V is a set of pages
- *E* is a set of edges
- Each edge $(u, v) \in E$ represents that page u points/references to page v

Adjacent List

- A data structure for a graph
- $Adj[u] = \{v: (u, v) \in E\}$ contains each vertex v being adjacent to u
- Example: Adj[2] = {3, 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" r_i for page j

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

 d_i ... out-degree of node i



"Flow" equations:

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

Solving the Flow Equations

Flow equations:

 $r_v = r_v / 2 + r_a / 2$

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

 $r_m = r_a /2$

 3 equations, 3 unknowns, no constants

- No unique solution
- All solutions equivalent modulo the scale factor
- Additional constraint forces uniqueness:

$$\bullet r_y + r_a + r_m = 1$$

• Solution:
$$r_y = \frac{2}{5}$$
, $r_a = \frac{2}{5}$, $r_m = \frac{1}{5}$

But, we need a better method for large web-size graphs

PageRank: Matrix Formulation

Stochastic adjacency matrix M

Let page *i* has *d_i* out-links

• If
$$i \to j$$
, then $M_{ji} = \frac{1}{d_i}$ else $M_{ji} = 0$

- *M* is a column stochastic matrix
 - 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

$$r = M \cdot r$$



Example

- Remember the flow equation:
- Flow equation in the matrix form

$$\boldsymbol{M}\cdot\boldsymbol{r}=\boldsymbol{r}$$

Suppose page *i* links to 3 pages, including *j*





Example: Flow Equations & M



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

 $r = M \cdot r$

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

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

Power Iteration Method

- Given a web graph with N nodes, where the nodes are pages and edges are hyperlinks
- **Power iteration:** a simple iterative scheme
 - Suppose there are N web pages
 - Initialize: $\mathbf{r}^{(0)} = [1/N, ..., 1/N]^{T}$
 - Iterate: $\mathbf{r}^{(t+1)} = \mathbf{M} \cdot \mathbf{r}^{(t)}$
 - Stop when $|\mathbf{r}^{(t+1)} \mathbf{r}^{(t)}|_1 < \varepsilon$



```
d_i \, \dots \, out-degree of node i
```

 $|\mathbf{x}|_1 = \sum_{1 \le i \le N} |x_i|$ is the L₁ norm Can use any other vector norm, e.g., Euclidean

PageRank: How to solve?

Power Iteration:

- Set $r_j = 1/N$
- **1**: $r'_j = \sum_{i \to j} \frac{r_i}{d_i}$
- **2**: r = r'
- Go to **1**





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

Example:

(r _y)	1/3	1/3	5/12	9/24	6/15
$ \mathbf{r}_a =$	1/3	3/6	1/3	11/24	6/15
r _m	1/3	1/6	3/12	1/6	3/15

Iteration 0, 1, 2, ...

Why Power Iteration works? (1)



Power iteration:

A method for finding dominant eigenvector (the vector corresponding to the largest eigenvalue)

• $r^{(1)} = M \cdot r^{(0)}$

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

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

Claim:

Sequence $M \cdot r^{(0)}$, $M^2 \cdot r^{(0)}$, ... $M^k \cdot r^{(0)}$, ... approaches the dominant eigenvector of M (M is stochastic/Markov matrix)

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

Optimal r is the first or principal eigenvector of M, with corresponding eigenvalue 1

Why Power Iteration works? (2)

Claim: Sequence $M \cdot r^{(0)}, M^2 \cdot r^{(0)}, \dots M^k \cdot r^{(0)}, \dots$ approaches the dominant eigenvector of MNOTE: v is an eigenvector with the sequence of the sequ

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

Proof:

- Assume **M** has **n** linearly independent eigenvectors, $x_1, x_2, ..., x_n$ with corresponding eigenvalues $\lambda_1, \lambda_2, ..., \lambda_n$, where $\lambda_1 > \lambda_2 > \cdots > \lambda_n$
- Vectors $x_1, x_2, ..., x_n$ form a basis and thus we can write: $r^{(0)} = c_1 x_1 + c_2 x_2 + \dots + c_n x_n$
- $Mr^{(0)} = M(c_1 x_1 + c_2 x_2 + \dots + c_n x_n)$ = $c_1(Mx_1) + c_2(Mx_2) + \dots + c_n(Mx_n)$ = $c_1(\lambda_1 x_1) + c_2(\lambda_2 x_2) + \dots + c_n(\lambda_n x_n)$
- **Repeated multiplication on both sides produces** $M^{k}r^{(0)} = c_{1}(\lambda_{1}^{k}x_{1}) + c_{2}(\lambda_{2}^{k}x_{2}) + \dots + c_{n}(\lambda_{n}^{k}x_{n})$

Details!

Why Power Iteration works? (3)

- Claim: Sequence $M \cdot r^{(0)}, M^2 \cdot r^{(0)}, \dots M^k \cdot r^{(0)}, \dots$ approaches the dominant eigenvector of M
- Proof (continued):
 - Repeated multiplication on both sides produces $M^{k}r^{(0)} = c_{1}(\lambda_{1}^{k}x_{1}) + c_{2}(\lambda_{2}^{k}x_{2}) + \dots + c_{n}(\lambda_{n}^{k}x_{n})$

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

- Since $\lambda_1 > \lambda_2$ then fractions $\frac{\lambda_2}{\lambda_1}, \frac{\lambda_3}{\lambda_1} \dots < 1$ and so $\left(\frac{\lambda_i}{\lambda_1}\right)^k = 0$ as $k \to \infty$ (for all $i = 2 \dots n$).
- Thus: $M^k r^{(0)} \approx c_1(\lambda_1^k x_1)$
 - Note if c₁ = 0 then the method won't converge
- The largest eigenvalue of a stochastic matrix is always 1.

Details!

PageRank: Three Questions



- Does this converge?
- Does it converge to what we want?
- Are results reasonable?

Does this converge?



Iteration 0, 1, 2, ...

Does it converge to what we want?



Iteration 0, 1, 2, ...

PageRank: Problems

2 problems:

- (1) Some pages are dead ends (have no out-links)
 - "Vote" has "nowhere" to go to
 - Such pages cause importance to "leak out"

(2) Spider traps:

(all out-links are within the group)

- "Vote" gets "stuck" in a trap
- And eventually spider traps absorb all importance



Problem: Spider Traps

Power Iteration:

• Set
$$r_j = 1$$

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

And iterate



m is a spider trap

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

• Example:

$\left(r_{y} \right)$	1/3	2/6	3/12	5/24	0
$ \mathbf{r}_a =$	1/3	1/6	2/12	3/24	0
(r _m)	1/3	3/6	7/12	16/24	1

Iteration 0, 1, 2, ...

All the PageRank score gets "trapped" in node m.

Solution: Teleports!

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



Problem: Dead Ends

Power Iteration:

• Set
$$r_j = 1$$

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

And iterate



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

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

Example:

r _y	1/3	2/6	3/12	5/24		0
$ \mathbf{r}_a =$	1/3	1/6	2/12	3/24	•••	0
r _m	1/3	1/6	1/12	2/24		0
	Iterati	on 0 1	2			

Here the PageRank "leaks" out since the matrix is not stochastic.

Solution: Always Teleport!

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



Why Teleports Solve the Problem?

- Spider-traps are not a problem, but with traps PageRank scores are not what we want
 - Solution: Never get stuck in a spider trap by teleporting out of it in a finite number of steps
- Dead-ends are a problem
 - The matrix is not column stochastic so our initial assumptions are not met
 - Solution: Make matrix column stochastic by always teleporting when there is nowhere else to go

Solution: Random Teleports

- Google's solution that does it all: At each step, random surfer has two options:
 - With probability β , follow a link at random
 - With probability 1-β, jump to some random page
- PageRank equation [Larry Page and Sergey Brin 1998]

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

d_i ... out-degree of node i

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:

$$A = \beta M + (1 - \beta) \left[\frac{1}{N}\right]_{N \times N}$$

[1/N]_{NxN}...N by N matrix where all entries are 1/N

- We have a recursive problem: $r = A \cdot r$ And the Power method still works!
- What is β ?
 - In practice $\beta = 0.8 \sim 0.9$ (make 5 steps on avg., jump)

Random Teleports (β = 0.8)



У	1/3	0.33	0.24	0.26		7/33
a =	1/3	0.20	0.20	0.18	•••	5/33
m	1/3	0.46	0.52	0.56		21/33

MapReduce Program for PageRank

```
Map(key, value) {
    // key: a page,
    // value: page rank of the page
    For each page in Adj[key]
        emit(page, PR(key)/sizeof(Adj[key]);
}
```

```
Reduce(key, values) {
    // key: a page,
    // values: a list of page ranks from all its incoming pages
    PR(key)=1-β;
    For each pagerank in values
        PR(key) = PR(key) + β*pagerank;
    emit(key, PR(key));
}
```

MapReduce Program for PageRank



Web Search Engines

- Indexer
 - Process the retrieved pages/documents and represents them in efficient search data structures (inverted files)
- Query Server
 - Accept the query from the user and return the result pages by consulting the search data structure

