1/20

Scalable Data-driven PageRank: Algorithms, System Issues, and Lessons Learned

Xinxuan Li ¹

¹University of Maryland Baltimore County

November 13, 2015

Data-driven PageRank

Outline



- 2 Topology-driven PageRank
 - Power Method

3 Data-driven PageRank

- Data-driven pull-based PageRank
- Data-driven push-based PageRank
- Data-driven pull-push-based PageRank

4 References

Motivation	Topology-driven PageRank 0000	Data-driven PageRank 0000000	References

- **PageRank** is a numeric value that represents the importance of a page present on the web.
- when one page links to another page, it is effectively casting a vote for the other page.
- More votes implies more importance.
- Importance of each vote is taken into account when a page's PageRank is calculated.
- PageRank is Google's way of deciding a page's importance.
- It matters because it is one of the factors that determines a page's ranking in the search results.

4 / 20

Proposed by Sergey Brin and Larry Page

The PageRank Citation Ranking: Bringing Order to the Web

January 29, 1998

Abstract

The importance of a Web page is an interently subjective matter, which depends on the readers interests, knowledge and attitudes. But there is still much that can be said objectively about the relative importance of web pages. This paper describes PageRank, a method for rating Web pages objectively and mechanically, effectively measuring the human interest and attention deroted to them.

We compare PageRank to an idealized random Web surfer. We show how to efficiently compute PageRank for largenumbers of pages. And, we show how to apply PageRank to search and to user mavigation.

1 Introduction and Motivation

The Wirdd Wole Web create many new challenge for information retrieval. It is very harge and heterogeneous. Current suitants are that there are over 150 million web papes with a doubled life of low than one year. More importantly, the web papes are extremely diverse, ranging from What is be having for hand bady? It is pursual should information retrieval. In addition to there may be challenges, reach engines on the Web must also contend with inexperienced users and page engineered to mainplate search engineer ranking functions.

However, unlike "fait' document collections, the World Wide Web is hypertext and provides considerable auxiliary information on top of the text of the web pages, such as faik structure and link text. In this paper, we take advantage of the link structure of the Web to produce a global "importance" ranking of every web page. This ranking, called PageRank, heby search engines and users quickly make sense of the varia the teregoneiro of the Werld Web.

1.1 Diversity of Web Pages

Although there is a level y_{i} is given between a scalar is obtained and the scalar is present which are regulationally received, we do pare problem for each other is presented with the scalar is presented with the



Data-driven PageRank 00000000

Example

Website	PageRank
http://www.google.com/	9 / 10
http://www.youtube.com/	9 / 10
http://www.harvard.edu/	8 / 10
http://www.math.umbc.edu/ kogan/	4 / 10

Table: PageRank Table. Results are from http://www.seocentro.com/tools/search-engines/pagerank.html

Data-driven PageRank 00000000

Calculate PageRank



 $x_A = x_B + x_C + x_D = 0.75$ $x_A = \frac{x_B}{3} + \frac{x_C}{2} + \frac{x_D}{2} = 0.458$

Hence, the general form for any node v is $x_v = \Sigma_{\omega \in S_v} \frac{x_\omega{}^k}{|\Gamma_\omega|}\,.$

< □ ▶ < □ ▶ < ■ ▶ < ■ ▶ < ■ ▶ = のへで 6/20

Teleportation Parameter

Teleportation Parameter $\boldsymbol{\alpha}$ is the probability that a person will continue.

$$x_{\mathsf{v}} = \alpha \Sigma_{\omega \in S_{\mathsf{v}}} \frac{x_{\omega}}{|\Gamma_{\omega}|} + (1 - \alpha)$$

The sum of of all PageRank is N.

$$x_{\nu} = \alpha \Sigma_{\omega \in S_{\nu}} \frac{x_{\omega}}{|\Gamma_{\omega}|} + \frac{(1-\alpha)}{N}$$

The sum of of all PageRank is 1. However, we are taking $x = \frac{x}{\|x\|_1}$. Both formats have same results.

Motivation	Topology-driven PageRank ●000	Data-driven PageRank 00000000	References
Power Method			
Power Me	ethod		

For every iteration,

$$x_{\nu} = \alpha \sum_{\omega \in S_{\nu}} \frac{x_{\omega}}{|\Gamma_{\omega}|} + \frac{(1-\alpha)}{N}$$

$$\implies \mathbf{x}(\mathbf{t}+\mathbf{1}) = \alpha \mathbf{P}^{\mathsf{T}} \mathbf{x}(\mathbf{t}) + (1-\alpha) \mathbf{e}.$$

P is a transition probability. when $t \to \infty$, $\mathbf{x}(t+1) \to \mathbf{x}(t)$. $\|\mathbf{x}\|_1 = 1$ so $\mathbf{E}\mathbf{x} = \mathbf{1}$, where \mathbf{E} is a matrix with all entries as 1.

$$\mathbf{x} = (\alpha \mathbf{P}^T + (1 - \alpha) \mathbf{E}) \mathbf{x}$$

Motivation	Topology-driven PageRank 0●00	Data-driven PageRank 00000000	References
Power Method			
Example			

 $v = \{Dr.Kogan, Maria, Zois, Math students, Engineering students\}$

$$A = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 \\ 1 & 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 \end{bmatrix} \qquad D = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & 4 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 2 \end{bmatrix}$$
$$P = D^{-1}A^{T}$$
$$\mathbf{r} = (1 - \alpha)\mathbf{e} - (\mathbf{I} - \alpha\mathbf{P}^{T})\mathbf{x}$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Topology-driven PageRank

Data-driven PageRank

Power Method

Topology-driven PageRank

- Given a graph G = (V, ξ) where V is a vertex set, ξ is an edge set.
- The PageRank vector x = (1 α)e. Where α is a teleportation parameter between 0 and 1;
 e is the vector of all 1's.
- x_v^{k+1} = αΣ_{ω∈S_v} x_ω^k/|Γ_ω| + (1 − α) where x_vk is the node v's PageRank at kth-iteration; S_v is the set of incoming neighbors of node v; Γ_v is the set of outgoing neighbors of node v.
- The PageRank values are repeatedly computed until the difference between $x_v{}^k$ and $x_v{}^{k+1}$ is smaller than ϵ for all nodes.

Topology-driven PageRank

Data-driven PageRank 00000000

Power Method

Topology-driven PageRank

Define a row-stochastic matrix **P**, such that $\mathbf{P} = \mathbf{D}^{-1}\mathbf{A}$ where **A** is an adjacency matrix and **D** is the degree diagonal matrix.

The **PageRank** problem requires solving the linear system:

$$(\mathbf{I} - \alpha \mathbf{P}^{\mathsf{T}})\mathbf{x} = (1 - \alpha)\mathbf{e}.$$

The residual is defined to be

$$\mathbf{r} = (1 - \alpha)\mathbf{e} - (\mathbf{I} - \alpha \mathbf{P}^{\mathsf{T}})\mathbf{x}.$$

NЛ	Otr	vat	101	
		vuu	101	

Data-driven PageRank

Data-driven pull-based PageRank

Data-driven PageRank

- Given a graph G = (V, ξ) where V is a vertex set, ξ is an edge set.
- The PageRank vector $\mathbf{x} = (1 \alpha)\mathbf{e}$.
- Pick a node in the working list V, computing the node's PageRank and adding its outgoing neighbors to the worklist.

Topology-driven PageRank

Data-driven pull-based PageRank

Comparison

Algorithm 1. Topology-driven PageRank Input: graph $G = (\mathcal{V}, \mathcal{E}), \alpha, \epsilon$ Output: PageRank x 1: Initialize $\mathbf{x} = (1 - \alpha)\mathbf{e}$ 2: while true do 3: for $v \in \mathcal{V}$ do $x_v^{(k+1)} = \alpha \sum_{w \in \mathcal{S}_v} \frac{x_w^{(k)}}{|\mathcal{T}_w|} + (1 - \alpha)$ 4: $\delta_v = |x_v^{(k+1)} - x_v^{(k)}|$ 5: 6:7:8: end for if $\|\delta\|_{\infty} < \epsilon$ then break; 9: end if 10: end while 11: $\mathbf{x} = \frac{1}{\|\mathbf{x}\|_1}$

Algorithm 2. Data-driven PageRank **Input:** graph $G = (\mathcal{V}, \mathcal{E}), \alpha, \epsilon$ Output: PageRank x 1: Initialize $\mathbf{x} = (1 - \alpha)\mathbf{e}$ 2: for $v \in V$ do 3: worklist.push(v)4: end for 5: while !worklist.isEmpty do 6: v = worklist.pop() $x_v^{new} = \alpha \sum_{w \in S_v} \frac{x_w}{|\mathcal{T}_w|} + (1 - \alpha)$ 7: 8: $\begin{array}{ll} \text{if} & |x_v^{new} - x_v| \geq \epsilon & \text{then} \\ & x_v = x_v^{new} \end{array}$ 9: 10: for $w \in \mathcal{T}_v$ do 11: if w is not in worklist then 12: worklist.push(w) 13: end if 14: end for 15: end if 16: end while 17: $\mathbf{x} = \frac{1}{\|\mathbf{x}\|_1}$

(日) (同) (三) (三)

13/20

3

Motivation	Topology-driven PageRank 0000	Data-driven PageRank 00●00000	References
Data-driven pull-bas	ed PageRank		
Pull and	push		

Each active nodes v reads(pulls) its incoming neighbor's PageRank values and updates(pushes) residuals to its outgoing neighbors.

Pull

$$x_{v}^{k+1} = \alpha \Sigma_{\omega \in S_{v}} \frac{x_{\omega}^{k}}{|\Gamma_{\omega}|} + (1 - \alpha)$$

• Push

for
$$\omega \in \Gamma_{\omega}$$

 $r_{\omega} = r_{\omega} + \frac{r_{\nu}\alpha}{|\Gamma_{\omega}|}$

イロン イヨン イヨン イヨン

14/20

Topology-driven PageRank

Data-driven PageRank

Data-driven push-based PageRank

Data-driven push-based PageRank

In the push-based algorithms, an active node updates its own value, and only pushes its neighbor's values.

1 Initialize
$$\mathbf{x} = (1 - \alpha)\mathbf{e}$$

2
$$\mathbf{r}^{0} = (1 - \alpha)\alpha \mathbf{P}^{T} \mathbf{e}$$

 $\mathbf{r} = (1 - \alpha) - (\mathbf{I} - \alpha \mathbf{P}^{T})(1 - \alpha)\mathbf{e}$

Opdate all the active nodes v

$$x_v^{\mathsf{new}} = x_v^{\mathsf{old}} + r_v$$

$$r_{\omega}^{\mathsf{new}} = r_{\omega}^{\mathsf{old}} + \frac{\alpha v}{|\Gamma_{\omega}|}$$

If r_ω^{new} > ε or r_ω^{old} < ε then send ω to push list(repeat 3).
Set r_ν = 0.

Topology-driven PageRank

Data-driven PageRank

References

Data-driven pull-push-based PageRank

Data-driven pull-push-based PageRank

- Given a graph G = (V, ξ) where V is a vertex set, ξ is an edge set.
- **2** The PageRank vector $\mathbf{x} = (1 \alpha)\mathbf{e}$.
- Pick a active node in the working list V, computing the node's PageRank by reading its incoming neighbors.
- Updating its out-coming neighbors follow the rule by push-based PageRank.

 $r_{\omega}^{\text{new}} > \epsilon$ or $r_{\omega}^{\text{old}} < \epsilon$ then send ω to push list(repeat 3).

Topology-driven PageRank

Data-driven PageRank

(日) (同) (三) (三)

3

17/20

Data-driven pull-push-based PageRank

Comparison

Algorithm 3. Pull-Push-based PageRank Input: graph $G = (\mathcal{V}, \mathcal{E}), \alpha, \epsilon$ Output: PageRank x 1: Initialize $\mathbf{x} = (1 - \alpha)\mathbf{e}$ 2: Initialize r = 0 3: for $v \in \mathcal{V}$ do 4: for $w \in S_v$ do $r_v = r_v + \frac{1}{|\mathcal{T}_w|}$ 5: 6: end for 7: $r_n = (1 - \alpha)\alpha r_n$ 8: end for 9: for $v \in \mathcal{V}$ do 10: worklist.push(v) 11: end for 12: while !worklist.isEmpty do 13: v = worklist.pop() $x_v = \alpha \sum_{w \in S_v} \frac{x_w}{|\mathcal{T}_w|} + (1 - \alpha)$ 14: 15: for $w \in \mathcal{T}_v$ do $r_w^{old} = r_w$ 16: $r_w = r_w + \frac{r_v \alpha}{|\mathcal{T}_v|}$ 17: 18: $\text{ if } r_w \geq \epsilon \text{ and } r_w^{old} < \epsilon \text{ then } \\$ 19: worklist.push(w) 20: end if 21: end for 22: $r_{v} = 0$ 23: end while 24: $\mathbf{x} = \frac{1}{\|\mathbf{x}\|_1}$

Algorithm 4. Push-based PageRanl Input: graph $G = (\mathcal{V}, \mathcal{E}), \alpha, \epsilon$ Output: PageRank x 1: Initialize $\mathbf{x} = (1 - \alpha)\mathbf{e}$ 2: Initialize r = 03: for $v \in V$ do for $w \in S_v$ do 4: $r_v = r_v + \frac{1}{|\mathcal{T}_w|}$ 5: 6: end for 7: $r_v = (1 - \alpha)\alpha r_v$ 8: end for 9: for $v \in \mathcal{V}$ do 10: worklist.push(v) 11: end for 12: while !worklist.isEmpty do 13: v = worklist.pop() 14: $x_v^{new} = x_v + r_v$ 15: for $w \in \mathcal{T}_v$ do 16: $r_w^{old} = r_w$ $r_w = r_w + \frac{r_v \alpha}{|\mathcal{T}_v|}$ 17: $\text{ if } r_w \geq \epsilon \text{ and } r_w^{old} < \epsilon \text{ then } \\$ 18: 19: worklist.push(w) 20: end if 21: end for 22: $r_v = 0$ 23: end while 24: $x = \frac{x}{1}$ $\|\mathbf{x}\|_{1}$

Motivation	Topology-driven PageRank 0000	Data-driven PageRank ○○○○○○●○	References
Data-driven pull-pu	sh-based PageRank		
Schedulir	າຍ		

- Schedule process the node has largest residual first. Define: a node v's priority p_v to be the residual per unit work. Let $p_v = \frac{p_v}{d_v}$ For pull-push-based PageRank: $d_v = |\Gamma_v| + |S_v|$. For push-based PageRank: $d_v = |\Gamma_v|$
- NUMA-aware OBIM priority scheduler(stealing) This scheduler uses an approximate priority consensus protocol to inform a per-thread choice to search for stealable high-priority work or to operate on local near-high-priority work.
- Bulk-synchronous priority scheduler. This scheduler operates in rounds. Each rounds, all items with priority above a threshold are executed. Generated tasks and unexecuted items are placed in the next round. Threshold are updated each round based on the distribution of priorities observed recomputed every round.

ΝЛ	Δt_1	1/2	+10	h in
	ou	va	uις	211

Data-driven PageRank

References

Data-driven pull-push-based PageRank

Results



19/20

20 / 20

References

- Whang, J., Lenharth, A., Dhillon, I., Pingali, K., Larsson Traff, J., Hunold, S., and Versaci, F. (2015). Scalable Data-Driven PageRank: Algorithms, System Issues, and Lessons Learned. doi:10.1007/978-3-662-48096-034
- Frank McSherry, A Uniform Approach to Accelerated PageRank Computation, 2005