

Fast Asynchronous Anti-TrustRank for Web Spam Detection

Joyce Jiyoung Whang^{*1}, Yeonsung Jung¹, Inderjit S. Dhillon², Seongwoo Kang³, and Jungmin Lee³

¹ Sungkyunkwan University (SKKU), ² The University of Texas at Austin, ³ Naver Corporation.

ACM International Conference on Web Search and Data Mining (WSDM)

Workshop on MIS2: Misinformation and Misbehavior Mining on the Web, 2018

Main Contributions

- Asynchronous Anti-TrustRank algorithms
 - Significantly **reduce the number of arithmetic operations** compared to the traditional synchronous Anti-TrustRank algorithm
 - Without degrading the performance in **detecting Web spams**
- Convergence** of the **asynchronous Anti-TrustRank algorithms**
- Experiments on a real-world Web graph indexed by **NAVER** which is the most popular search engine in Korea.

Notation

- $G' = (\mathcal{V}, \mathcal{E}')$: a graph with **reverse edges**, i.e., if an edge $\{i, j\} \in \mathcal{E}$ then $\{j, i\} \in \mathcal{E}'$. Also, let \mathbf{A} denote the adjacency matrix of G' .
- $\mathbf{P} \equiv \mathbf{D}^{-1} \mathbf{A}$ (\mathbf{D} is the degree diagonal matrix)
- \mathcal{Q}_i : the set of incoming neighbors of node i on G'
- \mathcal{T}_i : the set of outgoing neighbors of node i on G'
- \mathbf{x} : a vector of the ATR scores, \mathbf{r} : a vector of the residuals
- \mathbf{e}_s : a vector with ones for the positions of **the seed spam documents** and zeros for other positions

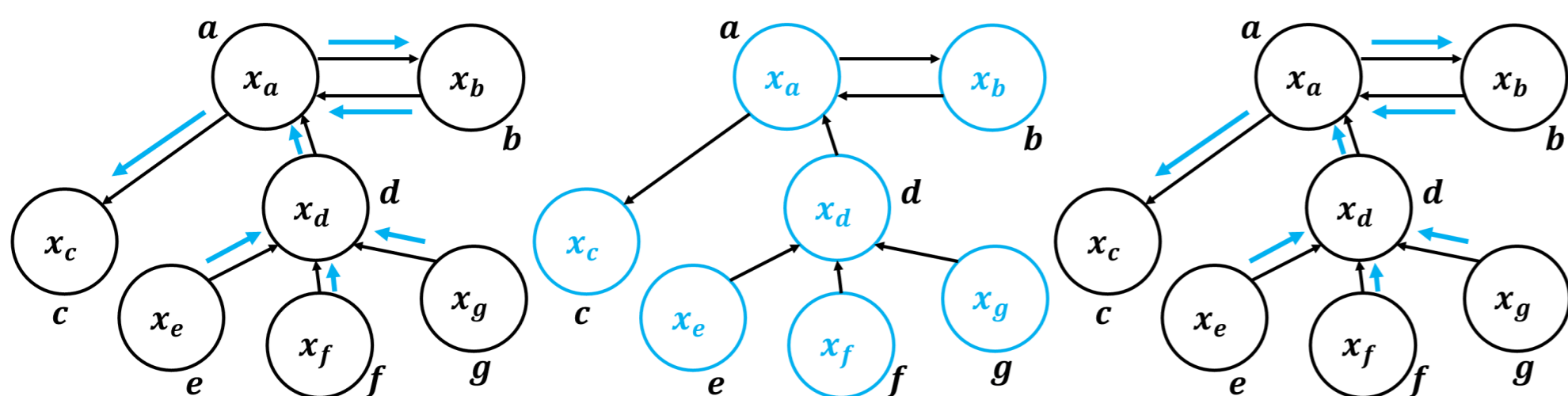
Anti-TrustRank*

- Spam pages are likely to be **referred by other spam pages**.
- Documents with high **Anti-TrustRank (ATR) score** \rightarrow spam pages
- From **spam seeds**, the ATR scores are propagated to **incoming neighbors** of the nodes so that the documents having links to the spam documents end up with having high ATR scores.

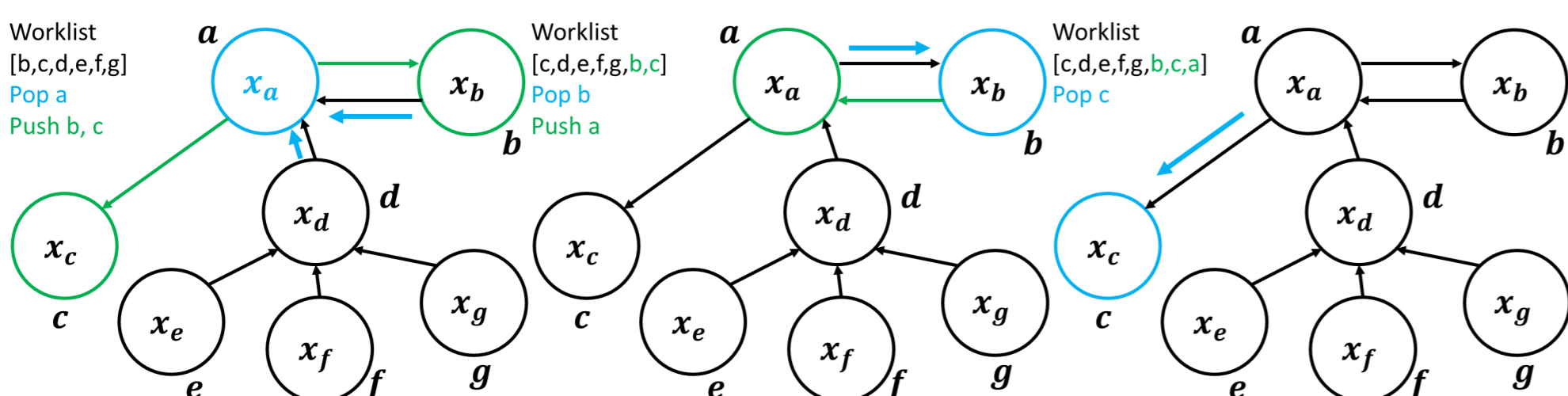
* V. Krishnan et al., Web spam detection with anti-trust rank. *AIRWeb*, 2006.

Algorithms

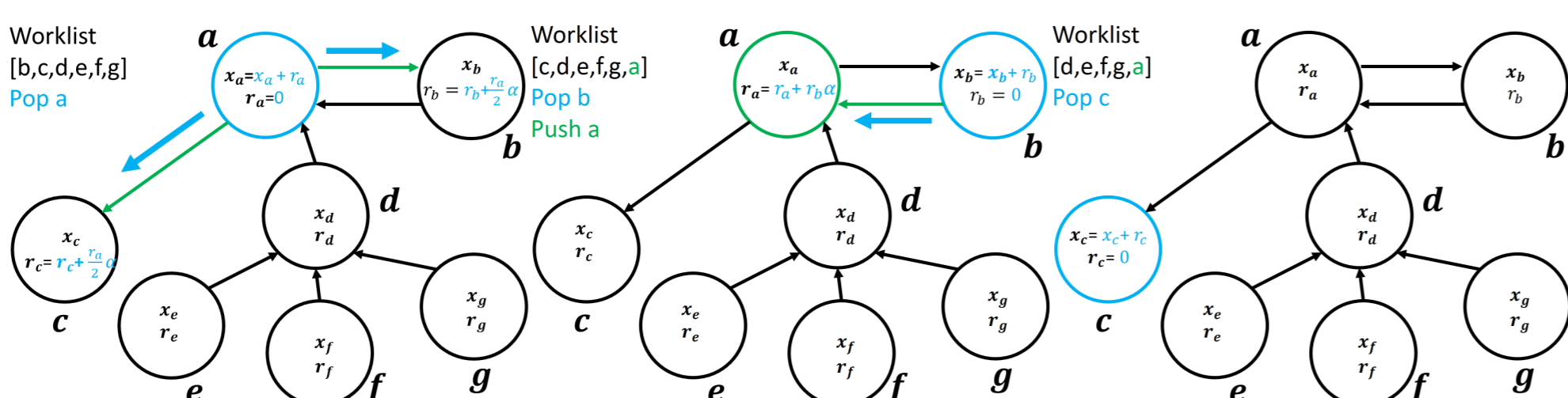
- Synchronous Anti-TrustRank (SYNC ATR)
 - The scores are updated after all the nodes re-compute the scores.



- Asynchronous Anti-TrustRank (ASYNC ATR)
 - worklist**: a set of nodes whose ATR scores need to be updated.



- Residual-based Asynchronous Anti-TrustRank (RASYNC ATR)
 - new ATR = current ATR + current residual (explicitly maintain the **residual** of each node)
 - Filtering out unnecessary work** in the worklist.

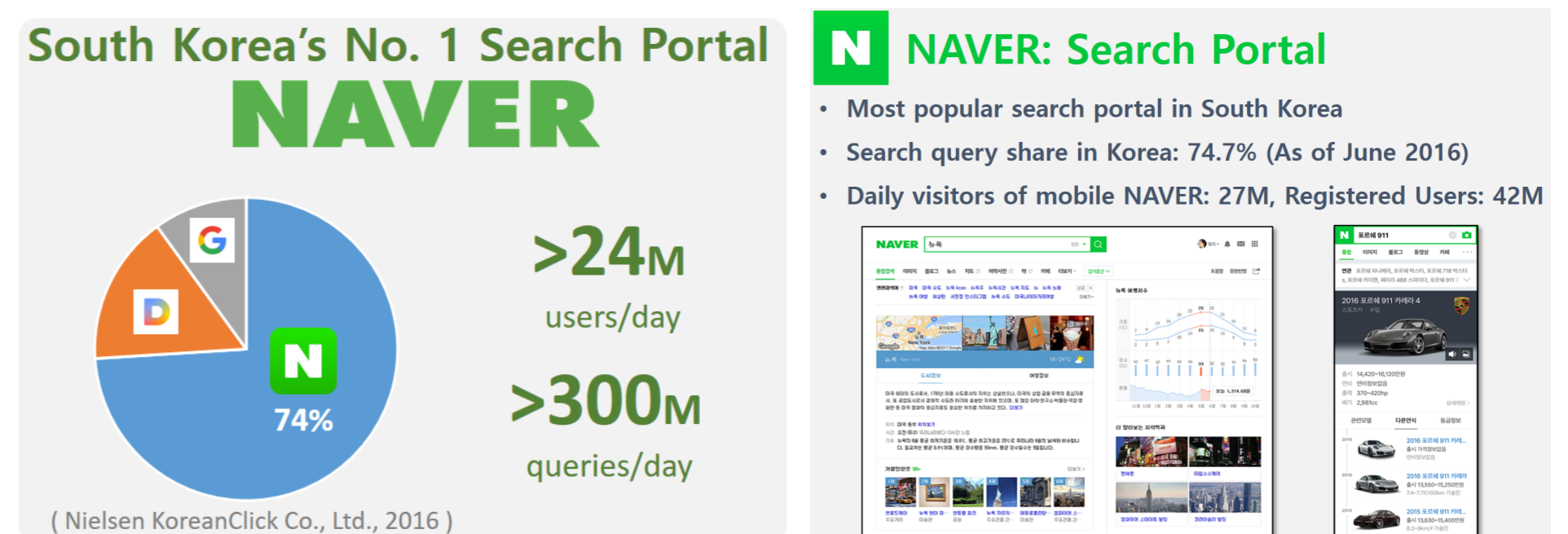


Pseudocodes

Algorithm: SYNC ATR	Algorithm: ASYNC ATR	Algorithm: RASYNC ATR
Input: $G' = (\mathcal{V}, \mathcal{E}'), \mathcal{S}, \alpha, \epsilon$ Output: ATR vector \mathbf{x} 1: Initialize $\mathbf{x} = (1 - \alpha)\mathbf{e}_s$ 2: while true do 3: for $i \in \mathcal{V}$ do 4: if $i \in \mathcal{S}$ then 5: $x_i^{new} = \alpha \sum_{j \in \mathcal{Q}_i} \frac{x_j}{ \mathcal{T}_j } + (1 - \alpha)$ 6: else 7: $x_i^{new} = \alpha \sum_{j \in \mathcal{Q}_i} \frac{x_j}{ \mathcal{T}_j }$ 8: end if 9: $\delta_i = x_i^{new} - x_i $ 10: end for 11: $\mathbf{x} = \mathbf{x}^{new}$ 12: if $\ \delta\ _\infty < \epsilon$ then 13: break ; 14: end if 15: end while 16: $\mathbf{x} = \frac{\mathbf{x}}{\ \mathbf{x}\ _1}$	Input: $G' = (\mathcal{V}, \mathcal{E}'), \mathcal{S}, \alpha, \epsilon$ Output: ATR vector \mathbf{x} 1: Initialize $\mathbf{x} = (1 - \alpha)\mathbf{e}_s$ 2: for $i \in \mathcal{V}$ do 3: $wlist.push(i)$ 4: end for 5: while !wlist.empty do 6: $i = wlist.pop()$ 7: if $i \in \mathcal{S}$ then 8: $x_i^{new} = \alpha \sum_{j \in \mathcal{Q}_i} \frac{x_j}{ \mathcal{T}_j } + (1 - \alpha)$ 9: else 10: $x_i^{new} = \alpha \sum_{j \in \mathcal{Q}_i} \frac{x_j}{ \mathcal{T}_j }$ 11: end if 12: if $ x_i^{new} - x_i \geq \epsilon$ then 13: $x_i = x_i^{new}$ 14: for $j \in \mathcal{T}_i$ do 15: if j is not in wlist then 16: $wlist.push(j)$ 17: end if 18: end for 19: end if 20: end while 21: $\mathbf{x} = \frac{\mathbf{x}}{\ \mathbf{x}\ _1}$	Input: $G' = (\mathcal{V}, \mathcal{E}'), \mathcal{S}, \alpha, \epsilon$ Output: ATR vector \mathbf{x} 1: Initialize $\mathbf{x} = (1 - \alpha)\mathbf{e}_s$ 2: Initialize $\mathbf{r} = (1 - \alpha)\alpha \mathbf{P}^T \mathbf{e}_s$ 3: for $i \in \mathcal{V}$ do 4: $wlist.push(i)$ 5: end for 6: while !wlist.empty do 7: $i = wlist.pop()$ 8: $x_i^{new} = x_i + r_i$ 9: for $j \in \mathcal{T}_i$ do 10: $r_j^{old} = r_j$ 11: $r_j = r_j + \frac{r_i \alpha}{ \mathcal{T}_i }$ 12: if $r_j \geq \epsilon$ and $r_j^{old} < \epsilon$ then 13: $wlist.push(j)$ 14: end if 15: end for 16: $r_i = 0$ 17: end while 18: $\mathbf{x} = \frac{\mathbf{x}}{\ \mathbf{x}\ _1}$

Experimental Results

- Real-world Web graph from **NAVER** corporation
 - 584,092 documents and 2,470,557 edges
 - 437,386 (74.88%) normal docs and 45,641 (7.81%) spam docs
 - 101,065 (17.30%) documents are unlabeled.



- Most of the retrieved documents are **correctly classified into spam**.
 - $|\mathcal{L}| = p|\mathcal{V}|$ where \mathcal{L} denotes the set of labeled documents
 - Pick top $q|\mathcal{S}|$ documents where \mathcal{S} denotes the set of spam seeds

Table: Accuracy of the retrieved documents

		$q = 1$	$q = 3$	$q = 5$
$p = 0.01$	spam docs	1,367 (100%)	4,099 (99.951%)	6,833 (99.971%)
	normal docs	0 (0%)	0 (0%)	0 (0%)
	unlabeled docs	0 (0%)	2 (0.049%)	2 (0.029%)
$p = 0.02$	spam docs	3,083 (100%)	9,113 (98.530%)	15,279 (99.117%)
	normal docs	0 (0%)	107 (1.157%)	107 (0.694%)
	unlabeled docs	0 (0%)	29 (0.314%)	29 (0.188%)
$p = 0.03$	spam docs	3,910 (100%)	11,593 (98.832%)	19,413 (99.299%)
	normal docs	0 (0%)	107 (0.912%)	107 (0.547%)
	unlabeled docs	0 (0%)	30 (0.256%)	30 (0.154%)

- The asynchronous algorithms, ASYNC and RASYNC, make **much fewer ATR updates** than the synchronous algorithm, SYNC.
- RASYNC significantly reduces the number of **arithmetic computations**.

Table: No. of ATR updates and arithmetic operations

p	ϵ		SYNC	ASYNC	RASYNC	
0.01	10^{-8}	No. of ATR updates	2,336,368	20,361	20,361	
		No. of arithmetics	24,442,660	8,424,970	2,516,097	
		No. of ATR updates	2,920,460	39,483	39,483	
	10^{-12}	No. of arithmetics	30,553,325	17,845,604	3,284,207	
		0.03	No. of ATR updates	2,336,368	20,628	20,628
			No. of arithmetics	24,452,832	10,716,065	2,703,630
No. of ATR updates	2,920,460		39,804	39,804		
10^{-12}	No. of arithmetics	30,566,040	25,817,600	3,932,405		