

Fast Pessimistic PageRank for dynamic directed graphs

Rohith Parjanya, Suman Kundu

Cognitive and Social Analytics Lab,
Department of Computer Science and Engineering,
Indian Institute of Technology Jodhpur, India
{parjanya.1, suman}@iitj.ac.in

Abstract. The paper presents a new algorithm for PageRank on dynamic directed graphs. The algorithm considers the addition of nodes and recalculates PageRanks. It takes cues from Monte Carlo methods for calculating Google’s PageRank. The proposed algorithm computes PageRank of a newly added node in constant time. Extensive experiments have been performed on different synthetic and real world networks. The experimental result shows that the ranks generated by the proposed method is the closest to the rank of the benchmark algorithm of Power Iteration.

Keywords: Dynamic Network, Monte Carlo, Link Sensitivity Index

1 Introduction

Calculating PageRank [2] in a dynamic network is an important and challenging research problem. A network is dynamic if its topology changes with time, either by adding nodes and edges or by deleting nodes and edges. Dynamic networks are ubiquitous in the age of information technology, e.g., Twitter mentions network changes with each Tweet post or product purchasing network changes for every product purchased, or the friendship of Facebook (or professional connection in LinkedIn) network changes over time with the addition of new friends (or contacts). In such cases, fast recalculation of the PageRank plays a vital role in providing real-time services to its users. A node, if added to a network, brings several edges along with it. This disrupts the existing PageRank scores. Furthermore, the score of the newly added node needs to be determined. The trivial solution is to recalculate the PageRank for the whole network. However, it is time-consuming and not sustainable when nodes are added to the network rapidly. In this connection, the research problem is to recalculate the PageRanks for the affected nodes only instead of the whole network. Can we recalculate fast without affecting the actual rank of the nodes?

PageRank, first proposed in [2], is conventionally used for ranking web pages. It is based on the idea of Eigenvector centrality. Since its inception in 1999 [2], many different variants of PageRank algorithms [1, 3, 4, 12, 8, 6, 7] are proposed. Many of these algorithms work only on static networks. While the deterministic

algorithms are slow for large networks, many randomized Monte Carlo based approaches [4, 7, 8, 6, 5] are developed to provide fast solutions. Monte Carlo based algorithms provide a good estimate of the PageRanks. However, the modern network changes over time at a rapid speed. Many dynamic PageRank [10–15] approaches are proposed to deal with the dynamic network. However, they cannot cope with the pace of the streaming graph, especially for Big Data Social Networks. A simple and efficient algorithm to find the PageRanks for topological changes in dynamic streaming graphs is the need of the hour.

In this paper, we propose a simple algorithm to update the PageRanks in a directed network after adding a new node in the network. The proposed algorithm, namely, Fast Pessimistic PageRank (FPPR), uses the expected value of a random surfer to calculate the score of the new node by adding the scores contributed by in-links and estimating to what extent the chain may contribute to the score of the node. The same method is used to update existing nodes through the out-links of the newly added node. The proposed algorithm executes $\mathcal{O}(k \times d_{avg})$ time complexity, where k is the number of new nodes added and d_{avg} is the average degree of the k nodes added to the graph. Further, FPPR takes only $\mathcal{O}(|V|)$ additional space, specifically, $4 \times |V|$ space. The experiment performed over several synthetic and real-world networks shows that the ranking of the proposed method is highly correlated with the benchmark Power Iteration based recalculation and better than Fast Incremental PageRank(FIPR) [13].

The paper is organized as: Section 2 provides a brief literature review, the proposed FPPR method and its rationale is presented in Section 3, experiments performed and corresponding results are reported in Section 4 and finally, Section 5 describes the conclusions of the research work.

2 Related Work

PageRank [2] defines the importance of a node based on the quality of the neighboring nodes and the degree of the node. It is traditionally used for the hyperlink network of WWW. In simple terms, it uses a random surfer model where the random surfer clicks out-links with probability $1 - \rho$ and terminates its walk to start a new walk at a random page with probability ρ .

There are many Monte Carlo based approaches [4, 7, 8, 6, 5] for approximating the PageRank. Monte Carlo methods majorly follow four different strategies. These are (i) Monte Carlo end-point with a random start, (ii) Monte Carlo end-point with cyclic start, (iii) Monte Carlo Complete Path, and (iv) Monte Carlo complete path stopping at the dangling nodes. Out of these methods, the last one shows better performance in terms of execution time compared to the former three approaches [7]. In this method, R number of random walks are simulated starting from each node, and the random walk terminates at dangling nodes. The final rank vector is defined as $\pi = \frac{\text{visits of each node}}{\text{total visits}}$.

Many algorithms for the dynamic network were proposed in the literature. These are broadly classified into two categories viz, (i) aggregation algorithms [10–12] and (ii) Monte Carlo based algorithms [13, 15]. Although aggregate meth-

ods do not recompute the PageRanks of the whole graph, it aggregates the vicinity of the disturbed area imposed by the new changes in the graph. Thus, these methods recompute the PageRanks locally. The disadvantage of these methods is the aggregation overhead. In addition, accuracy highly depends on the size of the subset selected for the aggregation or the area of the vicinity [9].

On the other hand, Monte Carlo based approaches use the theory of Markov chains [14]. In order to recalculate PageRanks on network changes, all the random walks which pass through the updated edges (or nodes) are required to be adjusted [13, 15]. Hence, it requires keeping track of all the random walks it performed earlier, and the same reason adds the overhead of additional space. For example, in FIPR [13], at each node, the unique segment IDs pass through the node, and the total number of visits to that node is stored. In addition, a hashtable is kept in each node, with unique segment IDs and the nodes that each segment covers whenever the graph is updated, visited, and segments are modified for particular to that updated node or edge. Hence, these methods are not suitable for large-scale streaming graphs due to the significant memory overhead.

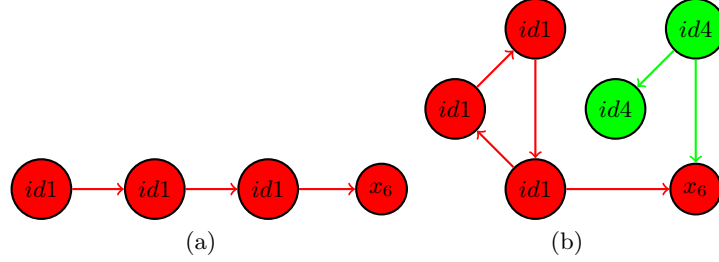


Fig. 1. Examples

3 Proposed Fast Pessimistic PageRank (FPPR) Algorithm for Dynamic Directed Graphs

The real-time use cases like finding the top- k popular products in an online shopping cart or finding top- k spreaders of news considering reshare network etc. will need high-speed computation of PageRank. In this section, we will present the FPPR algorithm, which is capable of recalculating the PageRanks of a directed network upon the addition of new nodes. In order to achieve this, FPPR calculates the expected score of random surfer considering the static graph, i.e., if the Monte Carlo method is executed on the graph after the addition of the node, what would be the expected score generated by the random walks. For example, in Fig. 1(a), let 'x6' be a new node. Intuitively, the expected value of 80% (considering $\rho = 0.2$) of random walks through out-link reaching 'x6' must have visited the new node 'x6' considering a static Monte Carlo simulation. This is what we would like to calculate in the first phase (Approximate Visits,

Definition 2). We assume the incoming link is linear in the second phase (Link Sensitivity Index, Definition 1), as shown in Fig. 1(a). For example, R random walks starting from the node of ‘id1’ must have reached ‘x6’ in 3 hops, i.e., the expected number of visits is $R * (0.8)^3$. Similarly, the next node must have gone ‘x6’ in 2 hops (expected visit is $R * (0.8)^2$), and so on. All these contribute to the score of the new node. Note that the same formulation may not be valid for Fig. 1(b) ‘id4’. In such a case algorithm is prone to error. However, the algorithm thinks pessimistic that “the ‘id4’ has a loop”. Then all the random walks made on those nodes must have reached the node ‘x6’, and our linear assumption is valid. Accordingly, we are calculating the PageRank scores. A similar kind of process is adopted in addressing the out link from the new node as well.

Algorithm 1 Calculate approximate PageRank of the newly created node

```

1: Input: new node  $u$ ,  $R = 1000$ 
2: for all  $v \in \Gamma_{in}(u)$  do
3:   Assign linkID to the node as described in text
4:    $linkID[v].length \leftarrow linkID[v].length + 1$ 
5:    $temporary\_Var \leftarrow R \times LSI(n = linkID[v].length)$  {LSI calculated by Eqn. 1}
6:    $linkSensitivityIndex \leftarrow \max(linkSensitivityIndex, temporary\_Var)$ 
7:    $approxVisits \leftarrow AV(v)$  {Using Eqn. 2}
8: end for
9:  $approxVisits \leftarrow approxVisits + R + linkSensitivityIndex$ 
10:  $finalResult[u] \leftarrow approxVisits$ ,  $totalVisits \leftarrow totalVisits + approxVisits$ 
11:  $approxVisits, linkSensitivityIndex \leftarrow 0$ 
12: for all  $v \in \Gamma_{out}(u)$  do
13:   if  $u$  has bidirectional edge with  $v$  then
14:     replace both nodes with maximum visits
15:   else
16:     Assign link ID to the node as described in text
17:      $linkID[u].length \leftarrow linkID[u].length + 1$ 
18:      $temporary\_Var \leftarrow R \times LSI(n = linkID[u].length)$  {LSI calculated by Eqn. 1}
19:      $approxVisits \leftarrow AV(u)$  {Using Eqn. 2}
20:      $finalResult[v] \leftarrow approxVisits + linkSensitivityIndex$ 
21:      $totalVisits \leftarrow totalVisits + approxVisits + linkSensitivityIndex$ 
22:   end if
23: end for
24: return  $finalResult \div totalVisits$ 

```

Definition 1 (Link Sensitivity Index (LSI)). *LSI defines the extent one node affects the whole link chain it belongs. We assume the link chain is linear and the LSI is mathematically defined by:*

$$c + c(c)^1 + c(c)^2 + c(c)^3 + \dots \rightarrow \frac{c(1 - (c)^n)}{(1 - c)} \quad (1)$$

Here, $c = 1 - \rho$ is the probability that a random surfer moves forward and n is the length of the chain. In our experiment we took $c = 0.8$.

Definition 2 (Approximate Visits (AV)). *AV estimates the contribution of the incoming links to the newly added node, that a random walk might have used if the random walk is executed from the scratch. If nodes have bidirectional edges, a max of approximate visits is considered for both nodes. It is calculated as:*

$$AV(v) \leftarrow \sum_{u \in \Gamma_{in}(v)} \frac{1}{d_{out}(u)} * (1 - \rho) * (AV(u)) \quad (2)$$

Here, $\Gamma_{in}(\cdot)$ and $d_{out}(\cdot)$ return the set of incoming neighbors and out degree respectively.

FPPR does not keep track of all the random walk segments or aggregations. FPPR takes only $4|V|$ space. That is, space complexity is $\mathcal{O}(|V|)$. For a network with billions of nodes, this provides a significant improvement.

The Algorithm 1 shows the steps of FPPR algorithm. The *linkID* of a node is the ID given to a node to identify the chain to which it belongs. In other words, *linkID* keeps track of all the links or distinct chains in the graph. One may note that a lower *linkID* is assigned to a node if multiple links exist through that node. This convention help resolving the conflicts. As we are keeping track of every link(chain) in the graph, we store the length of each *linkID*.

4 Experiments and Results

Experiments have been conducted on synthetically generated networks as well as real world networks. Both the Random Network and Barabasi-Albert [16] networks are used. The real world networks are Wiebo reshare network of [17]. The salient features of these graphs are presented in Table 1. In order to get the dynamic character, we started with one node and added each node according to their creation.

Table 1. Dataset Description

| Name | Nodes | Edges | Min In-Deg | Max In-Deg | Min Out-Deg | Max Out-Deg |
|-----------|-------|-------|------------|------------|-------------|-------------|
| RW1 [17] | 40 | 66 | 0 | 12 | 0 | 6 |
| RW2 [17] | 206 | 206 | 0 | 45 | 0 | 2 |
| RW3 [17] | 759 | 780 | 0 | 201 | 0 | 3 |
| RW4 [17] | 817 | 901 | 0 | 28 | 0 | 297 |
| Random51 | 51 | 93 | 0 | 5 | 0 | 7 |
| Random126 | 126 | 233 | 0 | 5 | 0 | 8 |
| Random201 | 201 | 386 | 0 | 7 | 0 | 9 |
| Random276 | 276 | 543 | 0 | 10 | 0 | 7 |
| BA1 [16] | 55 | 154 | 0 | 35 | 0 | 3 |
| BA2 [16] | 130 | 376 | 0 | 61 | 0 | 3 |
| BA3 [16] | 205 | 604 | 0 | 101 | 0 | 3 |
| BA4 [16] | 280 | 829 | 0 | 135 | 0 | 3 |

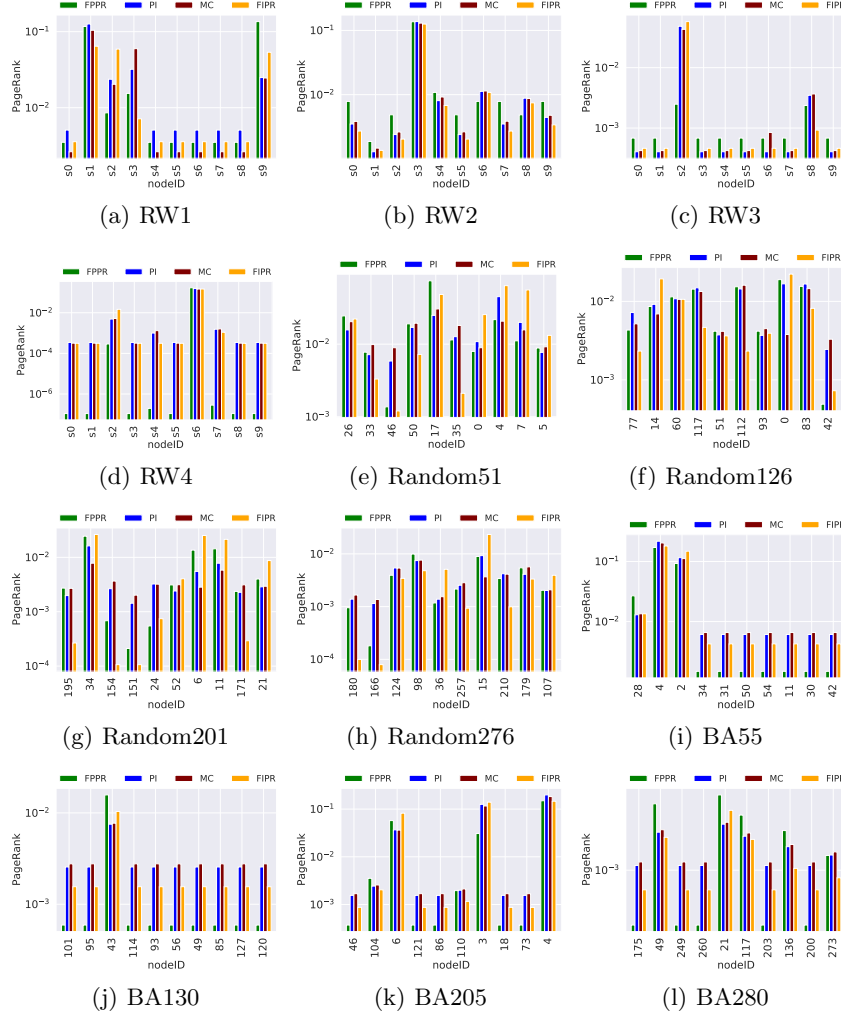


Fig. 2. Comparing plots of Randomly Sampled 10 nodes for different algorithms

Comparing Methods The proposed methods has been compared with the following methods.

- **PowerIteration (PI)** [2]: We consider this method as a benchmark in our experiments. In the experiment, we restarted PI algorithm and recalculated PageRank for the whole graph on each node addition.
- **Static Monte Carlo (MC)** [7]: Static Monte-Carlo method is the method of approximate PageRanks of the network using Monte Carlo method. We implemented the version of complete path with dangling nodes. Similar to PI,

we recalculated the PageRanks every time a node is added to the network. The number of random walk is considered in the experiment is 1000.

- **Fast Incremental PageRank on Dynamic Networks (FIPR)** [13]: The method is designed for dynamic network and proposed in 2019. As our algorithm is designed for dynamic networks, we included this method as a related research. Parameter R is set to 16 in the experiments.

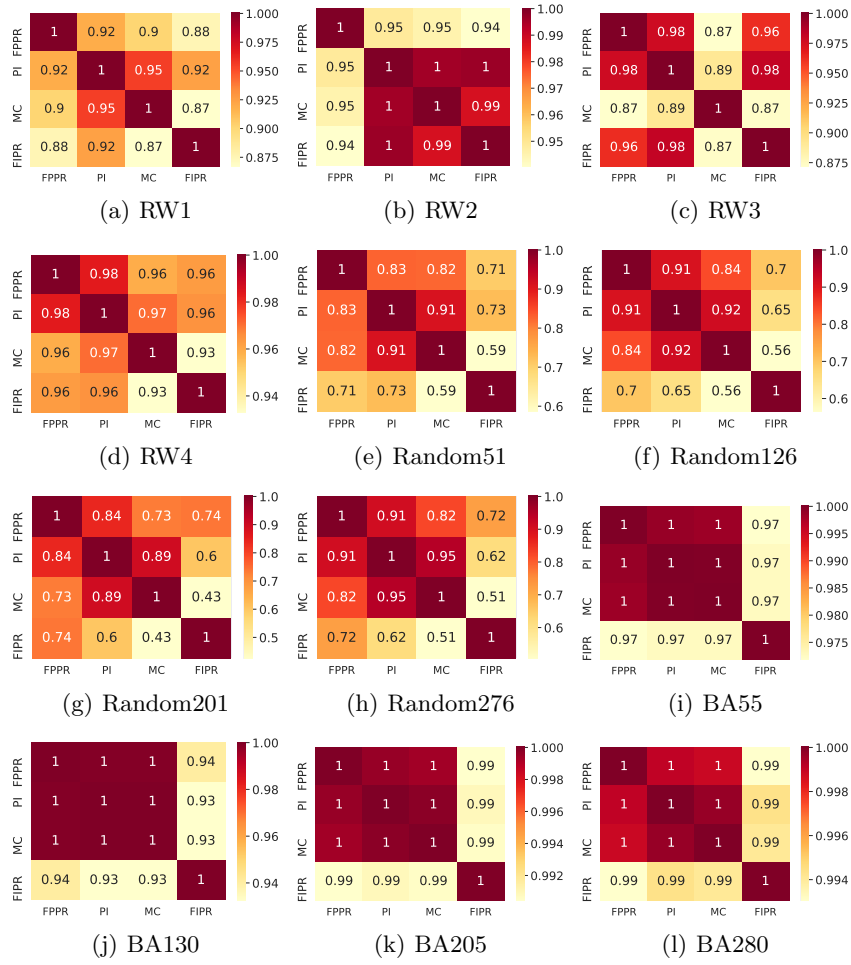


Fig. 3. Spearman Correlation between all four approaches over the different data sets

4.1 Results

Accuracy All the comparing algorithms were executed on different graphs. As expected, the absolute values of PageRanks by different algorithms of a node are different. The results for 10 random nodes for all the dataset are shown in 2 for reference. Hence, comparing different algorithms in terms of absolute values of the PageRank is not fare. So the Spearman’s rank correlation coefficient[18] is used to compare the ranking of the proposed FPPR with that of the comparing methods. The Spearman correlation between two vectors will be high when observations have a similar rank between the two variables and it will be lower otherwise. A value of it in between 0.8 to 1 is considered to be strongly correlated[19]. The results are shown in the Figure 3. It is evident from the result that the proposed FPPR performed equal or better than FIPR for all the network except RW2 while comparing with the PI and MC as benchmarks. Note that MC and PI are highly correlated to each other as both of these are recalculated over the full graph once a new node is added to the graph.

Execution Time: We checked the overall execution time of different comparing algorithms and found that the proposed algorithm is magnitude faster than all the methods we have experimented with for all the data sets. The comparison is presented in the Figure 4. As expected FIPR takes second lowest time as it is designed for dynamic network and only updates the PageRank locally.

4.2 Computation Complexity

Computation complexity of different algorithms are shown in Table 2. The proposed algorithm take constant time for each new node added to the network. So the worth case complexity of the proposed algorithm is $\mathcal{O}(k \times d_{avg})$, where k is the number of node added to the network.

Table 2. Computation complexities of different PageRank algorithms in Dynamic setup, k is number of updates, n is number of nodes, E is edgeset, R is the no of random walk simulations.

| Algorithm | Time complexity |
|--------------------|---------------------------------------|
| Power Iteration | $\Omega_{\frac{1}{1-\rho}}(kn^2)$ |
| Monte Carlo method | $\Omega_{(\rho)}(\frac{knR}{\rho})$ |
| FIPR | $\mathcal{O}_{(E)}(\frac{kn}{ E })$ |
| Our Method | $\mathcal{O}(k \times d_{avg})$ |

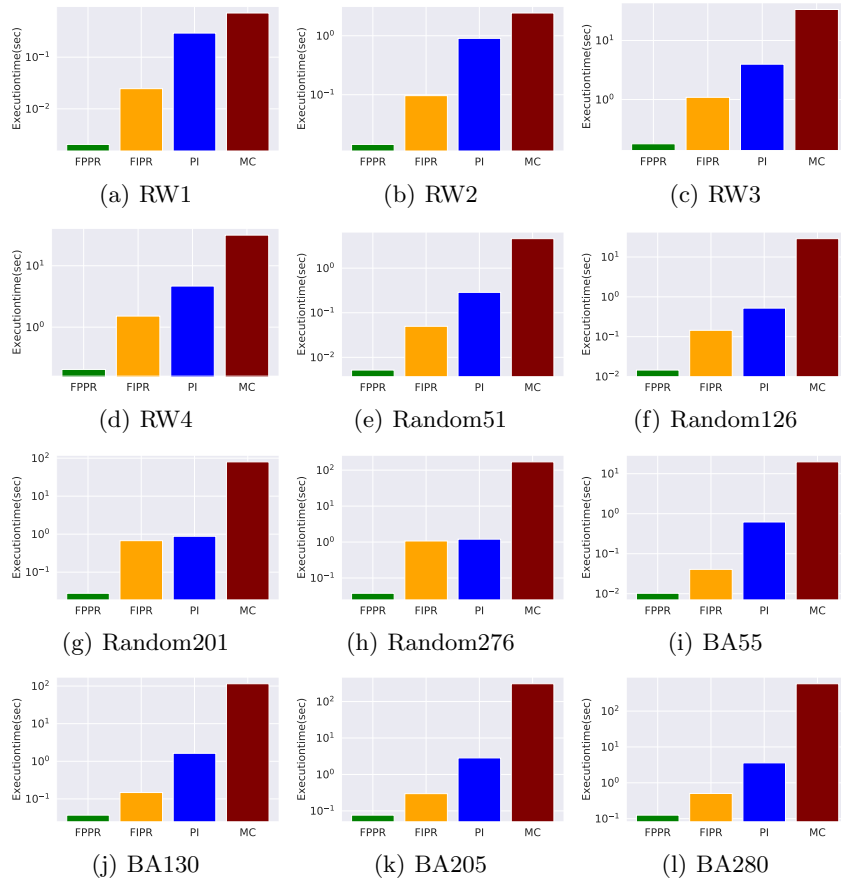


Fig. 4. Comparing plots of execution time of different algorithms

5 Conclusion

In the present paper we proposed a new algorithm for calculating PageRank for dynamic directed graph. We consider the node addition for the same. We showed thorough experimental result that the proposed algorithm perform better in terms of ranking from the existing FIPR algorithm. The execution time is much faster while providing better results. Power Iteration algorithm is considered as benchmark.

Experiments were performed on wide range of networks including random network, Barabasi-Albert network and real world network. It is also shown that the PageRank values differs between algorithms while the ranking shows high correlation in terms of Spearman ranking correlation index.

Although the algorithm is executed for node addition, we believe this algorithm can be modified and applied to the link addition as well. Similarly, node

and edge deletion may also be employed. These are kept as the future works in the current research.

References

1. Kleinberg, J. & Others Authoritative sources in a hyperlinked environment.. *SODA*. **98** pp. 668-677 (1998)
2. Page, L., Brin, S., Motwani, R. & Winograd, T. The PageRank Citation Ranking: Bringing Order to the Web.. (Stanford InfoLab,1999,11), Previous number = SIDL-WP-1999-0120
3. Lempel, R. & Moran, S. The stochastic approach for link-structure analysis (SALSA) and the TKC effect. *Computer Networks*. **33**, 387-401 (2000)
4. Breyer, L. Markovian page ranking distributions: some theory and simulations. *Preprint*. (2002)
5. Avrachenkov, K., Litvak, N., Nemirovsky, D. & Osipova, N. Monte Carlo methods in PageRank computation: When one iteration is sufficient. *SIAM Journal On Numerical Analysis*. **45**, 890-904 (2007)
6. Zhou, Z. Evaluation of Monte Carlo Method in PageRank. (University of Missouri-Columbia,2015)
7. Vargas, B. Exploring PageRank Algorithms: Power Iteration & Monte Carlo Methods. (2020)
8. Salehi, O. PageRank algorithm and Monte Carlo methods in PageRank Computation. *Semantic Scholar*. (2007)
9. Bahmani, B., Chowdhury, A. & Goel, A. Fast incremental and personalized pagerank. *ArXiv Preprint ArXiv:1006.2880*. (2010)
10. Desikan, P., Pathak, N., Srivastava, J. & Kumar, V. Incremental page rank computation on evolving graphs. *Special Interest Tracks And Posters Of The 14th International Conference On World Wide Web*. pp. 1094-1095 (2005)
11. Chien, S., Dwork, C., Kumar, R., Simon, D. & Sivakumar, D. Link evolution: Analysis and algorithms. *Internet Mathematics*. **1**, 277-304 (2004)
12. Langville, A. Meyer, C. Updating pagerank with iterative aggregation. *Proceedings Of The 13th International World Wide Web Conference On Alternate Track Papers & Posters*. pp. 392-393 (2004)
13. Zhan, Z., Hu, R., Gao, X. & Huai, N. Fast incremental pagerank on dynamic networks. *International Conference On Web Engineering*. pp. 154-168 (2019)
14. Seneta, E. Non-negative matrices and Markov chains. (Springer Science & Business Media,2006)
15. Liao, Q., Jiang, S., Yu, M., Yang, Y. & Li, T. Monte Carlo based incremental PageRank on evolving graphs. *Pacific-Asia Conference On Knowledge Discovery And Data Mining*. pp. 356-367 (2017)
16. Albert, R. & Barabási, A. Statistical mechanics of complex networks. *Reviews Of Modern Physics*. **74**, 47 (2002)
17. Chuai, Y. & Zhao, J. Anger makes fake news viral online. *ArXiv Preprint ArXiv:2004.10399*. (2020)
18. Spearman, C. The proof and measurement of association between two things. *The American Journal Of Psychology*. **100**, 441-471 (1987)
19. Zar, J. Spearman rank correlation: Overview. *Wiley StatsRef: Statistics Reference Online*. (2014)