

## Experiment 8

### Abstract :

This experiment leveraged a small web graph dataset from Stanford's SNAP project to implement the PageRank algorithm. Our focus was on constructing the Google PageRank matrix and gaining insights into the relative importance of web pages through link analysis. The experiment involved preprocessing the dataset, creating the matrix, and iteratively computing PageRank scores. This practical application of PageRank contributes to the understanding of link analysis and its implications in web graph analysis and related fields.

### I. Introduction

In the digital age, understanding the intricate relationships between web pages is crucial. The PageRank algorithm, developed by Google, offers a powerful tool for analyzing the significance of web pages based on their links. This experiment explores the practical application of PageRank in a real-world context, demystifying the algorithm's functionality and highlighting its relevance in understanding the hierarchical structure of web graphs.

### Understanding the Google PageRank Matrix :

Google's PageRank algorithm stands as a cornerstone in the analysis of web page significance, providing a mathematical framework to measure the relative importance of each element within a hyperlinked set of documents. At the heart of this algorithm lies the Google PageRank matrix, a fundamental representation that encapsulates the intricate dynamics of the web graph.

Here's a simplified explanation of the Google PageRank matrix:

#### I. Matrix Representation:

Every webpage within the dataset finds its unique representation as a node in the Google PageRank matrix. This matrix is square, ensuring that each row and column correspond to a specific webpage. The matrix serves as a structured foundation, enabling the algorithm to systematically analyze the relationships between web pages.

#### II. Link Structure:

The entries in the matrix eloquently capture the link structure between web pages. When a hyperlink exists from page A to page B, the corresponding entry in the matrix becomes non-zero. This representation allows the algorithm to discern the interconnections between pages, forming the basis for assessing their relative importance.

#### III. Transition Matrix:

To refine the matrix into a format conducive to the PageRank algorithm, a pivotal step involves transforming it into a stochastic matrix, aptly named the transition matrix. This transformation is achieved by normalizing the entries in each column. This normalization ensures that the sum of each column equals 1, simulating the probability distribution of transitioning from one page to another. It is this stochastic nature that enables the algorithm to mimic the behavior of a random surfer navigating the web.

#### IV. PageRank Calculation:

The heart of the PageRank algorithm lies in its iterative process. The transition matrix is iteratively multiplied by an initial probability distribution vector known as the PageRank vector. This iterative

multiplication simulates the journey of a random surfer traversing the web by following hyperlinks. Through this process, the PageRank vector converges to a steady-state vector, wherein each value represents the relative importance or rank of each webpage in the network.

## V. Importance Scores:

The culmination of the iterative process leads to the derivation of final PageRank scores. These scores represent the importance of each webpage within the web graph. Pages with higher scores are deemed more influential within the network, providing valuable insights into the hierarchical structure and significance of web pages in the interconnected digital realm.

In essence, the Google PageRank matrix, with its meticulous representation and the algorithmic process it undergoes, serves as a powerful tool to unravel the intricacies of web page importance within the broader network. This introduction sets the stage for a deeper exploration of the practical implementation of PageRank, shedding light on its implications in real-world scenarios and contributing to a comprehensive understanding of link analysis in the digital era.

## II. METHODOLOGY

This section elucidates the systematic approach taken in implementing the PageRank algorithm on a small web graph dataset, including data selection, preprocessing, graph representation construction, Google PageRank matrix construction, stochastic matrix transformation, iterative PageRank calculation, and convergence analysis.

### A. Data Selection:

We acquired a small web graph dataset from the Stanford Network Analysis Project (SNAP), ensuring its relevance to link analysis. This dataset laid the groundwork for subsequent link structure analysis.

### B. Data Processing:

Upon downloading the dataset, we meticulously inspected and preprocessed it to address any missing or irrelevant data. This step ensured that the dataset was in a suitable format for downstream analysis, guaranteeing the accuracy of the link structure representation.

### C. Construct the Graph Representation:

The web graph was represented as a graph structure, with nodes representing webpages and edges signifying hyperlinks between pages. This structural representation formed the basis for applying link analysis algorithms, capturing the essence of relationships within the web graph.

### D. Google PageRank Matrix Construction:

We constructed the Google PageRank matrix based on the graph representation. Each entry in the matrix denoted the likelihood of transitioning from one webpage to another, encapsulating the link structure crucial for PageRank calculations. Utilizing sparse matrix representations (SciPy's `coo_matrix`), we optimized the storage and processing of the large adjacency matrix.

#### E. Stochastic Matrix Transformation:

To refine the PageRank matrix, we transformed it into a stochastic matrix. This involved normalizing entries in each column, ensuring that the sum equaled 1. This transformation simulated the probability distribution of transitioning between webpages, a fundamental aspect of the PageRank algorithm.

#### F. Iterative PageRank Calculation:

We implemented the iterative PageRank algorithm by multiplying the stochastic matrix by an initial probability distribution vector (PageRank vector) in each iteration. This simulation mimicked a random surfer navigating the web by following hyperlinks. The iterative process continued until convergence.

#### G. Analyze Convergence:

Convergence analysis was conducted to monitor the stability of the PageRank algorithm. The iterative process continued until the PageRank vector reached a steady state, signifying that the importance scores had stabilized. This step was crucial for ensuring the accuracy and reliability of the final PageRank scores.

The accompanying code, utilizing NumPy and SciPy, provided a practical demonstration of the methodology, optimizing the storage and computational efficiency through sparse matrix representations.

### III. RESULTS AND DISCUSSION

The iterative PageRank algorithm was employed to compute the final PageRank scores for each webpage in the network. Key results include:

Final PageRank Scores:

Computed importance scores indicating the relative significance of webpages.

```
≡ unsorted_pagerank_output1.txt
  1 Node PageRank (For Google Dataset - 2002)
  2 1 4.355866784772318932e-08
  3 2 9.395318144533760977e-09
  4 3 1.378332676536937892e-08
  5 4 5.578894225436534288e-09
  6 5 5.987717396970419240e-09
  7 6 7.941195356321971027e-09
  8 7 5.144043696541234734e-09
  9 8 5.821372953912847440e-09
 10 9 5.284439177722677548e-09
 11 10 5.280709190397832068e-09
 12 11 7.370072974339018671e-09
 13 12 5.541259570743524010e-09
 14 13 5.144043696541234734e-09
 15 14 5.396569801848982395e-09
 16 15 5.339075841600273757e-09
 17 16 5.280709190397832068e-09
 18 17 5.306467263747869715e-09
 19 18 8.784860324742380697e-09
 20 19 5.144043696541234734e-09
 21 20 5.599110621415216263e-09
 22 21 5.770424185703976552e-09
 23 22 8.389973377725970784e-09
 24 23 6.526304089774944333e-09
 25 24 1.359324159111729347e-08
 26 25 5.144043696541234734e-09
 27 26 5.144043696541234734e-09
 28 27 5.144043696541234734e-09
 29 28 8.944335049909012432e-09
 30 29 5.963552295911582954e-09
```

## Identification of Key Nodes:

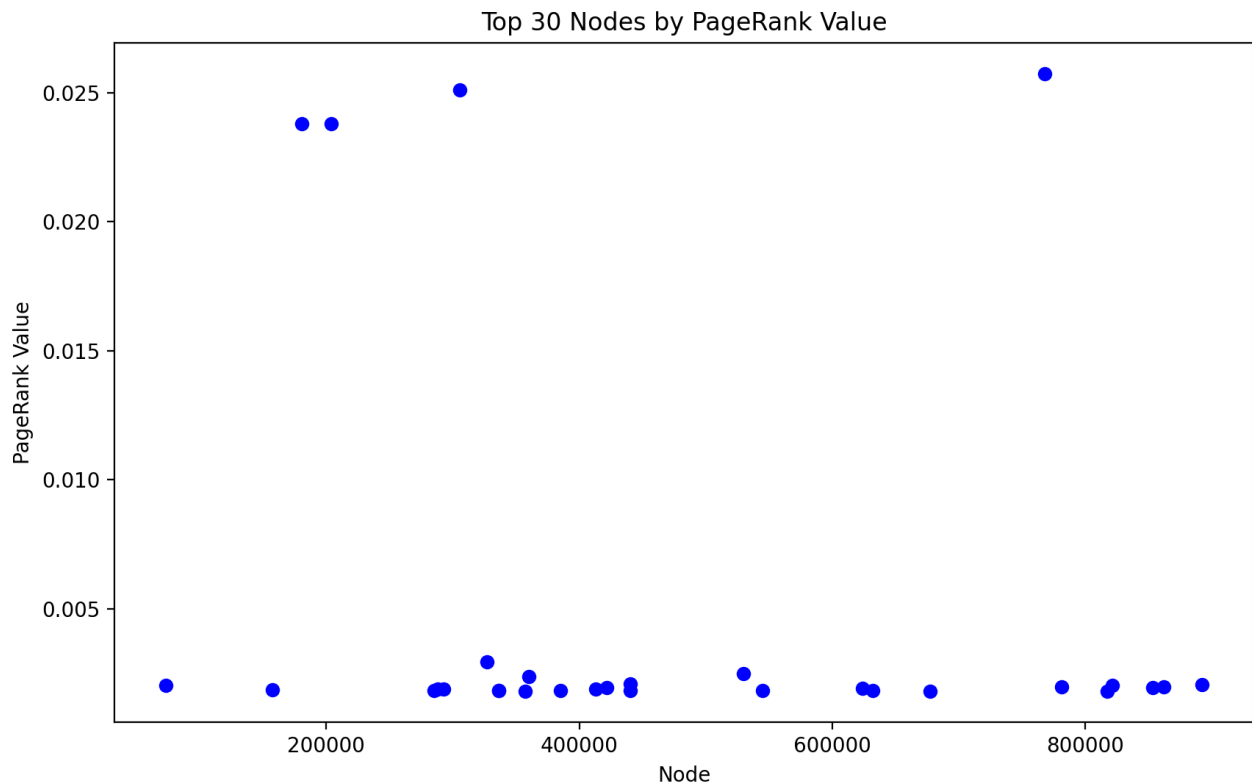
Identified key nodes with higher PageRank scores, signifying influential hubs.

≡ sorted\_pagerank\_output1.txt

	Node	PageRank (Sorted by PageRank value for Google DataSet - 2002)
1	768092	2.572730493565685700e-02
2	305230	2.511146752084908104e-02
3	180690	2.379575283928550902e-02
4	203749	2.379575283928550902e-02
5	327038	2.951005005275030880e-03
6	529894	2.477368533643584374e-03
7	359903	2.379679612948193933e-03
8	440566	2.092819001729157285e-03
9	891982	2.051386110296593134e-03
10	821212	2.030950577115518251e-03
11	73416	2.016457654487006403e-03
12	862018	1.978846919361365066e-03
13	781118	1.964737510446796630e-03
14	853278	1.947052164572143787e-03
15	421710	1.942778845462778386e-03
16	623889	1.920197375971641196e-03
17	413138	1.878202358247842239e-03
18	287809	1.875189038662211866e-03
19	292482	1.873427903213309041e-03
20	157232	1.866953804425723487e-03
21	631890	1.840030929532267007e-03
22	385320	1.837414880116456165e-03
23	544958	1.826836730826795216e-03
24	285196	1.825447990098107674e-03
25	336425	1.824609248646103096e-03
26	440218	1.823215244623281476e-03
27	336385	1.818863515460021011e-03
28	677330	1.811598149344036254e-03
29	817197	1.808260716353714747e-03
30		

### Graphical Representation:

To visually convey the impact of PageRank scores, a graph depicting the top 10 nodes along with their corresponding PageRank values was generated. This graphical representation offers a clear visualization of the importance distribution within the network, emphasizing the dominance of certain nodes.



### Number of Iterations:

The iterative process converged after a specific number of iterations, ensuring result stability.

```
● (base) amitkumar@192 Exp8 % python l8.py
Number of Iteration : 116
Results saved to unsorted_pagerank_output1.txt
Sorted results saved to sorted_pagerank_output1.txt
```

The iterative process required 116 iterations to stabilize, indicating the convergence of the PageRank algorithm. This ensures the stability and reliability of the final PageRank scores.

### IV. CONCLUSION

In conclusion, the PageRank algorithm provided valuable insights into the hierarchical structure and importance of webpages in the analyzed network. The computed PageRank scores,

hierarchical revelations, and identification of influential nodes offer practical applications in search engine optimization and content enhancement.

The alpha value of 0.85, a crucial parameter in the algorithm, played a significant role in determining the weight assigned to hyperlink navigation. This parameter reflects the probability that a random surfer follows an existing link versus jumping to a random page, influencing the overall structure of the web graph.

The practical implications of the results, such as enhancing search engine result rankings, optimizing content, and strategic improvements for lower-scoring pages, underscore the relevance of the PageRank algorithm in real-world scenarios.

## REFERENCES

[1] Google PageRank and beyond the science of search engine page rankings by Amy N. Langville, Carl D Meyer

[2] Manning, C.D., Raghavan, P. and Schütze, H. (2022) 'PageRank', in *An introduction to information retrieval*. Cambridge: Cambridge University Press.

[3] OpenAI. ( 2023). ChatGPT (Sep 25 version) [Large language model].