# Accelerated Pagerank Algorithm

20337025 Canming Cui and 20337138 Zhanhao Xu

Department of Computer Science, Sun Yat-Sen University.

**Abstract**

This article is the final major assignment for the Principles and Techniques of Big Data class. We gave a brief overview of the PageRank algorithm and performed some optimizations and implemented three techniques for accelerating convergence, namely initial value design, Quadratic Extrapolation, and fix converged nodes. We give mathematical derivations and show in detail the experimental results in my report, which demonstrate the correctness and efficiency of the optimisation algorithm. In addition, We have implemented the PageRank algorithm in MapReduce form.

**Keywords:** PageRank, Quadratic Extrapolation, fix converged nodes, MapReduce

## 1   Introduction

The PageRank algorithm, developed by Larry Page and Sergey Brin[1], is used to evaluate and rank web pages based on their importance. It plays a vital role in search engine optimization and web page ranking.

The algorithm assesses the importance of web pages by considering the quantity and quality of links. Each web page is treated as a node, and links are viewed as directed edges between nodes(as shown in Fig. 1). By analyzing these links, the algorithm determines the significance of web pages.The core idea behind the algorithm is "link voting." When a web page is linked by other pages, it receives more "link votes," increasing its importance. The importance of a page is also influenced by the importance of the pages that link to it.

To calculate the importance of web pages, the PageRank algorithm uses an iterative approach. In each iteration, the algorithm redistributes link weights and performs calculations based on the updated weights. This process continues until convergence is achieved. The final result is a score assigned to each web page, reflecting its importance for ranking purposes.

Apart from search engine rankings, the PageRank algorithm finds applications in social network analysis and recommendation systems. It provides an effective way to evaluate node importance in networks and has significantly impacted information retrieval and network analysis

There have been many improvements to the PageRank algorithm over the last few decades. Our accelerated algorithm includes the Quadratic Extrapolation algorithm. It was proposed by Kamvar
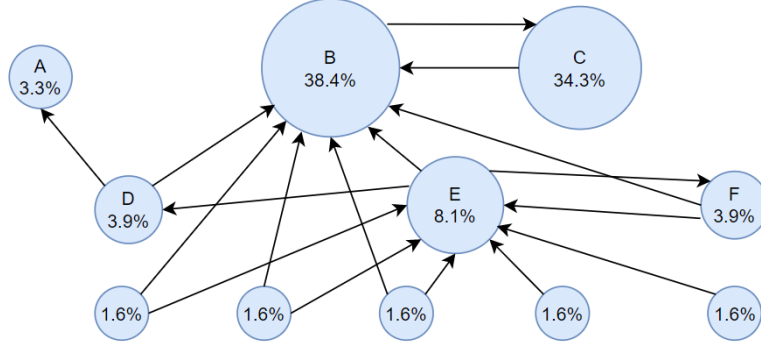
Figure 1: A simple illustration of the Pagerank algorithm. The percentage shows the perceived importance, and the arrows represent hyperlinks.

et al.[2], which aimed to expedite the convergence rate of the PageRank vector. Through their work, they achieved a significant improvement in convergence, reducing the convergence rate from 25% to 300%. W.Xing et al.[3] proposed and developed the weighted PageRank algorithm (WPR) as an extension to the standard PageRank algorithm. They introduced the concept of incorporating the importance of both inlinks and outlinks of web pages into the ranking process. Ispen and Selee[4] presented a fundamental algorithm for computing the PageRank vector using the Google matrix. They introduced a novel approach where all dangling nodes are combined into a single node. The algorithm calculates the PageRank of non-dangling nodes separately and then calculates the PageRank for the dangling nodes.

The innovation of our accelerated Pagerank algorithm is the combination of various techniques: sparse matrices, initial value design, Quadratic Extrapolation and fix converged nodes. And we combine it with MapReduce.

## 2   PageRank Review

The PageRank algorithm simulates user click behavior. Suppose there are users browsing web pages, and at a given moment, they randomly click on a link on the current page and navigate to another page. As time goes on, the number of users on each page should stabilize, representing the page's PageRank value.

Let $num_nodes$ represents the total number of nodes, $Out[j]$ represents the number of outLinks of page j, $In[j]$ represents the number of inLinks of page j, $\beta(j)$ is the set of pages pointing into page j, $PageRank[j]$ represents the PageRank values of page j, $A[\ ][\ ]$ represents the adjacency matrix of pages, $p$ represents the probability of random walk, and $\epsilon$ represents the convergence threshold. The formula for the aforementioned PageRank algorithm is expressed as **??**:

$$PageRank[i] = \sum_{j \in \beta(j)} \frac{PageRank[j]}{Out[j]} \tag{1}$$

To give precise and manageable meaning to such a roughly stated definition let us denote by A the adjacency matrix with the elements:

$$A[j][i] = a_{ij} = \begin{cases} \frac{1}{Out[j]} & \text{if } j \in \beta(i) \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

It can be shown that each PageRank value converges (i.e. the magnitude of change tends to 0) when the adjacency matrix A between pages satisfies these three conditions, independent of the initial value: (i)A is a Markov matrix. (ii)A corresponds to a graph that is strongly connected. (iii)A is non-periodic. However, this is only the ideal case; in real web pages, the matrix A will not usually satisfy the above conditions. Therefore, to ensure that the algorithm converges correctly, we need to modify the matrix A slightly.

A Markov matrix is that all lements are greater than or equal to 0 and the sum of the elements of any row is equal to 1. Why we need A to be a Markov matrix? Becaues the equivalent condition for matrix A to be a Markov matrix is that all nodes have outLinks. If page i has no outLinks, meaning the i-th row of matrix A is all zeros, then the PageRank values that transition to this page will be lost and cannot be transferred to other pages. As a result, the sum of PageRank values in the entire system will continuously decrease until it reaches 0. Finally the PageRank values of all pages converge to 0. So we need to make matrix A to be a Markov matrix, i.e. we need to ensure that all pages without outLinks have outLinks. So, link the pages without outLinks to all pages:

$$A[j][i] = a_{ij} = \begin{cases} \frac{1}{Out[j]} & \text{if } j \in \beta(i) \\ \frac{1}{num\_nodes} & \text{if page j hasn't outLinks} \end{cases} \quad (3)$$

Strong connectivity ensures that every node can be reached from any other node through directed edges. Periodicity refers to the property of a matrix where the result of matrix multiplication repeats itself after a certain number of multiplications. It implies that the matrix follows a cyclic pattern in its multiplication process. If the matrix A does not satisfy strong connectivity, the value of PageRank at convergence can only be 0 or 1. If the matrix A does not satisfy non-periodicity, it may never converge.

A random wandering strategy is used to solve both problems, i.e. the user closes the current page and a new page is opened at random:

$$A[j][i] = (1 - p) \cdot A[j][i] + p \quad (4)$$

The equation for the transfer from the (k-1)th PageRank to the kth PageRank can be expressed as follows(by matrix multiplication):

$$PageRank_k = A \cdot PageRank_{k-1} \quad (5)$$

The mainstream PageRank algorithm primarily employs iterative calculation methods to quickly obtain the approximate convergence value of the PageRank, so the complexity of the PageRank algorithm lies in iterative computation. Next we will discuss iterative methods and acceleration.

3

# 3    Accelerated Algorithm

In this section we describe in detail the four accelerated iteration algorithms of the PageRank algorithm.

## 3.1    Sparse Matrix

The above mentioned PageRank algorithm based on matrix power implementation is as follows. Suppose we have obtained the adjacency matrix A from Eq. (3) and Eq. (4), and then initialise the values of PageRank, and the sum of PageRank is 1. Then start iterating by letting A multiply by PageRank until PageRank converges(the amount of change is less than threshold $\epsilon$), and the algorithm would be if sparse matrices are not used:

---

**Algorithm 1:** Computing PageRank by the power method

An adjacency matrix $A$ has been obtained;
**for** *all page i* **do**
 $\quad$ $PageRank_0[i] = \frac{1}{num\_nodes}$;
**end**
k=1;
**while** $\delta < \epsilon$ **do**
 $\quad$ $PageRank_k = A \cdot PageRank_{k-1}$;
 $\quad$ $\delta = max_i(\mid PageRank_k[i] - PageRank_{k-1}[i] \mid)$;
 $\quad$ $k+ = 1$;
**end**

---

The complexity of the algorithm is $O(num\_nodes^2 \cdot n)$, where $n$ is the number of iterations. Since each element of the matrix A is non-zero, the multiplication time occupies most of the runtime. However, we can observe that most of the elements in A are the same, so they can be extracted as constants, and removing these elements from A will turn it into a sparse matrix. Sparse matrix can reduce the number of multiplications significantly compared to the original matrix. For elements in the matrix that are obtained by random wandering or no outLinks equally divided among all pages, they can be extracted as constants. Modify the update formula of PageRank as follows(where $sum_{k-1} = \sum_{i \text{ hasn't outLinks}} PageRank_{k-1}[i]$):

$$PageRank_k = A' \cdot PageRank_{k-1} + sum_{k-1} \cdot (1 - p) + p \qquad (6)$$

The sparse matrix power iteration algorithm is as Algorithm 2, where $A'[i][j]$(j¡number of inLinks for page i) denotes the page from which the jth inLink of page i originates. The complexity of the algorithm can be analysed as $O(n \cdot (num\_nodes + e))$(e is the total number of links). Because in the sparse matrix:$num\_nodes + e < num\_nodes^2$, so the sparse matrix power algorithm complexity is lower than the original matrix power algorithm.

---

**Algorithm 2:** Computing PageRank by the power method using sparse matrix

---

An sparse matrix $A'$ has been obtained;
**for** *all page i* **do**
   | $PageRank_0[i] = \frac{1}{num\_nodes}$;
**end**
k=1;
**while** $\delta < \epsilon$ **do**
   $sum = 0$;
   **for** *page i that hasm't outLinks* **do**
     | $sum = sum + PageRank_{k-1}[i]$;
   **end**
   **for** *all page i* **do**
     **for** $1 \le j \le In[i]$ **do**
       $PageRank_k[i]+ = PageRank_{k-1}[\frac{A'[i][j]\cdot(1-p)}{Out[[A'[i][k]]]}]$;
       $PageRank_k[i]+ = \frac{sum(1-p)}{num\_nodes}$;
       $PageRank_k[i]+ = p$;
     **end**
   **end**
   $\delta = max_i(|\ PageRank_k[i] - PageRank_{k-1}[i]\ |)$;
   $k+ = 1$;
**end**

---

## 3.2  Initial Value Design

The initial PageRank value can be equally distributed or determined based on some heuristic. However, the algorithm will converge faster if the initial value can be closer to the final converged state. So we set the initial value of PageRank according to the number of inLinks, and the size of the initial value is proportional to the number of inLinks, which is shown in Eq. (7)(*num_nodes* is added to prevent the denominator from being 0).

$$PageRank_0[i] = \frac{In[i] + 1}{\sum_{\le j \le num\_nodes} In[j] + num\_nodes} \tag{7}$$

The code implementation is as follows, see Appendix A for full codes:

```
1  //init the value of matrix
2  void Init_value(int Init){
3          if(Init==0){
4                  for(int i=1;i<=num_total;i++){
5                          val[i]=1/double(num_total);
6                  }
7          }
8          if(Init==1){
```

```
9              int inlinks_sum=0;
10             for(int i=1;i<=num_total;i++){
11                     inlinks_sum+=count_inlinks[i]+1 ;
12             }
13             for(int i=1;i<=num_total;i++){
14                     val[i]=(count_inlinks[i]+1)/double(inlinks_sum);
15             }
16        }
17 }
```

## 3.3 Quadratic Extrapolation

The quadratic extrapolation method comes from the work of Kamvar et al[2]. and we have applied it based on the guidance of their paper. The quadratic extrapolation method is a technique that speeds up the iterative process of the PageRank algorithm. It leverages the convergence properties of PageRank values. This method assumes that the change in PageRank values follows a quadratic convergence pattern, meaning that the change in PageRank is proportional to the square of the iteration step. Based on this assumption, the quadratic extrapolation method can use the results of the early iterations to estimate the values of the subsequent iterations, thereby accelerating the convergence of the algorithm.

The steps of the quadratic extrapolation method are as follows:

1. Assume that $PageRank_{k-2}$ is a linear combination of the feature vectors corresponding to the first and second features:

$$PageRank_{k-2} = u_1 + \alpha_2 u_2$$

2. From the mathematical analysis we have that the maximum eigenvalue $\lambda = 1$, so for k-1 iterations and k iterations :

$$PageRank_{k-1} = A \cdot PageRank_{k-2} = u_1 + \alpha_2 \lambda_2 u_2$$

$$PageRank_k = A \cdot PageRank_{k-1} = u_1 + \alpha_2 \lambda_2^2 u_2$$

3. Let the variables g,h:

$$g = PageRank_{k-1} - PageRank_{k-2} = \alpha_2(\lambda_2 - 1)u_2$$

$$h = PageRank_k - PageRank_{k-1} = \alpha_2(\lambda_2 - 1)\lambda_2 u_2$$

4. The formula for calculating $\lambda_2$ is:
$$\lambda_2 = \frac{h^T \cdot g}{g^T \cdot g}$$

5. Convergence can be accelerated by updating the original PageRank with the new value of

6

PageRank as follows:

$$PageRank_k = u_1 = \frac{\lambda_2 PR_{k-1} - PR_k}{\lambda_2 - 1}$$

The code implementation is as follows, see Appendix A for full codes:

```
1   //Quadratic Extrapolation method
2   void Quadratic_Extrapolation(){
3         double lamda;
4         double sum1=0;
5         double sum2=0;
6         //Calculated as shown in the formula
7         for(int i=1;i<=num_total;i++){
8               sum1+=last_difference[i]*last_difference[i];
9               sum2+=last_difference[i]*dif[i];
10        }
11        lamda=sum2/sum1;
12        for(int i=1;i<=num_total;i++){
13              newval[i]=(lamda*val[i]-newval[i])/(lamda-1);
14        }
15  }
```

## 3.4   Fix Converged nodes

In the iterative calculation of PageRank algorithm, the convergence speed of different nodes is different, some nodes converge faster, but because other nodes have not converged, then the converging nodes have to follow along with the iterative calculation, which will cost a lot of unnecessary calculations. Therefore, we can fix the value of the converged nodes to reduce the amount of computation in the iterative process.

We can monitor $\delta$(The amount of change in PageRank for the two iterations before and after) of a node during the iteration process, and if it is less than a certain threshold, it will be judged to have converged. And this threshold needs many experiments to select the best threshold.

The code implementation is as follows, see Appendix A for full codes:

```
1   ...
2   ...
3   for (int iter=1;iter<max_iteration;iter++){
4         double sum_val0=0;
5         // Calculate the sum of values for nodes with no outgoing links
6         for (int j=0;j<num_of_nooutlinks;j++)
7         sum_val0+=val[no_outlinks[j]];
8         // Calculate the new PageRank values for each node
9         for (int i=1;i<=num_total;i++){
10            if(Fix&&is_Converge[i]) continue;
```

```
11              // Nodes with outgoing links
12              // If the node has outgoing links, assign it a value based
                     on the sum of values for nodes with no outgoing links
13              if(count_outlinks[i]!=0)newval[i]=Poblity/double(num_total)
                     +(1-Poblity)*sum_val0/double(num_total-1);
14              // If the node has no outgoing links, assign it a value
                     based on the sum of values for nodes with no outgoing
                     links
15              else  newval[i]=Poblity/double(num_total)+(1-Poblity)*(
                     sum_val0-val[i])/double(num_total-1);
16              // The second part of the PageRank value is obtained by "
                     absorbing" from all incoming edges to the node
17              for(int  j=0;j<count_inlinks[i];j++){
18                  // The new PageRank value of node i is updated by adding
                         the probability of not performing a random walk,
19                  // the weight of the j-th incoming edge of node i, and
                         the PageRank value of the incoming node
20                  newval[i]+=(1-Poblity)*val[spares_matrix_cite[i][j]]*
                         weight[i][j];
21              }
22          }
23          ...
24          ...
25 }
26 ...
27 ...
```

# 4 MapReduce

The idea for the section comes from the work of ST Wierzchoń et al[5]. The MapReduce programming framework was proposed by Dean and Ghemawat[6]. It is used for handling and producing vast amounts of data. It involves two main functions: the map function and the reduce function. The map function takes a key/value pair and generates a collection of intermediate key/value pairs. The reduce function then combines all the intermediate values associated with the same intermediate key.

The process can be summarized as follows:

1. Data Splitting: The input dataset is divided into multiple equally sized data blocks, which are processed in parallel.

2. Mapping: Each data block is processed by a set of mappers, applying a user-defined map function to produce intermediate key-value pairs.

3. Shuffling and sorting: Intermediate key-value pairs are sorted and grouped by keys to aggregate values with the same key.

4. Reducing: Each reducer processes intermediate keys and their associated values using a user-defined reduce function.

5. Output generation: Results from reducers are written to output files or storage systems.

The "mapper" and "reducer" can be represented symbolically in the following manner($< k1, v1 >$ represents a key-value pair, where 'k1' is the key and 'v1' is the corresponding value. Additionally, the use of $[\cdot]$ denotes a list containing multiple elements):

$$\text{map:} \quad < k_1, v_1 > \rightarrow [< k_2, v_2 >]$$

$$\text{reduce:} \quad < k\_2, [v_2] > \rightarrow [< k_3, v_3 >]$$

In order to implement Algorithm 1 within the MapReduce framework, a dual set of map-reduce functions needs to be applied. The first pair of map-reduce functions is utilized to perform matrix-vector multiplication. First, knowing page i, its current PageRank $PageRank[i]$, and the list of outlinks $OutLinks_i$, the mapper returns the set of pairs $< j, w_j >$, where j is a page linked from i and $w_j = PageRank[i]/out_i$, with $out_i$ being the length of $OutLinks_i$. The reducer is responsible for aggregating the information:

$$\text{map:} \quad < i; (PageRank[i], OutLinks_i) > \rightarrow [< j, PageRank[i]/out_i >]$$

$$\text{reduce:} \quad < j, [w_i] > \rightarrow [< j, \sum_i w_i >]$$

The second pair finishes the job, i.e. the pairs $< j, \sum_i w_i >$ are transmitted by the mapper to the reducer that computes their sum. We can also implement the optimised algorithm by a similar way.

## 5   Experiment

### 5.1   Data Set

We chose the Wikipedia voting dataset for our experiments(https://snap.stanford.edu/data/wiki-Vote.html). The network contains all the Wikipedia voting data from the inception of Wikipedia till January 2008. Nodes in the network represent wikipedia users and a directed edge from node i to node j represents that user i voted on user j. The network has 7115 nodes and 103689 directed edges.

### 5.2   Results

In order to compare the performance of various versions of the PageRank algorithm, we compare the execution time of each version of the algorithm to reach the convergence threshold in terms of the order of magnitude of the convergence threshold, and the resulting table is shown in Tab. 1.

Figure 2: Line graph of execution time for different versions.



Figure 3: Line graph of speed for different versions.

The meaning of the symbols in the table is as follows:

- origin: The basic Pagerank implemented with a dense adjacency matrix.

- $v_1$: Sparse matrix optimisation.

- $v_2$: Sparse matrix + Initial value design.

- $v_3$: Sparse matrix + Quadratic Extrapolation.

- $v_4$: Sparse matrix + Fix Converged nodes

- $v_5$: Accelerated Pagerank (Sparse matrix + Quadratic Extrapolation + Fix Converged nodes)

In order to show the experimental results more intuitively, they are plotted as a line graph, as shown in Fig. 2 and Fig. 3(the pagerank of dense matrices takes longer and is omitted here).

| $\epsilon$ | $origin$ | $v_1$ | $v_2$ | $v_3$ | $v_4$ | $v_5$ |
|---|---|---|---|---|---|---|
| $10^{-3}$ | 144.555 | 8.867 | 20.917 | 9.983 | 10.814 | 7.495 |
| $10^{-4}$ | 319.118 | 19.717 | 26.011 | 13.921 | 15.522 | 12.478 |
| $10^{-5}$ | 461.762 | 24.092 | 34.380 | 24.401 | 18.905 | 17.000 |
| $10^{-6}$ | 614.731 | 31.214 | 38.623 | 29.051 | 25.189 | 24.486 |
| $10^{-7}$ | 791.330 | 35.597 | 42.651 | 32.875 | 29.588 | 29.016 |
| $10^{-8}$ | 945.651 | 39.312 | 44.803 | 39.572 | 34.172 | 33.737 |

Table 1: Table of execution times for different convergence thresholds, table entry represents execution time (ms)

From the experimental results, it can be learnt that the execution time of all algorithms rise monotonically with decreasing threshold $\epsilon$, which is approximately linear. Each optimisation plays a role, i.e., speeds up the execution of the Pagerank algorithm. And the Accelerated PageRank algorithm has the best performance among all versions, which proves that our accelerated algorithm is effective.

# 6    Summary

In this article, we give a brief overview of the PageRank algorithm and improve it. In Sec. 1, we introduced some concept, background and related works of the PageRank algorithm. In Sec. 2, we described the principles of the PageRank algorithm. In Sec. 3, we present the accelerated PageRank algorithm, applying four optimisations. In Sec. 4 we described the MapReduce framework and how to apply it to PageRank algorithm. In Sec. 5, we show and analyse the results of the experiments, which demonstrate the effectiveness of the accelerated algorithm. In the end, we learnt a lot through this final assignment and our abilities were honed.

# References

[1]    Sergey Brin and Lawrence Page. "The anatomy of a large-scale hypertextual web search engine". In: *Computer networks and ISDN systems* 30.1-7 (1998), pp. 107–117.

[2]    Sepandar D Kamvar et al. "Extrapolation methods for accelerating PageRank computations". In: *Proceedings of the 12th international conference on World Wide Web.* 2003, pp. 261–270.

[3]    Wenpu Xing and Ali Ghorbani. "Weighted pagerank algorithm". In: *Proceedings. Second Annual Conference on Communication Networks and Services Research, 2004.* IEEE. 2004, pp. 305–314.

[4]    Ilse CF Ipsen and Teresa M Selee. "PageRank computation, with special attention to dangling nodes". In: *SIAM Journal on Matrix Analysis and Applications* 29.4 (2008), pp. 1281–1296.

[5] Sławomir T Wierzchoń et al. "Accelerating PageRank computations". In: *Control and Cybernetics* 40.2 (2011), pp. 259–274.

[6] Jeffrey Dean and Sanjay Ghemawat. "MapReduce: simplified data processing on large clusters". In: *Communications of the ACM* 51.1 (2008), pp. 107–113.

# A   Appendix

Run the main function of the programme for conducting experiments.

main.c:

```c
1   //main function
2   int main(int argc, char *argv[]) {
3       bool N=0;
4           bool I=0;
5           bool E=0;
6           bool F=0;
7           bool W=0;
8           bool S=0;
9           bool H=0;
10      string name ;
11      if(argc<2){
12                  printf("##########  error  #########\n");
13          printf("-n : Origin pagerank\n") ;
14          printf("-s : Sparse Matrix\n") ;
15          printf("-i : Initial Value\n") ;
16          printf("-e : Quadratic Extrapolation\n") ;
17          printf("-f : Fix Converged nodes\n") ;
18                  printf("##########  error  #########\n");
19          return 0 ;
20          }
21          for(argc--,argv++;argc>0;argc-=1,argv+=1){
22          if(!strcmp(*argv,"-n")){
23              name+="naive";
24                          N=1;
25                      }
26          else if(!strcmp(*argv,"-s")){
27              name+="sparse";
28                          S=1;
29          }
30          if(!strcmp(*argv,"-i")){
31              name+="_init";
32                          I=1;
```

```cpp
33                }
34                if (!strcmp(*argv,"-e")){
35                    name+="_extra";
36                            E=1;
37                }
38                if (!strcmp(*argv,"-f")){
39                    name+="_is_Converge";
40                            F=1;
41                }
42        }
43            QueryPerformanceFrequency(&frequ);
44            create_idx();
45            build_statistics(W);
46            QueryPerformanceCounter(&start_time);
47            Init_value(I);
48            if(N)PageRank_origin();
49            if(S)PageRank_sparse(E,F);
50            QueryPerformanceCounter(&end_time);
51        cout<<(double)(end_time.QuadPart-start_time.QuadPart)/frequ.QuadPart
            <<endl;
52            output_to_txt(name+".txt");
53            return 0;
54  }
```

Header file for the accelerated pagerank algorithm, defining some variables and functions.
accelerated_pagerank.h:

```cpp
1   //pagerank
2   #include <iostream>
3   #include <map>
4   #include <algorithm>
5   #include <vector>
6   #include <windows.h>
7   #include <string>
8   #include <fstream>
9   #include <cstring>
10  using namespace std;
11  #define max_iteration 100000
12  #define epslion 1e-10
13  #define epslionIndex 10
14  #define maxnodes 8000
15  #define maxedges 403689
16  #define maxindegree 10000
```

```cpp
17  #define Poblity 0.15
18  #define maxlength 20
19  LARGE_INTEGER frequ;
20  LARGE_INTEGER start_time, end_time;
21  map<int, string> id_2_name;//id to name
22  map<string, int> name_2_id;//name to id
23  typedef pair<int, double> weigh_edges;//edges have weights
24  vector<weigh_edges> wedges_vec;
25  unsigned short spares_matrix_cite[maxnodes][maxindegree];
26  float matrix[maxnodes][maxnodes];
27  int no_outlinks[maxnodes];
28  int num_total;
29  int num_of_edges;
30  int num_of_nooutlinks;
31  int count_outlinks[maxnodes];//cnt of outlinks of node
32  int count_inlinks[maxnodes];//cnt of inlinks of node
33  int count_outlinks_in[maxnodes];
34  bool is_Converge[maxnodes];
35  float weight[maxnodes][maxindegree];
36  double val[maxnodes];
37  double newval[maxnodes];
38  //The difference between the two iterations
39  double last_difference[maxnodes];
40  int dif[maxnodes];
41  double the_maxchange[max_iteration];//Maximum change in weights per
        iteration
42  //the number of the largest difference in each generation
43  int number[max_iteration];
44  //the ovalues of the point
45  double v[max_iteration];
46  double vn[max_iteration];
47  int iterate[epslionIndex+1];//number of iteration
48  double time[epslionIndex+1];//using time
49
50  //compare
51  struct cmp_fun {
52          bool operator()(const weigh_edges& the_left, const weigh_edges&
                the_right) {
53          return the_left.second > the_right.second;
54          }
55  };
56  //read data
```

```
57  void create_idx();
58  void build_statistics(int weighted);
59  //init the value of matrix
60  void Init_value(int Init);
61  //check if converged
62  int check(int iter);
63  //Quadratic Extrapolation method
64  void Quadratic_Extrapolation();
65  //the origin pagerank alg
66  void PageRank_origin();
67  //pagerank alg using sparse matrix
68  void PageRank_sparse(int Extra, int Fix);
69  // Function to generate the result and save it to a file
70  void output_to_txt(string name);
```

accelerated_pagerank.cpp implements the functions in the header file.

accelerated_pagerank.cpp:

```
1   #include "accelerated_pagerank.h"
2
3   //read data
4   void create_idx(){
5           FILE *file1=fopen("../Wiki-Vote-id.txt", "r");
6           FILE *file2=fopen("../Wiki-Vote.txt", "r");
7           if(file1==NULL||file2==NULL){
8                   cout<<"open input file error!!!!!!"<<endl ;
9                   return;
10      }
11          int i;
12          for(i=1;i<=maxnodes;i++){
13                  char name[maxlength];
14                  if(fgets(name,20,file1)==0)break;
15                  name[strlen(name)-1]=0;
16                  id_2_name[i]=name;
17                  name_2_id[name]=i;
18      }
19      num_total=i-1 ;
20          for(i=1;i<=maxedges;i++){
21                  char temp[2*maxlength+5];
22                  char name1[maxlength];
23                  char name2[maxlength];
24                  int j=0;
25                  if(fgets(temp,2*maxlength+5,file2)==0)break ;
```

```cpp
26                        for(j=0;temp[j+1]!='\t';j++){
27                                name1[j]=temp[j];
28                        }
29                        name1[j]=0;
30                        int j0=j+2;
31                        for(j=j0;temp[j-1]!=0;j++){
32                                name2[j-j0]=temp[j];
33                        }
34                        name2[j-j0-2]=0;
35                        if(name_2_id.find(name1)==name_2_id.end())continue ;
36                        if(name_2_id.find(name2)==name_2_id.end())continue ;
37                        int index1=name_2_id[name1];
38                        int index2=name_2_id[name2];
39                        //Record this citation relationship
40            spares_matrix_cite[index2][count_inlinks[index2]]=index1;
41            count_inlinks[index2]++;
42            count_outlinks[index1]++;
43            //index1 link to index2
44            matrix[index2][index1]+=1;
45        }
46        num_of_edges=i-1;
47        fclose(file1);
48        fclose(file2);
49 }
50
51 void build_statistics(int weighted)
52 {
53            cout<<"#############################################\n";
54            cout<<"the number of pages: "<<num_total<<endl;
55            cout<<"the number of links: "<<num_of_edges<<endl;
56            cout<<"#############################################\n";
57            int the_max_num=0;
58            //Initial state, each point has the same weight. All points have
                    weights that sum to 1
59            for(int i=1;i<=num_total;i++){
60                    if(count_inlinks[i]>the_max_num)the_max_num=
                        count_inlinks[i];
61            }
62            for(int i=1;i<=num_total;i++){
63                    for(int j=0;j<count_inlinks[i];j++){
64                            count_outlinks_in[spares_matrix_cite[i][j]]+=
                                count_inlinks[i];
```

```
65                          }
66              }
67              if ( weighted==0){
68                      for ( int i =1; i<=num_total ; i++){
69                              for ( int j =0; j<count_inlinks [ i ] ; j++){
70                                      weight [ i ] [ j]=(double)1/ count_outlinks [
                                            spares_matrix_cite [ i ] [ j ] ]  ;
71                              }
72                      }
73              }
74              else {
75                      for ( int i =1; i<=num_total ; i++){
76                              for ( int j =0; j<count_inlinks [ i ] ; j++)weight [ i ] [ j
                                    ]=(double) count_inlinks [ i ] / count_outlinks_in
                                    [ spares_matrix_cite [ i ] [ j ] ]  ;
77                      }
78              }
79          // Calculate  the  transfer  matrix
80      for ( int i =1; i<=num_total ; i++){
81          for ( int j =1; j<=num_total ; j++){
82                              if ( matrix [ i ] [ j]!=0){
83                                      if ( weighted==0)matrix [ i ] [ j]=matrix [ i ] [ j
                                            ]*(1−Poblity ) / count_outlinks [ j ];
84                                      else  matrix [ i ] [ j]=matrix [ i ] [ j]*(1−
                                            Poblity )* count_inlinks [ i ] /
                                            count_outlinks_in [ j ];
85                              }
86                              matrix [ i ] [ j]+=(double) Poblity / num_total ;
87              }
88          }
89      // Record  points  that  have  not  outlinks
90      for ( int i =1;  i<=num_total ; i++){
91          if ( count_outlinks [ i]==0){
92              no_outlinks [ num_of_nooutlinks]=i ;
93              num_of_nooutlinks++;
94                          //  Update  the  transfer  matrix  by  connecting  the
                                points  that  do  not  have  an  edge  to  all  the
                                points .
95              for ( int j =1; j<=num_total ; j++)matrix [ j ] [ i]+=(double)(1−
                    Poblity ) / num_total ;
96          }
97      }
```

```cpp
98   }
99
100  //init the value of matrix
101  void Init_value(int Init){
102          if(Init==0){
103                  for(int i=1;i<=num_total;i++){
104                          val[i]=1/double(num_total);
105                  }
106          }
107          if(Init==1){
108                  int inlinks_sum=0;
109                  for(int i=1;i<=num_total;i++){
110                          inlinks_sum+=count_inlinks[i]+1 ;
111                  }
112                  for(int i=1;i<=num_total;i++){
113                          val[i]=(count_inlinks[i]+1)/double(inlinks_sum);
114                  }
115          }
116  }
117
118  //check if converged
119  int check(int iter){
120          static float threshold=0.001;    //Threshold of convergence
121          static int power=3;        //the power of Threshold
122          the_maxchange[iter]=0;   //max change of pagerank
123          //Calculate the maximum change
124          for(int i=1;i<=num_total;i++){
125                  if(val[i]-newval[i]>the_maxchange[iter]){
126                          number[iter]=i;
127                          v[iter]=val[i];
128                  vn[iter]=newval[i];
129                          the_maxchange[iter]=val[i]-newval[i];
130                  }
131                  else if(newval[i]-val[i]>the_maxchange[iter]){
132                  number[iter]=i;
133                          v[iter]=val[i];
134                          vn[iter]=newval[i];
135                          the_maxchange[iter]=newval[i]-val[i];
136                  }
137                  //some node Converge
138                  if(newval[i]-val[i]<epslion/10&&-newval[i]+val[i]<
                     epslion/10)is_Converge[i]=1;
```

```
139            val[i] = newval[i] ;
140        }
141            //start record time
142            if(the_maxchange[iter]<threshold){
143                    threshold/=10;
144                    iterate[power]=iter;
145                    QueryPerformanceCounter(&end_time);
146            time[power++]=(double)(end_time.QuadPart-start_time.QuadPart)/
                   frequ.QuadPart;
147            }
148        //Determine whether the convergence condition is satisfied
149        if(the_maxchange[iter]<epslion)return 1;
150        else return 0;
151  }
152
153  //Quadratic Extrapolation method
154  void Quadratic_Extrapolation(){
155            double lamda;
156            double sum1=0;
157            double sum2=0;
158            //Calculated as shown in the formula
159            for(int i=1;i<=num_total;i++){
160                    sum1+=last_difference[i]*last_difference[i];
161                    sum2+=last_difference[i]*dif[i];
162            }
163            lamda=sum2/sum1;
164            for(int i=1;i<=num_total;i++){
165                    newval[i]=(lamda*val[i]-newval[i])/(lamda-1);
166            }
167  }
168
169  //the origin pagerank alg
170  void PageRank_origin(){
171            //Calculated as shown in the formula
172            for(int iter=1;iter<=max_iteration/4;iter++){
173                    for(int i=1;i<=num_total;i++){
174                            newval[i]=0;
175                            for(int j=1;j<=num_total;j++)newval[i]+=matrix[i
                                ][j]*val[j];
176                    }
177                    if(check(iter)==1)break;
178        }
```

```
179  }
180
181  //pagerank alg using sparse matrix
182  void PageRank_sparse(int Extra, int Fix){
183      for (int iter=1;iter<max_iteration;iter++){
184          double sum_val0=0;
185          // Calculate the sum of values for nodes with no outgoing links
186          for (int j=0;j<num_of_nooutlinks;j++)sum_val0+=val[no_outlinks[j
                  ]];
187          // Calculate the new PageRank values for each node
188          for (int i=1;i<=num_total;i++){
189              if(Fix&&is_Converge[i])continue;
190              // Nodes with outgoing links
191              // If the node has outgoing links, assign it a value based
                      on the sum of values for nodes with no outgoing links
192              if(count_outlinks[i]!=0)newval[i]=Poblity/double(num_total)
                      +(1-Poblity)*sum_val0/double(num_total-1);
193              // If the node has no outgoing links, assign it a value
                      based on the sum of values for nodes with no outgoing
                      links
194              else newval[i]=Poblity/double(num_total)+(1-Poblity)*(
                      sum_val0-val[i])/double(num_total-1);
195              // The second part of the PageRank value is obtained by "
                      absorbing" from all incoming edges to the node
196              for(int j=0;j<count_inlinks[i];j++){
197                  // The new PageRank value of node i is updated by adding
                          the probability of not performing a random walk,
198                  // the weight of the j-th incoming edge of node i, and
                          the PageRank value of the incoming node
199                  newval[i]+=(1-Poblity)*val[spares_matrix_cite[i][j]]*
                          weight[i][j];
200              }
201          }
202          // Use quadratic extrapolation to accelerate the convergence (if
                  enabled)
203          if(Extra==1){
204              if(iter==1){
205                  for(int i=1;i<=num_total;i++)
206                      dif[i]=newval[i]-val[i];
207              }
208                      else{
209                  for(int i=1;i<=num_total;i++){
```

```
210                        last_difference[i]=dif[i];
211                        dif[i]=newval[i]-val[i];
212                    }
213                    Quadratic_Extrapolation();
214                }
215            }
216            // Check if the PageRank values have converged
217            if(check(iter)==1)break;
218        }
219    }
220
221    // Function to generate the result and save it to a file
222    void output_to_txt(string name) {
223        // Calculate the sum of all node values. If computed correctly, the
                sum of all node values should be 1.
224        double allsum=0;
225        for(int i=1;i<=num_total;i++) {
226            // Insert the node and its value into a vector for sorting
227            wedges_vec.push_back(make_pair(i,val[i]));
228            allsum+=val[i];
229        }
230        // Sort the nodes based on their values
231        sort(wedges_vec.begin(),wedges_vec.end(),cmp_fun());
232        // Open the result file
233        cout << "############# saving result ################" << endl;
234        FILE *file3;
235        FILE *file4;
236        file3=fopen((string("output/PageRank_with_")+name).c_str(),"w");
237        // Convert the node IDs back to author names and output the names
                and PageRank values
238        for(int i=0;i<num_total;i++)fprintf(file3,"%s\t%f\n",id_2_name[
                wedges_vec[i].first].c_str(),wedges_vec[i].second);
239        fclose(file3);
240        // Output the convergence information for each iteration
241        file4=fopen((string("output/statistics_")+name).c_str(),"w");
242        fprintf(file4, "acc\ttimes\n");
243        for(int power=3;power<=epslionIndex;power++)fprintf(file4,"%d\t%llf\
                n",power,time[power]);
244        fprintf(file4,"acc\titeration\n");
245        for(int power=3;power<=epslionIndex;power++)fprintf(file4,"%d\t%d\n"
                ,power,iterate[power]);
246        fprintf(file4,"iteration\tmax change node\tlast val\tnew val\tchange
```

```
            \n");
247     for(int  i=1;i<=max_iteration;i++){
248         fprintf(file4,"%d\t%d\t%.30llf\t%.30llf\t%.30llf\n",i,number[i],
                v[i],vn[i],the_maxchange[i]);
249         if(the_maxchange[i]<epslion)break;
250     }
251     fclose(file4);
252     wedges_vec.clear();
253 }
```