

暨南大学本科实验报告专用纸

课程名称	数值计算实验		成绩评定	
实验项目名称	Computing Problems		指导老师	Liangda Fang
实验项目编号	04	实验项目类型	验证型	实验地点
学生姓名	学号			
学院	国际学院	系	专业	计算机科学与技术
实验时间	2023 年 12 月 1 日上午 10:30 ~ 12:10			

I. Problem

1. Implement a simple web crawler;
2. Acquire the google matrix of the first 500 web pages from the main webpage of any university via the above web crawler, and give their adjacency matrix;
3. Implement the Power Method;
4. Compute the dominant eigenvector of the google matrix;
5. List the top 20 web pages.

II. Algorithm Summary

1. PageRank Algorithm

The PageRank algorithm, devised by Larry Page and Sergey Brin, co-founders of Google, is a method used to determine the importance of a webpage based on its relevance and connectivity within the web. It serves as a fundamental algorithm for Google’s search engine, aiding in the sorting and ranking of webpages.

Key concepts of the PageRank algorithm include:

1. Link Graph: The internet is visualized as a vast graph where webpages are represented as nodes, and hyperlinks between them form directed edges.
2. Importance Transmission: PageRank measures the importance of a webpage by evaluating both the quantity and quality of incoming links. A webpage is considered more important if it is linked by other pages, especially those of high importance.
3. Markov Process: PageRank employs the concept of a Markov chain to model the transition between webpages. Webpages are treated as states, and the hyperlinks between them represent the probabilities of transitioning between states.
4. Iterative Calculation: The algorithm uses iterative methods to compute the weight (importance) of each webpage. Initially, all webpages are assigned equal weight, and through multiple iterations based on linking relationships, the weight of each webpage is updated until convergence.

5. Stability and Convergence: PageRank continues its iterative calculations to update web-page weights until they reach a stable state where further iterations do not significantly change the weights.

In short, the algorithm is implemented based on Markov process and power iteration method. Let's look at the following example flow:

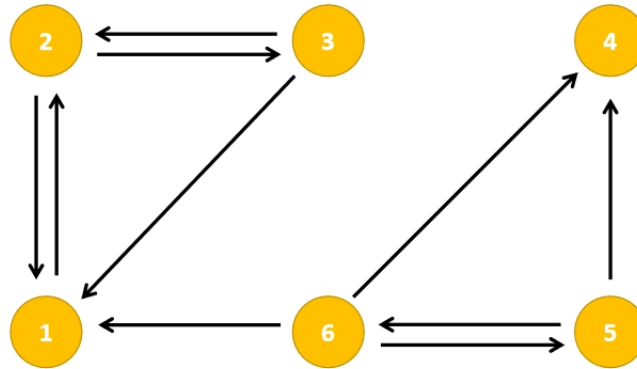


Figure 1: Page links diagram

This is a directed graph where outgoing links from a website refer to other pages. In this case, Site 1 refers to both Site 2 and Site 3. In computer science, directed graphs are commonly represented using an adjacency matrix \mathbf{A} . For instance:

$$\mathbf{A} = \begin{bmatrix} 0 & 1 & 1 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

Now, The adjacency matrix \mathbf{A} , we set $\mathbf{A}_{ij} = 1$ which indicates that Site i refers to Site j , while $\mathbf{A}_{ij} = 0$ signifies no reference. This representation helps depict the directed connections within the graph of web pages.

Typically, web pages are assigned priorities, and displaying high-priority pages at the top aims to facilitate easier access for users. However, determining the importance of these pages poses a challenge. We primarily rely on the relevance between web pages as a guide. If a webpage is linked to by many other pages, it implies higher importance. Similarly, pages referenced by high-priority pages might also carry significance. Common intuition suggests that users tend to spend more time on pages perceived as highly important.

Suppose there's a student named **Hyq** currently browsing the internet. **Hyq** has two choices while surfing:

暨南大学本科实验报告专用纸 (附页)

1. Moving to a Random Site: **Hyq** might decide to move to a random site across the internet with a probability denoted as α .
2. Clicking a Link on the Current Page: Alternatively, **Hyq** may choose to click on a link present on the current page they're visiting with a probability of $(1 - \alpha)$.

It's assumed that the probability of **Hyq** clicking on a link on any given site is equal, and there exist 'n' total websites available. Consequently, the probability of **Hyq** transitioning from Site i to Site j can be expressed as:

$$p_{ji} = \begin{cases} \frac{\alpha}{n} + (1 - \alpha) * \frac{A_{ji}}{\sum_{k=1}^n A_{ki}}, & \text{other} \\ \frac{1}{n}, & \sum_{k=1}^n A_{ki} = 0 \end{cases}$$

p_{ji} means Rosen goes from site i to j .

By this time, let's transform our adjacency matrix **A** into a "surfing" matrix **M**:

$$\mathbf{M} = \begin{bmatrix} \frac{a}{6} & \frac{a}{6} + \frac{(1-a)}{2} & \frac{a}{6} + \frac{(1-a)}{2} & \frac{1}{6} & \frac{a}{6} & \frac{a}{6} + \frac{(1-a)}{3} \\ \frac{a}{6} + (1-a) & \frac{a}{6} & \frac{a}{6} + \frac{(1-a)}{2} & \frac{1}{6} & \frac{a}{6} & \frac{a}{6} \\ \frac{a}{6} & \frac{a}{6} + \frac{(1-a)}{2} & \frac{a}{6} & \frac{1}{6} & \frac{a}{6} & \frac{a}{6} \\ \frac{a}{6} & \frac{a}{6} & \frac{a}{6} & \frac{1}{6} & \frac{a}{6} + \frac{(1-a)}{2} & \frac{a}{6} + \frac{(1-a)}{3} \\ \frac{a}{6} & \frac{a}{6} & \frac{a}{6} & \frac{1}{6} & \frac{a}{6} & \frac{a}{6} + \frac{(1-a)}{3} \\ \frac{a}{6} & \frac{a}{6} & \frac{a}{6} & \frac{1}{6} & \frac{a}{6} + \frac{(1-a)}{2} & \frac{a}{6} \end{bmatrix}$$

The **M** refers to the "Google matrix/Stochastic Matrix/Markov Matrix." An interesting fact to note is that all elements of the matrix are positive, and the sum of each column of the matrix equals 1. In other words, there exists a left eigenvector $\mathbf{v} = (1, 1, \dots, 1)$:

$$\mathbf{vM} = \mathbf{v}$$

That means the matrix **M** has an eigenvalue of 1. As per the Gershgorin circle theorem, it's established that the eigenvalues of **M** cannot exceed 1. Given the positivity of **M**, according to the Perron-Frobenius theorem, it follows that the absolute values of all eigenvalues of the matrix are less than 1, excluding the aforementioned eigenvalue of 1.

Now, when we say matrix **M** possesses a steady-state distribution, what does it signify? Let's assume a vector **p**, and \mathbf{p}_i , represents the number of people staying on page i . Assuming:

$$\mathbf{P} = (100, 0, 0, \dots, 0)$$

This implies that initially, there were 100 people staying on page 1. After a considerable duration, we're interested in determining the number of people currently staying on the various web pages:

$$\mathbf{P}_{\text{now}} = \lim_{n \rightarrow \infty} \mathbf{G}^n \mathbf{P}$$

You'll be amazed to observe the convergence of the matrix power \mathbf{P}_{now} ! This incredible outcome is due to the fact that 1 is the largest absolute eigenvalue of **M** ! Take a look at this:

Assume $c_1, c_2 \neq 0$ and v_i is eigenvector of **M**

And $|\lambda_1| > |\lambda_2| \geq |\lambda_3| \geq \dots \geq |\lambda_n|$

$$\mathbf{P} = c_1 v_1 + \dots + c_n v_n$$

$$\mathbf{P}_k = \mathbf{M}^k \mathbf{P} = c_1 \lambda_1^k v_1 + \dots + c_n \lambda_n^k v_n$$

$$\frac{\mathbf{P}_k}{\lambda_1^k} = c_1 v_1 + c_2 \left(\frac{\lambda_2}{\lambda_1}\right)^k v_2 + \dots + c_n \left(\frac{\lambda_n}{\lambda_1}\right)^k v_n$$

$$\lim_{k \rightarrow \infty} \left(\frac{\lambda_j}{\lambda_1}\right)^k \rightarrow 0$$

That's precisely why \mathbf{P}_{now} will converge! It will gradually converge towards the eigenvector of the greatest absolute eigenvalue of \mathbf{M} , eventually reaching a multiple of that eigenvalue. This process is referred to as the power iteration method.

By sorting the vector \mathbf{P}_{now} from largest to smallest, we've successfully obtained the ranking of pages! We've accomplished it gracefully.

III. Experimental procedures

Step1: Define the Web_Crawler class to implement the method of crawling pages and building an adjacency matrix between pages.

Step2: Start crawling pages from "https://www.jnu.edu.cn" and continue crawling until 500 web pages have been crawled.

Step3: Draw the adjacency matrix of these pages and convert the matrix to a google matrix.

Step4: Implement the Power Iteration Method and calculate the dominant eigenvector of the google matrix.

Step5: List the top 20 web pages with the highest rank value.

暨南大学本科实验报告专用纸 (附页)

IV. Result Analysis

First we crawled 500 web pages:

```
E:\Miniconda3\python.exe D:\icloud\Drive\Documents\文\暨大JNU_Course\Junior\NumericalComputation\Lab\Lab4\H3A\1\main.py
Crawling Page 1 : https://www.jnu.edu.cn , Status: 200
Crawling Page 2 : https://info.jnu.edu.cn , Status: 200
Crawling Page 3 : https://mail.jnu.edu.cn , Status: No-HTML content
Crawling Page 4 : https://ehall.jnu.edu.cn , Status: 200
Crawling Page 5 : https://lib.jnu.edu.cn , Status: 200
Crawling Page 6 : https://english.jnu.edu.cn , Status: 200
Crawling Page 7 : https://qsh.jnu.edu.cn , Status: 200
Crawling Page 8 : https://pamy.jnu.edu.cn , Status: 200
Crawling Page 9 : https://zh.jnu.edu.cn , Status: 200
Crawling Page 10 : https://sz.jnu.edu.cn , Status: 200
Crawling Page 11 : https://zsh.jnu.edu.cn , Status: 200
Crawling Page 12 : https://yz.jnu.edu.cn , Status: 200
Crawling Page 13 : https://hwy.jnu.edu.cn/13405/main.htm , Status: 200
Crawling Page 14 : https://jyxy.jnu.edu.cn , Status: 200
Crawling Page 15 : https://career.jnu.edu.cn , Status: 200
Crawling Page 16 : https://fseulty.jnu.edu.cn , Status: 200
Crawling Page 17 : https://jwc.jnu.edu.cn , Status: 200
Crawling Page 18 : https://hwy.jnu.edu.cn , Status: 200
Crawling Page 19 : https://kjc.jnu.edu.cn , Status: 200
Crawling Page 20 : https://skc.jnu.edu.cn , Status: 200
Crawling Page 21 : https://www.jnu.edu.cn/jyxy/list.htm , Status: 200
Crawling Page 22 : https://wsc.jnu.edu.cn , Status: 200
Crawling Page 23 : https://gccgl.jnu.edu.cn/jmzpf/ajzp/base/index , Status: 200
Crawling Page 24 : https://hrdam.jnu.edu.cn/345/list.htm , Status: 200
Crawling Page 25 : https://zhongh.jnu.edu.cn , Status: 200
Crawling Page 26 : https://jwc.jnu.edu.cn , Status: 200
Crawling Page 27 : https://netc.jnu.edu.cn , Status: 200
Crawling Page 28 : https://www.jnu.edu.cn/tz/list.asp , Status: 200
Crawling Page 29 : https://www.jnu.edu.cn/gg/list.asp , Status: 200
Crawling Page 30 : https://news.jnu.edu.cn/jnyw/yw/2023/11/10/18371652545.html , Status: 404
Crawling Page 31 : https://mp.weixin.qq.com/s/sfhARiab97U1A7Wc13PQg , Status: 200
Crawling Page 32 : https://news.jnu.edu.cn/cq12.html , Status: 200
Crawling Page 33 : https://ddesd.jnu.edu.cn , Status: 200
Crawling Page 34 : https://news.jnu.edu.cn/cq11.html , Status: 200
Crawling Page 35 : https://news.jnu.edu.cn/cq13.html , Status: 200
Crawling Page 36 : https://news.jnu.edu.cn , Status: 200
Crawling Page 37 : https://jwc.jnu.edu.cn/2023/1201/c8320a783087/page.htm , Status: 200
Crawling Page 38 : https://jwc.jnu.edu.cn/2023/0823/c8320a761121/page.htm , Status: 200
Crawling Page 39 : https://news.jnu.edu.cn/cq18.html , Status: 200
Crawling Page 40 : https://news.jnu.edu.cn/cq19.html , Status: 200
Crawling Page 41 : https://news.jnu.edu.cn/cq15.html , Status: 200
```

Figure 2: Crawling websites

```
Crawling Page 468 : https://hwy.jnu.edu.cn/2022/062c/c13398a707387/page.htm , Status: 200
Crawling Page 469 : https://hwy.jnu.edu.cn/13399/list.htm , Status: 200
Crawling Page 470 : https://hwy.jnu.edu.cn/2022/062c/c13438a704655/page.htm , Status: 200
Crawling Page 471 : https://hwy.jnu.edu.cn/13397/list.htm , Status: 200
Crawling Page 472 : https://hwy.jnu.edu.cn/13402/list.htm , Status: 200
Crawling Page 473 : https://hwy.jnu.edu.cn/2023/0831/c13397a761843/page.htm , Status: 200
Crawling Page 474 : https://hwy.jnu.edu.cn/xsjy/main.htm , Status: 200
Crawling Page 475 : https://ehall.jnu.edu.cn/infoplus/form/KYXMYST750/start , Status: 200
Crawling Page 476 : https://kjc.jnu.edu.cn/2016/1206/c3501a93205/page.htm , Status: 200
Crawling Page 477 : https://kjc.jnu.edu.cn/web/onlineSurvey/api/lookOnlineSurvey/new.cat7_p=XYWNH1m0D0xHTmc0Dx3m09T1V_banner-1&tt=0_8846275260812899&request_locale=zh_CN&localeChinese=256check=2 , Status: 200
Crawling Page 478 : https://kjc.jnu.edu.cn/2023/0705/c917a635537/page.htm , Status: 200
Crawling Page 479 : https://jwc.jnu.edu.cn/11/31/c6234a758355/page.asp , Status: 200
Crawling Page 480 : https://csc.jnu.edu.cn , Status: 200
Crawling Page 481 : https://hrdam.jnu.edu.cn , Status: 200
Crawling Page 482 : https://kjc.jnu.edu.cn/33611/list.htm , Status: 200
Crawling Page 483 : https://kjc.jnu.edu.cn/2023/1207/c3403a783665/page.htm , Status: 200
Crawling Page 484 : https://kjc.jnu.edu.cn/33211/list.htm , Status: 200
Crawling Page 485 : https://kjc.jnu.edu.cn/2023/1122/c3403a782055/page.htm , Status: 200
Crawling Page 486 : https://kjc.jnu.edu.cn/33205/list.htm , Status: 200
Crawling Page 487 : https://kjc.jnu.edu.cn/2023/1205/c3403a783403/page.htm , Status: 200
Crawling Page 488 : https://kjc.jnu.edu.cn/33645/list.htm , Status: 200
Crawling Page 489 : https://kjc.jnu.edu.cn/33629/list.htm , Status: 200
Crawling Page 490 : https://kjc.jnu.edu.cn/33607/list.htm , Status: 200
Crawling Page 491 : https://kjc.jnu.edu.cn/33215/list.htm , Status: 200
Crawling Page 492 : https://kjc.jnu.edu.cn/2022/0516/c33581a699277/page.htm , Status: 200
Crawling Page 493 : https://kjc.jnu.edu.cn/784/list.htm , Status: 200
Crawling Page 494 : https://kjc.jnu.edu.cn/33215/list.htm , Status: 200
Crawling Page 495 : https://kjc.jnu.edu.cn/33625/list.htm , Status: 200
Crawling Page 496 : https://kjc.jnu.edu.cn/33173/list.htm , Status: 200
Crawling Page 497 : https://kjc.jnu.edu.cn/33603/list.htm , Status: 200
Crawling Page 498 : https://kjc.jnu.edu.cn/33207/list.htm , Status: 200
Crawling Page 499 : https://kjc.jnu.edu.cn/33617/list.htm , Status: 200
Crawling Page 500 : https://kjc.jnu.edu.cn/9360/list.htm , Status: 200
Crawling Page 501 : https://kjc.jnu.edu.cn/9355/list.htm , Status: 200
Crawling Page 502 : https://xsc.jnu.edu.cn/5840/list.htm , Status: 200
Crawling Page 503 : https://www.jnu.edu.cn/2022/0829/c2593a714477/page.htm , Status: 200
Crawling Page 504 : https://english.jnu.edu.cn/2016/1031/c2490a47917/page.htm , Status: No-HTML content
Crawling Page 505 : https://lkye.southcn.com , Status: 200
Crawling Page 506 : https://jyxy.jnu.edu.cn/28101/list.htm , Status: 200
Crawling Page 507 : https://qddi.southcn.com , Status: 200
Crawling completed. Total sites: 500
errs: 500
```

Figure 3: Crawling websites(cont.)

Based on the link relationships between these 500 pages, an adjacency matrix diagram can be constructed as follows:

暨南大学本科实验报告专用纸 (附页)

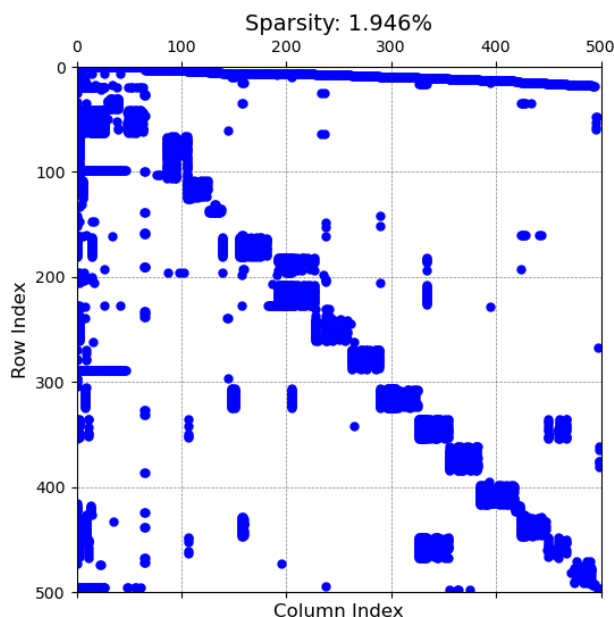


Figure 4: Sparsity matrix of these 500 pages

According to the PageRank algorithm, we get the top 20 web pages with the highest rank values:

```
TOP 1: 暨南大学图书馆
URL: https://lib.jnu.edu.cn
Rank Value: 0.036824973012285524
TOP 2: 暨南大学-华侨最高学府-国家“双一流”建设高校
URL: https://www.jnu.edu.cn
Rank Value: 0.03621698036050209
TOP 3: 暨南大学图书馆
URL: https://lib.jnu.edu.cn/department/index/11
Rank Value: 0.0240492226913659
TOP 4: 图书馆概况-暨南大学图书馆
URL: https://lib.jnu.edu.cn/home/service/0/2774
Rank Value: 0.020536069120696315
TOP 5: 暨南大学统一身份认证登录指南
URL: https://netc.jnu.edu.cn/2020/1124/c10374a565499/page.htm
Rank Value: 0.02033873219076459
TOP 6: Mynet校园网用户自助服务平台
URL: https://mynet.jnu.edu.cn/app/forgetPassword.html
Rank Value: 0.02033873219076459
TOP 7: 暨南大学图书馆
URL: https://lib.jnu.edu.cn/department/index/36
Rank Value: 0.01763508439933019
TOP 8: 暨南大学图书馆
URL: https://lib.jnu.edu.cn/department/index/22
Rank Value: 0.017047205259862395
TOP 9: 暨南大学图书馆
URL: https://lib.jnu.edu.cn/department/index/37
Rank Value: 0.01685197034276468
TOP 10: 校地合作
URL: https://kjc.jnu.edu.cn/33625/list.htm
Rank Value: 0.009215398454833084
```

Figure 5: Top 20 web pages(1-10)

```
TOP 11: 其他
URL: https://kjc.jnu.edu.cn/33611/list.htm
Rank Value: 0.008737289959345666
TOP 12: 统一身份认证平台
URL: https://siteadmin.jnu.edu.cn
Rank Value: 0.008130817200861949
TOP 13: 创新平台科
URL: https://kjc.jnu.edu.cn/33215/list.htm
Rank Value: 0.007815244223543832
TOP 14: JiNan University Undergraduate Admission
URL: https://admission.jnu.edu.cn
Rank Value: 0.007740225664963853
TOP 15: 首页 - 暨南大学新闻网
URL: https://news.jnu.edu.cn
Rank Value: 0.007602933214455526
TOP 16: 机构简介
URL: https://kjc.jnu.edu.cn/33205/list.htm
Rank Value: 0.007337932383741568
TOP 17: 支部建设
URL: https://skc.jnu.edu.cn/9317/list.htm
Rank Value: 0.007305802612076512
TOP 18: 国侨办项目
URL: https://skc.jnu.edu.cn/9340/list.htm
Rank Value: 0.007305802612076512
TOP 19: 国家级
URL: https://skc.jnu.edu.cn/9355/list.htm
Rank Value: 0.007305802612076512
TOP 20: 教学科研 - 暨南大学新闻网
URL: https://news.jnu.edu.cn/col9.html
Rank Value: 0.0072411138069266735
```

Figure 6: Top 20 web pages(11-20)

Finally, we can plot the website rank value distribution:

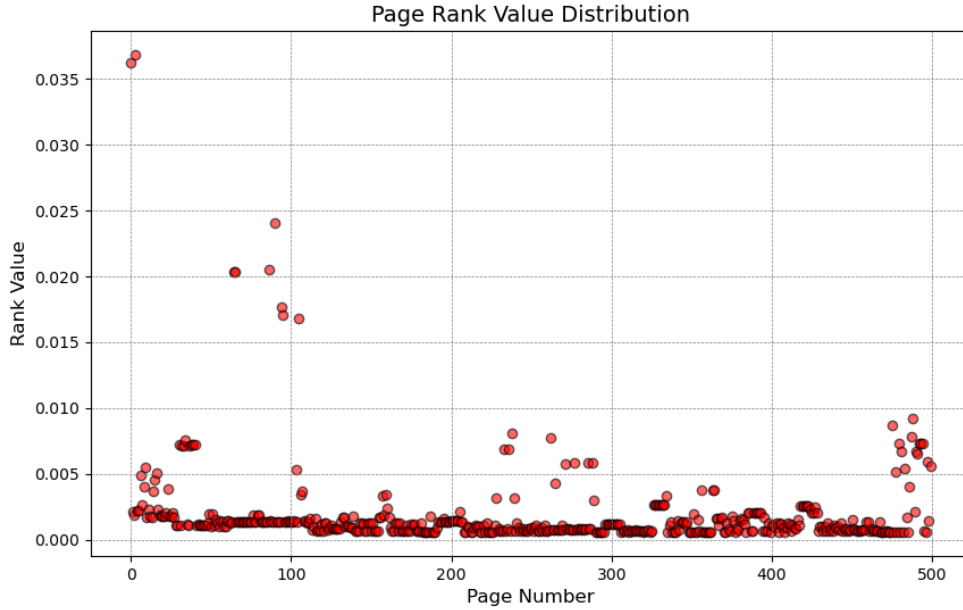


Figure 7: Rank value distribution

V. Experimental Summary

In this lab, we extracted data from 500 distinct web pages linked to Jinan University's official site, "<https://www.jnu.edu.cn>", forming an adjacency matrix \mathbf{A} . Utilizing \mathbf{A} , we constructed the Google matrix \mathbf{G} and employed the power iteration algorithm to determine the relative importance and rankings of these sites.

It's important to note that for the adjacency matrix \mathbf{A} , we ensured the sum of each column equals 1, contrary to the row sum. Additionally, we imposed a limit on the number of relative URLs per page to a maximum of 20. This was particularly relevant for pages like "<https://www.jnu.edu.cn>" to prevent overrepresentation of nodes from the same website's sub-pages. Moreover, we encountered numerous URLs referring to the same site (for instance, "<http://www.jnu.edu.cn>", "<https://www.jnu.edu.cn>", and "<https://www.jnu.edu.cn/>"). To address this, we focused exclusively on URLs starting with 'https' and eliminated any trailing '/' to avoid redundancy.

VI. Appendix: Source Code

In this experiment, the code implementation part uses Python and a series of related third-party libraries, including: **numpy**, **matplotlib**, **requests** and **beautifulsoup4**.

1. web_crawler.py

```
1 from requests import get, head
2 from bs4 import BeautifulSoup
3 from random import sample
4 from urllib.parse import urljoin
5
6 # Constants for crawler configuration
7 MAX_SITES = 500
8 MAX_RELATIVE_LINKS = 20
9
10
11 class Web_Crawler:
12     HEADERS = {
13         'user-agent': 'Mozilla/5.0 (Macintosh; Intel Mac OS X 10_14_5) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/87.0.4280.88 Safari/537.36'
14     }
15
16     def __init__(self, base_url, verbose=False):
17         self.to_crawl = [base_url]
18         self.site_count = 1
19         self.crawled = set()
20         self.unreachable = set()
21         self.link_graph = {}
22         self.verbose = verbose
23
24     def extract_links(self, url):
25         try:
26             response = head(url)
27             if 'text/html' not in response.headers.get('Content-Type', ''):
28                 raise ValueError("No-HTML content")
29             response = get(url, headers=self.HEADERS, timeout=10)
30         except:
31             self.crawled.remove(url)
32             self.unreachable.add(url)
33             self.site_count = self.site_count - 1
34             return -1
35
36         response.encoding = 'utf-8'
37         soup = BeautifulSoup(response.text, 'lxml')
38
39         # Collect and normalize links
40         links = [link['href'] for link in soup.find_all('a', href=True)]
41         links = self._normalize_links(url, links)
42
43         for link in links:
44             if link not in self.to_crawl and \
45                 link not in self.crawled and \
46                 link not in self.unreachable and \
47                 self.site_count < MAX_SITES:
48                 self.to_crawl.append(link)
49                 self.site_count += 1
```



```
50
51     self.link_graph[url] = links
52     response.close()
53     return response.status_code
54
55     def _normalize_links(self, base_url, links):
56         relative = [link for link in links if link.startswith('/')]
57         absolute = [link for link in links if link.startswith('https')]
58         relative = sample(relative, min(len(relative), MAX_RELATIVE_LINKS))
59         links = [link.rstrip('/') for link in absolute + [urljoin(base_url,
60             rel) for rel in relative]]
61         return links
62
63     def run(self):
64         index = 1
65         while self.to_crawl:
66             current_url = self.to_crawl.pop(0)
67             self.crawled.add(current_url)
68             status = self.extract_links(current_url)
69             if self.verbose:
70                 print('Crawling Page', index, ':', current_url, ', Status:',
71                     'No-HTML content' if status == -1 else status)
72             index += 1
73         if self.verbose:
74             print('Crawling completed. Total sites:', self.site_count)
75             print('keys:', len(self.link_graph.keys()))
76
77     def get_adjacent_matrix(self):
78         keys = tuple(self.link_graph.keys())
79
80         for key in keys:
81             for site in self.link_graph[key]:
82                 if site not in keys or site == key:
83                     self.link_graph[key] = [x for x in self.link_graph[key]
84                         if x != site]
85
86         matrix = [[0] * self.site_count for _ in range(self.site_count)]
87
88         for site, links in self.link_graph.items():
89             site_index = keys.index(site)
90             for link in links:
91                 link_index = keys.index(link)
92                 matrix[site_index][link_index] = 1
93
94         return matrix
95
96     def get_page_title(url):
97         try:
98             response = get(url, headers=Web_Crawler.HEADERS)
99             response.encoding = 'utf-8'
100             soup = BeautifulSoup(response.text, 'lxml')
101             title = soup.find('title').text
102         finally:
103             response.close()
104         return title
```

2. algorithm.py

```
1 import numpy as np
2
3
4 def power_iteration(matrix):
5     num_nodes = matrix.shape[0]
6     # Initialize vector: [1, 0, 0, ..., 0]
7     rank_vector = np.zeros(num_nodes)
8     rank_vector[0] = 1
9     while True:
10         new_rank_vector = np.dot(matrix, rank_vector)
11         change = np.linalg.norm(new_rank_vector - rank_vector)
12         # Convergence check (threshold set arbitrarily)
13         if change < 1e-20:
14             break
15         else:
16             rank_vector = new_rank_vector
17
18     return new_rank_vector
19
20
21 def convert_to_google_matrix(adj_matrix, num_sites, alpha):
22     for row in range(num_sites):
23         row_sum = sum(adj_matrix[row])
24         if row_sum == 0:
25             for col in range(num_sites):
26                 adj_matrix[row][col] = 1 / num_sites
27         else:
28             for col in range(num_sites):
29                 random_surf = alpha * (1 / num_sites)
30                 outlink_surf = (1 - alpha) * (adj_matrix[row][col] / row_sum)
31                 adj_matrix[row][col] = random_surf + outlink_surf
32     return adj_matrix
```

3. utils.py

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3 from json import dump, load
4
5 # Constants
6 DATA_FILE = 'crawled_data.json'
7 ALPHA = 0.15
8 NUM_TOP_PAGES = 20
9
10
11 def save_crawl_data(crawler, filename):
12     data = {
13         'adjacent_matrix': crawler.get_adjacent_matrix(),
14         'directed_graph': crawler.link_graph,
15         'num_site': crawler.site_count
16     }
17     with open(filename, 'w', encoding='utf-8') as file:
18         dump(data, file, ensure_ascii=False, indent=4)
19
20
21 def load_data(filename):
22     with open(filename, 'r') as file:
23         return load(file)
24
25
26 def plot_sparsity(matrix):
27     size = len(matrix)
28     plt.figure(figsize=(10, 6)) # Set a larger figure size
29     plt.spy(matrix, markersize=5, marker='o', color='blue') # Increase
30     # marker size and change color
31     plt.xlim([0, size])
32     plt.ylim([size, 0])
33
34     # Calculate sparsity
35     non_zero_count = np.count_nonzero(matrix)
36     sparsity_percentage = (non_zero_count / (size ** 2)) * 100
37     plt.title(f'Sparsity: {sparsity_percentage:.3f}%', fontsize=14)
38     plt.xlabel('Column Index', fontsize=12)
39     plt.ylabel('Row Index', fontsize=12)
40     plt.grid(color='gray', linestyle='--', linewidth=0.5)
41     plt.show()
42
43 def plot_rank_distribution(rank_values):
44     plt.figure(figsize=(10, 6)) # Set a larger figure size
45     plt.scatter(range(len(rank_values)), rank_values, color='red', alpha=0.6,
46               edgecolor='black') # Adjust color, transparency, and edge
47     plt.title('Page Rank Value Distribution', fontsize=14)
48     plt.xlabel('Page Number', fontsize=12)
49     plt.ylabel('Rank Value', fontsize=12)
50     plt.grid(True, color='gray', linestyle='--', linewidth=0.5)
51     plt.show()
```

4. main.py

```
1 from web_crawler import Web_Crawler, get_page_title
```

暨南大学本科实验报告专用纸 (附页)

```
2 from algorithm import power_iteration, convert_to_google_matrix
3 from utils import *
4 from os.path import exists
5 from datetime import datetime
6 import numpy as np
7
8
9 def crawl_website(base_url, verbose=True):
10     crawler = Web_Crawler(base_url=base_url, verbose=verbose)
11     start_time = datetime.now()
12     crawler.run()
13     end_time = datetime.now()
14     print('Crawling_time_used:', (end_time - start_time).total_seconds(), '
15         seconds')
16     return crawler
17
18 def display_top_pages(page_rank, num_pages):
19     for i, (url, rank) in enumerate(page_rank[:num_pages]):
20         title = get_page_title(url)
21         print(f'TOP_{i+1}:_{title}\nURL:_{url}\nRank_Value:_{rank}')
22
23
24 def main():
25     # Crawl and save data if not already done
26     if not exists(DATA_FILE):
27         crawler = crawl_website('https://www.jnu.edu.cn')
28         save_crawl_data(crawler, DATA_FILE)
29
30     # Load data from file
31     data = load_data(DATA_FILE)
32
33     # Plotting and processing
34     plot_sparsity(data['adjacent_matrix'])
35     google_matrix = convert_to_google_matrix(np.array(data['adjacent_matrix'
36     ], dtype=float), data['num_site'], ALPHA)
37
38     # Power iteration
39     rank_values = power_iteration(google_matrix.transpose())
40     plot_rank_distribution(rank_values)
41
42     # Ranking and displaying top pages
43     ranked_pages = [(key, rank_values[i]) for i, key in enumerate(data['
44         directed_graph'].keys())]
45     page_rank = sorted(ranked_pages, key=lambda x: x[1], reverse=True)
46     display_top_pages(page_rank, NUM_TOP_PAGES)
47
48 if __name__ == "__main__":
49     main()
```