

Marisabel Chang Chan, Ocean Lu, Richard Pham, Derek Yee  
CS4990 - Web Search  
Dr. Ben Steichen  
9 May 2019

## Project 2 - Report

### Introduction

The purpose of this project is to work with two important components of web search and search engines: Web Crawler and PageRank Calculation. The Web Crawler finds and downloads web pages by following links that are contained in the visited pages. The PageRank Calculation shows the rank of the most important web pages, people, citations, etc. in a collection.

### Method (Code Analysis)

We separated our code into three different sections: The utility methods contains the functions that we use in the other two sections (Main driver and Print PageRank). In the main driver, we use the function to crawl and update the outlink table and PageRank list.

For the utility methods, we used different functions to parse the HTML text, update the outlink table, and check if the link is contained in the queue. To parse the HTML text, we use the `get_html()` function that returns the downloaded HTML texts only if the URL is a "text/html" type and the HTTP status is OK. The `parse_html()` function gets the links that are contained in the HTML document and returns a set that contains unique links. To update the outlink table, we created the `updates_outlink_table` method that adds a row and a column of zeros to the outlink table when a new page is visited and set 1 to the specific row (codomain) and column (domain) in the table. Then, it increments the number of outlinks to the page that the outlink comes from. Also, we used the `divide_outlinks()` function that divides all items of each column with its respective number of outlinks, so that the sum of all items of each column is equal to 1. For the checking functions, we used the `has_link()` function that checks if the queue contains the link already.

Finally, for the utility function, we have `calculate_page_rank()` that calculates the PageRank of each page contained in the graph. In this function, we give the same weight to each page in the PageRank list. Then, we do a matrix multiplication of the outlink table and the PageRank list until the previous PageRank list is equal to the actual PageRank.

In the Main Driver section, we request the HTML document, get the outlink list, and update the outlink table. Then, we check if each outlink is from the same base URL of the seed and if it has been crawled. If so, we increment one to the specific position in the outlink table. If not, we add the outlink to the queue or update the outlink list for the specific link. Finally, we divide each column a the specific number of outlinks.

In the Print PageRank section, we calculate the PageRank and display the top 100 pages that have higher PageRank.

## **Challenges**

We encountered different challenges, such as retrieving the correct links to be parsed, updating the outlink table, and optimizing the runtime performance of the program.

In terms of getting the correct links, we had a problem that some links used alias names, causing outlink weights to be split. For example, `cpp.edu` and `cpp.edu/index.shtml` were considered two different pages, even though both links are linked to the same page. We fixed this problem by searching for the string “index” and removing it and everything after it. Also, we had a problem that we were concatenating a path of a URL to a URL that already contains the same path (ex. `cpp.edu/~aboutcpp/~aboutcpp/`). We fixed this problem by checking if the path was already contained in the URL before concatenating.

We encountered an issue when attempting to update the outlink table, where our table retrieved incorrect results from the web crawl. We solved this problem by creating a list that contained all the pages that land to the same page. Thus, every time a page is removed from the queue, the outlink table is updated with the changes.

Finally, we ran into another problem where some pages were taking more time checking if their outlinks were already visited or in the queue. Then, we realized that the list of the links that was returned in the `parse_html()` method was not unique. We solved this problem by changing the data structure of the list to a set. This optimized the amount of time it took to check outlinks, since duplicates do not exist in sets. Also, we changed the data structure of the `visited_links` variable to a set, too, for the same reason.

## PageRanks

SEED: '<https://www.cpp.edu>'

Number Pages Downloaded: 2500

Lambda: 0.20

Rank	Link	PageRank
1	<a href="https://www.cpp.edu">https://www.cpp.edu</a>	0.04525294055587124331
2	<a href="https://www.cpp.edu/~aboutcpp">https://www.cpp.edu/~aboutcpp</a>	0.04525294055587124331
3	<a href="https://www.cpp.edu/file-viewers.shtml">https://www.cpp.edu/file-viewers.shtml</a>	0.04521133778471275733
4	<a href="https://www.cpp.edu/~library">https://www.cpp.edu/~library</a>	0.02367918239143323653
5	<a href="https://www.cpp.edu/accessibility.shtml">https://www.cpp.edu/accessibility.shtml</a>	0.02236785926170458197
6	<a href="https://www.cpp.edu/website-feedback.php">https://www.cpp.edu/website-feedback.php</a>	0.02236785926170458197
7	<a href="https://www.cpp.edu/privacy.shtml">https://www.cpp.edu/privacy.shtml</a>	0.02236785926170458197
8	<a href="https://www.cpp.edu/events/day.php?date=04/19/2019">https://www.cpp.edu/events/day.php?date=04/19/2019</a>	0.00792685254681482765
9	<a href="https://www.cpp.edu/~stratcomm/for-the-media.shtml">https://www.cpp.edu/~stratcomm/for-the-media.shtml</a>	0.00373862069775227068
10	<a href="https://www.cpp.edu/~class">https://www.cpp.edu/~class</a>	0.00346597183863431172
11	<a href="https://www.cpp.edu/~aboutcpp/why-cpp/well-ranked.shtml">https://www.cpp.edu/~aboutcpp/why-cpp/well-ranked.shtml</a>	0.00306330151733182586
12	<a href="https://www.cpp.edu/~aboutcpp/calpolypomona-overview/mission-and-values.shtml">https://www.cpp.edu/~aboutcpp/calpolypomona-overview/mission-and-values.shtml</a>	0.00306277065477273888
13	<a href="https://www.cpp.edu/~aboutcpp/calpolypomona-overview/facts-and-figures.shtml">https://www.cpp.edu/~aboutcpp/calpolypomona-overview/facts-and-figures.shtml</a>	0.00305932128978270981
14	<a href="https://www.cpp.edu/~aboutcpp/why-cpp/affordable-and-accessible.shtml">https://www.cpp.edu/~aboutcpp/why-cpp/affordable-and-accessible.shtml</a>	0.00305932128978270981
15	<a href="https://www.cpp.edu/~aboutcpp/why-cpp/quality-learning-experiences.shtml">https://www.cpp.edu/~aboutcpp/why-cpp/quality-learning-experiences.shtml</a>	0.00305932128978270981
16	<a href="https://www.cpp.edu/~aboutcpp/calpolypomona-overview/public-good.shtml">https://www.cpp.edu/~aboutcpp/calpolypomona-overview/public-good.shtml</a>	0.00305932128978270981
17	<a href="https://www.cpp.edu/~aboutcpp/heritage/history-archives.shtml">https://www.cpp.edu/~aboutcpp/heritage/history-archives.shtml</a>	0.00305932128978270981
18	<a href="https://www.cpp.edu/~aboutcpp/calpolypomona-overview/campus-community.shtml">https://www.cpp.edu/~aboutcpp/calpolypomona-overview/campus-community.shtml</a>	0.00305932128978270981
19	<a href="https://www.cpp.edu/~aboutcpp/heritage/kellogg-house-pomona.shtml">https://www.cpp.edu/~aboutcpp/heritage/kellogg-house-pomona.shtml</a>	0.00305932128978270981
20	<a href="https://www.cpp.edu/~aboutcpp/calpolypomona-overview/presidents-message.shtml">https://www.cpp.edu/~aboutcpp/calpolypomona-overview/presidents-message.shtml</a>	0.00305932128978270981
21	<a href="https://www.cpp.edu/~aboutcpp/heritage/kellogg-connection.shtml">https://www.cpp.edu/~aboutcpp/heritage/kellogg-connection.shtml</a>	0.00305932128978270981
22	<a href="https://www.cpp.edu/~agri">https://www.cpp.edu/~agri</a>	0.00275032687609386628
23	<a href="https://www.cpp.edu/~stratcomm/sharing-news.shtml">https://www.cpp.edu/~stratcomm/sharing-news.shtml</a>	0.00266080444849505979
24	<a href="https://www.cpp.edu/~stratcomm/getting-news.shtml">https://www.cpp.edu/~stratcomm/getting-news.shtml</a>	0.00266080444849505979
25	<a href="https://www.cpp.edu/~stratcomm/contact-us.shtml">https://www.cpp.edu/~stratcomm/contact-us.shtml</a>	0.00262119628939905345
26	<a href="https://www.cpp.edu/~cba">https://www.cpp.edu/~cba</a>	0.00259770633231651196
27	<a href="https://www.cpp.edu/~class/news">https://www.cpp.edu/~class/news</a>	0.00243586586860777311
28	<a href="https://www.cpp.edu/~ceis">https://www.cpp.edu/~ceis</a>	0.00242422467303862223
29	<a href="https://www.cpp.edu/~alumni">https://www.cpp.edu/~alumni</a>	0.00240611008942580827
30	<a href="https://www.cpp.edu/~class/departments.shtml">https://www.cpp.edu/~class/departments.shtml</a>	0.00234286470377767725
31	<a href="https://www.cpp.edu/~class/centers.shtml">https://www.cpp.edu/~class/centers.shtml</a>	0.00233857709872082199
32	<a href="https://www.cpp.edu/~class/contact-us">https://www.cpp.edu/~class/contact-us</a>	0.00233857709872082199
33	<a href="https://www.cpp.edu/~class/giving">https://www.cpp.edu/~class/giving</a>	0.00233857709872082199
34	<a href="https://www.cpp.edu/~class/alumni">https://www.cpp.edu/~class/alumni</a>	0.00233857709872082199
35	<a href="https://www.cpp.edu/~class/about">https://www.cpp.edu/~class/about</a>	0.00233857709872082199
36	<a href="https://www.cpp.edu/~admissions">https://www.cpp.edu/~admissions</a>	0.00230084433245097411
37	<a href="https://www.cpp.edu/events/day.php?date=04/19/2019&amp;audience=all">https://www.cpp.edu/events/day.php?date=04/19/2019&amp;audience=all</a>	0.00228168339442923153
38	<a href="https://www.cpp.edu/~extended-education">https://www.cpp.edu/~extended-education</a>	0.00228150164806902684
39	<a href="https://www.cpp.edu/~giving">https://www.cpp.edu/~giving</a>	0.00214335725244028334
40	<a href="https://www.cpp.edu/~outreach/tours.shtml">https://www.cpp.edu/~outreach/tours.shtml</a>	0.00205053938013871834
41	<a href="https://www.cpp.edu/~aboutcpp/administration.shtml">https://www.cpp.edu/~aboutcpp/administration.shtml</a>	0.00195504803411726893
42	<a href="https://www.cpp.edu/~sci">https://www.cpp.edu/~sci</a>	0.00194510332999734683
43	<a href="https://www.cpp.edu/~collins">https://www.cpp.edu/~collins</a>	0.00181985533679030398
44	<a href="https://www.cpp.edu/~alumni/outstanding-alumni/distinguished-alumni">https://www.cpp.edu/~alumni/outstanding-alumni/distinguished-alumni</a>	0.00180132409527504890
45	<a href="https://www.cpp.edu/~commencement">https://www.cpp.edu/~commencement</a>	0.00170074912942810982
46	<a href="https://www.cpp.edu/~agri/student-life">https://www.cpp.edu/~agri/student-life</a>	0.00163290246774064687
47	<a href="https://www.cpp.edu/~agri/giving">https://www.cpp.edu/~agri/giving</a>	0.00163290246774064687
48	<a href="https://www.cpp.edu/~agri/about">https://www.cpp.edu/~agri/about</a>	0.00163290246774064687
49	<a href="https://www.cpp.edu/~agri/degrees">https://www.cpp.edu/~agri/degrees</a>	0.00163290246774064687
50	<a href="https://www.cpp.edu/~agri/contact-us">https://www.cpp.edu/~agri/contact-us</a>	0.00163290246774064687



51	<a href="https://www.cpp.edu/~agri/alumni">https://www.cpp.edu/~agri/alumni</a>	0.00163290246774064687
52	<a href="https://www.cpp.edu/~agri/faculty-staff">https://www.cpp.edu/~agri/faculty-staff</a>	0.00163290246774064687
53	<a href="https://www.cpp.edu/~agri/departments">https://www.cpp.edu/~agri/departments</a>	0.00163290246774064687
54	<a href="https://www.cpp.edu/~officeofequity/titleIX">https://www.cpp.edu/~officeofequity/titleIX</a>	0.00161765111346531538
55	<a href="https://www.cpp.edu/~aboutcpp/visitor-information/dining.shtml">https://www.cpp.edu/~aboutcpp/visitor-information/dining.shtml</a>	0.00152961545132126119
56	<a href="https://www.cpp.edu/~aboutcpp/visitor-information/parking.shtml">https://www.cpp.edu/~aboutcpp/visitor-information/parking.shtml</a>	0.00150509600981396912
57	<a href="https://www.cpp.edu/~aboutcpp/visitor-information/attractions.shtml">https://www.cpp.edu/~aboutcpp/visitor-information/attractions.shtml</a>	0.00149706247936096234
58	<a href="https://www.cpp.edu/~aboutcpp/news.shtml">https://www.cpp.edu/~aboutcpp/news.shtml</a>	0.00148530596610023206
59	<a href="https://www.cpp.edu/~aboutcpp/visitor-information/shopping.shtml">https://www.cpp.edu/~aboutcpp/visitor-information/shopping.shtml</a>	0.00148530596610023206
60	<a href="https://www.cpp.edu/~aboutcpp/visitor-information/lodging.shtml">https://www.cpp.edu/~aboutcpp/visitor-information/lodging.shtml</a>	0.00148530596610023206
61	<a href="https://www.cpp.edu/~aboutcpp/visitor-information/directions.shtml">https://www.cpp.edu/~aboutcpp/visitor-information/directions.shtml</a>	0.00148530596610023206
62	<a href="https://www.cpp.edu/~aboutcpp/why-cpp/diverse.shtml">https://www.cpp.edu/~aboutcpp/why-cpp/diverse.shtml</a>	0.00148530596610023206
63	<a href="https://www.cpp.edu/~aboutcpp/visitor-information/getting-around.shtml">https://www.cpp.edu/~aboutcpp/visitor-information/getting-around.shtml</a>	0.00148530596610023206
64	<a href="https://www.cpp.edu/~cba/student-success-center">https://www.cpp.edu/~cba/student-success-center</a>	0.00146489223613554729
65	<a href="https://www.cpp.edu/~cba/academics">https://www.cpp.edu/~cba/academics</a>	0.00145488777720749799
66	<a href="https://www.cpp.edu/~cba/college-giving">https://www.cpp.edu/~cba/college-giving</a>	0.00145488777720749799
67	<a href="https://www.cpp.edu/~cba/alumni-friends">https://www.cpp.edu/~cba/alumni-friends</a>	0.00145488777720749799
68	<a href="https://www.cpp.edu/~cba/about">https://www.cpp.edu/~cba/about</a>	0.00145488777720749799
69	<a href="https://www.cpp.edu/~cba/our-faculty">https://www.cpp.edu/~cba/our-faculty</a>	0.00145488777720749799
70	<a href="https://www.cpp.edu/~ceis/about">https://www.cpp.edu/~ceis/about</a>	0.00138116806810778124
71	<a href="https://www.cpp.edu/~safety">https://www.cpp.edu/~safety</a>	0.00137333260761528524
72	<a href="https://www.cpp.edu/alpha-index">https://www.cpp.edu/alpha-index</a>	0.00135193078362459999
73	<a href="https://www.cpp.edu/campus-safety-plan.shtml">https://www.cpp.edu/campus-safety-plan.shtml</a>	0.00135193078362459999
74	<a href="https://www.cpp.edu/annual-security-report.shtml">https://www.cpp.edu/annual-security-report.shtml</a>	0.00135193078362459999
75	<a href="https://www.cpp.edu/~jobs">https://www.cpp.edu/~jobs</a>	0.00133655625459262045
76	<a href="https://www.cpp.edu/contact.shtml">https://www.cpp.edu/contact.shtml</a>	0.00130271035419350927
77	<a href="https://www.cpp.edu/~ceis/news-events">https://www.cpp.edu/~ceis/news-events</a>	0.00129682993312513510
78	<a href="https://www.cpp.edu/~ceis/alumni-and-friends">https://www.cpp.edu/~ceis/alumni-and-friends</a>	0.00129682993312513510
79	<a href="https://www.cpp.edu/~ceis/contact-us">https://www.cpp.edu/~ceis/contact-us</a>	0.00129682993312513510
80	<a href="https://www.cpp.edu/~ceis/giving">https://www.cpp.edu/~ceis/giving</a>	0.00129682993312513510
81	<a href="https://www.cpp.edu/~ceis/departments.shtml">https://www.cpp.edu/~ceis/departments.shtml</a>	0.00129682993312513510
82	<a href="https://www.cpp.edu/~alumni/outstanding-alumni/spotlight/andrew-kopp.shtml">https://www.cpp.edu/~alumni/outstanding-alumni/spotlight/andrew-kopp.shtml</a>	0.00129571614489005062
83	<a href="https://www.cpp.edu/events/index.php?cat=Cultural+and+Performing+Arts">https://www.cpp.edu/events/index.php?cat=Cultural+and+Performing+Arts</a>	0.00128579891403741184
84	<a href="https://www.cpp.edu/student-gateway">https://www.cpp.edu/student-gateway</a>	0.00128483358806129452
85	<a href="https://www.cpp.edu/~compass">https://www.cpp.edu/~compass</a>	0.00128306481886802026
86	<a href="https://www.cpp.edu/~alumni/about">https://www.cpp.edu/~alumni/about</a>	0.00128290936058151933
87	<a href="https://www.cpp.edu/~alumni/contact">https://www.cpp.edu/~alumni/contact</a>	0.00128290936058151933
88	<a href="https://www.cpp.edu/~alumni/programs">https://www.cpp.edu/~alumni/programs</a>	0.00128290936058151933
89	<a href="https://www.cpp.edu/~alumni/student-outreach">https://www.cpp.edu/~alumni/student-outreach</a>	0.00128290936058151933
90	<a href="https://www.cpp.edu/~alumni/chapters">https://www.cpp.edu/~alumni/chapters</a>	0.00128290936058151933
91	<a href="https://www.cpp.edu/~aboutcpp/visitor-information?C=D;O=A">https://www.cpp.edu/~aboutcpp/visitor-information?C=D;O=A</a>	0.00125042494102697464
92	<a href="https://www.cpp.edu/~aboutcpp/visitor-information?C=M;O=A">https://www.cpp.edu/~aboutcpp/visitor-information?C=M;O=A</a>	0.00125042494102697464
93	<a href="https://www.cpp.edu/~aboutcpp/visitor-information?C=S;O=A">https://www.cpp.edu/~aboutcpp/visitor-information?C=S;O=A</a>	0.00125042494102697464
94	<a href="https://www.cpp.edu/~engineering">https://www.cpp.edu/~engineering</a>	0.00123921686645503144
95	<a href="https://www.cpp.edu/~financial-aid">https://www.cpp.edu/~financial-aid</a>	0.00121732662330244538
96	<a href="https://www.cpp.edu/~admissions/international?C=D;O=A">https://www.cpp.edu/~admissions/international?C=D;O=A</a>	0.00119772809755955322
97	<a href="https://www.cpp.edu/~admissions/international?C=S;O=A">https://www.cpp.edu/~admissions/international?C=S;O=A</a>	0.00119772809755955322
98	<a href="https://www.cpp.edu/~admissions/international?C=M;O=A">https://www.cpp.edu/~admissions/international?C=M;O=A</a>	0.00119772809755955322
99	<a href="https://www.cpp.edu/~env">https://www.cpp.edu/~env</a>	0.00119742440835321701
100	<a href="https://www.cpp.edu/faculty-staff-gateway">https://www.cpp.edu/faculty-staff-gateway</a>	0.00119742440835321701

## Contributions

Everybody contributed to writing the report, testing the calculations of PageRank by using a different number of pages and fixing bugs. Marisabel wrote the code to crawl and calculated the PageRank.