

# Assignment 4

Information Retrieval and Web Search

Winter 2016

Total points: 100

Issued: 03/14/2016 Due: 03/28/2016

All the code has to be your own (exceptions to this rule are specifically noted below). The code must run on the CAEN environment without additional installation or additional files (except for the data files specified in the assignment).

You can discuss the assignment with others, but the code is to be written individually. You are to abide by the University of Michigan/Engineering honor code; violations will be reported to the Honor Council.

## **[50 points] Web crawler.**

Write a Web crawler that collects the URL of webpages from the umich domain. Your crawler will have to perform the following tasks:

- a. Start with `http://www.eecs.umich.edu`
- b. Perform a Web traversal using a breadth-first strategy. You should only crawl html webpages.
- c. Keep track of the traversed URLs, making sure:
  - they are part of the `eecs.umich` domain
  - they were not already traversed (i.e., avoid duplicates, avoid cycles)
- d. Stop when you reach 2000 URLs.

### Programming guidelines:

Write a program called *crawler.py* that implements the Web crawler.. The program will receive two arguments on the command line: the first one consists of the name of a file containing all the seed URLs (for the purpose of this assignment, this file will only contain one seed, namely `http://www.eecs.umich.edu`); the second one consists of the maximum number of URLs to be crawled (for the purpose of this assignment, the value of this parameter will be 2000)

The *crawler.py* program should be run using a command like this:

```
% python crawler.py myseedURLs.txt 2000
```

It should produce a list of all the URLs being identified by the crawler, in the order in which they are crawled. E.g.:

`http://www.eecs.umich.edu`

`http://eecs.umich.edu/eecs/academics/academics.html`

`http://eecs.umich.edu/eecs/about/why-eecs.html`

Etc.

Save the output of your program in a file called URLs.txt

*[If necessary, you can add extra arguments. E.g., you may choose to also output a file with all the (source\_URL, URL) pairs that you identify in your crawl, which will essentially serve as the directed edges in the graph representation for the next problem]*

## **[50 points] PageRank.**

Implement the PageRank algorithm and apply it to determine the PageRank score for each of the 2000 URLs you crawled. The PageRank implementation should assume:

- An initial value of 0.25 for all the URLs
- A value of 0.85 for  $d$
- Convergence when the difference between the scores obtained with two iterations of PageRank for each of the URLs falls below 0.001.

### Programming guidelines:

Write a program called *pagerank.py* that implements the PageRank algorithm. The program will receive two arguments on the command line: the first one consists of the name of a file containing all the URLs (for the purpose of this assignment, this file will contain all the URLs you identified in the previous problem); the second one consists of the convergence threshold for PageRank (for the purpose of this assignment, the value of this parameter will be 0.001)

The *pagerank.py* program should be run using a command like this:

```
% python pagerank.py URLs.txt 0.001
```

It should produce a list of all the URLs in the *URLs.txt* file, along with their PageRank score. E.g.:

http://www.eecs.umich.edu 0.9

http://eecs.umich.edu/eecs/academics/academics.html 0.6

http://eecs.umich.edu/eecs/about/why-eecs.html 0.2

Etc.

Save the output of your program in a file called PageRankURLs.txt

*[If necessary, you can add extra arguments. E.g., you can pass as input a file with all the (source\_URL, URL) pairs that was (optionally) generated by the crawler.py program.]*

### Write-up guidelines:

Create a file called answers.txt. Include in answers.txt the following information:

1. Amount of time (in seconds) your crawler needed to download the 2000 URLs.
2. Number of iterations your PageRank implementation performed to convergence.

General Canvas submission instructions:

- Include all the files for this assignment in a folder called *[your-username].Assignment4/*  
**Do not** include the content of the pages you crawled.

For instance, lahiri.Assignment4/ will contain crawler.py, pagerank.py, URLs.txt, PageRankURLs.txt, answers.txt.

- Archive the folder using tgz or zip and submit on Canvas by the due date.
- Make sure you include your name and username in each program and in the answers.txt file.
- Make sure all your programs run correctly on the CAEN machines.