# Assignment 3

Information Retrieval and Web Search
Winter 2022
Total points: 100 points
Issued: 02/11/2022 Due: 02/24/2022

All the code has to be your own (exceptions to this rule are specifically noted below). The code must run on the CAEN environment. For the purpose of this assignment, you can use the following libraries: requests, bs4, numpy. If you need other libraries, please check with the instructors.

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.

Note: requests, bs4, numpy are not available by default on CAEN. To use them please install them in a virtual environment. You can use the following sequence of commands:

> *python3 -m venv env/*
> *source env/bin/activate*
> *pip3 install bs4,requests,numpy*

Use a try-except block to prevent network issues that may crash your code.

## [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 https://eecs.engin.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.engin.umich domain (or the old eecs.umich domain). This will include any pages under eecs.umich, eecs.engin.umich, ece.engin.umich, cse.engin.umich

   - 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 contains one seed, namely http://eecs.engin.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 two files:
- A file called *crawler.output* including a list of all the URLs being identified by the crawler, in the order in which they are crawled. E.g.:
  http://eecs.engin.umich.edu
  http://eecs.engin.umich.edu/eecs/academics/academics.html
  http://eecs.engin.umich.edu/eecs/about/why-eecs.html
  Etc.
- A file called *links.output*, including all the links (in the form of *(source_URL, URL))* that you identify in your crawl, which will serve as the directed edges in the graph representation for the next problem. E.g.,
  http://eecs.engin.umich.edu http://cse.engin.umich.edu
  http://eecs.engin.umich.edu http://ece.engin.umich.edu
  http://eecs.engin.umich.edu https://cse.engin.umich.edu/news/
  Etc.


# [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 three arguments on the command line:
- the first one consists of the name of the file containing all the URLs (for the purpose of this assignment, this file will contain all the URLs you identified in the previous problem, i.e., crawler.output);
- the second argument is the name of the file containing all the links in the form of *(source_URL, URL)* pairs (for the purpose of this assignment, this file will contain all the URLs you identified in the previous problem, i.e., links.output);
- the third 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  crawler.output links.output 0.001*

It should produce a file called *pagerank.output* including a list of all the URLs in the *crawler.output* file, along with their PageRank score. The list should be sorted in reverse order of 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.

Notes:
- The number of iterations is not important as long as your implementation converges
- If you identify links not in the domain, you can disregard them

<u>Write-up guidelines:</u>

Create a file called *crawler.pagerank.answers*. Include in this file 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.

<u>General Canvas submission instructions:</u>
- Include all the files for this assignment in a folder called *[your-uniqname].Assignment3/* **Do not** include the content of the pages you crawled*.*
  For instance, mihalcea.Assignment3/ will contain crawler.py, crawler.output, links.output, pagerank.py, pagerank.output, crawler.pagerank.answers
- Archive the folder using tgz or zip and submit on Canvas by the due date.
- Make sure you include your name and uniqname in each program and in the *crawler.pagerank.answers* file.
- Make sure all your programs run correctly on the CAEN machines.