

# 311 PA1 Report

Pseudocode:

crawl() {

My graph webgraph;

Queue is linked list of type UnAnndDepth  
list discovered

int pageCounter; webhitcounter;

while (Q is not empty & < max pages

- get first element in Q

- check for maxDepth

- ~~check~~ check for 50 searches

↳ webcounter == 50

↳ if 50 → wait 3 seconds

- search first element if valid

- ++webhitcounter

- for all links (less than num pages)

- if (not in discovered)

- update Queue

- update discovered

- ++pagecounter

- add edge to webgraph

- remove element from Q

- return webgraph

## makeIndex()

- webhits = 0
  - for (every entry URL in URL map)
    - new hash map wordmap
    - Check for 50 webhits
      - ↳ if so wait 3 seconds
    - J soup on body
    - scanner on body
    - webhits ++
    - while (scanner.hasNext())
      - <sup>scanner.next()</sup> check for punctuation & stop word
        - ↳ then put word in wordmap
    - close scanner
    - ~~put~~ put wordmap into index map hashmap (hashmap of hashmaps)
- 

## search()

- Array list list of type URL and ranking
  - for (every entry in index map)
    - if (word w is in entry)
      - ranking = indegree + #occurrences
      - if (ranking > 0)
        - add to list;
  - sort list
-

## Search cont.

→ ArrayList of type Tagged vertex

- create finalist of url & rankings from list.
- return finalist.

## Search with And()

- ArrayList list of type URL and ranking
- for (every entry in index map)
  - if (both words,  $(w_1, \& w_2)$  are in entry)
    - ranking = indegree \*  $\left( \frac{\# \text{ of } w_1}{\text{occurrences}} + \frac{\# \text{ of } w_2}{\text{occurrences}} \right)$
    - if (ranking > 0)
      - add to list
- sort list
- create finalist ArrayList of type Tagged vertex
- include entry urls & rankings from list finalist
- return finalist

## Search with Or()

- ArrayList of type url and ranking called list
- for (every entry in index map)
  - if (either  $(w_1 \text{ or } w_2)$  is in entry) <sup>at least one</sup>
    - ranking = indegree \*  $\left( \frac{w_1}{\text{occurrence}} + \frac{w_2}{\text{occurrence}} \right)$
    - if (ranking > 0)
      - add to list
- sort list
- ArrayList finalist of type Tagged vertex
- include entry url & corresponding ranking
- return finalist

adjust for 0 value

## Search with $N \log C$

- ArrayList list of type UrlAndRanking
  - for (each entry in indexmap)
    - if ( $w_1$  is in entry &  $w_2$  is not in entry)
      - ranking = in degree + # of  $w_1$  occurrences
      - if (ranking > 0)
        - add to list
  - sort list
  - ArrayList finalList of type ~~Tagged~~ ranked
  - include entry urls & rankings from list in finalList
  - return finalList
- 

## Choice of data structures :

for my Graph :

I decided to use a hash map so I could store the URL's & their corresponding links. The URL's are the keys. So crawler implemented this hash map & added to it. for the indexing, I chose to use 2 hash maps. 1 hash map holds the word map, which is a hash map of words. & the next hash map holds the URL's & their corresponding word hash maps.

## Run times of algorithms

Crawl():  $O(\text{maxPages} * \text{maxdepth})$   
// depend on Jsoup

Make Index():  $O(\text{maxPages} * \text{maxDepth} * \text{size of url})$   
// depends on Jsoup

Search():  $O(\text{num URLs}) = O(\text{num URLs})$

Search with And():  $O(\text{num URLs})$

Search with Or():  $O(\text{num URLs})$

Search with NOT():  $O(\text{num URLs})$

↑  
# of URLs in entry list

add vertex():  $O(1)$

add Edge():  $O(1)$

vertex Data():  $O(1)$

Vertex Data with Incoming Edges():  $O(\text{Vertex Data Size})$

get Neighbors():  $O(\text{URL neighbors})$

↑  
how many neighbors it has

get Incoming():  $O(\text{size of map})$