

## **Programming Assignment 3**

Nathan Friend

COM S 535: Algorithms for Large Data Sets:  
Theory and Practice

2018/11/16

### Data structure used to implement Weighted Q:

Internally, my `WeightedQ` class uses a simple `ArrayList<WeightedItem>`. The `WeightedItem` class is a simple container structure that holds the item (a `String` link, in this project) and the item's weight. This `WeightedItem` class implements the `Comparable` interface, which allows the items to be sorted (by weight) using the `Collections.sort()` method.

When the `add()` method is called, items are simply appended to the end of the internal list, regardless of their weight. When the `extract()` method is called, the `WeightedQ` class first sorts the list and then removes and returns the first item in the list. This allows `add()`s to be very fast, which is important for this use case, since `add()` is called many more times than `extract()`.

**Pseudo code of your crawling algorithm. Please describe how your procedure will ensure that the graph formed will have exactly max many vertices:**

```
// variable declarations/initializations
Create a list of edges that will eventually be written to fileName;
Create a new WeightedQ instance (wq) and add seedUrl;
Create an int to track the number of URLs we've visited (visitedCount);
Create a list of Strings to track which URLs we've visited (visited);

// crawling algorithm
While wq is not empty and visitedCount < max:

    String url = wq.extract();

    If we've already visited this URL, continue;
    If this URL is disallowed by robots.txt, continue;

    visitedCount++;
    visited.add(url);

    Make a GET request to get the URL's HTML;
    Extract all the links from the HTML;
    For each link:
        Add link to wq;
        Add the edge (URL -> link) to the final list of edges;

    If we've made 10 requests, pause for 2 seconds;

Finally, remove any edges from the list of all edges that include URLs
that we didn't crawl;

Write the list of edges to the file;
```

(The real implementation of this pseudo-code can be viewed in *WikiCrawler.java*, in the *crawl()* method.)

This procedure ensures that the final graph will only have **max** many vertices by:

1. ensuring that only **max** many webpages are downloaded (see the **while** loop condition), and
2. removing any edges that include webpages that weren't downloaded.

### Output of the program MyWikiRanker (top 20 page rank, in degree and outdegree pages and Jaccard Similarities):

Seed URL: */wiki/Cello*

Keywords: [*"cello", "bow", "neck", "fingerboard", "pegbox", "scroll", "pegs", "endpin", "f-holes", "f-hole", "pizzicato", "vibrato", "thumb", "string", "strings", "rosin", "orchestra", "symphony", "Stradivarius", "Stradivari", "Yo-Yo", "Rostropovich", "Casals", "Maisky", "Isserlis", "Starker", "violoncello"*]

Output:

A: Top 20 links based on outdegree:

1. */wiki/Bows\_for\_Musical\_Instruments*
2. */wiki/Cello\_Concerto\_in\_E\_major\_(Cassado-Tchaikovsky)*
3. */wiki/Scottish\_National\_Orchestra*
4. */wiki/String\_Quartet,\_Op.\_3\_(Berg)*
5. */wiki/Symphony\_No.\_2\_(Honegger)*
6. */wiki/Ariel\_String\_Quartet*
7. */wiki/The\_Regent\_String\_Quartet*
8. */wiki/Concerto\_for\_Double\_String\_Orchestra\_(Tippett)*
9. */wiki/Cello\_Sonata\_No.\_1\_(Reger)*
10. */wiki/National\_Taiwan\_Normal\_University*
11. */wiki/Kutcher\_String\_Quartet*
12. */wiki/Albert\_Augustine\_Strings*
13. */wiki/Aviv\_String\_Quartet*
14. */wiki/London\_Festival\_Orchestra*
15. */wiki/Benedetto\_Marcello*
16. */wiki/New\_Orford\_String\_Quartet*
17. */wiki/Stanford\_String\_Quartet*
18. */wiki/Tyburn\_String\_Quartet*
19. */wiki/Oslo\_String\_Quartet*
20. */wiki/String\_Quartets\_(Schumann)*

B: Top 20 links based on indegree:

1. */wiki/BWV\_1011*
2. */wiki/String\_Quartet\_No.\_12\_(Schubert)*
3. */wiki/Frankfurt\_Radio\_Symphony\_Orchestra*
4. */wiki/Bows\_for\_Musical\_Instruments*
5. */wiki/Solo\_Cello\_Sonata\_(Ligeti)*

6. /wiki/Cello\_Concerto\_in\_E\_major\_(Cassado-Tchaikovsky)
7. /wiki/Cello\_Symphony
8. /wiki/Sound-hole
9. /wiki/Cello\_suites
10. /wiki/String\_Quartet\_No.\_2\_(Schoenberg)
11. /wiki/String\_Quartet,\_Op.\_3\_(Berg)
12. /wiki/Cello\_Concerto\_in\_D\_minor\_(Cassado)
13. /wiki/String\_basses
14. /wiki/Cello\_Rock
15. /wiki/Autopista\_de\_Pau\_Casals
16. /wiki/String\_quintets
17. /wiki/Baroque\_bow
18. /wiki/Symphony\_No.\_7\_(Williamson)
19. /wiki/Electric\_cello
20. /wiki/Cello\_Concerto\_(Carter)

C: Top 20 links based on page rank with approximation = 0.01 and teleportation = 0.85:

1. /wiki/Cello
2. /wiki/String\_quartet
3. /wiki/Orchestra
4. /wiki/String\_instrument
5. /wiki/Bow\_(music)
6. /wiki/String\_section
7. /wiki/Symphony
8. /wiki/Mstislav\_Rostropovich
9. /wiki/Pizzicato
10. /wiki/Yo-Yo\_Ma
11. /wiki/Fingerboard
12. /wiki/Chicago\_Symphony\_Orchestra
13. /wiki/Philadelphia\_Orchestra
14. /wiki/Boston\_Symphony\_Orchestra
15. /wiki/Pablo\_Casals
16. /wiki/Rosin
17. /wiki/String\_orchestra
18. /wiki/Tuning\_peg
19. /wiki/Orchestral\_enhancement
20. /wiki/Sympathetic\_string

D: Top 20 links based on page rank with approximation = 0.005 and teleporation = 0.85:

1. /wiki/Cello
2. /wiki/String\_quartet
3. /wiki/Orchestra
4. /wiki/String\_instrument
5. /wiki/Bow\_(music)
6. /wiki/String\_section
7. /wiki/Symphony
8. /wiki/Mstislav\_Rostropovich
9. /wiki/Pizzicato

10. /wiki/Yo-Yo\_Ma
11. /wiki/Fingerboard
12. /wiki/Chicago\_Symphony\_Orchestra
13. /wiki/Philadelphia\_Orchestra
14. /wiki/Boston\_Symphony\_Orchestra
15. /wiki/Pablo\_Casals
16. /wiki/Rosin
17. /wiki/String\_orchestra
18. /wiki/Tuning\_peg
19. /wiki/Orchestral\_enhancement
20. /wiki/Sympathetic\_string

E: Top 20 links based on page rank with approximation = 0.001 and teleporation = 0.85:

1. /wiki/Cello
2. /wiki/String\_quartet
3. /wiki/Orchestra
4. /wiki/String\_instrument
5. /wiki/Bow\_(music)
6. /wiki/String\_section
7. /wiki/Symphony
8. /wiki/Mstislav\_Rostropovich
9. /wiki/Pizzicato
10. /wiki/Yo-Yo\_Ma
11. /wiki/Fingerboard
12. /wiki/Chicago\_Symphony\_Orchestra
13. /wiki/Philadelphia\_Orchestra
14. /wiki/Boston\_Symphony\_Orchestra
15. /wiki/Pablo\_Casals
16. /wiki/Rosin
17. /wiki/String\_orchestra
18. /wiki/Tuning\_peg
19. /wiki/Orchestral\_enhancement
20. /wiki/Sympathetic\_string

Exact Jaccard similarity between lists A and B: 0.0811  
Exact Jaccard similarity between lists A and C: 0.0000  
Exact Jaccard similarity between lists A and D: 0.0000  
Exact Jaccard similarity between lists A and E: 0.0000  
Exact Jaccard similarity between lists B and C: 0.0000  
Exact Jaccard similarity between lists B and D: 0.0000  
Exact Jaccard similarity between lists B and E: 0.0000  
Exact Jaccard similarity between lists C and D: 1.0000  
Exact Jaccard similarity between lists C and E: 1.0000  
Exact Jaccard similarity between lists D and E: 1.0000

Number of iterations for your page rank algorithm to converge (within  $\epsilon$ ) on the graph `wikiTennis.txt`, for all three choices of epsilon, when  $\beta = 0.85$  and  $\beta = 0.25$ :

Epsilon ( $\epsilon$ )	Teleportation ( $\beta$ )	Number of iterations
0.01	0.85	25
0.005	0.85	29
0.001	0.85	39
0.01	0.25	5
0.005	0.25	5
0.001	0.25	6

**A note on documentation:** the documentation for this programming assignment can be found in Javadoc format in the `/doc` directory.