

Homework: Web Crawling

1. Objective

In this assignment, you will work with a simple web crawler to measure aspects of a crawl, study the characteristics of the crawl, download web pages from the crawl and gather webpage metadata, all from pre-selected news websites.

2. Preliminaries

To begin we will make use of an existing open source Java web crawler called crawler4j. This crawler is built upon the open source crawler4j library which is located on github. For complete details on downloading and compiling see

<https://github.com/yasserg/crawler4j>

Also see the following document for help installing Eclipse and crawler4j

<http://www-scf.usc.edu/~csci572/2017Spring/hw2/Crawler4jinstallation.pdf>

3. Crawling

Your task is to configure and compile the crawler and then have it crawl a news website. In the interest of distributing the load evenly and not overloading the news servers, we have pre-assigned the news sites to be crawled according to your USC ID number, given in the table below.

The maximum pages to fetch can be set in crawler4j and it should be set to **20,000** to ensure a reasonable execution time for this exercise. Also, maximum depth should be set to **16** to ensure that we limit the crawling. You must crawl only the news site assigned to you.

You should crawl only the news websites assigned to you, and your crawler should be configured so that it does not visit pages outside of the given news website!

USC ID ends with	News Sites to Crawl	Root URL
01~20	LA_Times	http://www.latimes.com/
21~40	CNN	http://www.cnn.com/
41~60	NY_Times	http://www.nytimes.com/
61~80	ABC_News	http://abcnews.go.com/
81~00	NBC_News	http://www.nbcnews.com/

Limit your crawler so it only visits HTML, doc, pdf and different image format URLs and record the meta data for those file types

4. Collecting Statistics

Your first task is to enhance the crawler so it collects information about:

1. *the URLs it attempts to fetch*, a two column spreadsheet, column 1 containing the URL and column 2 containing the HTTP status code received; name the file **fetch_NewsSite.csv** (where the name “NewsSite” is replaced by the news website name in the table above that you are crawling). The number of rows should be no more than 20,000 as that is our pre-set limit.
2. *the files it successfully downloads*, a four column spreadsheet, column 1 containing the URLs successfully downloaded, column 2 containing the size of the downloaded file, column 3 containing the # of outlinks found, and column 4 containing the resulting content-type; name the file **visit_NewsSite.csv**; clearly the number of rows will be less than the number of rows in **fetch_NewsSite.csv**
3. *all of the URLs (including repeats) that were discovered and processed in some way*; a two column spreadsheet where column 1 contains the encountered URL and column two an indicator of whether the URL **a.** resides in the website (**OK**), or **b.** points outside of the website (**N_OK**). (A file points out of the website if its URL does not start with the initial domain name, e.g. when crawling ABC news all inside URLs must start with www.abcnews.go.com.) Name the file **urls_NewsSite.csv**. This file will be much larger than **fetch_*.csv** and **visit_*.csv**.

Note1: you should modify the crawler so it outputs the above data into three separate csv files; you may use them for processing later;

Note2: all uses of NewsSite should be replaced by the name given in the table on page 1.

Based on the information recorded by the crawler in its output files, you are to collate the following statistics for a crawl of your designated news website:

- Fetch statistics:
 - # fetches attempted:
The total number of URLs that the crawler attempted to fetch. This is usually equal to the `MAXPAGES` setting if the crawler reached that limit; less if the website is smaller than that.
 - # fetches succeeded:
The number of URLs that were successfully downloaded in their entirety, i.e. returning a HTTP status code of 2XX.
 - # fetches failed or aborted:
The number of fetches that failed for whatever reason, including, but not limited to: HTTP redirections (3XX), client errors (4XX), server errors (5XX) and other network-related errors.¹
- Outgoing URLs: statistics about URLs extracted from visited HTML pages
 - Total URLs extracted:
The grand total number of URLs extracted (including repeats) from all visited pages

¹ Based purely on the success/failure of the fetching process. Do not include errors caused by difficulty in parsing content *after* it has already been successfully downloaded.

- o # unique URLs extracted:
The number of *unique* URLs encountered by the crawler
- o # unique URLs within a news website:
The number of *unique* URLs encountered that are associated with the news website, i.e. the URL begins with the given root URL of the news website.
- o # unique URLs outside the news website:
The number of *unique* URLs encountered that were *not* from the news website.
- Status codes: number of times various HTTP status codes were encountered during crawling, including (but not limited to): 200, 301, 401, 402, 404, etc.
- File sizes: statistics about file sizes of visited URLs – the number of files in each size range (See Appendix A).
 - o 1KB = 1024B; 1MB = 1024KB
- Content Type: a list of the different content-types encountered

These statistics should be collated and submitted as a plain text file whose name is CrawlReport_domain.txt, following the format given in Appendix A at the end of this document.

Make sure you understand the crawler code and outputs before you commence collating these statistics. All the information that you are required to collect can be derived by processing the crawler output.

5. FAQ

Q: For the purposes of counting unique URLs, how to handle URLs that differ only in the query string? For example: `http://www.cnn.com/page?q=0` and `http://www.cnn.com/page?q=1`

A: These can be treated as different URLs.

Q: URL case sensitivity: are these the same, or different URLs?
`http://www.cnn.com/foo` and `http://www.cnn.com/FOO`

A: The path component of a URL is considered to be case-sensitive, so the crawler behavior is correct according to RFC3986. Therefore, these are different URLs.

The page served may be the same because:

- that particular web server implementation treats path as case-insensitive (some server implementations do this, especially windows-based implementations)
- the web server implementation treats path as case-sensitive, but aliasing or redirect is being used.

This is one of the reasons why deduplication is necessary in practice.

Q: Attempting to compile the crawler results in syntax errors.

A: Make sure that you have included crawler4j as well as all its dependencies.

Also check your Java version; the code includes more recent Java constructs such as the typed collection `List<String>` which requires at least Java 1.5.0.

Q: I get the following warnings when trying to run the crawler:

```
log4j: WARN No appenders could be found for logger
log4j: WARN Please initialize the log4j system properly.
```

A: You failed to include the `log4j.properties` file that comes with crawler4j.

Q: On Windows, I am encountering the error: `Exception_Access_Violation`

A: This is a Java issue. See: http://java.com/en/download/help/exception_access.xml

Q: I am encountering multiple instances of this info message:

```
INFO [Crawler 1] I/O exception (org.apache.http.NoHttpResponseException)
caught when processing request: The target server failed to respond
INFO [Crawler 1] Retrying request
```

A: If you're working off an unsteady wireless link, you may be battling network issues such as packet losses – try to use a better connection. If not, the web server may be struggling to keep up with the frequency of your requests.

As indicated by the info message, the crawler will retry the fetch, so a few isolated occurrences of this message are not an issue. However, if the problem repeats persistently, the situation is not likely to improve if you continue hammering the server at the same frequency. Try giving the server more room to breathe:

```
/*
 * Be polite: Make sure that we don't send more than
 * 1 request per second (1000 milliseconds between requests).
 */
config.setPolitenessDelay(2500);
/*
 * READ ROBOTS.TXT of the website - Crawl-Delay: 10
 * Multiply that value by 1000 for millisecond value
 */
```

Q: The crawler seems to choke on some of the downloaded files, for example:

```
java.lang.StringIndexOutOfBoundsException: String index out of range: -2
```

```
java.lang.NullPointerException: charsetName
```

A: Safely ignore those. We are using a fairly simple, rudimentary crawler and it is not necessarily robust enough to handle all the possible quirks of heavy-duty crawling and parsing. These problems are few in number (compared to the entire crawl size), and for this exercise we're okay with it as long as it skips the few problem cases and keeps crawling everything else, and terminates properly – as opposed to exiting with fatal errors.

Q: While running the crawler, you may get the following error:

SLF4J: Failed to load class "org.slf4j.impl.StaticLoggerBinder".

SLF4J: Defaulting to no-operation (NOP) logger implementation

SLF4J: See <http://www.slf4j.org/codes.html#StaticLoggerBinder> for further details.

A: The problem is described here -

<https://groups.google.com/forum/#!topic/crawler4j/urroQ5BRsWM>. A simple fix is to add more JARs that can be found contained as a .zip/.tar.gz in this link -

<http://logback.qos.ch/download.html>. Adding all these external JARs to the project in the same way as the crawler-4j JAR will make the crawler display logs now.

Q: What should we do with URL if it contains comma ?

A: Replace the comma with "-" or "_", so that it doesn't throw an error.

Q: What's the difference between aborted fetches and failed fetches?

A: failed: Can be due to HTTP errors and other network related errors

aborted: Client decided to stop the fetching. (ex: Taking too much time to fetch)

Q: For some reason my crawler attempts 19,999 fetches, even though max pages is set to 20,000, does this matter?

A: It can be possible because 20,000 is the limit that you will try to fetch (it may contain successful status code like 200 and other like 301). But the visit.csv will contain only the URL's for which you are able to successfully download the files.

Q: How to differentiate fetched pages and downloaded pages?

A: In this assignment we do not ask you save any of the files to the disk. Visiting a page means crawler4j processing a page (it will parse the page and extract relevant information like outgoing URLs). That means all visited pages are downloaded.

6. Submission Instructions

- Save your statistics report as a plain text file and name it based on the domain names assigned below:

USC ID ends with	Site
01~20	CrawlReport_LATimes.txt
21~40	CrawlReport_CNN.txt
41~60	CrawlReport_NYTimes.txt
61~80	CrawlReport_ABCNews.txt
81~00	CrawlReport_NBCNews.txt

- Also include the output files generated from your crawler run, using the extensions as shown above:
 - fetch_NewSite.csv
 - visit_NewSite.csv
- Do NOT include the output files
 - urls_NewSite.csvwhere _NewSite should be replaced by the name from the table above.
- Do **not** submit Java code or compiled programs; it is not required.
- Compress all of the above into a single zip archive and name it:

crawl.zip

Use only standard zip format. Do **NOT** use other formats such as zipx, rar, ace, etc. For example the zip file might contain the following three files:

1. CrawlReport_LATimes.txt, (the statistics file)
 2. fetch_LATimes.csv
 3. visit_LATimes.csv
- To submit your file electronically to the csci572 account enter the following command from your UNIX prompt:

```
$ submit -user csci572 -tag hw2 crawl.zip
```

Appendix A

Use the following format to tabulate the statistics that you collated based on the crawler outputs.

Note: The status codes and content types shown are only a sample. The status codes and content types that you encounter may vary, and should all be listed and reflected in your report. Do **NOT** lump everything else that is not in this sample under an “Other” heading. You may, however, exclude status codes and types for which you have a count of zero.

CrawlReport_NewsSite.txt

Name: Tommy Trojan
USC ID: 1234567890
News site crawled: latimes.com

Fetch Statistics

=====

fetches attempted:
fetches succeeded:
fetches aborted:
fetches failed:

Outgoing URLs:

=====

Total URLs extracted:
unique URLs extracted:
unique URLs within News Site:
unique URLs outside News Site:

Status Codes:

=====

200 OK:
301 Moved Permanently:
401 Unauthorized:
403 Forbidden:
404 Not Found:

File Sizes:

=====

< 1KB:
1KB ~ <10KB:
10KB ~ <100KB:
100KB ~ <1MB:
>= 1MB:

Content Types:

=====

text/html:
image/gif:
image/jpeg:
image/png:
application/pdf: