

ReadME

Description of used software

We use the crawler we used in project 1 which is crawler4j. It is an open source web crawler for Java which provides a simple interface for crawling the web pages.

Installation

NOTE:

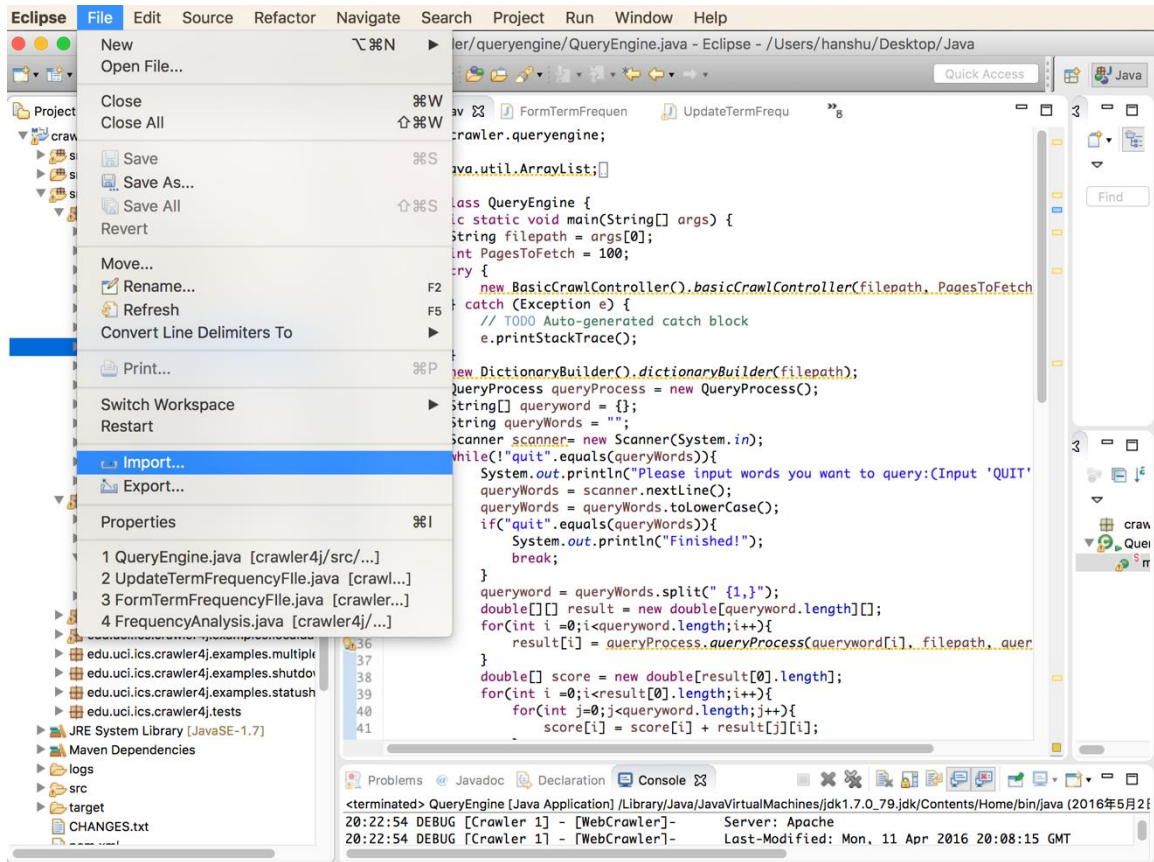
Using Maven

To use the latest release of crawler4j, please add the following snippet in the “pom.xml” file which is contained in “crawler4j-crawler4j-4.1” project

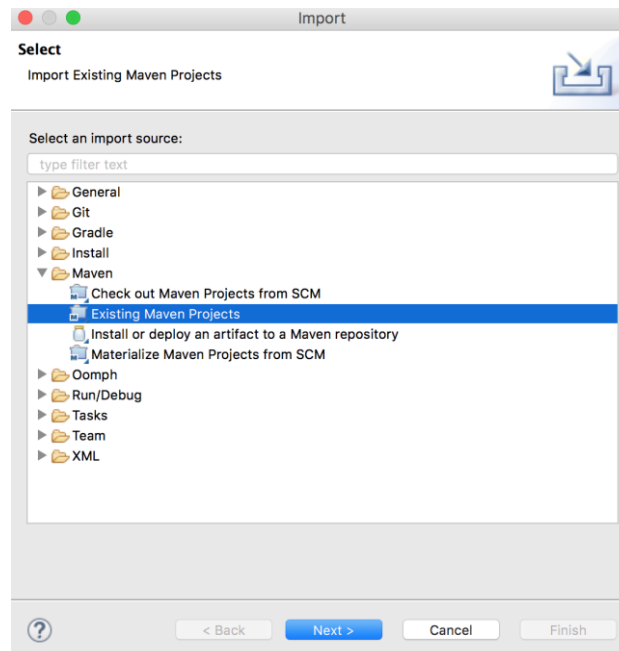
```
<dependency>
    <groupId>edu.uci.ics</groupId>
    <artifactId>crawler4j</artifactId>
    <version>4.2</version>
</dependency>
```

Step:

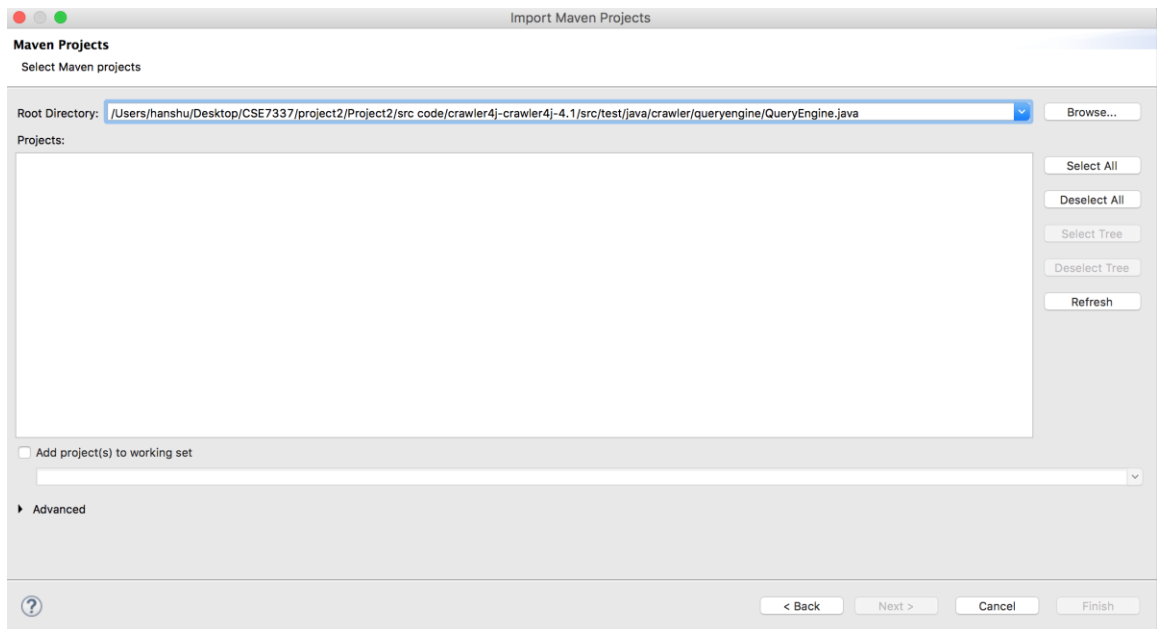
1. open Java IDE(Eclipse)--->click on “file” button--->select “Import” tab as show below:



2. select the “Existing Maven Projects” option under the Maven folder in the Import interface as show below:



3. select the “crawler4j-crawler4j-4.1” file as your “Existing Maven Projects”--->click on “next” which equals to “OK” as show below:



Then, set the right path.

4. as the notation content before, add the snippet

```
<dependency>

    <groupId>edu.uci.ics</groupId>

    <artifactId>crawler4j</artifactId>

    <version>4.2</version>

</dependency>
```

into the pom.xml file, then the import process started.

As we have added the dependency snippet and dependency package has already existed, therefore we cannot take a screen shot of the process, but the screen shot of the dependency files are listed below:

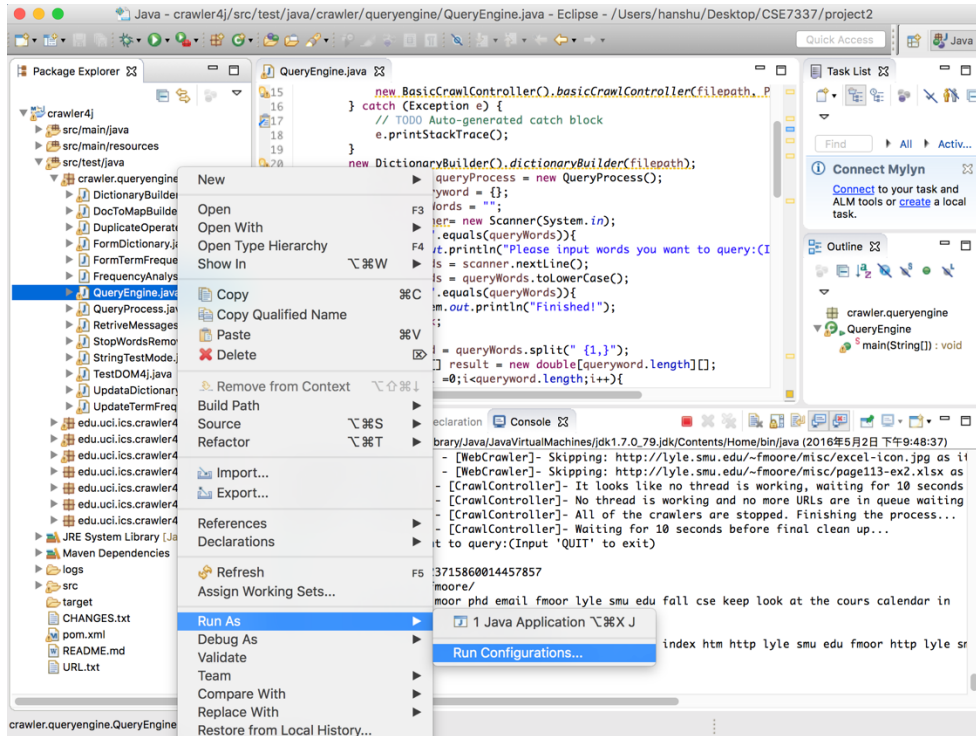
▼ Maven Dependencies

- ▶ slf4j-api-1.7.7.jar - /Users/hanshu/.m2/repository/org/slf4j/slf4j-api/1.7.7
- ▶ lidalia-slf4j-ext-1.0.0.jar - /Users/hanshu/.m2/repository/uk/org/lidalia/lidalia-slf4j-ext/1.0.0
- ▶ guava-14.0.1.jar - /Users/hanshu/.m2/repository/com/google/guava/guava/14.0.1
- ▶ logback-classic-1.1.2.jar - /Users/hanshu/.m2/repository/ch/qos/logback/logback-classic/1.1.2
- ▶ logback-core-1.1.2.jar - /Users/hanshu/.m2/repository/ch/qos/logback/logback-core/1.1.2
- ▶ httpclient-4.3.5.jar - /Users/hanshu/.m2/repository/org/apache/httpcomponents/httpclient/4.3.5
- ▶ httpcore-4.3.2.jar - /Users/hanshu/.m2/repository/org/apache/httpcomponents/httpcore/4.3.2
- ▶ commons-logging-1.1.3.jar - /Users/hanshu/.m2/repository/commons-logging/commons-logging/1.1.3
- ▶ commons-codec-1.6.jar - /Users/hanshu/.m2/repository/commons-codec/commons-codec/1.6
- ▶ je-5.0.73.jar - /Users/hanshu/.m2/repository/com/sleepycat/je/5.0.73
- ▶ tika-parsers-1.5.jar - /Users/hanshu/.m2/repository/org/apache/tika/tika-parsers/1.5
- ▶ tika-core-1.5.jar - /Users/hanshu/.m2/repository/org/apache/tika/tika-core/1.5
- ▶ vorbis-java-tika-0.1.jar - /Users/hanshu/.m2/repository/org/gagravarr/vorbis-java-tika/0.1
- ▶ vorbis-java-core-0.1-tests.jar - /Users/hanshu/.m2/repository/org/gagravarr/vorbis-java-core/0.1
- ▶ netcdf-4.2-min.jar - /Users/hanshu/.m2/repository/edu/ucar/netcdf/4.2-min
- ▶ apache-mime4j-core-0.7.2.jar - /Users/hanshu/.m2/repository/org/apache/james/apache-mime4j-core/0.7.2
- ▶ apache-mime4j-dom-0.7.2.jar - /Users/hanshu/.m2/repository/org/apache/james/apache-mime4j-dom/0.7.2
- ▶ commons-compress-1.5.jar - /Users/hanshu/.m2/repository/org/apache/commons/commons-compress/1.5
- ▶ xz-1.2.jar - /Users/hanshu/.m2/repository/org/tukaani/xz/1.2
- ▶ pdfbox-1.8.4.jar - /Users/hanshu/.m2/repository/org/apache/pdfbox/pdfbox/1.8.4
- ▶ fontbox-1.8.4.jar - /Users/hanshu/.m2/repository/org/apache/pdfbox/fontbox/1.8.4
- ▶ jempbox-1.8.4.jar - /Users/hanshu/.m2/repository/org/apache/pdfbox/jempbox/1.8.4
- ▶ bcmail-jdk15-1.45.jar - /Users/hanshu/.m2/repository/org/bouncycastle/bcmail-jdk15/1.45
- ▶ bcprov-jdk15-1.45.jar - /Users/hanshu/.m2/repository/org/bouncycastle/bcprov-jdk15/1.45
- ▶ poi-3.10-beta2.jar - /Users/hanshu/.m2/repository/org/apache/poi/poi/3.10-beta2
- ▶ poi-scratchpad-3.10-beta2.jar - /Users/hanshu/.m2/repository/org/apache/poi/poi-scratchpad/3.10-beta2
- ▶ poi-ooxml-3.10-beta2.jar - /Users/hanshu/.m2/repository/org/apache/poi/poi-ooxml/3.10-beta2
- ▶ poi-ooxml-schemas-3.10-beta2.jar - /Users/hanshu/.m2/repository/org/apache/poi/poi-ooxml-schemas/3.10-beta2
- ▶ xmlbeans-2.3.0.jar - /Users/hanshu/.m2/repository/org/apache/xmlbeans/xmlbeans/2.3.0
- ▶ dom4j-1.6.1.jar - /Users/hanshu/.m2/repository/dom4j/dom4j/1.6.1
- ▶ geronimo-stax-api-1.0_spec-1.0.1.jar - /Users/hanshu/.m2/repository/org/apache/geronimo/specs/geronimo-stax-api-1.0_spec/1.0.1
- ▶ tagsoup-1.2.1.jar - /Users/hanshu/.m2/repository/org/ccil/cowan/tagsoup/tagsoup/1.2.1
- ▶ asm-debug-all-4.1.jar - /Users/hanshu/.m2/repository/org/ow2/asm/asm-debug-all/4.1
- ▶ isoparser-1.0-RC-1.jar - /Users/hanshu/.m2/repository/com/googlecode/mp4parser/isoparser/1.0-RC-1
- ▶ aspectjrt-1.6.11.jar - /Users/hanshu/.m2/repository/org/aspectj/aspectjrt/1.6.11
- ▶ metadata-extractor-2.6.2.jar - /Users/hanshu/.m2/repository/com/drewnoakes/metadata-extractor/2.6.2
- ▶ xmpcore-5.1.2.jar - /Users/hanshu/.m2/repository/com/adobe/xmp/xmpcore/5.1.2
- ▶ xercesImpl-2.8.1.jar - /Users/hanshu/.m2/repository/xerces/xercesImpl/2.8.1
- ▶ xml-apis-1.3.03.jar - /Users/hanshu/.m2/repository/xml-apis/xml-apis/1.3.03

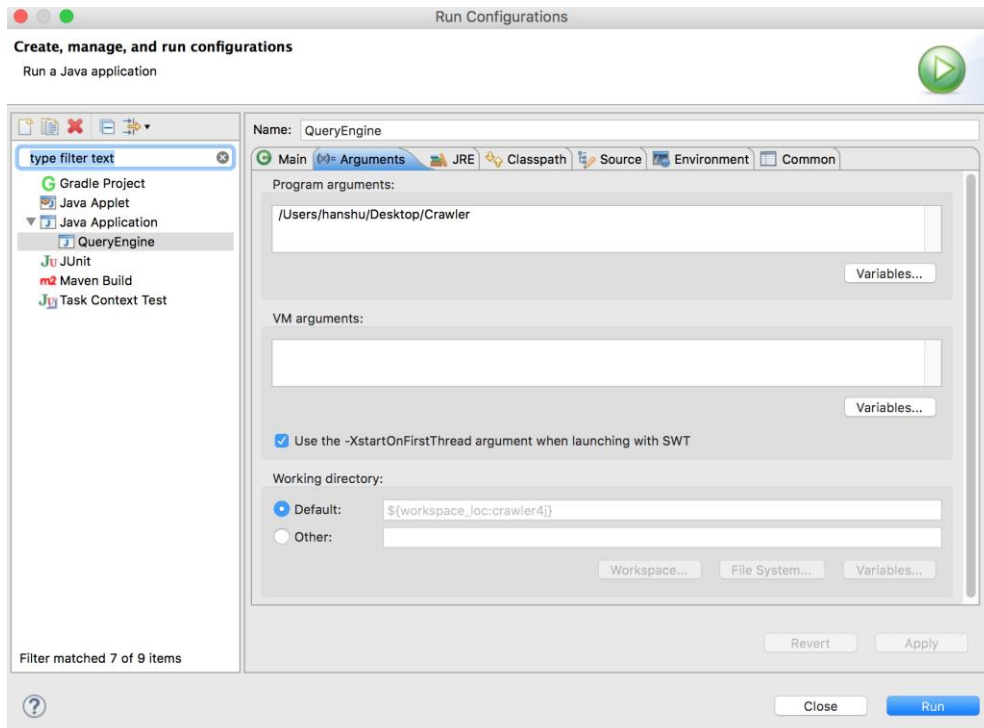
5. Completed!

Execution instructions

1. create a file to store the crawler result and create a stopwords.txt file under the file newly created to store the crawler result--->add the stop words as provided--->you have to run the java file “QueryEngine.java” as configuration



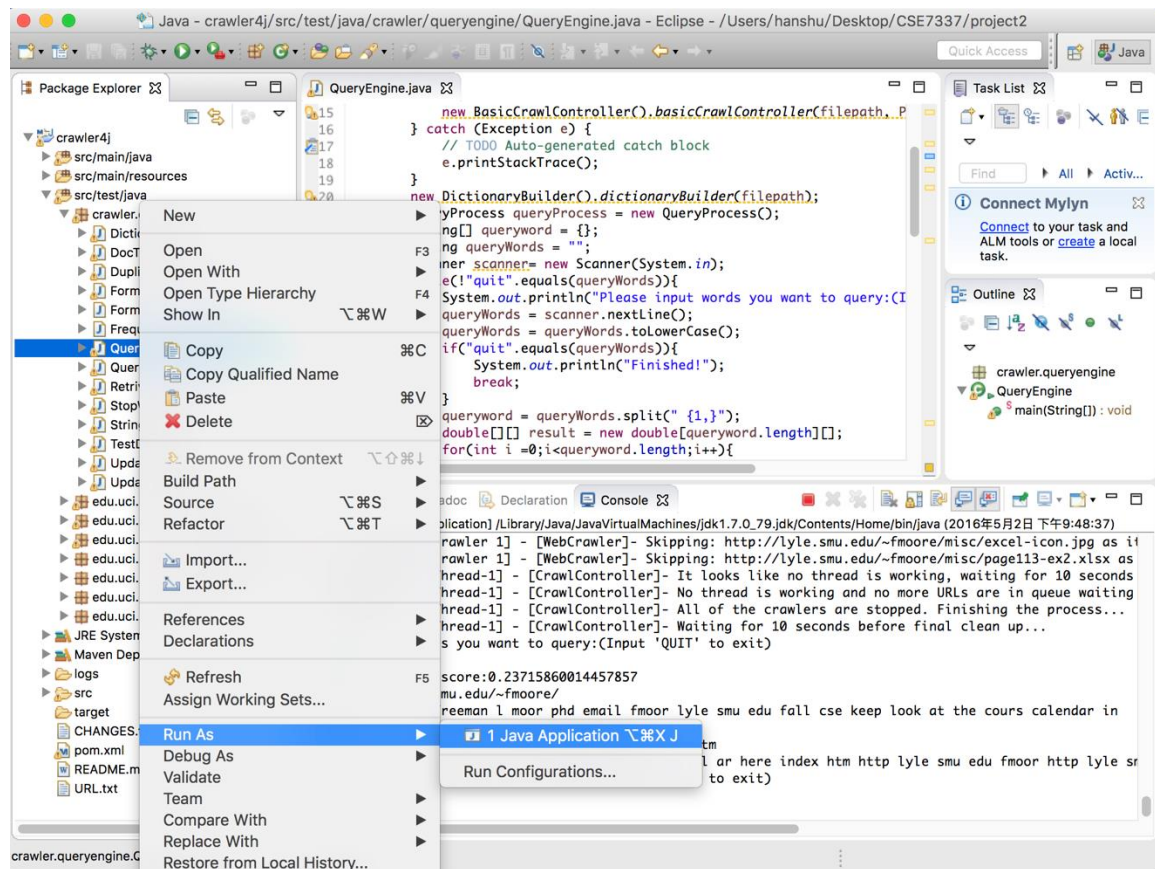
2. type in the Arguments---> press “run” button



set the right path to the folder where contains a stopwords.txt file, and it will be

used to save other result files.

3. run the QueryEngine.java as Java Application



4. type in the query words in the reminded district and the engine will start work

```
21:49:59 INFO [Thread-1] - [CrawlController]- It looks like no thread is working, waiting for 10 seconds to make sure...
21:50:09 INFO [Thread-1] - [CrawlController]- No thread is working and no more URLs are in queue waiting for another 10 seconds to make sure...
21:50:19 INFO [Thread-1] - [CrawlController]- All of the crawlers are stopped. Finishing the process...
21:50:19 INFO [Thread-1] - [CrawlController]- Waiting for 10 seconds before final clean up...
Please input words you want to query:(Input 'QUIT' to exit)
moore smu
docID:1----rank score:0.23715860014457857
URL:http://lyle.smu.edu/~fmoore/
Content: spring freeman l moor phd email fmoor lyle smu edu fall cse keep look at the cours calendar in
docID:27-----rank score:0.15811388300841894
URL:http://lyle.smu.edu/~fmoore/misc/urlexample1.htm
Content: exampl link februari how mani distinct url ar here index htm http lyle smu edu fmoor http lyle smu
Please input words you want to query:(Input 'QUIT' to exit)
```