Crawling Assignment Introduction to Information Retrieval Donald J. Patterson

Content adapted from Hinrich Schütze http://www.informationretrieval.org

Robust Crawling

A Robust Crawl Architecture



Setting up Eclipse

Create a new java project in Eclipse

- Create a new java project in Eclipse
- In your src folder, put your code from the Text Processing assignment

- Create a new java project in Eclipse
- In your src folder, put your code from the Text Processing assignment
- Download crawler4j's code (yasserg on github)
 - Maven
 - crawler4j and crawler4j's dependencies
 - crawler4j with dependences

Setting up Eclipse

• Create a "lib" folder as a sibling to your "src" folder

- Create a "lib" folder as a sibling to your "src" folder
- Put all the jars into the "lib" folder

- Create a "lib" folder as a sibling to your "src" folder
- Put all the jars into the "lib" folder
- Refresh Eclipse ("F5"), so that Eclipse becomes aware of the new files on your hard disk

- Right-click on your project name and select "Properties"
- Choose "Java Build Path"
- Select the "Libraries" Tab
- Select "Add jars" and add all the jars that you downloaded to your project. This allows your code to use the methods and objects that are defined in those jar files

type filter text Jav • Resource Builders Coverage FindBugs Java Build Path • Java Code Style • Java Code Style • Java Compiler • Java Editor Javadoc Location Project References Refactoring History Run/Debug Settings • Task Repository • Validation WikiText	Properties for newcrawler A Build Path Source Projects Libraries Order and Export JARs and class folders on the build path: Magache-mime4j-core-0.7.jar - newcrawler/lib Madd JARs	t tell
WIKITEXL	Add External JARs Add Variable Add Uariable Add Library Add Class Folder Add Class Folder Add External Class Folder	ve >n nit



- First you must have an account on wcpkneel
 - Prof. Patterson has to set that up for you

- Export as a runnable jar
 - Run your program on your local machine to create a run profile
 - Export your project as a "Runnable Jar"
 - A runnable jar packages all the required libraries into one big .jar file so you don't have to deal with moving them separately. That's convenient, but it probably violates licensing agreements. Don't do that with code you plan on distributing.
 - The launch configuration is asking what "main" you want the runnable jar to start from. If you tested your program in Eclipse, you should have an option in the dropdown



000	Runnable JAR File Export	
Runnable JAR File Specification 📃		
Select a 'Java Application' launch configuration to use to create a runnable JAR.		
Launch configuration	:	
Controller (1) - newcrawler \$		
Export destination:		
/Users/djp3/Downloads/crawler4openlab.jar Browse		
Library handling:	ibraries into generated IAR	
Package required libraries into generated IAR		
Copy required libraries into a sub-folder next to the generated JAR		
Save as ANT script		
ANT script location:	/Users/djp3/Documents/ClassResources/2014_01_INF141 Browse	
?	< Back Next > Cancel Finish	

Running on wcpkneel

 Make sure you update any hardcoded file path names to reflect the arrangement on wcpkneel

- Now you need to open a terminal window on wcpkneel
- I use "ssh" to login to the remote machine. If you are on Windows you will need to use PuTTY

- Now you have to move the jar you just made to wcpkneel
- On a Mac you can do that with "sftp" from a Terminal. On windows use WinSCP
- Now you need to open a terminal window on wcpkneel
- I use "ssh" to login to the remote machine. If you are on Windows you will need to use PuTTY

- Once there use "Is" to see if your jar made it okay
- Then make sure you are using the right version of java.
 - java -version

- Then use the command
 - nice -n 19 java -Xmx1024M -jar <myJar.jar> <args>
 - to start your crawler.
 - nice is a command that reduces resource usage
 - The -Xmx command tells java how much memory to use while crawling. You shouldn't need too much because the frontier is kept on disk.

- Once you are happy everything is going okay, detach from the screen with Ctrl-A then "d"
- This will send you back to the real terminal
- If you logout with the "exit" command then your crawler is still running away. You need to make sure you are monitoring it.

- To go back to it, log in to wcpkneel
- Type "screen -r" and you will reattach to the virtual screen that your crawler is running in. If you hit Ctrl-C it will kill your process, or you can just look at any output that you are generating. Hit Ctrl-A then "d" to detach again, as often as you like



What is it?

 crawler4j is an open source web crawler for Java which provides a simple interface for crawling the Web. Using it, you can setup a multi-threaded web crawler in few minutes.

How do you use it

 You need to create a crawler class that extends WebCrawler.
 This class decides which URLs should be crawled and handles the downloaded page.

public class TestCrawler extends WebCrawler

WebCrawler Class

- This method receives two parameters.
 - The page in which we have discovered this new url
 - The second parameter is the new url.
- You should implement this function to specify whether the given url should be crawled or not
 - Based on the destination URL
 - Based on the file type
 - Based on the referring page
 - Based on whatever logic you want to use.

```
@Override
public boolean shouldVisit(Page referringPage, WebURL url) {
    return true;
}
```

WebCrawler Class

• This function is called when a page is fetched and ready to be processed by your program.

```
@Override
public void visit(Page page) {
}
```

WebCrawler Class

- Useful methods in Page
 - getWebURL().getURL();
 - getParseData();
 - returns a type based on the data in the page
 - HtmlParseData
 - getText();
 - getHtml();
 - getOutgoingUrls();

Controller Class

Contains the main that you want to run

CrawlConfig config = new CrawlConfig();

Controller Class

 Specify where to keep the queues, the frontier, the robots files, etc.

config.setCrawlStorageFolder(crawlStorageFolder);

• "/share/bigspace/CS150-1-F15/<username>"

Controller Class

Set up the components of the crawl

PageFetcher pageFetcher = new PageFetcher(config); RobotstxtConfig robotstxtConfig = new RobotstxtConfig(); RobotstxtServer robotstxtServer = new RobotstxtServer(robotstxtConfig, pageFetcher); CrawlController controller = new CrawlController(config, pageFetcher, robotstxtServer);

Controller Class

- Add a seed URL
 - Repeat to add more

controller.addSeed("http://www.westmont.edu/");

Controller Class

- Start the crawl
 - This is a blocking operation

controller.start(MyCrawler.class, numberOfCrawlers);



Configuration Options

- Set a maximum depth of crawling
 - 0 crawls only the seed pages

config.setMaxDepthOfCrawling(maxDepthOfCrawling);

Configuration Options

• Set a maximum depth of crawling

config.setMaxPagesToFetch(maxPagesToFetch);

Configuration Options

• If you suspect that you might crash or have a long running crawl, you can inform crawler4j to maintain state and restart

config.setResumableCrawling(true);

Configuration Options

• User-agent string

config.setUserAgentString(userAgentString);

WESTMONT COMPUTER SCIENCE

