

Andrew Risse

USC ID: 5987029524

HW 5 Report

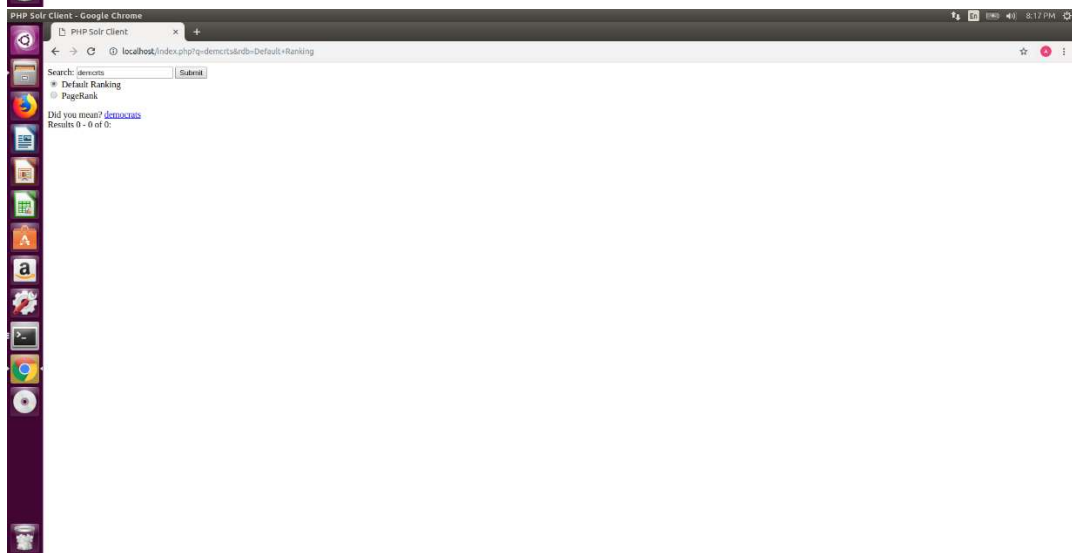
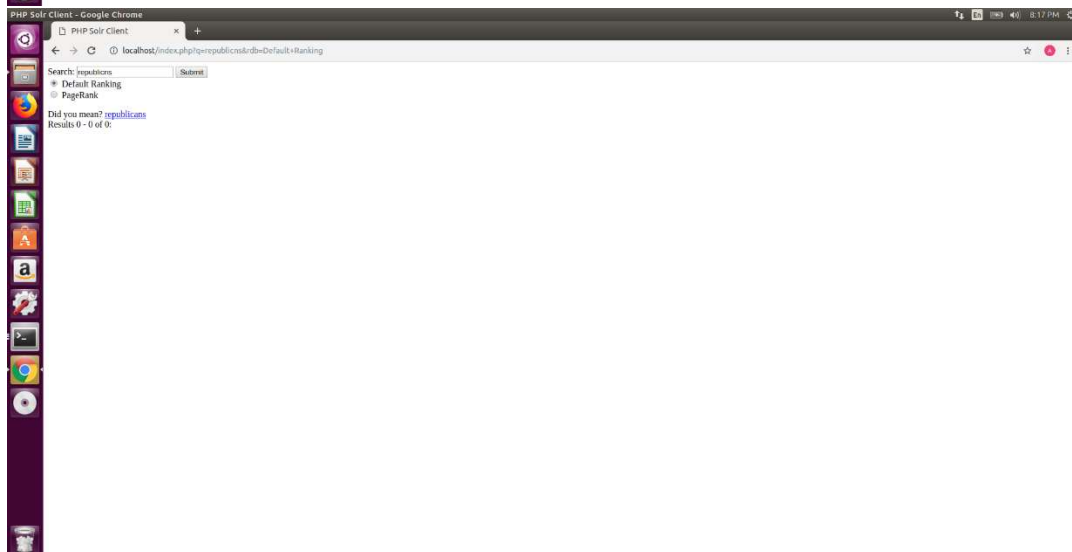
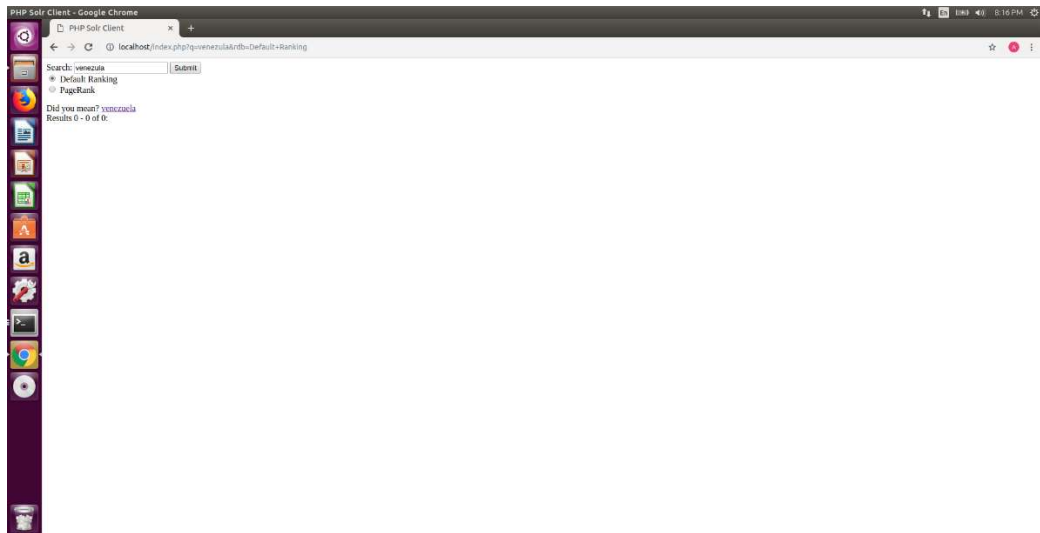
These are the steps I used to complete this project:

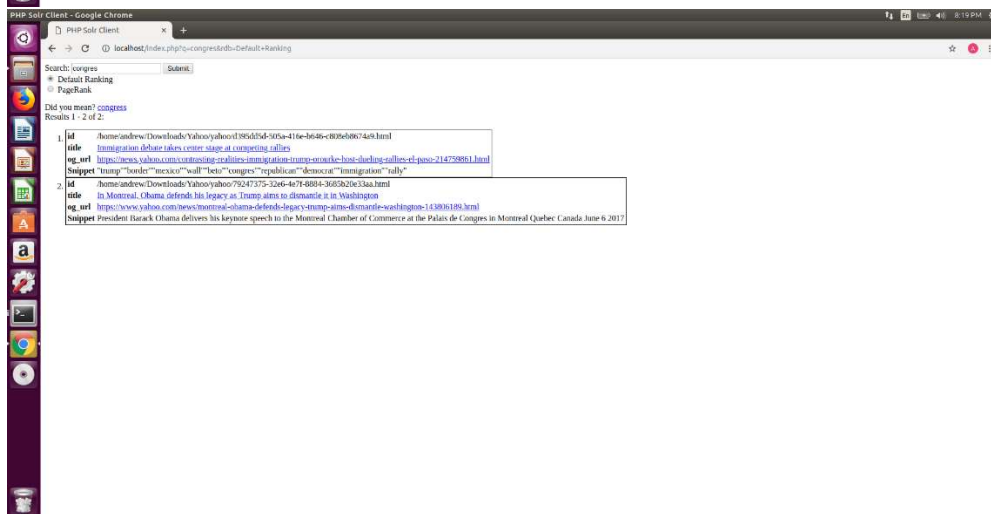
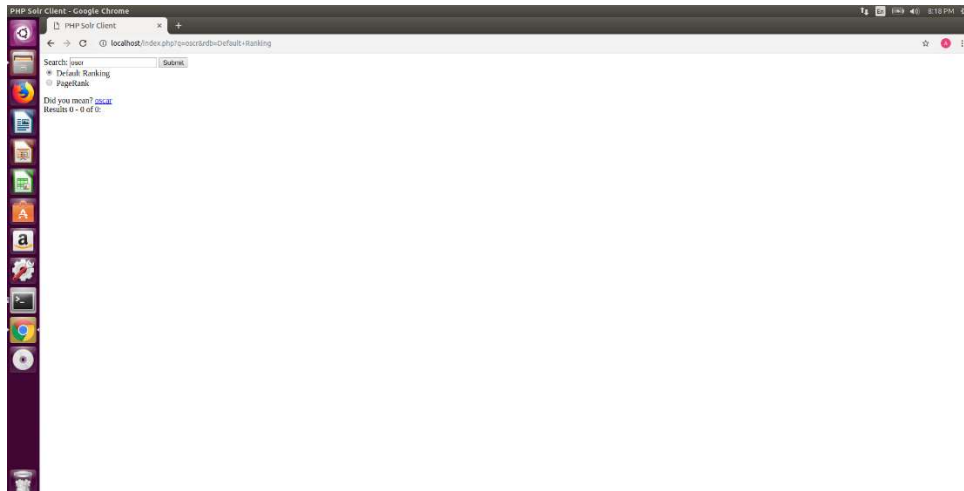
1. Write a script (tikascript.sh) that uses the command:

```
java -jar tika-app-1.20.jar --text-main "$file" >> big.txt
```

This script executed the command for all the "Yahoo" html files and put it's results in "big.txt"
It took several hours for this to complete.
2. I then followed the instructions for adding SOLR autocomplete. I wrote some PHP code that uses JQuery.ajax() to submit a /suggest query to my core and respond with auto-suggestions in JSON format. This format had to be parsed so only the terms were extracted and displayed in a drop down. I used this website as a reference: <https://examples.javacodegeeks.com/enterprise-java/apache-solr/solr-autocomplete-example/>
My version does not handle more than two words for the query when provided auto-complete, I ran out of time to continue messing with the solrconfig.xml.
3. I then tried to use my big.txt file with the PHP version of Norvig SpellCorrector. I was able to get it to spell correct with the Norvig big.txt file, but it would not work with mine. I tried everything recommended on the forums (permissions, looking for a serialized_dictionary.txt file, adding code to handle memory issues, etc) and nothing would work. I elected to add the query terms from HW4 to the beginning of the Norvig.txt file and move on. I had spent way too many hours troubleshooting. I decided to implement the spell correction in my HTML after the user submits a query so that any misspellings appear as a clickable link after the words "Did you mean?" on the results page.
4. The final step was to write PHP code that generates snippets (snippet.php). I tried using Tika to parse the HTML for text, but this caused snippet generation to be very slow even though the results were better. I ended up using file_get_contents and stripping the tags. I had to apply preg_replace several times with different regular expressions to clean up the output. There are still some issues with the snippets that are produced (leftover HTML stuff), but I ran out of time to perfect the regular expressions.

Screenshots of spell correct:





Screenshots of auto-complete:

