

Homework - 5

Enhancing Your Search Engine

Name: Anusha Mysore Swamy

USC ID: 5562476545

Steps followed for:

1. Autocomplete

Used Solr/Lucene's FuzzyLookupFactory to implement this feature, which automatically creates suggestions for misspelt words in fields. For this feature to work, I have added a search component and a request handler to solarconfig.xml, as suggested in the homework pdf "AutoCompleteInSolr". JQuery is used in the front end. An ajax request is made with the URL to Solr, whose response is the list of keywords matching the string input, and is later displayed to the user using the HTML5 dataList component.

2. Spell Correction

- Created a java program to parse all the HTML files for Fox news website, to create big.txt. Used Apache Tika library to extract the HTML content from each of the webpages
- Downloaded Norvig's SpellCorrector.php file from the link <http://www.phpclasses.org/package/4859-PHP-Suggest-corrected-spelling-text-in-pure-PHP.html#download>
- The input to Novig program's correct function is the big.txt. When the user types a word in the input, the program reads this from the php and sends the word to the SpellCorrector program. This program tries to check if the word exists in the dictionary(i.e. big.txt). If not, then it tries to find the nearest possible match and returns this as the corrected suggestion to the user and displays the results for the corrected word
- A link is also provided with the user's actual input to search for (if the user intends to continue with this input itself) and on clicking it, the obtained results are returned
- If the word does not exist in the webpages then the program returns an appropriate message to the user and asks him to try again with a different spelling

3. Snippets

The program tries to find the query term(s) in the html page contents and highlights them. If there are multiple terms, then it tries to find the first occurrence where all the terms occur together. If this is not found, then it tries to find the first sentence where at least one of the terms occurs. Once this is done, we extract the required amount of text and display to the user as a snippet.

Analysis of results:

Five examples of spell correction:

1.

The screenshot shows a web browser window titled "PHP Solr Client Example - Mozilla Firefox". The address bar displays "10.0.2.15/index.php?q=article&algo=luene&checkspelling=false". The search bar contains the word "article". Below the search bar, there are radio buttons for "Lucene Algorithm" (selected) and "PageRank Algorithm". A "Submit Query" button is present. The results section shows "Showing results for **article**" and "Search instead for **article**". It lists "Results 1 - 10 of 18560:". The first five results are:

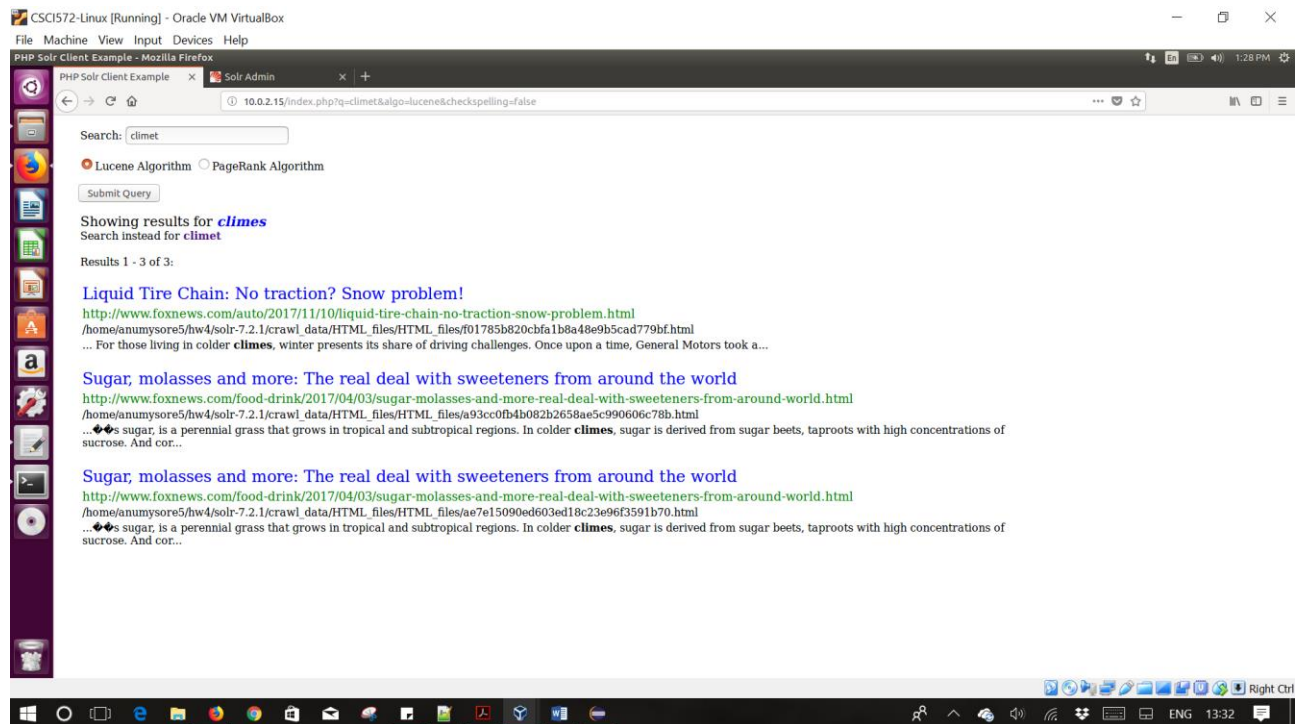
- CNN caught in family feud with ex-anchor Soledad O'Brien who called article 'mediocre,' 'facile'**
<http://www.foxnews.com/entertainment/2018/03/07/cnn-caught-in-family-feud-with-ex-anchor-soledad-obrien-who-called-article-mediocre-facile.html>
/home/anumysore5/hw4/solr-7.2.1/crawl_data/HTML_files/HTML_files/271b47c534c363a7b0f9966148603e7.html
... CNN caught in family feud with ex-anchor Soledad O'Brien who called **article** 'mediocre,' 'facile' ...
- Southern Poverty Law Center apologizes after painting journalists as fascists in retracted article**
<http://www.foxnews.com/entertainment/2018/03/16/southern-poverty-law-center-apologizes-after-painting-journalists-as-fascists-in-retracted-article.html>
/home/anumysore5/hw4/solr-7.2.1/crawl_data/HTML_files/HTML_files/09faa9d5530afa3c0be2f6235e991679.html
...outhern Poverty Law Center apologizes after painting journalists as fascists in retracted **article** By ...
- Southern Poverty Law Center apologizes after painting journalists as fascists in retracted article**
<http://www.foxnews.com/entertainment/2018/03/16/southern-poverty-law-center-apologizes-after-painting-journalists-as-fascists-in-retracted-article.html>
/home/anumysore5/hw4/solr-7.2.1/crawl_data/HTML_files/HTML_files/c7fe2a0a1f942ae7b0b0f671c9145c23.html
...outhern Poverty Law Center apologizes after painting journalists as fascists in retracted **article** By ...
- GoDaddy expels white supremacist site Daily Stormer after article on Charlottesville victim**
<http://www.foxnews.com/tech/2017/08/14/godaddy-expels-white-supremacist-site-daily-stormer-after-article-on-charlottesville-victim.html>
/home/anumysore5/hw4/solr-7.2.1/crawl_data/HTML_files/HTML_files/c0368b46cc750db4da05ef1f1ec00d51.html
... GoDaddy expels white supremacist site Daily Stormer after **article** on Charlottesville victim By ...
- Fox News**
<http://www.foxnews.com/entertainment/2018/03/07/cnn-caught-in-family-feud-with-ex-anchor-soledad-obrien-who-called-article-mediocre-facile.html>

2.

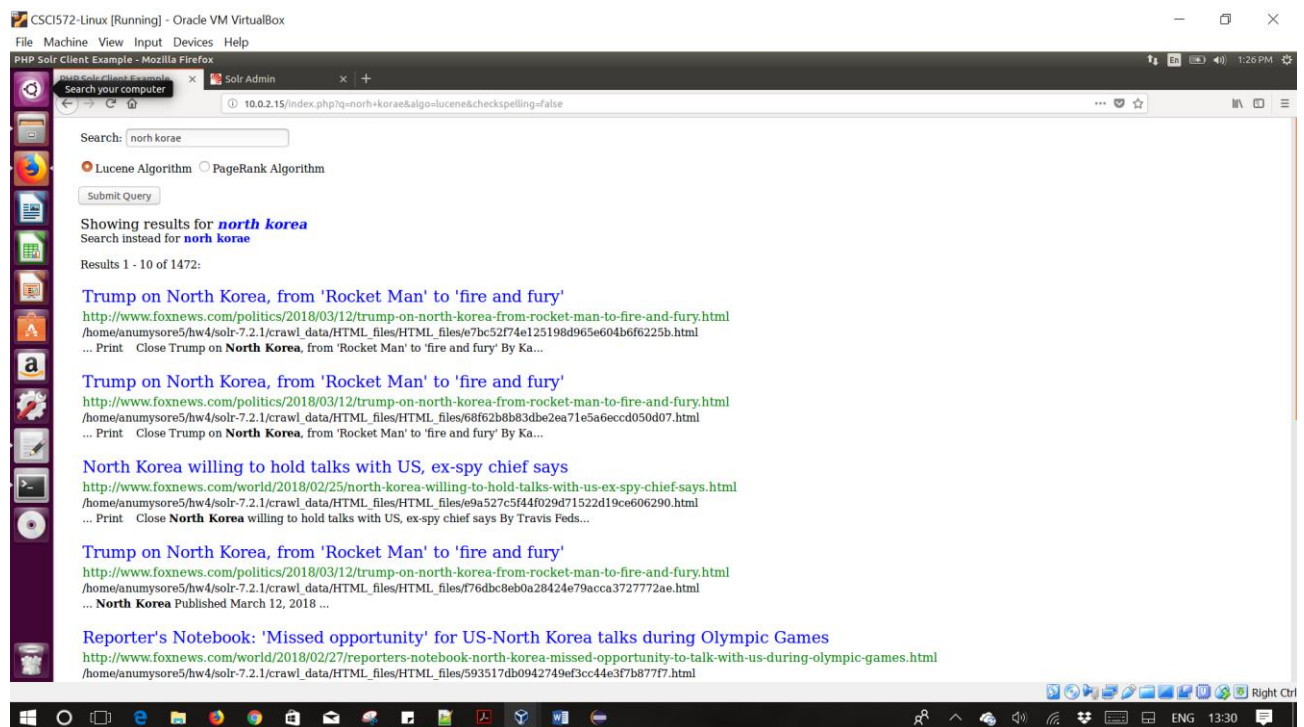
The screenshot shows a web browser window titled "PHP Solr Client Example - Mozilla Firefox". The address bar displays "10.0.2.15/index.php?q=bitcoin&algo=luene&checkspelling=false". The search bar contains the word "bitcoin". Below the search bar, there are radio buttons for "Lucene Algorithm" (selected) and "PageRank Algorithm". A "Submit Query" button is present. The results section shows "Showing results for **bitcoin**" and "Search instead for **bitcoin**". It lists "Results 1 - 10 of 110:". The first five results are:

- Lawsuit alleges bitcoin pioneer became a thief**
<http://www.foxnews.com/tech/2018/03/01/lawsuit-alleges-bitcoin-pioneer-became-thief.html>
/home/anumysore5/hw4/solr-7.2.1/crawl_data/HTML_files/HTML_files/015ce645f7726e2a99287ada6eea14c.html
... Print Close Lawsuit alleges **bitcoin** pioneer became a thief Published March 01, 2018 Newser ...
- Lawsuit alleges bitcoin pioneer became a thief**
<http://www.foxnews.com/tech/2018/03/01/lawsuit-alleges-bitcoin-pioneer-became-thief.html>
/home/anumysore5/hw4/solr-7.2.1/crawl_data/HTML_files/HTML_files/69ad9b81695a19c8cfbb7062cde80aa9.html
... Print Close Lawsuit alleges **bitcoin** pioneer became a thief Published March 01, 2018 Newser ...
- Q&A: How is the growth of bitcoin affecting the environment?**
<http://www.foxnews.com/us/2018/02/11/q-how-is-growth-bitcoin-affecting-environment.html>
/home/anumysore5/hw4/solr-7.2.1/crawl_data/HTML_files/HTML_files/b446969f3a836088f8e8c0652837739af.html
... Q&A: How is the growth of **bitcoin** affecting the environment? ...
- Lawsuit alleges bitcoin pioneer became a thief**
<http://www.foxnews.com/tech/2018/03/01/lawsuit-alleges-bitcoin-pioneer-became-thief.html>
/home/anumysore5/hw4/solr-7.2.1/crawl_data/HTML_files/HTML_files/8dddb24956d5a34a7f8503f0ea5897426.html
... Published March 01, 2018 Lawsuit alleges **bitcoin** pioneer became a thief By ...
- Seller puts lavish house on market for \$1.74 million, but also accepts bitcoin**
<http://www.foxnews.com/real-estate/2018/01/29/seller-puts-lavish-house-on-market-for-1-74-million-but-also-accepts-bitcoin.html>
/f43ae298c5b7e744e0618a7f34e0ac4.html

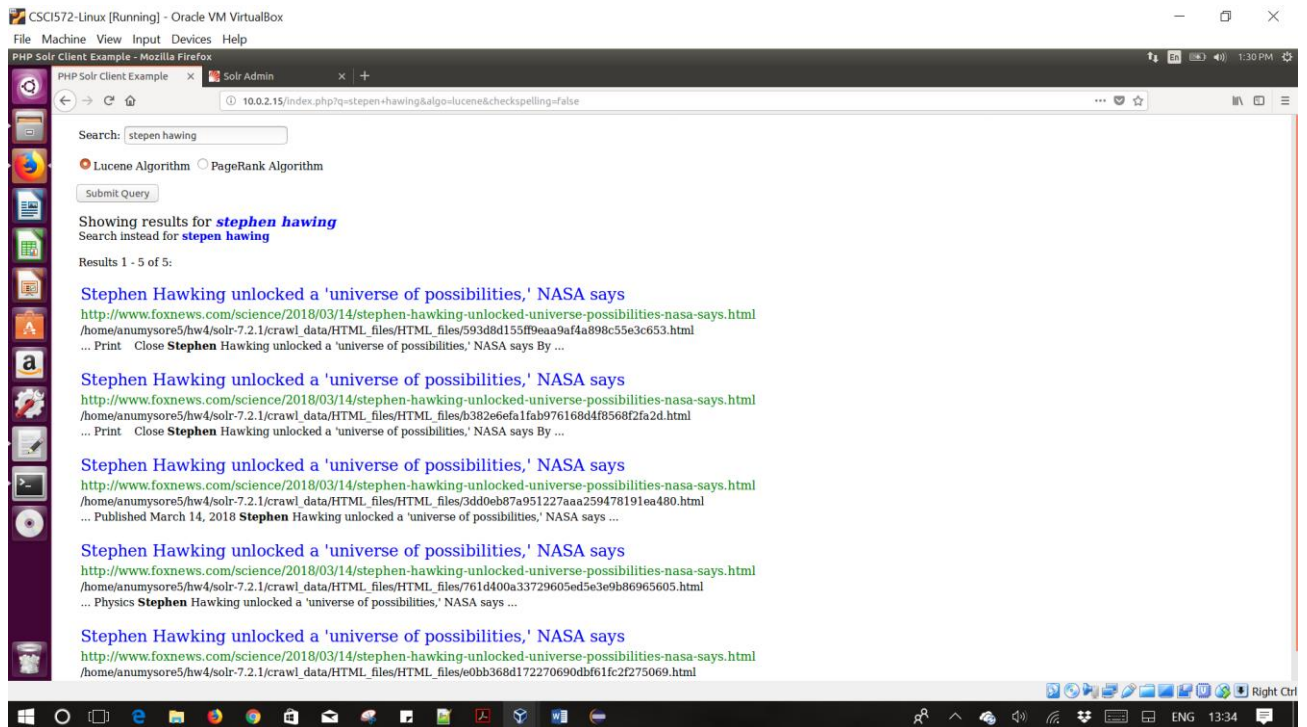
3.



4.

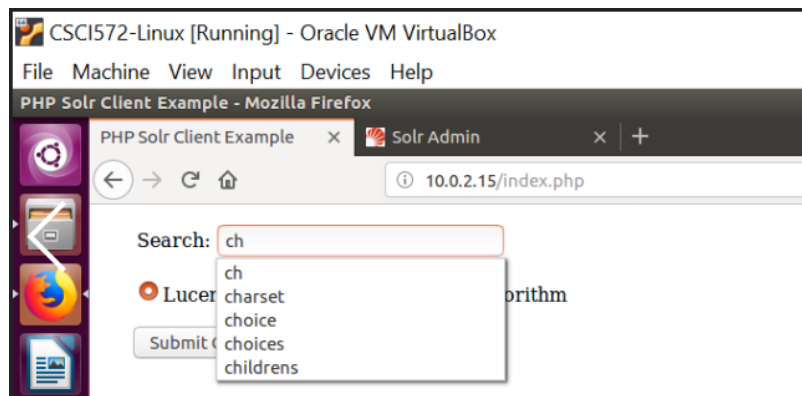


5.

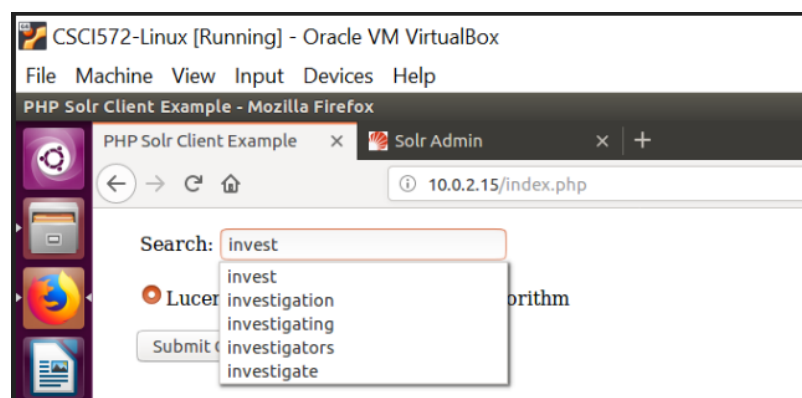


Five examples of auto-completion:

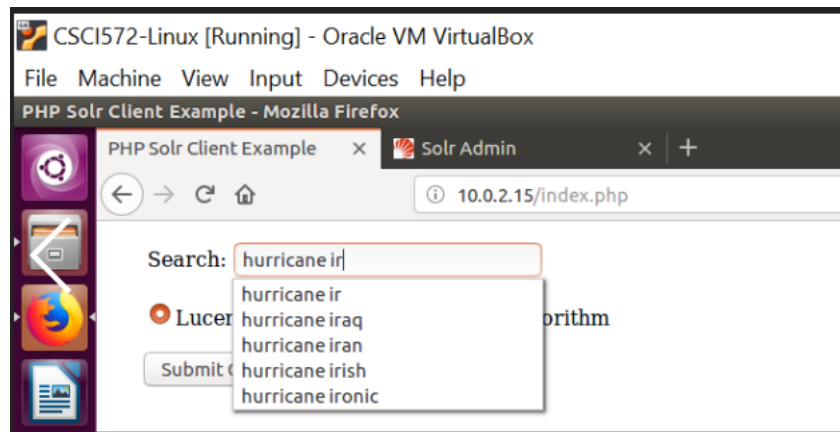
1.



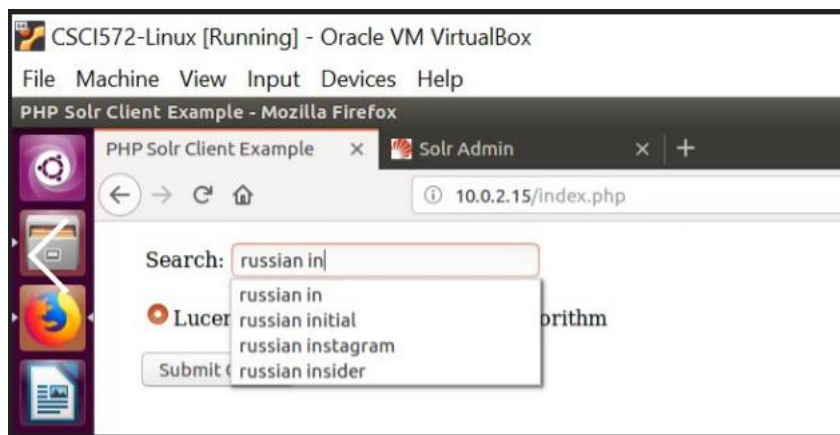
2.



3.



4.



5.

