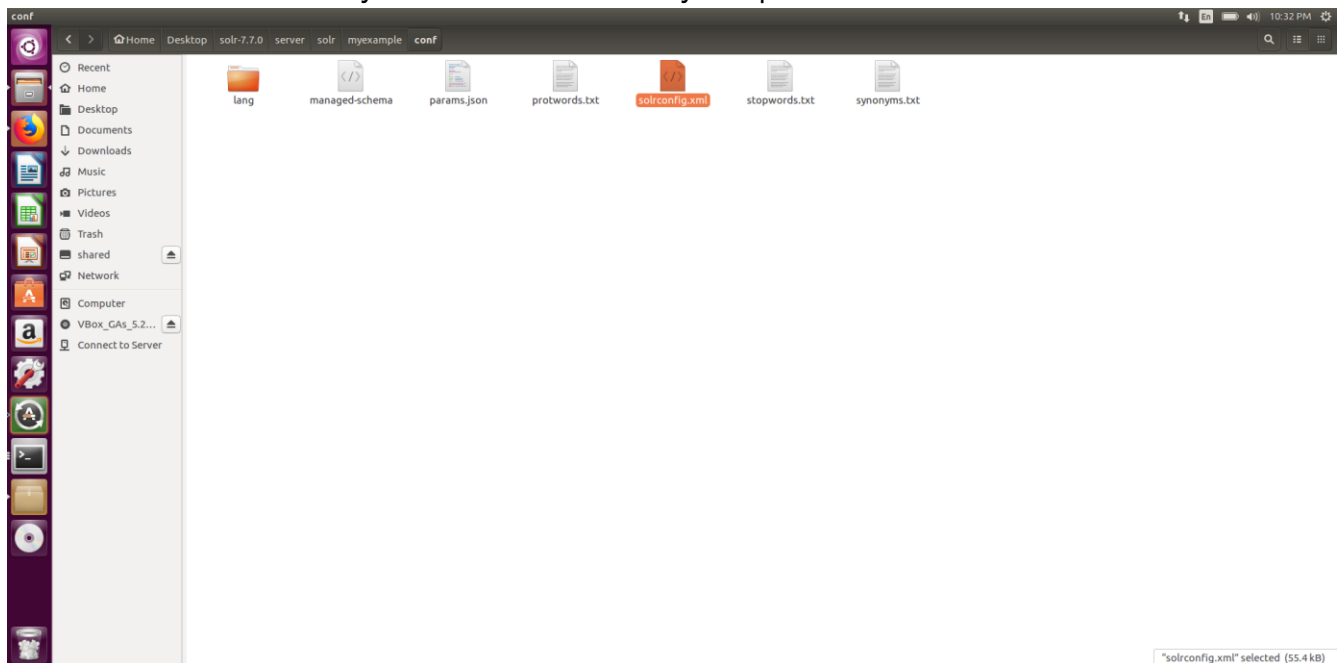Neekita Salvankar - USC ID: 8591-3366-93

Homework 4 Report (All steps and screenshots are explained in detailed manner)
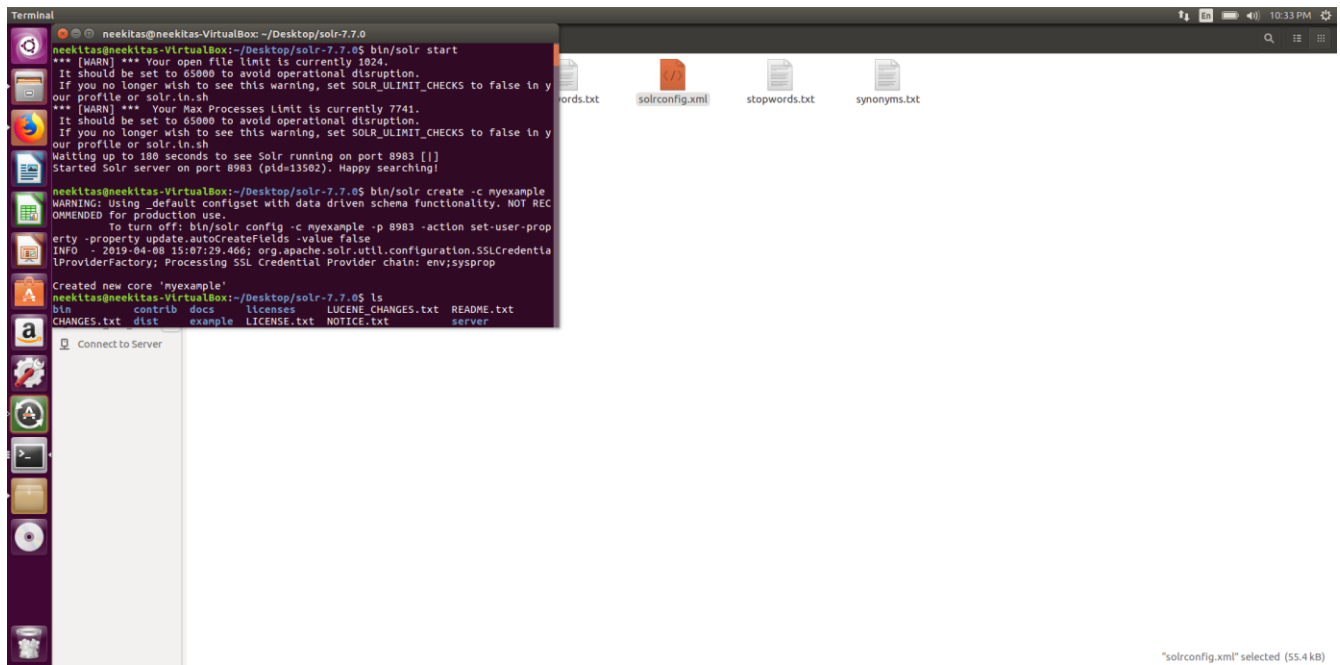
1.   Installation of Virtual Machine and Ubuntu
a. Installed and completed set up of Virtual Machine and Ubuntu using the tutorial provided on class website
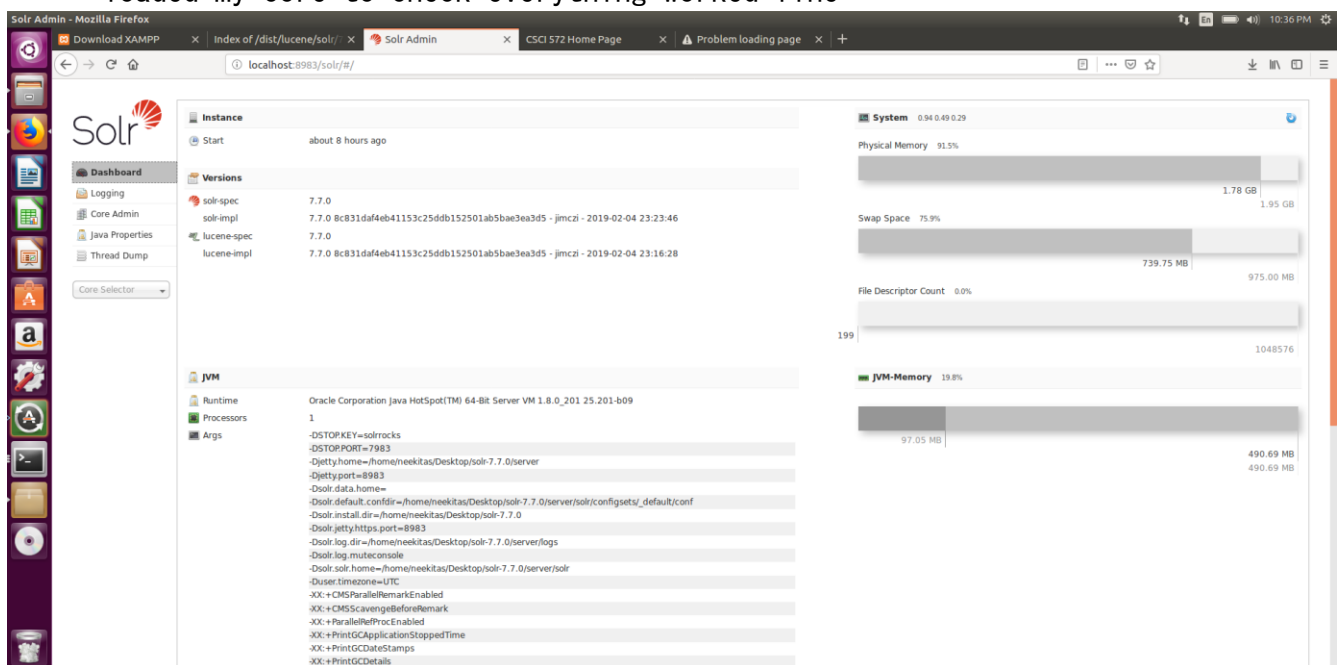
2. Downloaded and installed Solr

   a. Completed the download and installation of Solr using the tutorial provided on class website. Then I started running Solr using "bin/solr start" command . I then created my own core named "myexample".

b. I typed "localhost:8983" on my Mozilla browser to run the Solr admin and loaded my core to check everything worked fine



c. I made changes to 2 conf files :
1. managed-schema file : uncommented copyfield source tag
2. solrconfig.xml file: to default querying from the field <u>text</u> by defining default field in requestHandler and uncommented the str element with name "df", to specify the default query field to be "_text_"

d.  Downloaded the folder in Google Drive containing files for REUTER news as per
    my USC ID



e.   Indexed html files specified in the Reuter folder.



f.  Reloaded the core

## 3. Creating EdgeList file

I wrote a Java program and used the Jsoup library to create a file with edgelists. I used Eclipse to develop the same.

## 4. Creating External PageRank file

a. Used network library of python to create an external pageRank file. The input was the edgelist file created in the previous step.
b. I loaded the web graph using read_edgeList function with edgelist and digraph a parameters.
c. Used the pagerank function with following parameters:
   alpha=0.85, personalization=None, max_iter=30, tol=1e-06, nstart=None, weight = 'weight', dangling= None

## 5. I made changes to the conf files to add the external page rank files and reloaded the solr admin ( as mentioned in the tutorial)

a. Managed-schema file :
   Added the following:

```
<fieldType            name=" external"            keyField=" id"            defVal=" 0"
class=" solr.ExternalFileField" />
<field      name=" pageRankFile"      type=" external"      stored=" false"
indexed=" false" />
```

b. Solrconfig.xml file:
   Added the following:
   ```
   <listener event = "newSearcher" class =
   "org.apache.solr.schema.ExternalFileFieldReloader"/>
   <listener event = "firstSearcher" class =
   "org.apache.solr.schema.ExternalFileFieldReloader"/>
   ```

6. I reloaded the core and put in a sample query and got results both by default method and PageRank method
Below are the screenshots of both working:

Below are the results of 8 queries and also graph of the overlap


1. Venezuela

| No | Default | PageRank |
|---|---|---|
| 1 | https://www.reuters.com/investigates/special-report/venezuela-russia-rosneft/ | https://www.reuters.com/article/us-ethiopia-airplane-airlines-idUSKCN1QW2P8 |
| 2 | https://www.reuters.com/article/us-venezuela-russia-rosneft-special-repo-idUSKCN1QV1HN | https://www.reuters.com/article/us-barbie-anniversary-idUSKBN1QP2HG |
| 3 | https://www.reuters.com/article/us-usa-venezuela-sanctions-idUSKCN1QV2VS | https://www.reuters.com/news/archive/russia |

| 4 | https://www.reuters.com/article/us-venezuela-politics-usa-abrams-idUSKCN1QW2HM | https://www.reuters.com/article/aviva-ceo-idUSFWN20Q09G |
|---|---|---|
| 5 | https://www.reuters.com/article/us-venezuela-politics-tankers-idUSKCN1QW2L6 | https://www.reuters.com/journalists/julia-symmes-cobb |
| 6 | https://www.reuters.com/article/us-venezuela-russia-rosneft-money-idUSKCN1QV1HS | https://www.reuters.com/article/us-soccer-spain-usa-idUSKCN1N652I |
| 7 | https://www.reuters.com/journalists/nelson-bocanegra | https://www.reuters.com/article/us-northkorea-dissidents-idUSKCN1QW2Z |
| 8 | https://www.reuters.com/journalists/lesley-wroughton | https://www.reuters.com/finance/currencies/quote?destAmt=&srcAmt=1&srcCurr=USD&destCurr=CNY |
| 9 | https://www.reuters.com/journalists/roberta-rampton | https://www.reuters.com/article/us-usa-trump-russia-deripaska-idUSKCN1QW2CA |
| 10 | https://www.reuters.com/article/us-venezuela-politics-usa-banks-idUSKCN1QN2QE | https://www.reuters.com/journalists/nelson-bocanegra |

## 2. Senate

| No | Default | PageRank |
|---|---|---|
| 1 | https://www.reuters.com/article/us-usa-trump-congress-emergency-idUSKCN1QV1A8 | https://www.reuters.com/article/us-usa-trump-congress-emergency-idUSKCN1QW2AK |
| 2 | https://www.reuters.com/video/2019/03/04/senate-has-votes-to-block-border-emergen?videoId=521958334&videoChannel=-10465 | https://www.reuters.com/video/2019/03/06/breakingviews-tv-carry-trades?videoId=522657667&videoChannel=117766 |

| 3 | https://www.reuters.com/video/2019/03/05/mcconnell-senate-has-votes-to-block-bord?videoId=521939215&videoChannel=-10465 | https://www.reuters.com/video/2019/02/01/images-of-january?videoId=510237001&videoChannel=118069 |
|---|---|---|
| 4 | https://www.reuters.com/article/us-usa-politics-beto-orourke-special-rep-idUSKCN1QW26N | https://www.reuters.com/news/picture/photos-of-the-week-idUSRTX6R5YJ |
| 5 | https://www.reuters.com/article/us-usa-politics-beto-orourke-special-rep-idUSKCN1QW26N | https://www.reuters.com/article/us-usa-trump-russia-timing-explainer-idUSKBN1QU1E8 |
| 6 | https://www.reuters.com/article/us-usa-politics-beto-short-exclusive-idUSKCN1QW28H | https://www.reuters.com/article/us-myanmar-journalists-reaction-factbox-idUSKBN1FU030 |
| 7 | https://www.reuters.com/article/us-usa-politics-beto-short-exclusive-idUSKCN1QW28H | https://www.reuters.com/article/us-column-miller-socialsecurity-idUSKCN1PP1AN |
| 8 | https://www.reuters.com/politics | https://www.reuters.com/video/2019/03/15/tesla-launches-the-model-y-electric-cros?videoId=526236973&videoChannel=118169 |
| 9 | https://www.reuters.com/article/us-congo-politics-idUSKCN1QW28R | https://www.reuters.com/politics |
| 10 | https://www.reuters.com/article/us-usa-election-orourke-idUSKCN1QV00E | https://www.reuters.com/article/us-usa-trump-russia-cohen-idUSKCN1QH0K6 |

## 3. Democrats

| No | Default | PageRank |
|---|---|---|
| 1 | https://www.reuters.com/article/us-usa-election-medicare-idUSKBN1QU189 | https://www.reuters.com/article/us-usa-trump-congress-emergency-idUSKCN1QW2AK |
| 2 | https://www.reuters.com/article/us-usa-election-medicare-idUSKBN1QU189 | https://www.reuters.com/politics |

| | | |
|---|---|---|
| 3 | https://www.reuters.com/article/us-usa-congress-whitaker-idUSKCN1PX1LD | https://www.reuters.com/article/us-usa-trump-russia-cohen-idUSKCN1QH0K6 |
| 4 | https://www.reuters.com/video/2019/02/28/republicans-bash-cohen-democrats-at-cpac?videoId=520293110&videoChannel=13976 | https://www.reuters.com/news/archive/politicsNews |
| 5 | https://www.reuters.com/article/us-usa-election-biden-idUSKCN1QH1EX | https://www.reuters.com/article/us-usa-politics-trump-idUSKCN1QH29L |
| 6 | https://www.reuters.com/article/us-usa-trump-budget-idUSKBN1QS12D | https://www.reuters.com/article/usa-election-orourke-idINKCN1QV130 |
| 7 | https://www.reuters.com/video/2019/03/05/the-witch-hunt-continues-trump-on-democr?videoId=522293719&videoChannel=13976 | https://www.reuters.com/article/us-usa-census-ross-idUSKCN1QV27O |
| 8 | https://www.reuters.com/article/us-usa-trump-russia-congress-idUSKCN1QV23E | https://www.reuters.com/journalists/nathan-layne |
| 9 | https://www.reuters.com/politics | https://www.reuters.com/article/us-usa-election-booker-idUSKCN1QW2P4 |
| 10 | https://www.reuters.com/article/us-usa-trump-congress-emergency-idUSKCN1QW2AK | https://www.reuters.com/article/germany-saudi-arms-airbus-idUSL5N20M7OT |

## 4. Republicans

| No | Default | PageRank |
|---|---|---|
| 1 | https://www.reuters.com/article/us-usa-trump-congress-emergency-idUSKCN1QV1A8 | https://www.reuters.com/article/us-usa-trump-congress-emergency-idUSKCN1QW2AK |

| No | Default | PageRank |
|---|---|---|
| 2 | https://www.reuters.com/video/2019/02/28/republicans-bash-cohen-democrats-at-cpac?videoId=520293110&videoChannel=13976 | https://www.reuters.com/article/us-column-miller-socialsecurity-idUSKCN1PP1AN |
| 3 | https://www.reuters.com/article/us-usa-trump-congress-emergency-idUSKCN1QW2AK | https://www.reuters.com/politics |
| 4 | https://www.reuters.com/article/us-usa-trump-congress-emergency-idUSKCN1QW2AK | https://www.reuters.com/article/us-usa-trump-russia-cohen-idUSKCN1QH0K6 |
| 5 | https://www.reuters.com/article/us-usa-trump-congress-emergency-idUSKCN1QW2AK | https://www.reuters.com/news/archive/politicsNews |
| 6 | https://www.reuters.com/article/us-usa-trump-congress-emergency-idUSKCN1QW2AK | https://www.reuters.com/article/us-usa-politics-trump-idUSKCN1QH29L |
| 7 | https://www.reuters.com/article/us-usa-tax-survey-idUSKCN1QW1BY | https://www.reuters.com/article/us-usa-census-ross-idUSKCN1QV27O |
| 8 | https://www.reuters.com/article/us-california-fires-trump-idUSKCN1P31ND | https://www.reuters.com/video/2019/02/28/republicans-bash-cohen-democrats-at-cpac?videoId=520293110&videoChannel=13976 |
| 9 | https://www.reuters.com/article/us-usa-tax-survey-idUSKCN1QW1BY | https://www.reuters.com/article/us-usa-immigration-wall-insight-idUSKCN1QW1GR |
| 10 | https://www.reuters.com/article/us-usa-politics-trump-idUSKCN1QH29L | https://www.reuters.com/article/us-usa-court-census-idUSKCN1QW2MC |

## 5. Patriot Movement

| No | Default | PageRank |
|---|---|---|

| | | |
|---|---|---|
| 1 | https://www.reuters.com/article/us-usa-trump-turkey-idUSKCN1QV328 | "https://www.reuters.com/article/us-climate-change-youth-idUSKCN1QW01 |
| 2 | https://www.reuters.com/article/us-france-protests-idUSKCN1Q325A | https://www.reuters.com/journalists/fergal-smith |
| 3 | https://www.reuters.com/article/uk-people-rkelly-metoo-idUKKCN1QC0V7 | https://www.reuters.com/article/uk-britain-sterling-idUSKBN1QS0YM |
| 4 | https://www.reuters.com/article/us-britain-climatechange-education-idUSKCN1QW2QU | https://www.reuters.com/article/us-britain-lgbt-rights-idUSKCN1QW2QI |
| 5 | https://www.reuters.com/article/us-climate-change-youth-idUSKCN1QW01S | "https://www.reuters.com/article/us-taiwan-china-hongkong-idUSKCN1QW0NT |
| 6 | https://www.reuters.com/article/us-climate-change-youth-idUSKCN1QW01S | https://www.reuters.com/article/canada-forex-idUSL1N2121OB |
| 7 | https://www.reuters.com/article/us-climate-change-youth-idUSKCN1QW01 | https://www.reuters.com/journalists/crispian-balmer |
| 8 | https://www.reuters.com/article/us-usa-election-orourke-idUSKCN1QV13C | https://www.reuters.com/article/us-usa-directors-diversity-idUSKCN1Q20E8 |
| 9 | https://www.reuters.com/article/us-usa-election-orourke-idUSKCN1QV13C | https://www.reuters.com/journalists/jillian-kitchener |
| 10 | https://www.reuters.com/article/us-italy-china-analysis-idUSKCN1QW1E2 | https://www.reuters.com/article/us-myanmar-rakhine-investment-idUSKCN1QB0HO |

6.   Oscar 2019

| No | Default | PageRank |
|---|---|---|
| 1 | https://www.reuters.com/article/us-awards-oscars-women-idUSKCN1Q7130 | https://www.reuters.com/article/us-ethiopia-airplane-airlines-idUSKCN1QW2P8 |
| 2 | https://www.reuters.com/article/us-awards-oscars-diversity-idUSKCN1QE0N4 | https://www.reuters.com/article/us-barbie-anniversary-idUSKBN1QP2HG |
| 3 | https://www.reuters.com/article/us-aviancaholdings-fleet-idUSKCN1QW2XT | https://www.reuters.com/news/archive/russia |
| 4 | https://www.reuters.com/article/us-aviancaholdings-fleet-idUSKCN1QW2XT | https://www.reuters.com/article/aviva-ceo-idUSFWN20Q09G |
| 5 | https://www.reuters.com/article/us-aviancaholdings-fleet-idUSKCN1QW2XT | https://www.reuters.com/journalists/julia-symmes-cobb |
| 6 | https://www.reuters.com/journalists/pei-li | https://www.reuters.com/article/us-soccer-spain-usa-idUSKCN1N652I |
| 7 | https://www.reuters.com/article/us-britain-lgbt-celebrities-idUSKCN1QW2KY | https://www.reuters.com/article/us-northkorea-dissidents-idUSKCN1QW2ZL |
| 8 | https://www.reuters.com/theWire | https://www.reuters.com/finance/currencies/quote?destAmt=&srcAmt=1&srcCurr=USD&destCurr=CNY |
| 9 | https://www.reuters.com/article/us-britain-lgbt-rights-idUSKCN1QW2QI | https://www.reuters.com/news/archive/esgnews |
| 10 | https://www.reuters.com/article/us-britain-lgbt-rights-idUSKCN1QW2QI | https://www.reuters.com/news/archive/esgnews |

7.　Channel

| No | Default | PageRank |
|---|---|---|
| 1 | https://www.reuters.com/video/2019/03/06/breakingviews-tv-carry-trades?videoId=522657667&videoChannel=117766 | https://www.reuters.com/article/us-ethiopia-airplane-airlines-idUSKCN1QW2P8 |
| 2 | "https://www.reuters.com/video/2019/03/13/liverpool-result-at-bayern-wont-affect-p?videoId=525185681&videoChannel=79 | https://www.reuters.com/article/us-barbie-anniversary-idUSKBN1QP2HG |
| 3 | https://www.reuters.com/video/2019/03/14/boeing-fix-could-take-weeks-us-lawmakers?videoId=526002472&videoChannel=5 | https://www.reuters.com/article/aviva-ceo-idUSFWN20Q09G |
| 4 | https://www.reuters.com/video/2019/01/09/trump-threatens-california-wildfire-aid?videoId=501404962&videoChannel=118262 | https://www.reuters.com/article/us-soccer-spain-usa-idUSKCN1N652I |
| 5 | https://www.reuters.com/video/2019/03/10/trump-asks-congress-for-750b-to-fund-mil?videoId=524322587&videoChannel=-10465 | https://www.reuters.com/article/us-northkorea-dissidents-idUSKCN1QW2ZL |
| 6 | https://www.reuters.com/video/2019/03/10/brexit-in-peril-if-pm-mays-deal-is-rejec?videoId=524245642&videoChannel=118261 | https://www.reuters.com/article/us-usa-trump-russia-deripaska-idUSKCN1QW2CA |
| 7 | https://www.reuters.com/video/2019/03/15/malawi-floods-wreak-havoc-cyclone-idai-e?videoId=525916362&videoChannel=73 | https://www.reuters.com/article/us-soccer-germany-turkey-ozil-idUSKBN1KG1OX |
| 8 | https://www.reuters.com/video/2019/02/27/cohen-dishes-on-donald-jr-in-testimony?videoId=519943059&videoChannel=13976 | https://www.reuters.com/article/ecuador-imf-idUSL1N20Y1O6 |
| 9 | https://www.reuters.com/video/2019/03/13/woods-neck-not-painful-ahead-of-players?videoId=525458826&videoChannel=79 | https://www.reuters.com/article/china-consumerday-idUSL8N2102SG" |

| No | | |
|---|---|---|
| 10 | https://www.reuters.com/video/2019/02/15/sustainable-is-the-new-sexy-at-new-york?videoId=515332676&videoChannel=118285 | https://www.reuters.com/article/us-italy-china-mou-factbox-idUSKCN1QW1EB |

## 8. Wall

| No | Default | PageRank |
|---|---|---|
| 1. | http://www.reuters.com/investigates/section/wall-streets-way/" | https://www.reuters.com/news/archive/russia |
| 2. | https://www.reuters.com/article/us-usa-trump-budget-wall-exclusive-idUSKBN1QR0CW | https://www.reuters.com/news/archive/esgnews |
| 3. | https://www.reuters.com/article/us-usa-immigration-wall-insight-idUSKCN1QW1GR | https://www.reuters.com/finance/wealth |
| 4. | https://www.reuters.com/politics" | https://www.reuters.com/finance/funds |
| 5. | https://www.reuters.com/article/us-usa-trump-congress-emergency-idUSKCN1QW2AK | https://www.reuters.com/journalists/susan-mathew |
| 6. | https://www.reuters.com/article/us-usa-trump-congress-emergency-idUSKCN1QW2AK | https://www.reuters.com/article/us-usa-trump-congress-emergency-idUSKCN1QW2AK |
| 7. | https://www.reuters.com/article/us-usa-trump-congress-emergency-idUSKCN1QW2AK" | https://www.reuters.com/video/2019/03/06/breakingviews-tv-carry-trades?videoId=522657667&videoChannel=117766 |

| 8. | https://www.reuters.com/article/us-usa-trump-congress-emergency-idUSKCN1QW2AK" | https://www.reuters.com/video/2019/02/01/images-of-january?videoId=510237001&videoChannel=118069 |
|---|---|---|
| 9. | https://www.reuters.com/article/us-usa-trump-congress-emergency-idUSKCN1QV1A8 | "https://www.reuters.com/article/us-usa-trump-russia-timing-explainer-idUSKBN1QU1E8" |
| 10. | https://www.reuters.com/journalists/noel-randewich | https://www.reuters.com/news/archive/japan" |

## Number of overlaps (URL and ID were same)

| No | Query | No of overlaps |
|---|---|---|
| 1 | Venezuela | 1 |
| 2 | Senate | 5 |
| 3 | Democrats | 4 |
| 4 | Republicans | 5 |
| 5 | Patriot Movement | 3 |
| 6 | Oscar 2019 | 7 |
| 7 | Channel | 0 |
| 8 | Wall | 1 |

Graph of Overlaps

Overlap Graph

7. Installed Apache server and PHP on Ubuntu

8. I created the UI using PHP and below are the screenshots. Also, I cloned the git 'solr-php-repository' and placed it in the same folder as UI.

Added proper paths to Apache Sever and csv files with links. Added radio buttons for user to choose between Lucene and PageRank. Displayed appropriate values and set limit to 10.

CSCI-572-Information-R × | Neekita Assignment 4 × | Solr Admin × | Apache2 Ubuntu Default P: × | +

localhost/first.php

Search: Venezuela    Submit
○ Page Rank  ○ Lucene

---

CSCI-572-Information-R: × | Neekita Assignment 4 × | Solr Admin × | Apache2 Ubuntu Default P: × | hw 4 pics - neekita.salvi × | +

localhost/first.php?q=Venezuela&sort=Lucene

Page Rank ● Lucene

Results 1 - 10 of 1709 :

1. Title: How Russia sank billions of dollars into Venezuelan quicksand
   URL: https://www.reuters.com/article/us-usa-politics-beto-orourke-special-rep-idUSKCN1OW26N
   Description: A Russian oil company spent a fortune in Venezuela despite suspecting it was being ripped off. Sources say the Kremlin wanted to help its South American ally.
   ID: /home/neekitas/Downloads/Reuters/reutersnews/reutersnews/3ac186d9-6b06-4550-a2c1-3a83ffd9cb84.html

2. Title: Special Report: How Russia sank billions of dollars into Venezuela quicksand | Reuters
   URL: https://www.reuters.com/article/us-usa-politics-beto-orourke-special-rep-idUSKCN1OW26N
   Description: At the end of 2015, managers at Rosneft, the Russian state-controlled oil firm, sounded the alarm to their bosses about the company's investments in Venezuela. Rosneft's local partner, Venezuelan state oil company PDVSA, owed it hundreds of millions of dollars, according to internal documents, and there seemed no prospect things would get better.
   ID: /home/neekitas/Downloads/Reuters/reutersnews/reutersnews/2d57e926-d37b-4249-98be-1ff0d5b9d814.html

3. Title: U.S. considers sanctions to restrict Visa, Mastercard in Venezuela: official | Reuters
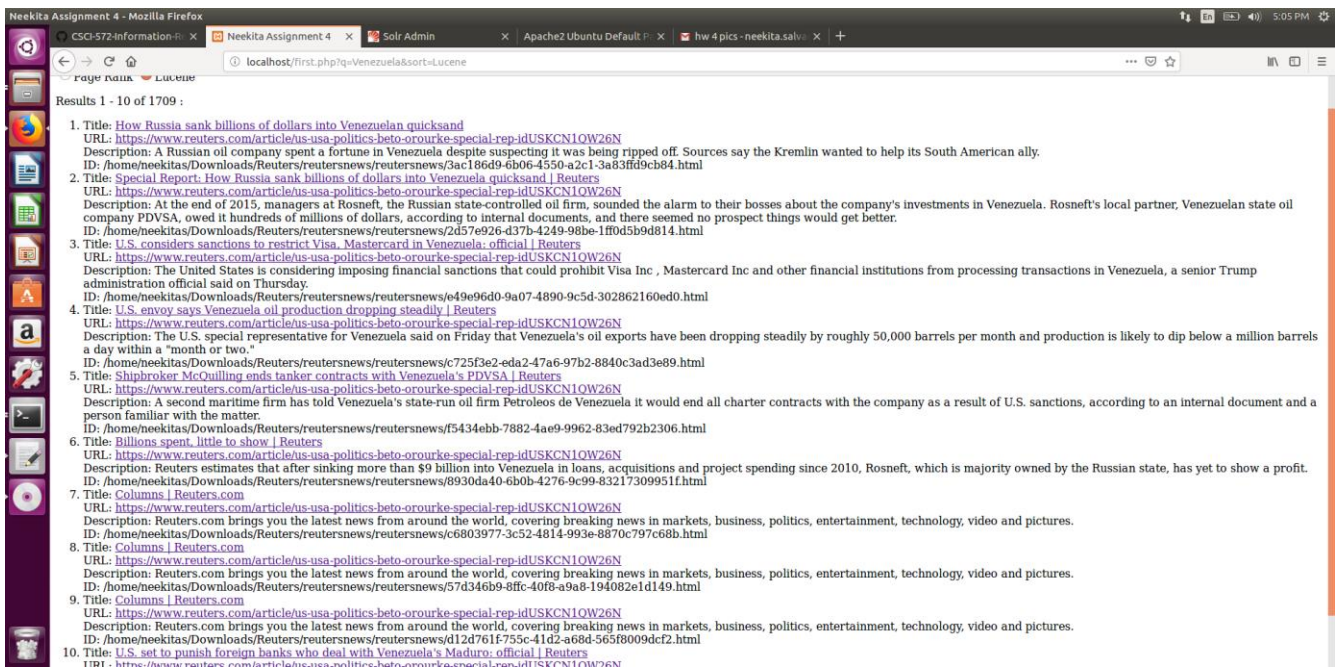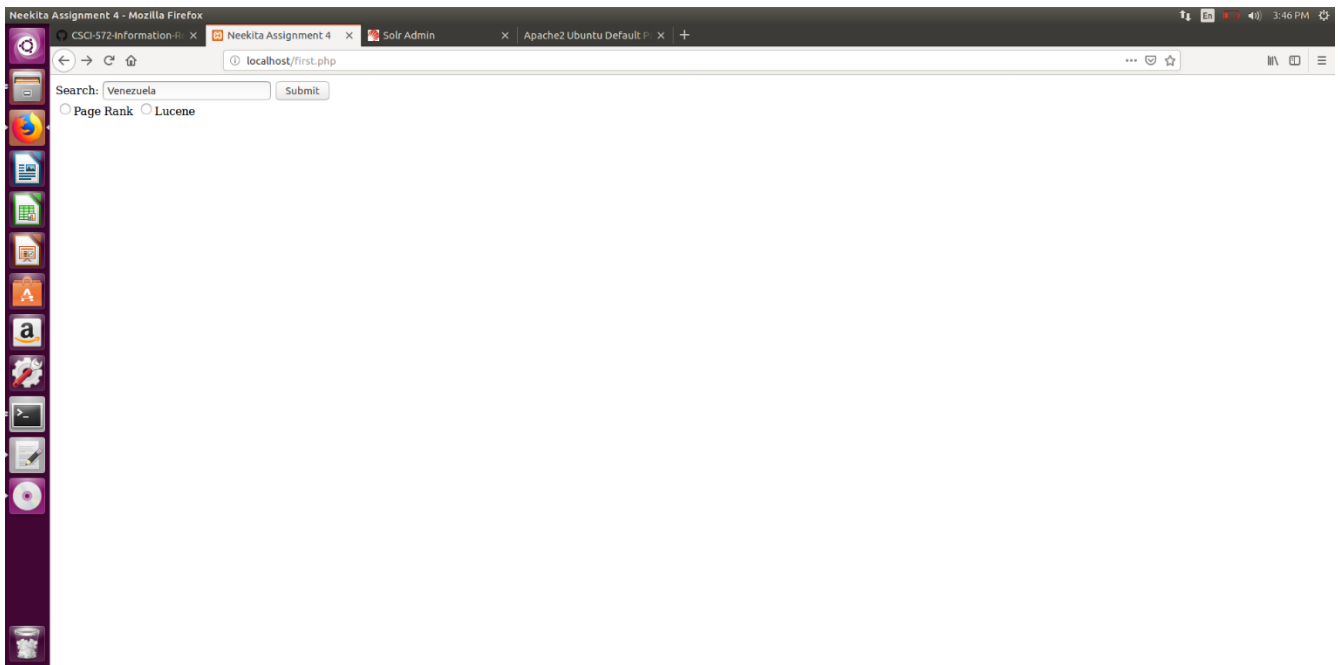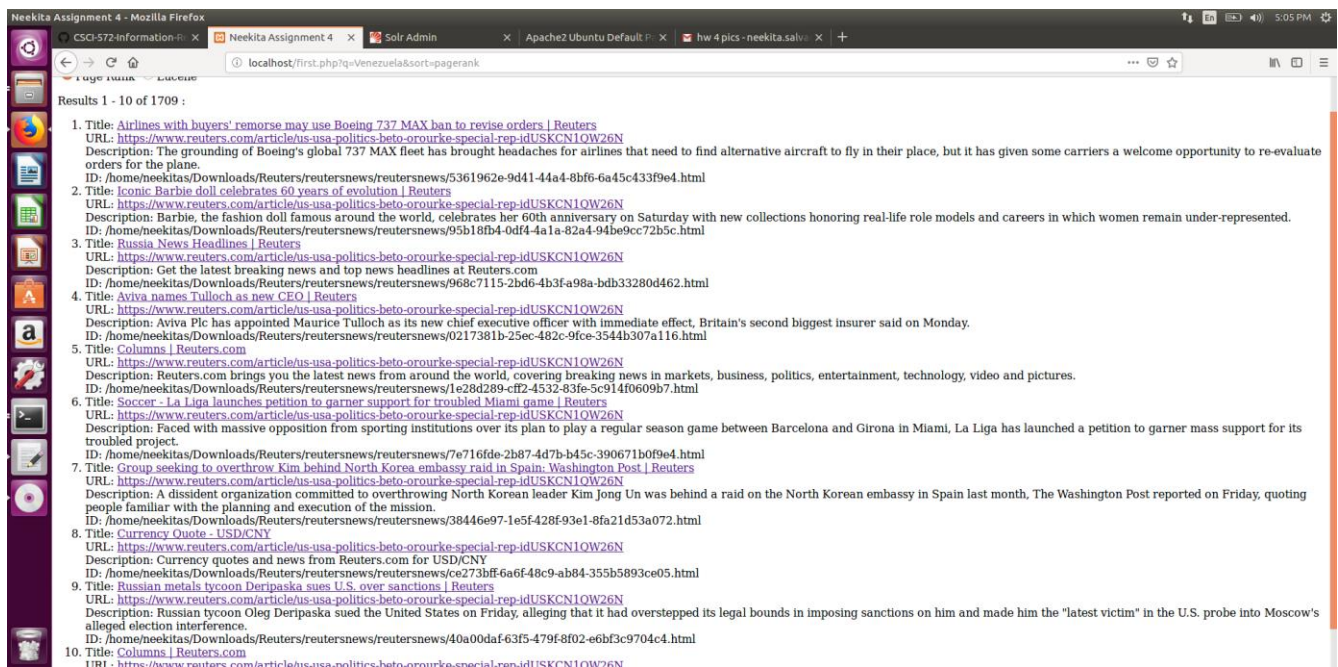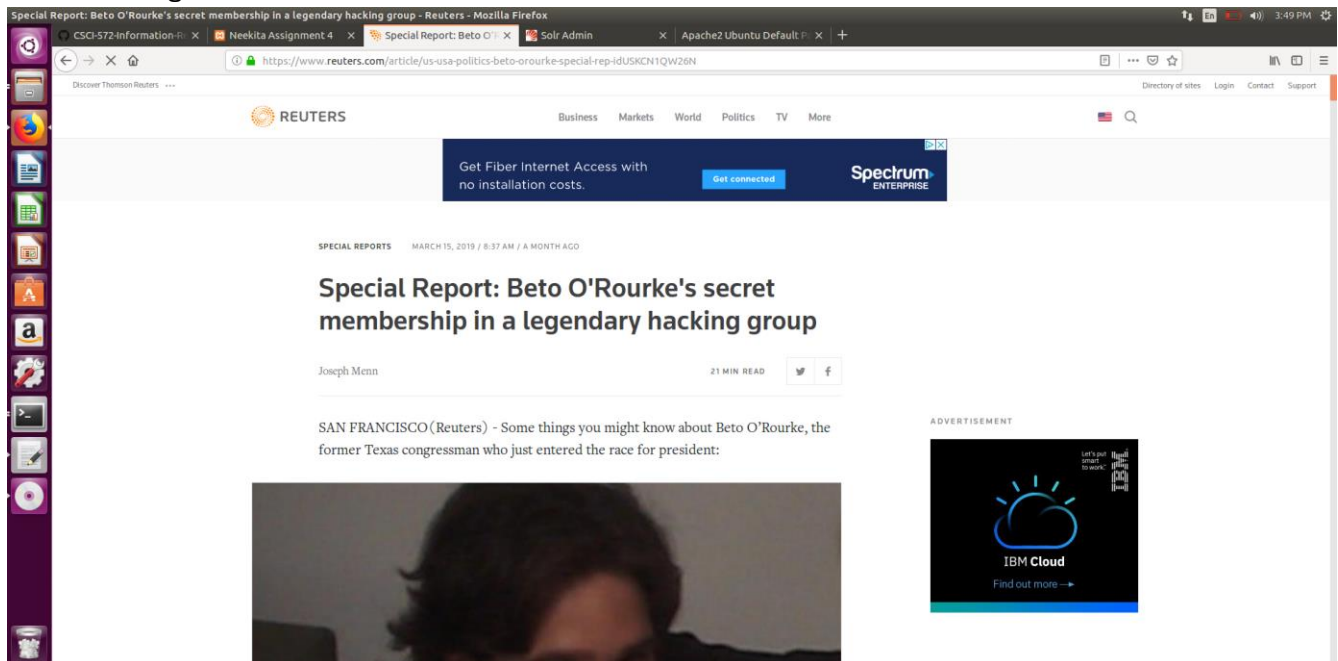   URL: https://www.reuters.com/article/us-usa-politics-beto-orourke-special-rep-idUSKCN1OW26N
   Description: The United States is considering imposing financial sanctions that could prohibit Visa Inc , Mastercard Inc and other financial institutions from processing transactions in Venezuela, a senior Trump administration official said on Thursday.
   ID: /home/neekitas/Downloads/Reuters/reutersnews/reutersnews/e49e96d0-9a07-4890-9c5d-302862160ed0.html

4. Title: U.S. envoy says Venezuela oil production dropping steadily | Reuters
   URL: https://www.reuters.com/article/us-usa-politics-beto-orourke-special-rep-idUSKCN1OW26N
   Description: The U.S. special representative for Venezuela said on Friday that Venezuela's oil exports have been dropping steadily by roughly 50,000 barrels per month and production is likely to dip below a million barrels a day within a "month or two."
   ID: /home/neekitas/Downloads/Reuters/reutersnews/reutersnews/c725f3e2-eda2-47a6-97b2-8840c3ad3e89.html

5. Title: Shipbroker McQuilling ends tanker contracts with Venezuela's PDVSA | Reuters
   URL: https://www.reuters.com/article/us-usa-politics-beto-orourke-special-rep-idUSKCN1OW26N
   Description: A second maritime firm has told Venezuela's state-run oil firm Petroleos de Venezuela it would end all charter contracts with the company as a result of U.S. sanctions, according to an internal document and a person familiar with the matter.
   ID: /home/neekitas/Downloads/Reuters/reutersnews/reutersnews/f5434ebb-7882-4ae9-9962-83ed792b2306.html

6. Title: Billions spent, little to show | Reuters
   URL: https://www.reuters.com/article/us-usa-politics-beto-orourke-special-rep-idUSKCN1OW26N
   Description: Reuters estimates that after sinking more than $9 billion into Venezuela in loans, acquisitions and project spending since 2010, Rosneft, which is majority owned by the Russian state, has yet to show a profit.
   ID: /home/neekitas/Downloads/Reuters/reutersnews/reutersnews/8930da40-6b0b-4276-9c99-83217309951f.html

7. Title: Columns | Reuters.com
   URL: https://www.reuters.com/article/us-usa-politics-beto-orourke-special-rep-idUSKCN1OW26N
   Description: Reuters.com brings you the latest news from around the world, covering breaking news in markets, business, politics, entertainment, technology, video and pictures.
   ID: /home/neekitas/Downloads/Reuters/reutersnews/reutersnews/c6803977-3c52-4814-993e-8870c797c68b.html

8. Title: Columns | Reuters.com
   URL: https://www.reuters.com/article/us-usa-politics-beto-orourke-special-rep-idUSKCN1OW26N
   Description: Reuters.com brings you the latest news from around the world, covering breaking news in markets, business, politics, entertainment, technology, video and pictures.
   ID: /home/neekitas/Downloads/Reuters/reutersnews/reutersnews/57d346b9-8ffc-40f8-a9a8-194082e1d149.html

9. Title: Columns | Reuters.com
   URL: https://www.reuters.com/article/us-usa-politics-beto-orourke-special-rep-idUSKCN1OW26N
   Description: Reuters.com brings you the latest news from around the world, covering breaking news in markets, business, politics, entertainment, technology, video and pictures.
   ID: /home/neekitas/Downloads/Reuters/reutersnews/reutersnews/d12d761f-755c-41d2-a68d-565f8009dcf2.html

10. Title: U.S. set to punish foreign banks who deal with Venezuela's Maduro: official | Reuters
    URL: https://www.reuters.com/article/us-usa-politics-beto-orourke-special-rep-idUSKCN1OW26N

On clicking the URL:



Question : Why PageRank of some pages is more than PageRank of others?

Answer: PageRank algorithm is based on how the web page is linked, that is , the number of in links and out links. If the page has many in links, it means the page is of relevance and is assigned a higher pagerank. Moreover, if these incoming links are also relevant, that is, they are pointed by many pages, then the PageRank value

is higher. In our implementation, we created the EdgeList and assigned it as an input parameter to python code to generate a graph. We computed the rank using this graph.