# INTRODUCTION TO INFORMATION RETRIEVAL CSE 535

## FALL-2019

## PROJECT-4 REPORT

**SUBMITTED BY:**

**BHARGAVA REDDY KUNCHA (50317932) (bkuncha)**
**VAMSHIDHAR REDDY KOMMIDI (50310287) (vamshidh)**
**ROHITH KUMAR GADALAY (50314899) (rohithku)**

# ANALYISING THE IMPACT OF POLITICAL RHETORIC IN TRADITIONAL AND SOCIAL MEDIA
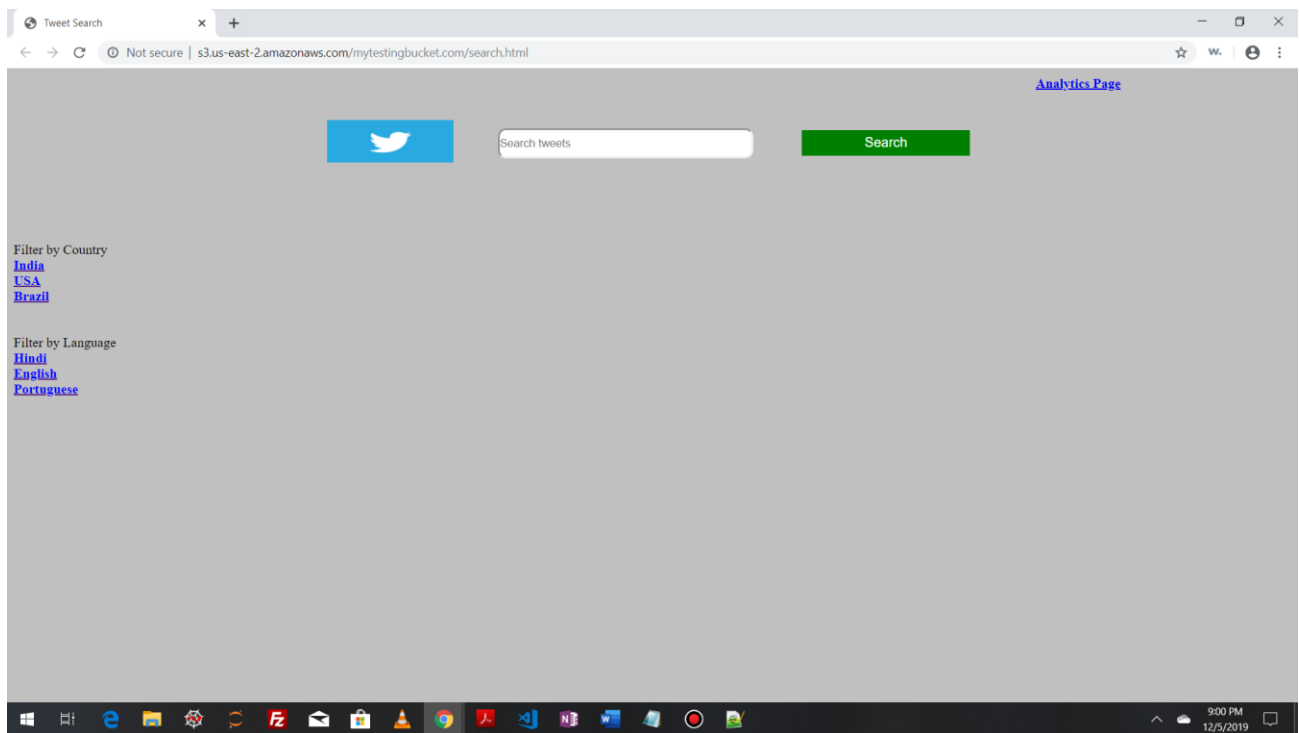
## OVERVIEW:

The data was crawled from Twitter using the twitter search API. The crawled data was then processed using python script to extract the required information such as text, hashtags, user names , language etc.. The processed tweet collection is then indexed in Solr. We developed a UI in which Solr acts as the backend of it and HTML and JavaScript is used for front end. We also included the filter search by language and search by country in-order to make a better interactive experience to the user.

The url of the search engine is:

http://s3.us-east-2.amazonaws.com/mytestingbucket.com/search.html

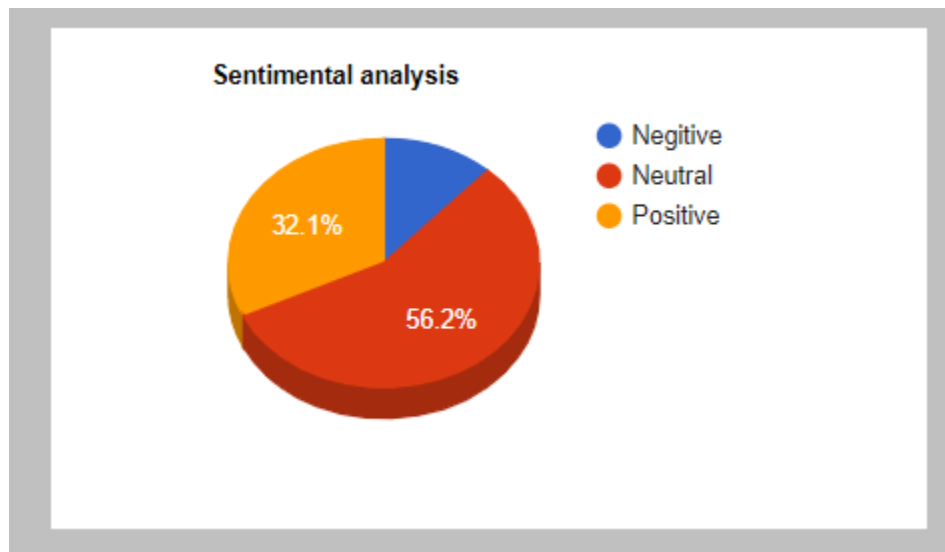The homepage of our search engine is as shown below:

## IMPLEMENTATION DETAILS:

### Preprocessing Data:

We utilized the twitter data retrieved of project 1 and implemented the preprocessing of it by using a python script. The data is processed in order to index the data effectively corresponding to the required fields. The preprocessing of data also included the addition of fields –topic obtained from topic analysis step and sentiment obtained from sentiment analysis step.
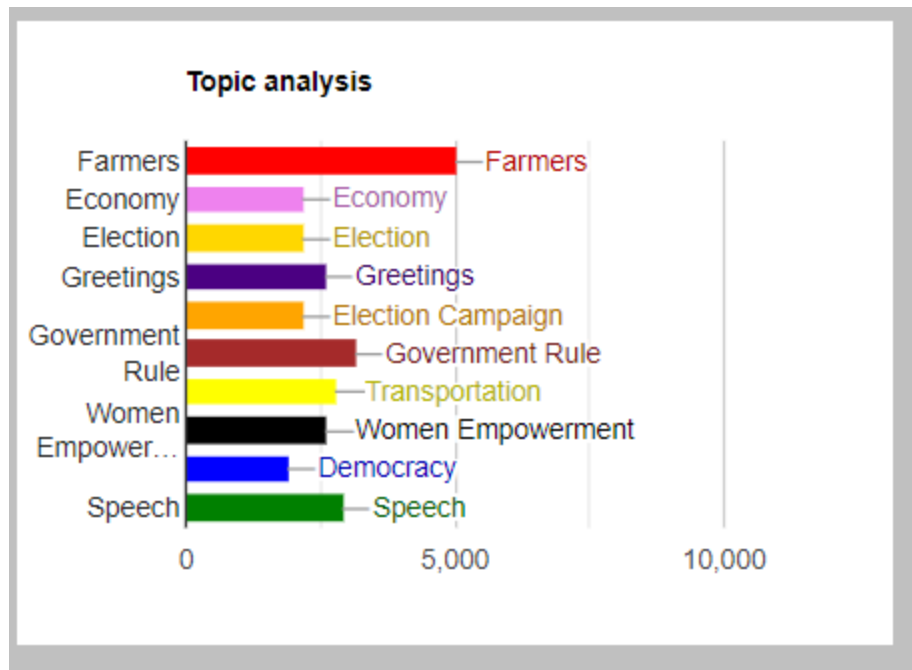
### Sentiment Analysis:

We have used VADER to perform sentimental analysis on the data retrieved. We used SentimentIntensityAnalyzer package and filtered the text of the retrieved data based on the language. The polarity scores are returned based on the tweet data and the sentiment values are classified as positive, negative and neutral respectively. The pictorial representation is done by pie-chart.
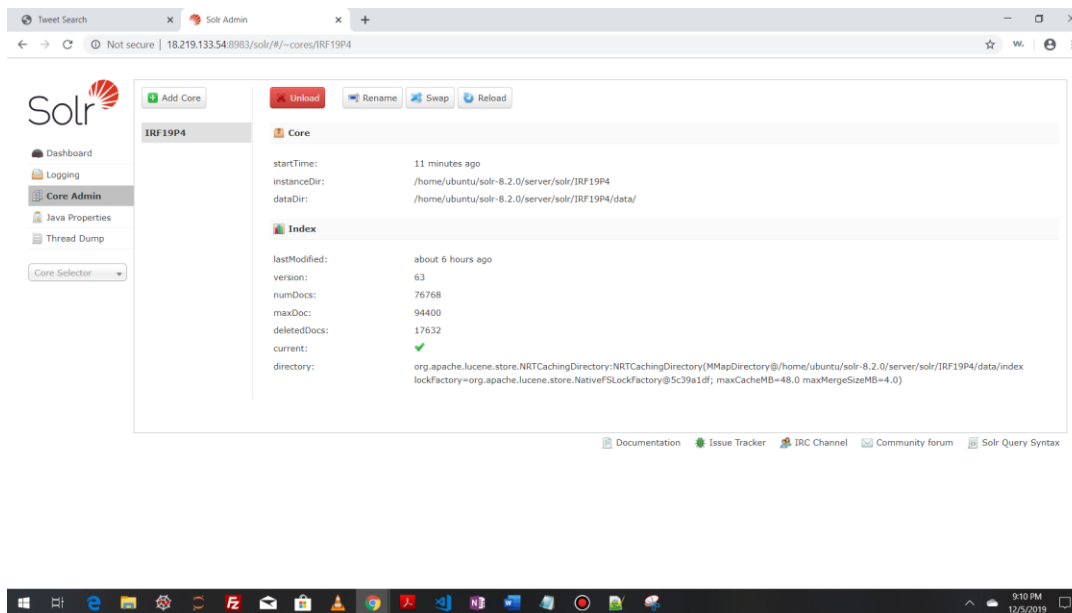


### Topic Analysis:

We have used unsupervised machine learning algorithm, LDA (Latent Dirichlet Allocation) to explore the topics on the retrieved data. We have used open source library spacy for text processing to load the language specific text. LDA'S approach to topic modeling is that it considers each document as a collection of topics in a certain proportion. And each topic as a collection of keywords, again in a certain proportion. Based on the scoring values of each word, we have identified the top 10 topics. We have performed topic modeling on the POI's tweet categorized by the country to show the impact of POI's tweets to show the impact on the society. Ten different topics from each country are chosen based on the POI's tweets. A sample of topic analysis for country India is represented in below figure.
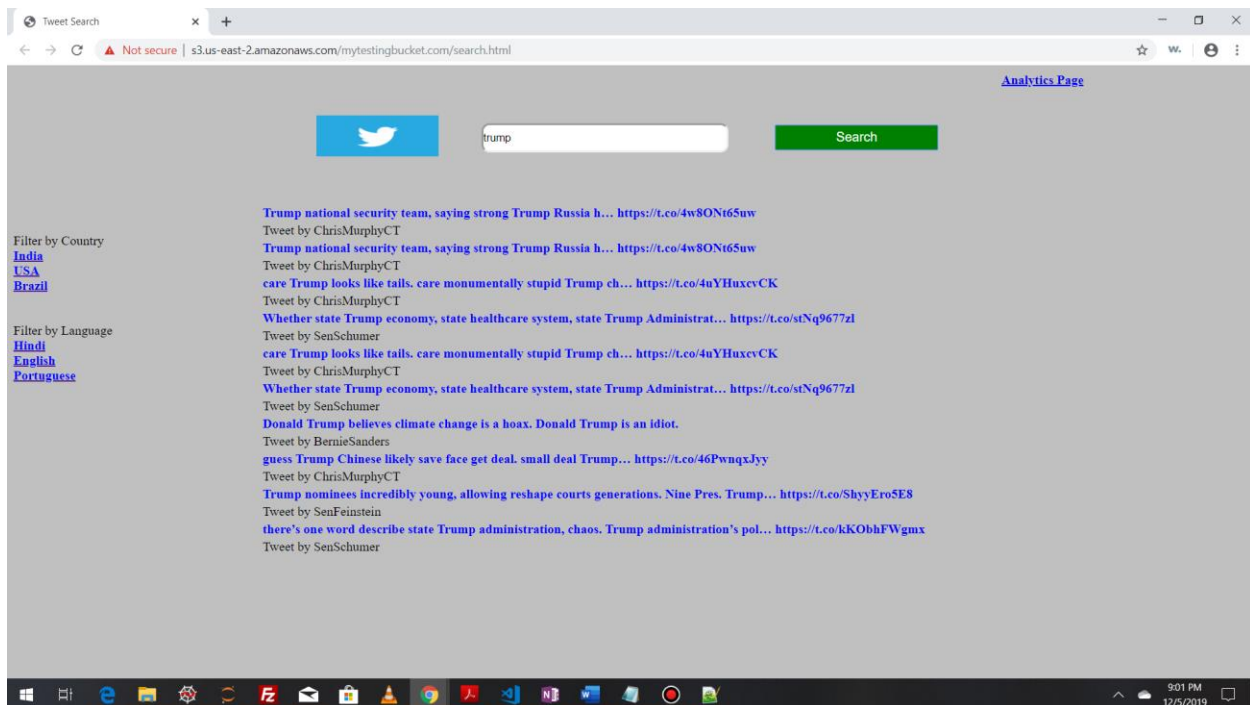
**Search Engine:**

We developed a UI for our application by using HTML, JavaScript for front end and backend. We used Google charts tool for making the pie chart and bar graph for the visualization of analytical data displayed on the webpage. The results page displays 10 results at a time and pagination is achieved. We have included the features search by language and search by country to make the UI more interactive. We have included an analytics page, when clicked will navigate to a new page which shows the visualizations implemented. We have used dismax query parser and boosted the query results in the Solr to achieve more relevant tweets based on the query. When a keyword is queried, the results are displayed accordingly and the person who has tweeted the tweet is also displayed.

The screenshot obtained from Solr after indexing is:

**RESULTS**:

When we search for trump keyword, the tweets are retrieved and the person who has done the original tweet is also represented. The resulting user interface is
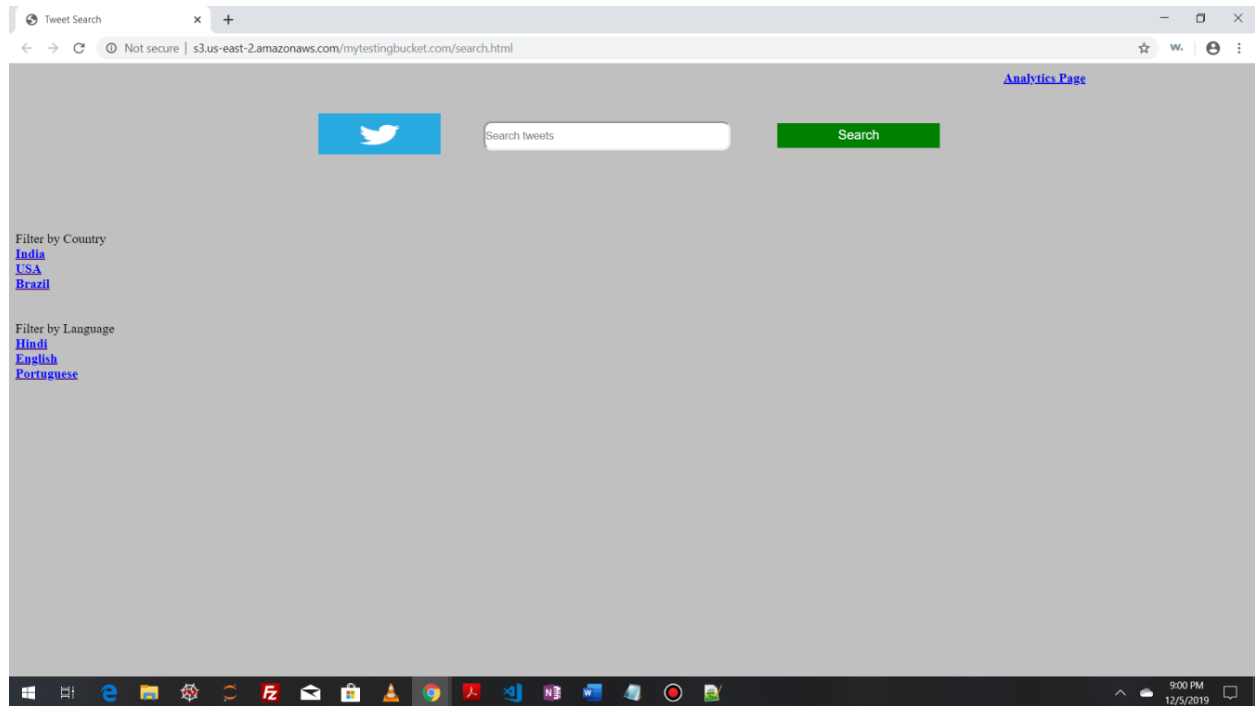
The Filter by Country and Filter by Language section is:



When the analytics page is clicked and one of search by country is selected, the resulting user interface is:

When the search page link in above screenshot is clicked, it navigates to the original page/home page



**VIDEO DEMONSTRATION**:

We demonstrated the functionality of our IR system in a video:

https://www.youtube.com/watch?v=A4oOCzHSeus&feature=youtu.be

In this video we queried for trump and retrieved the appropriate results from the tweets data.

**TEAM CONTRIBUTIONS**:

| TEAM MEMBER | UBIT NAME | UBID | TASKS |
|---|---|---|---|
| Bhargava Reddy Kuncha | bkuncha | 50317932 | UI Design, hosting and preprocessing tweets |
| Vamshidhar Reddy Kommidi | vamshidh | 50310287 | Sentimental Analysis and Topic Analysis |
| Rohith Kumar Gadalay | rohithku | 50314899 | Solr Indexing and writing report |

**REFERENCES:**

1)  Google Translator API – https://cloud.google.com/translate/docs/
2) Topic Analysis- https://www.machinelearningplus.com/nlp/topic-modeling-gensim-python/
3) Faceted searching – https://examples.javacodegeeks.com/enterprise-java/apache-solr/solr-faceted-search-example/
4) Sentiment Analysis - https://www.geeksforgeeks.org/python-sentiment-analysis-using-vader/
5)Google Chart tools - https://developers.google.com/chart/interactive/docs/