# PROJECT 4
# COMPLETE SEARCH AND ANALYTICS SOLUTION BASED ON DISSECTING TWITTER DATA
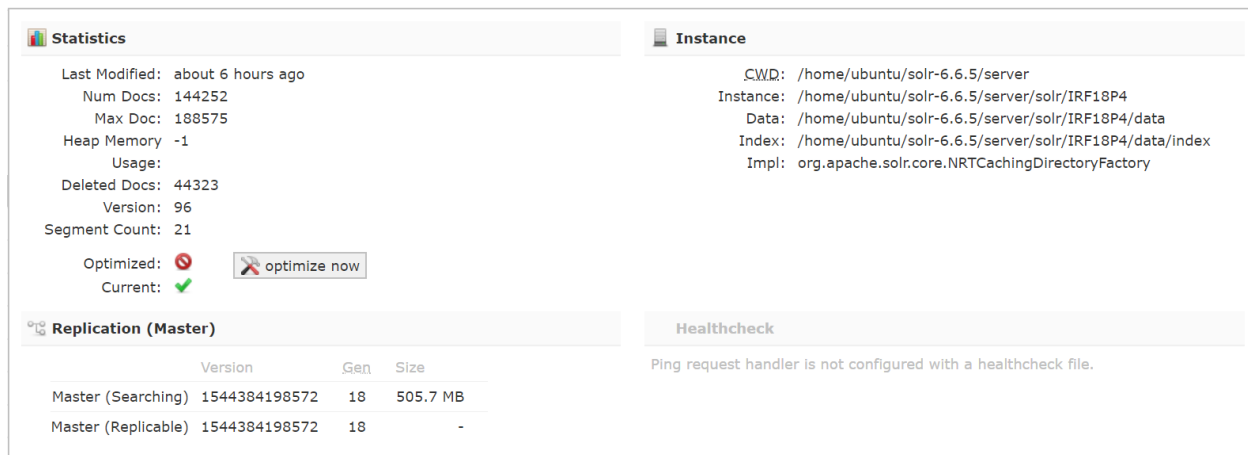
**Synopsis**:

The main goal of this project is to build an end to end IR system which enables the user to get the knowledge about the particular topic, relevant tweets and analysis results. In the first step, we collected the twitter data on few particular topics and the distribution of tweets was among the selected five cities and five languages. In the next step, we implemented the BM25 on Solr-6.6.5 schema and indexed the tweets. In the third step, we created a Web stack which would display the query results along with the analysis for the selected query and faceted search to build the IR system on these tweet topics.

**Twitter data:**

Collected data of over 2 lakh tweets using twitter API. Collection of tweets is based on five tweet topics namely social unrest, politics, environment, crime and infrastructure and these tweet topics are collected in five languages namely English, Hindi, French, Spanish and Thai and in five cities namely New Delhi, New York, Paris, Mexico, Bangkok.

**IR Model:**

We created a BM25 model schema on Solr and indexed around 1.5 lakh tweet



**Technology Stack:**

We used the java for the backend, the front end technology is JavaScript and javaservletpages to render the UI and server is Tomcat and Google cloud to host the IR system.

UI->SERVER->SOLR->SERVER(PROCESSING)->UI(DISPLAY)

**UI Display:**

The above image shows the retrieved tweets display. The middle part displays the tweets and the user who tweeted it and the text, the url directs to the twitter page where the tweet is present.

**Sentimental Analysis:**

Sentiment Analysis determines whether the writing is positive, negative or neutral. We have used Stanford.nlp package. Here in the below we have searched for trump and it shows the following sentimental analysis chart.

**Facet Analysis:**

Below is the graph showing the distribution of tweets on various dates.

Using the below, following url we have done the facet search which is inbuilt in solr

http://18.221.162.107:8984/solr/IRF18P4/select?facet.field=tweet_date&facet=on&fl=tweet_date&indent=on&wt=json&q=trump&rows=15

Negative 86.00%    CanvasJS.com

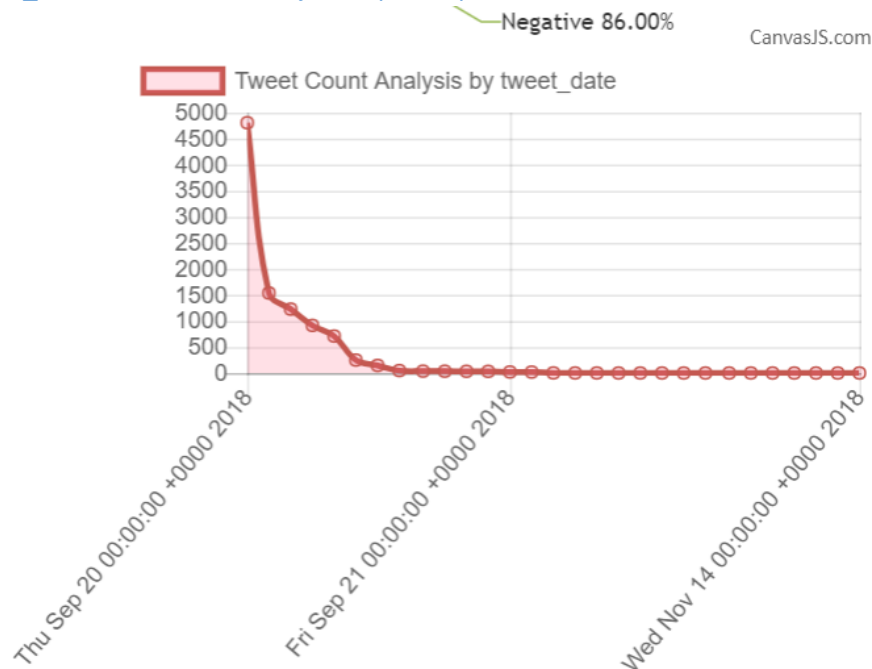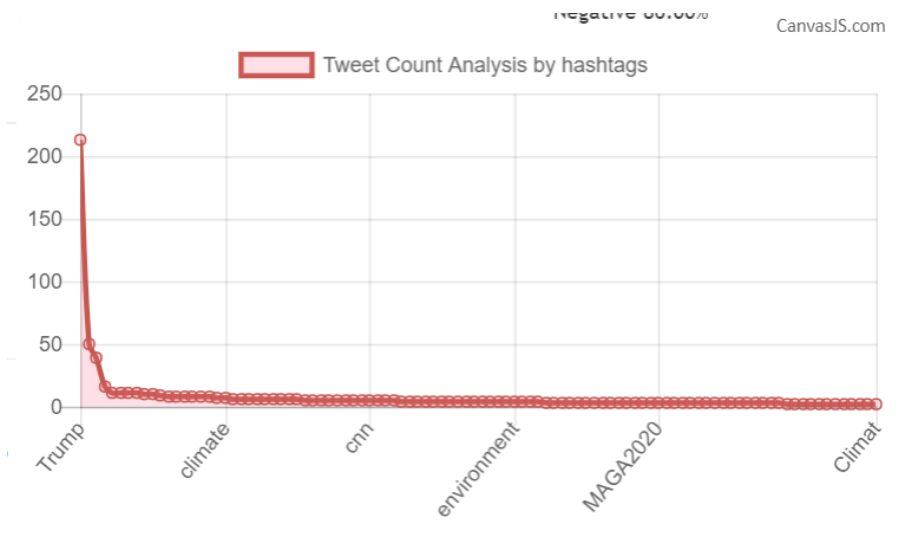Tweet Count Analysis by tweet_date

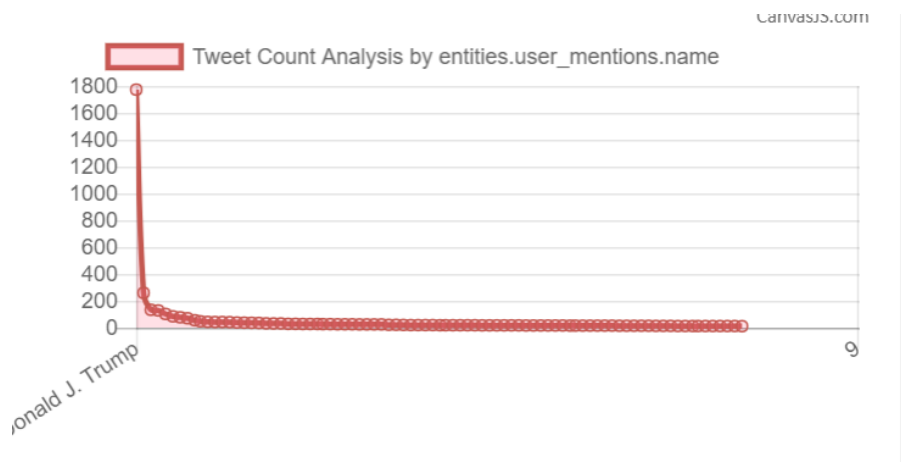Below is the hashtag count graph. This is being accessed using the following url:

http://18.221.162.107:8984/solr/IRF18P4/select?facet.field=hashtags &facet=on&fl=hashtags&indent=on&wt=json&q=trump&rows=15

This is query dependent and we can use filters to get hashtags across the cities, languages and topics.

This is for the mentions, url being used is:

http://18.221.162.107:8984/solr/IRF18P4/select?facet.field=entities.user_mentions.name&facet=on&fl
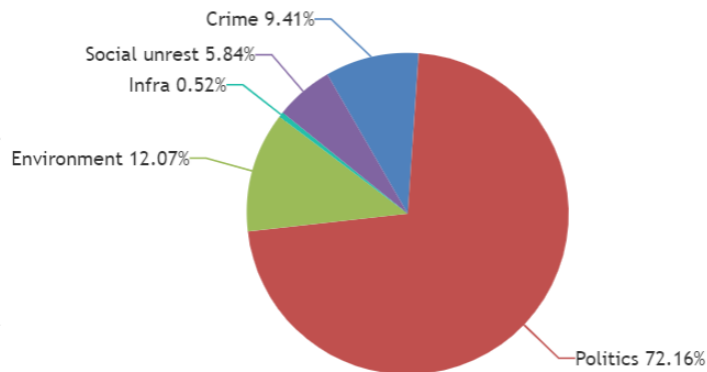= entities.user_mentions.name&indent=on&wt=json&q=trump&rows=15



The below is for the topic and the url is

http://18.221.162.107:8984/solr/IRF18P4/select?facet.field=topic&facet=on&fl=topic&indent=on&wt=j
son&q=trump&rows=15

**topic**

Crime 9.41%
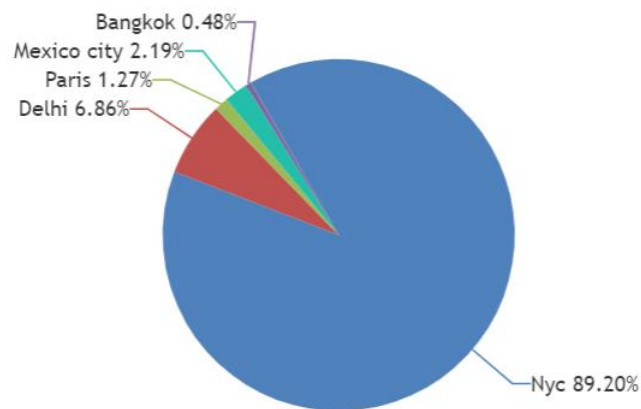Social unrest 5.84%
Infra 0.52%
Environment 12.07%
Politics 72.16%

The below is for the city, the url used is
http://18.221.162.107:8984/solr/IRF18P4/select?facet.field=city&facet=on&fl=city&indent=on&wt=json&q=trump&rows=15
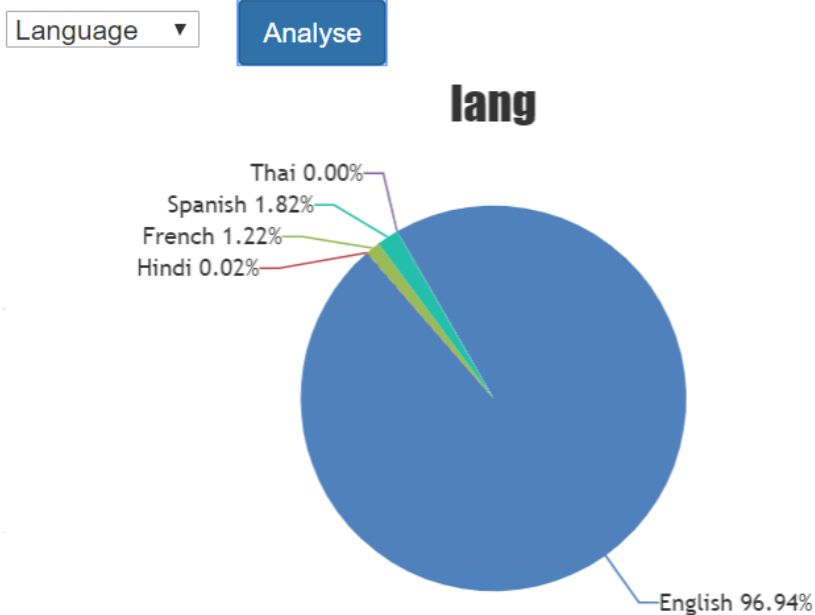


**city**

Bangkok 0.48%
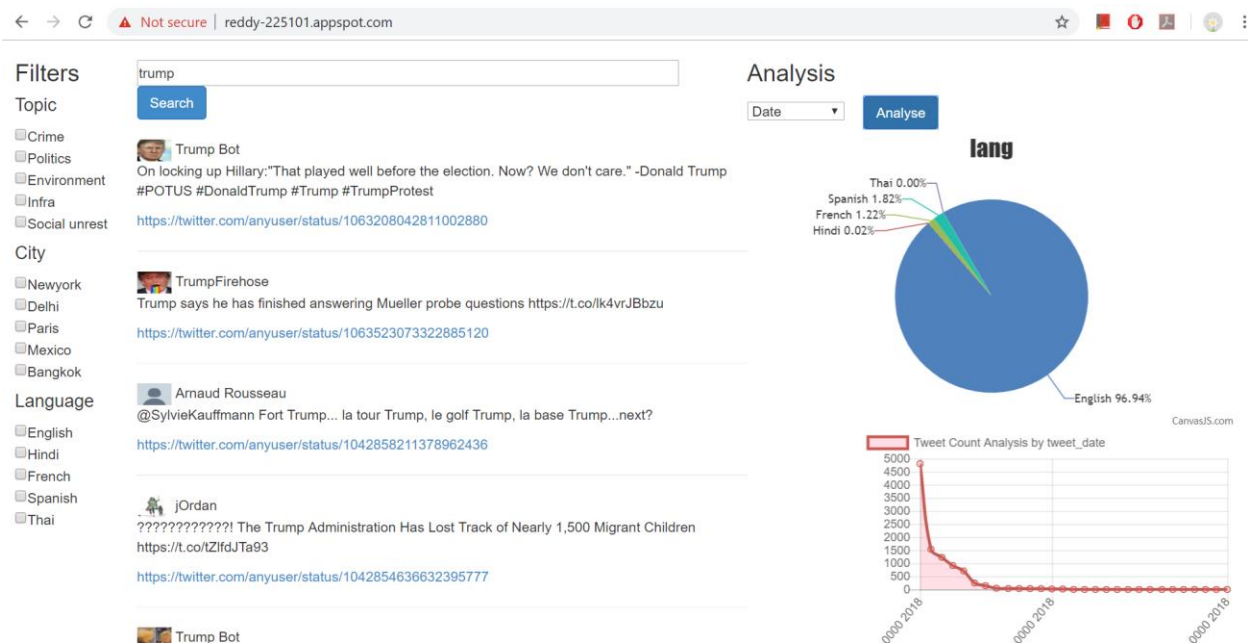Mexico city 2.19%
Paris 1.27%
Delhi 6.86%
Nyc 89.20%

The below is for the language, the url used is:

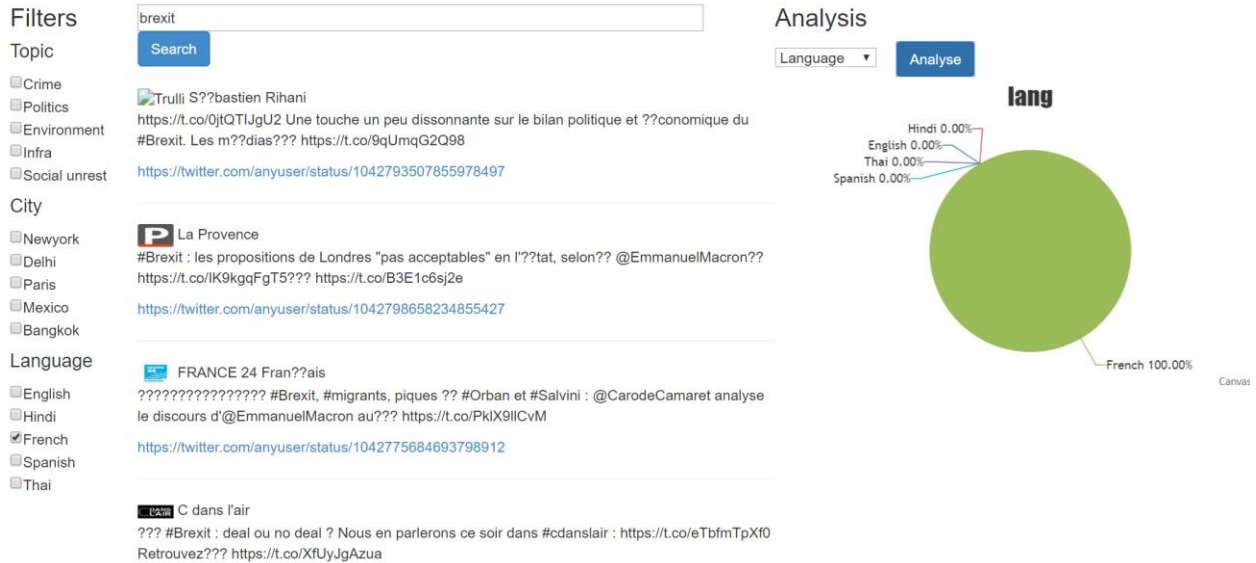http://18.221.162.107:8984/solr/IRF18P4/select?facet.field=lang&facet=on&fl=lang&indent=on&wt=json&q=trump&rows=15
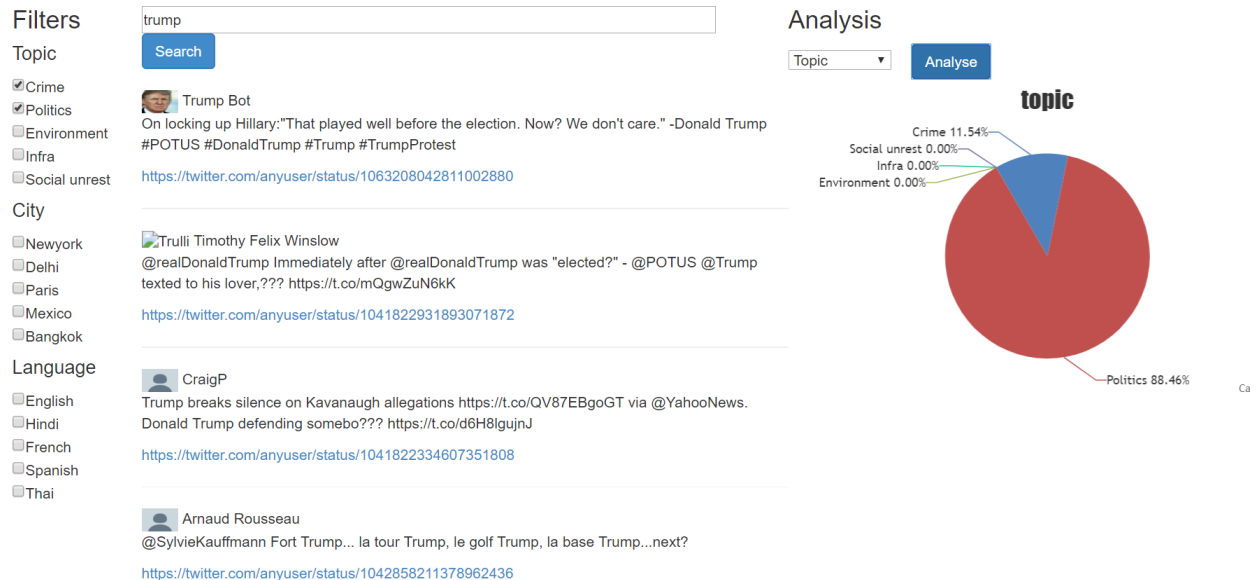
**Filters:**

The below image shows the Complete UI display.



After adding the language filter as French it displays the French language tweets as 100% in the chart:

Here after adding the filter in Topic as Crime and Politics, it retrieves all the tweets and to the right us shows the number of tweets related to each field.



**Video Link:**

Youtube video link: https://youtu.be/aq1fVDbMGYk

Google cloud link to access the project: https://reddy-225101.appspot.com/

**Contributions:**

*Tweets collection, Indexing, UI*: Mkalamda, sanand3.

*Server side hosting and backend web stack implementation*:  aalluri, lveerise, meharaje