



KNOWLEDGE DISCOVERY AND MANAGEMENT

PROJECT

SUMMER 2016

DONE BY:

- 1.SRICHARAN PAMURU(28)
- 2.PRUDHVI RAJ MUDUNURI (22)
- 3.SNEHAL VANTASHALA (41)
4. BHARGAV KRISHNA VELAGAPUDI (42)

## Contents

1. Project Proposal: .....	2
a. Motivation:.....	2
b. Objectives: .....	2
c. Expected Outcomes: .....	2
2. Domain:.....	3
3. Data Collection: .....	4
4. Report:.....	5
a. TF-IDF: .....	5
b. NER:.....	7
5. Design of features .....	13
a. Software Architecture: .....	13
i. Architecture of Model:.....	13
ii. UML Modeling: .....	14
b. Workflow: .....	18
c. Existing Services used: .....	19
d. New feature implemented: .....	20
6. Project Management: .....	21
a. Individual Contribution:.....	21
b. Zenhub/Github Screenshots: .....	22
i. Full timeline: .....	22
ii. Report 1:.....	23
iii. Report 2:.....	23
iv. Report 3:.....	24
v. Report 4:.....	24
vi. Burndown chart:.....	25
vii. Pulse:.....	25
viii. Traffic: .....	26
ix. Commits: .....	26
x. Code Frequency: .....	27
c. Concerns/Issue: .....	28
d. Future Work: .....	28
7. Bibliography: .....	29

## 1. Project Proposal:

### a. Motivation:

The motivation of building a recommendation system for articles is an aspect of providing some kind of intuitive outcome for people who read a specific post or blog. In our case, the thought of binding different sources altogether in order to obtain some sort of output which will yield to dynamic data fed in the form of RSS, has been an ideological aspect that we have considered significant to understanding how the data trends and topics of approach have to be brought together. This combination will enable a user to stream into topics of interest across more than one domain. The machine would be able to return a suggestive article in the list through various technologies that we plan on implementing in this project

### b. Objectives:

The goal of our project is to come up with a recommendation system for articles in different domain areas to bring the user a specific search blog or engine. It will help the user to understand how the scalability and necessity of the article is considered. The page rank would be used as a measure to understand the sequence of recommendations put forward. Based upon this, the articles returned would focus on terms that have a similarity index, matching the score of the previously iterated indices. There would be a relative rank associated that would be used to figure out the magnanimity of the article which would have to be returned.

Our recommendation engine would return the articles to be recommended in a list along with the associated hyperlink. This would help the user view topics of similar interest. Machine learning algorithms run in the background to bring a result yielding a knowledge graph where there is a proper connection between the elements, which in our case is articles, dynamically.

### c. Expected Outcomes:

The outcome of the project would be to bring out articles to be recommended in a list which will help the user view similarity index based articles. Our project is very useful for people who read blogs and articles. It would be an interactive web page where similar index based documents have a value returned which match the preexisting ones.

## 2. Domain:

Our project aims at fetching real time data in the form of RSS(Rich Site Summary) which will be dynamic. The recommendation engine will use collaborative filtering and it serves as a branch of deep learning. It will enable the machine to undergo a deep analysis and filter the data collection that is appropriate for mining. Technology has developed immensely and it is required for the machine to be able to understand each fragment in order to return a better perspective of how the inner mechanism works. Neural networks would basically iterate through every fragment and in our case the articles have to be searched and read through to understand how the term frequency and inverse document frequency matches. (TF-IDF) is a key constraint that will enable information retrieval to be scaled. We would obtain the respective important term that will differentiate a specific document when compared with another.

Recommendation engine is a platform for the machine to be trained and tested. Simulation of such technology actually occurs by rapidly searching the inner context. Big data mining tools focus on scalability and larger data sets are compared in order to retrieve information which is adaptive to the respective target, which is the articles. The NLP tools will help us extract a relation and naming the entities that is crucial to construction of feature vectors. They are combined with machine learning algorithms and tools to form an ontology. The ontologies learned will altogether contribute in constructing the recommendation system, which is the ultimate motivation of our project.

Moreover, the information retrieval can be characteristic of the artificial intelligence based approach that occurs in the learning of in gram to retrieve a continuous set of words which will bind the data and integrate the search. The recommendation system would then focus on setting up the link between the APIs and return the set of article.

### 3. Data Collection:

The data collection for this specific increment has been generated from the BBC news sites where the various articles sections have been formed. It has helped us understand how the articles are generated and their category mapping.

Sample article under Sports Category:



South Africa far too strong again

Second one-day international, Durban South Africa 329-6 beat Zimbabwe 198-7 by 131 runs

South Africa lead three-match series 2-0

Smith's 117 set up the best total at Durban in 25 one-day internationals. He was given fine support by

The same bowler struck soon afterwards as Rogers clipped a fullish delivery to wide mid-on. Nicky Boje

The spinners had to do plenty of work as a result and one of them, Prosper Utseya, went for 60 in six

Graeme Smith (captain), Adam Bacher, Jacques Rudolph, Herschelle Gibbs, Ashwell Prince, Justin Kemp, Stuart Matsikenyeri, Barney Rogers, Hamilton Masakadza, Alester Maregwede, Brendan Taylor, Tatenda Ta

Similarly, the data collection is static right now and we have used the BBC datasets as a training set. The overall outcome would be to dynamically obtain data in the form of RSS feed. The RSS feed would use the standard web feed formats in order to publish updates of information that can be categorized as blogs, headlines and also includes audio, video files. The metadata would explain about the specific data content and focus on publish dates, time and author information.

We intend to pull the required information and accordingly sort it as per the user's search. Dynamically, it would be feasible to categorize it based on the XML file format generated. The RSS reader would check how the user checks feed for information and the user interface will serve as a front-end for interpreting such information. The back end would include the process of fetching data and perform the NLP operations to understand the grammar and contextual meaning.

URL: <http://mlg.ucd.ie/datasets/bbc.html>



This was followed by determination of the term frequency and inverse document frequency as a measure of how the most important terms are obtained. The top terms are taken and there is an integration of top constraints to obtain a specific link.

```
((Aus),,24.78577651117808)
((Ind),,24.78577651117808)
(popular,23.63504822137096)
(Steyn,21.496620714543194)
(ruling,19.805727707708055)
(Al,19.805727707708055)
(tax,19.805727707708055)
(Division,19.805727707708055)
(overseas,18.58933238338356)
((SL),,18.58933238338356)
(Van,18.58933238338356)
(breaks,18.58933238338356)
(Sri,17.78693399089631)
(evidence,17.726286166028217)
(US,17.726286166028217)
(Cronje,17.726286166028217)
(fours,17.197296571634556)
(Rae,17.05685551208559)
(most,16.816055852417257)
(there's,16.509890841703726)
```

The files used here are the cricket directories that contain content and the top most words are filtered to return a set of terms. The words are returned. Stop words have been filtered and accordingly the TF-IDF table can be executed. The code is scalable

## 2. NER:

```
Gilchrist-8
NER
Gilchrist-8:PERSON

,-9

S-10

Katich-11
NER
Katich-11:PERSON

,-12

D-13

Martyn-14
NER
Martyn-14:PERSON

,-15

M-16

Clarke-17
NER
Clarke-17:PERSON

,-18

M-19

Hussey-20
NER
Hussey-20:LOCATION
```

```
30-year-old-2
NER
30-year-old-2:DURATION

will-3

move-4

to-5

England-6
NER
England-6:LOCATION

under-7

the-8

Kolpak-9
NER
Kolpak-9:PERSON

ruling-10
```

The Named Entity recognition module was run to return the specific naming conventions. The knowledge graph is extracted based on this and we can show the interlinking of an article through this. The conventional data format is used to show how the data can be superimposed can be seen as the manner for showing the format and class to which it belongs.



## II. Project Report 2

This increment focuses on detecting the topic of the document using LDA.

The recommendation systems have to be accurate on suggesting the content to the user. We have implemented the topic detection in our project as the recommendation system we build is going to suggest the user with online articles of his interest we have to collect the user interest based on the topics of the articles the user is reading.

Generally the another functioning of the recommendation system can be like recommending the articles with highest click through rate but the article is of no interest to the user. To avoid such un required suggestions in our model we have included the LDA to detect the user interest and get the model up and running such that all the recommendations emitted as the output will be of use to the user.

Machine learning algorithms such as collaborative filtering uses the ratings given by the user to find the similarity of choices and uses them as the knowledge base to recommend the user different online articles. In some cases there are latent values like one article which no one has ever rated before. This article might be a good article but no one has visited it till now. In collaborative filtering this article has the least chance of popping up as a recommendation.

Using our approach as discussed will populate the articles using collaborative filtering but with additional layer of filtering using the LDA for the topic detection. This way the user is given suggestions which contain articles with the same topic of interest and all the suggestions will be of the same topic.

The only drawback which we are considering is that, if all the suggestions belong to the same topic then we are just hindering the user from discovering other interests which is also damaging for the recommendation systems as user preferences keep changing from time to time.

In-efficient recommendation systems neglect the topic discovery and go with only collaborative filtering.

1. Results: -

```
----- Synonyms for Top words -----  
safa :  
[rahdi,0.9467719924022248]  
[slovaks,0.9141268978018406]  
[bridge,0.9137033667642116]  
[costina,0.9067222792972934]  
  
contepomi :  
[slot,0.9841768311248991]  
[shine,0.9484465903081952]  
[mignonus,0.9429226887138443]  
[simmer,0.9387021785030112]  
  
uganda :  
[road,0.959459576287833]  
[quit,0.9512214074110048]  
[beyond,0.9468891713055234]  
[iran,0.9435937062022871]  
  
fai :  
[management,0.9576998600685174]  
[southampton,0.9337444695776065]  
[grounds,0.9326236220442882]  
[lifelong,0.932017878375858]  
  
faster :  
[case,0.9873011790399964]  
[negotiate,0.9768225542234494]  
[small,0.9763763199839935]  
[hear,0.9731447115463905]
```

```
ind :
[anil,0.9929668927396992]
[earnshaw,0.9834188295976105]
[pak,0.9831599765827246]
[cockbain,0.9811122582708479]

newry :
[translation,0.9417721458200161]
[overlap,0.9341212573864177]
[injury-hit,0.9254891460342861]
[aftenposten,0.9206342292858737]

mirza :
[shortstop,0.9170902002403842]
[howley,0.9120305731203591]
[petulance,0.8972193056254024]
[cordoba,0.8956154094307746]

pennant :
[solskjaer,0.9784773923863354]
[well-known,0.9772468422933377]
[eriksson,0.9697201045302642]
[helguera,0.9660812391474295]

dechy :
[donny,0.9250336145732326]
[artell,0.9227017589180909]
[dusty,0.90607254643031]
[benhables,0.8956187244843185]
```

The specific topics are mined together and the frequency can be determined. This will help enhance the similarity and observe patterns which provide an interlinking of features. Feature vectors observe the relation. The lexical dictionary that has been used here is WordNet which provides correlation and the specific meaning is returned. The synonyms and exact referencing has been done over here. The synonyms provide us an insight into the overall mechanism of the algorithm. The words are searched through the lexical frame and the observation of top words will return the output associated with TF-IDF processing.

```

nce Dallaglio 's team have now get it all to do in the quest for a quarter-final place give tha
nce Dallaglio 's team have now get it all to do in the quest for a quarter-final place give tha
nce Dallaglio 's team have now get it all to do in the quest for a quarter-final place give tha
nce Dallaglio 's team have now get it all to do in the quest for a quarter-final place give tha
nce Dallaglio 's team have now get it all to do in the quest for a quarter-final place give tha
nother setback arrive in 2001 as a serious knee injury cut short Dallaglio 's involvement on th
rs begin to circulate that he career be over but , in typical Dallaglio style , he embark on a
rs begin to circulate that he career be over but , in typical Dallaglio style , he embark on a
isbergen ; Lewsey , Erinle , Abbott , Voyce ; King , Dawson ; Dowd , Greening , Green ; Shaw ,
cement : gotting , McKenzie , Lock , Hart , Biljon , Brooks , Hoadley .
y ; Rabeni , Smith , Gibson , Healey ; Goode , Ellis ; Rowntree , Chuter , White , m Johnson -l
nly player to play every minute of England 's World Cup triumph in Australia , Dallaglio could
nly player to play every minute of England 's World Cup triumph in Australia , Dallaglio could
nly player to play every minute of England 's World Cup triumph in Australia , Dallaglio could
nly player to play every minute of England 's World Cup triumph in Australia , Dallaglio could
cement -lr- from -rr- : Buckland/Cockerill , Morris , Kay , W Johnson/B Deacon , H Tuilagi ,
roblem now for England be how to replace the almost irreplaceable .
roblem now for England be how to replace the almost irreplaceable . the likes of Matt Dawson ,
roblem now for England be how to replace the almost irreplaceable . the likes of Matt Dawson ,

```

#### Corpus summary:

```

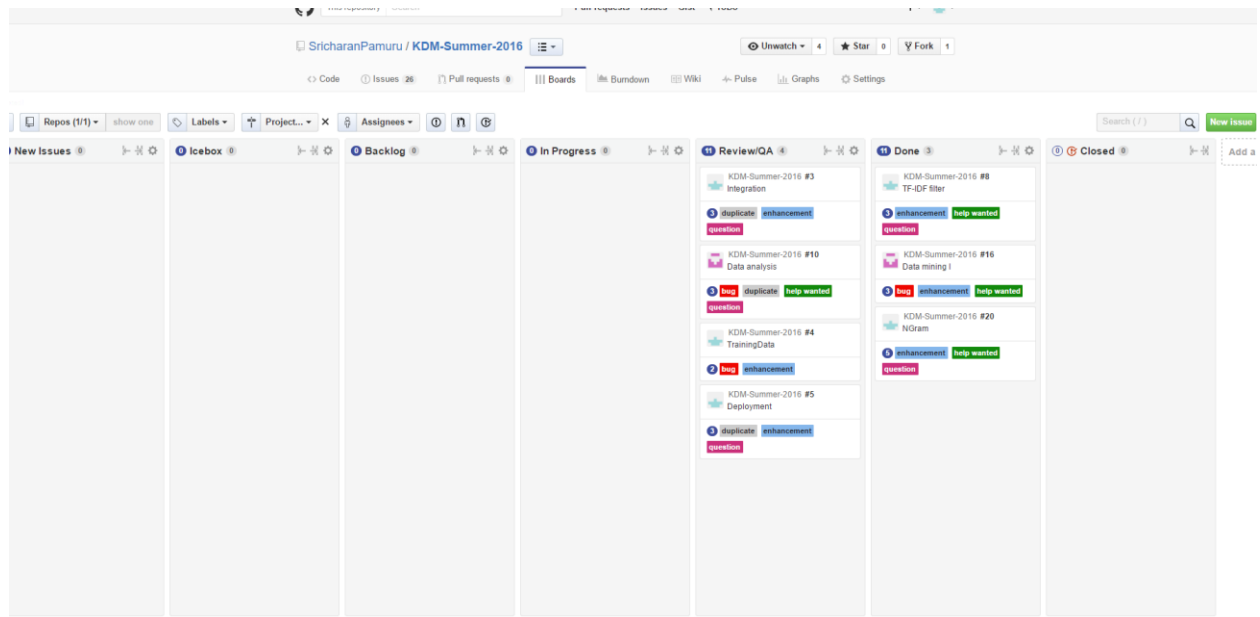
Training set size: 7485 documents
Vocabulary size: 12040 terms
Training set size: 122024 tokens
Preprocessing time: 145.254931418 sec

```

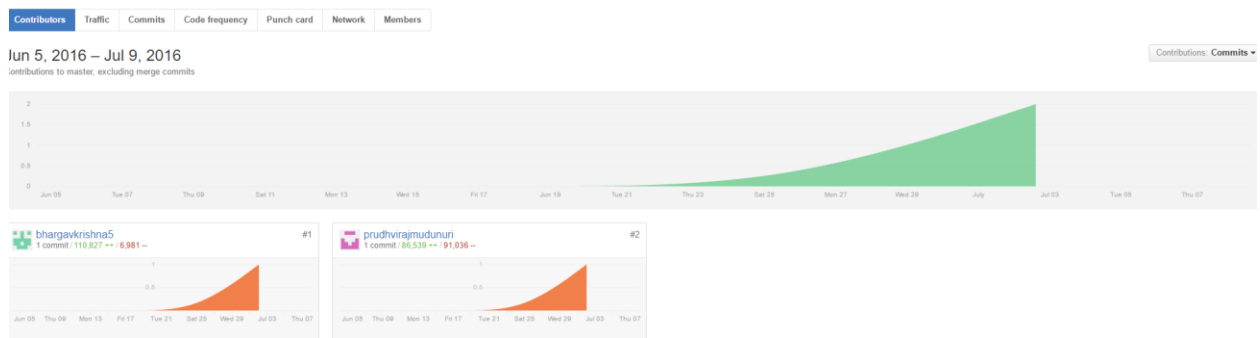
We have use Word2vec to get the synonyms for the top-idf words so that the synonyms can be matched against the topics obtained from LDA and suggest the user with articles which also contain the synonyms of the user interest like if the user likes soccer and is interested in reading articles on soccer then soccer is called football in other countries. This issue can be resolved using the synonyms.

## 2. ZenHub screenshots:

### a. Issues of Project development:



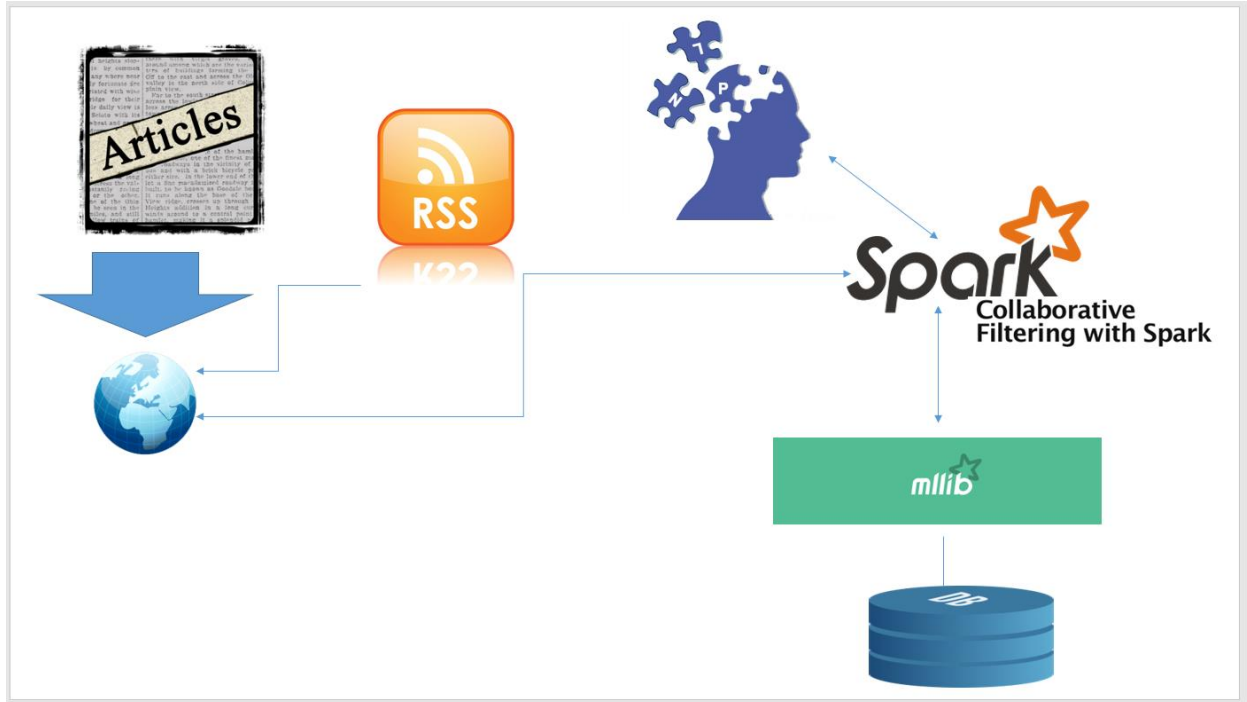
### b. Contribution:



## 5. Design of features

### a. Software Architecture:

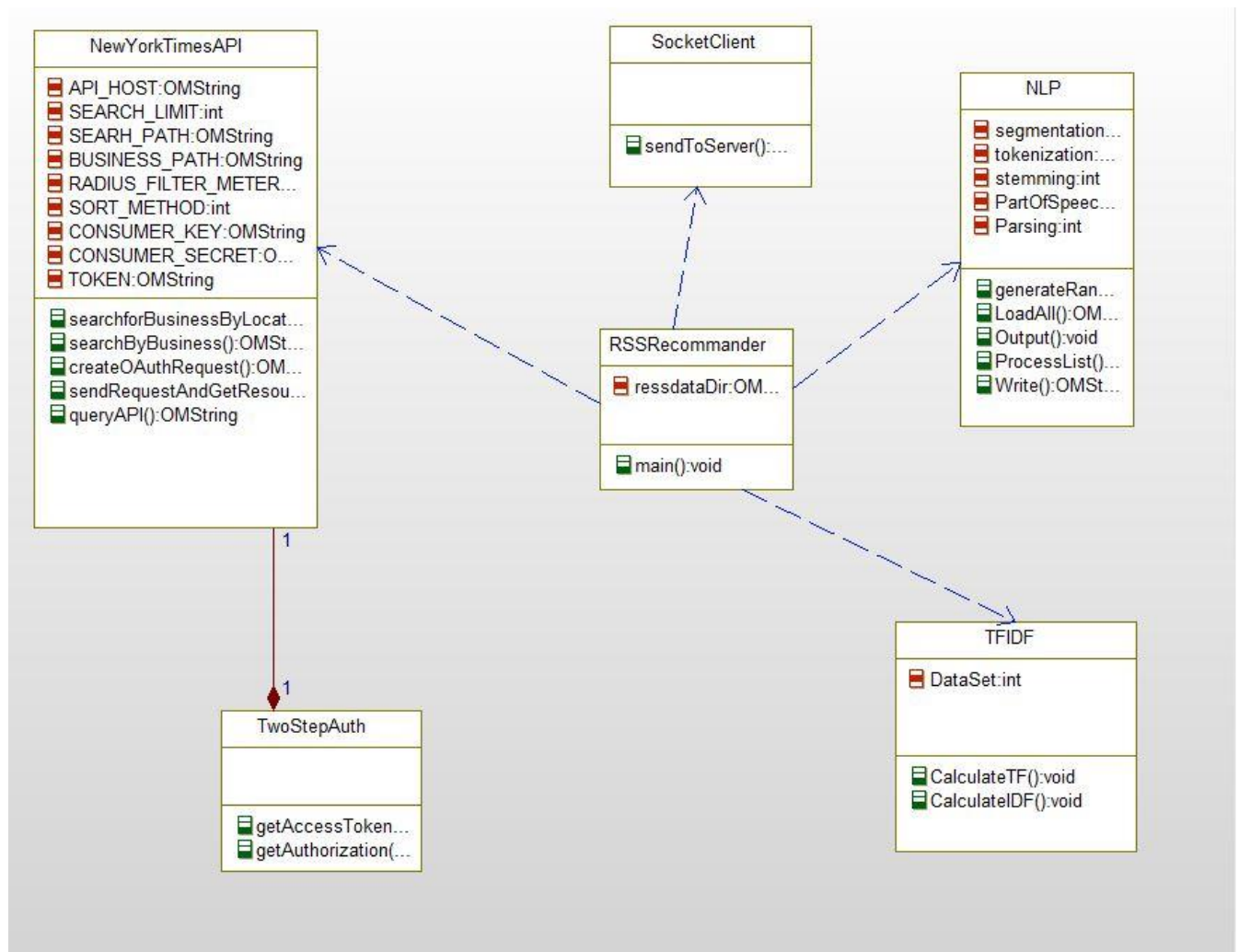
#### i. Architecture of Model:



The Articles are brought in online and the different RSS feed is present. The web would then combine with Spark where the collaborative filtering and NLP techniques are processed. The MLLib would serve as the library where the different algorithms are run. Machine learning techniques are being processed to indicate how the NLP parser are used. The articles have their key terms obtained and a recommendation is sent based on the user's search. This recommendation system serves as the main output for the project.

## ii. UML Modeling:

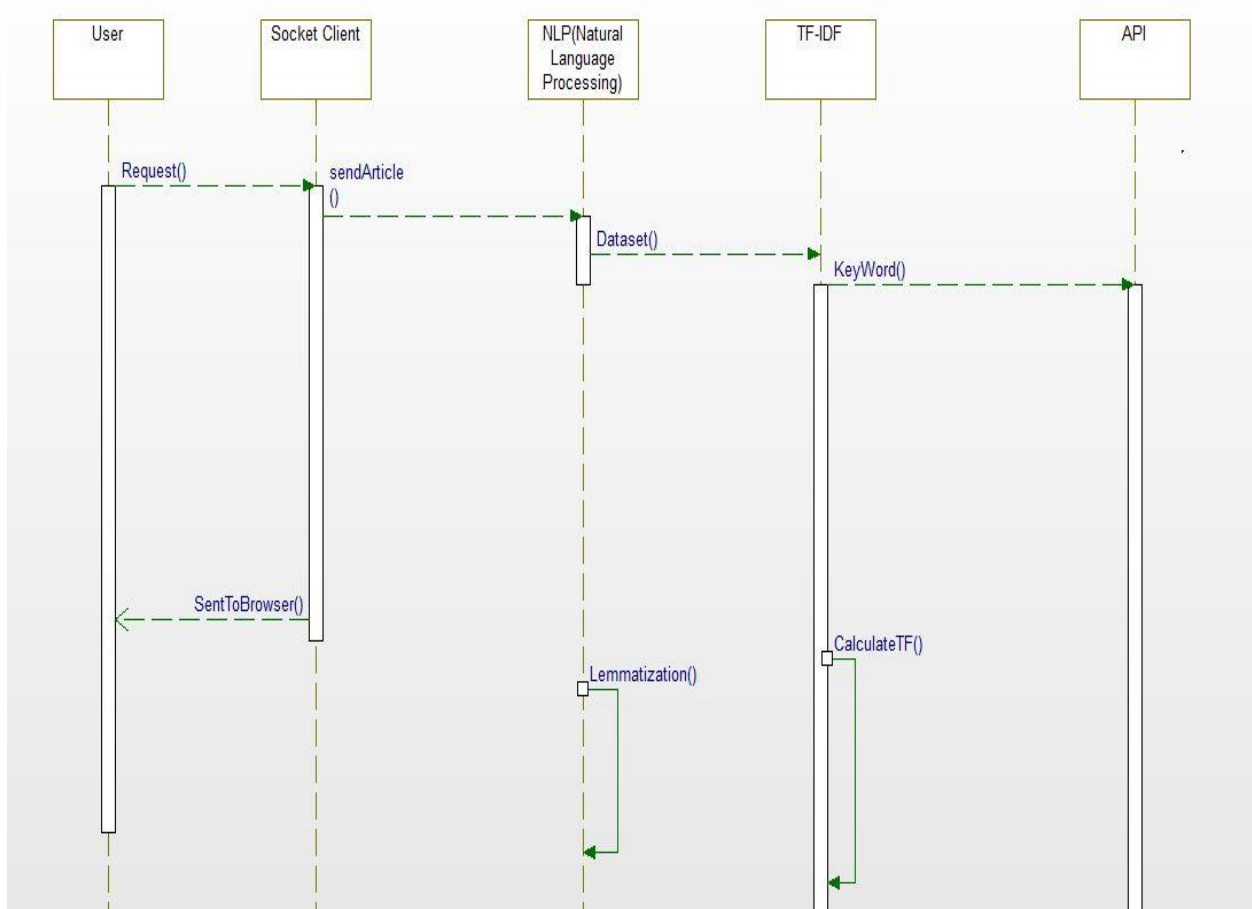
### 1. Class Diagram:



This diagram showcases the different classes that we have used in our program as a key component of the execution and the workflow can be depicted that will serve as base for the knowledge graph eventually processed.

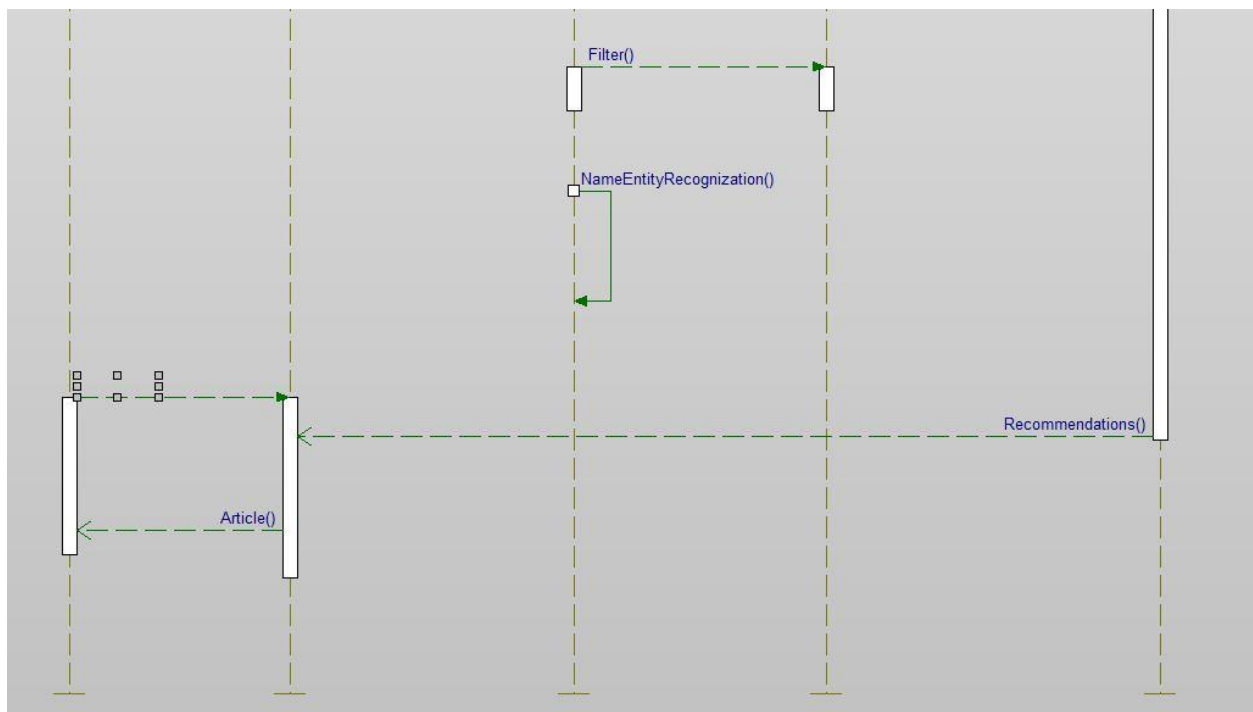
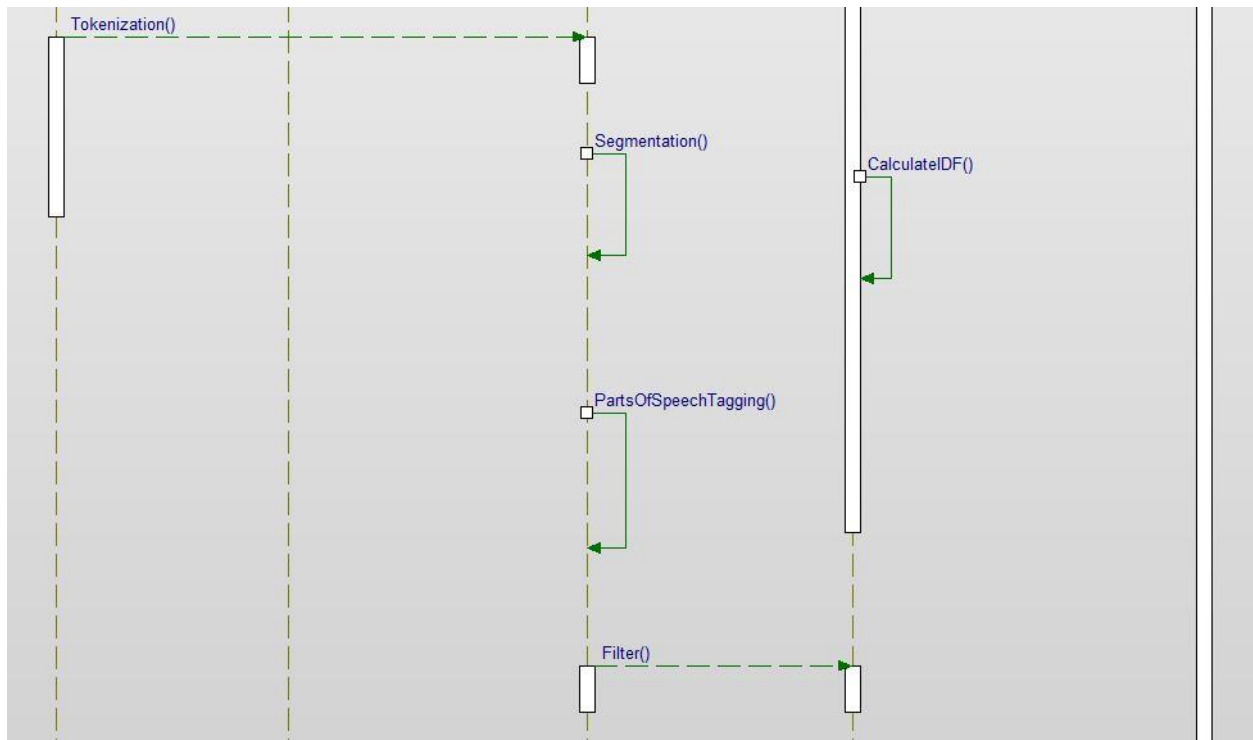
## 2. Sequence Diagram:

The sequence of execution for each of the classes is shown below. The user would use the present client side application to read articles. The client side program would then perform the various NLP operations on it to filter the text and observe a relative extraction mechanism.



The TF-IDF runs and the keywords are sent to the API. The NLP processes would include lemmatization, Named Entity recognition and stemming that forms the key part of our project. The fragment of code would then be run to check for the relative frequency of words present in the current article. After the program is executed, the recommendation of articles is returned based on process such as collaborative filtering and information retrieval occurs. The entire portion of the code would be executed dynamically using RSS feed.

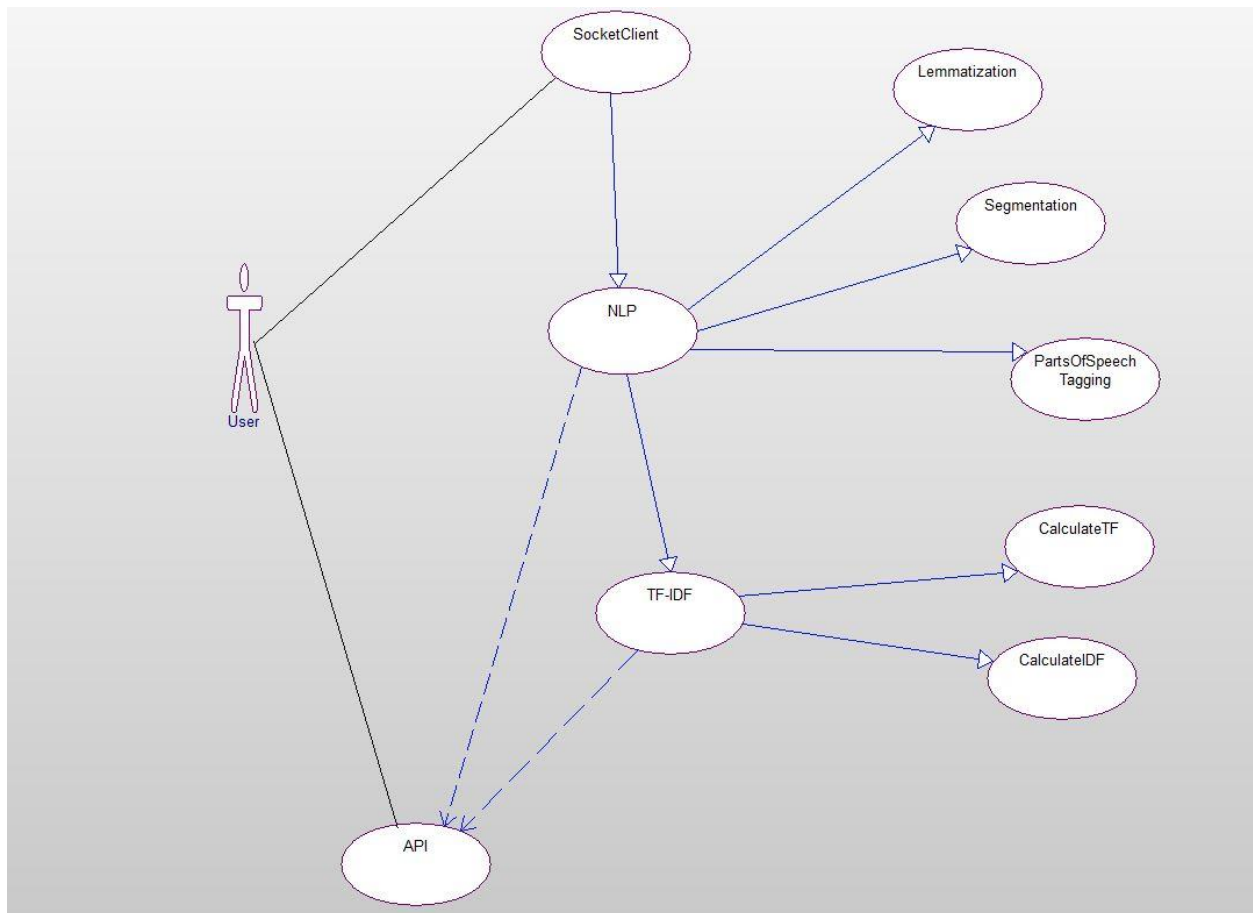




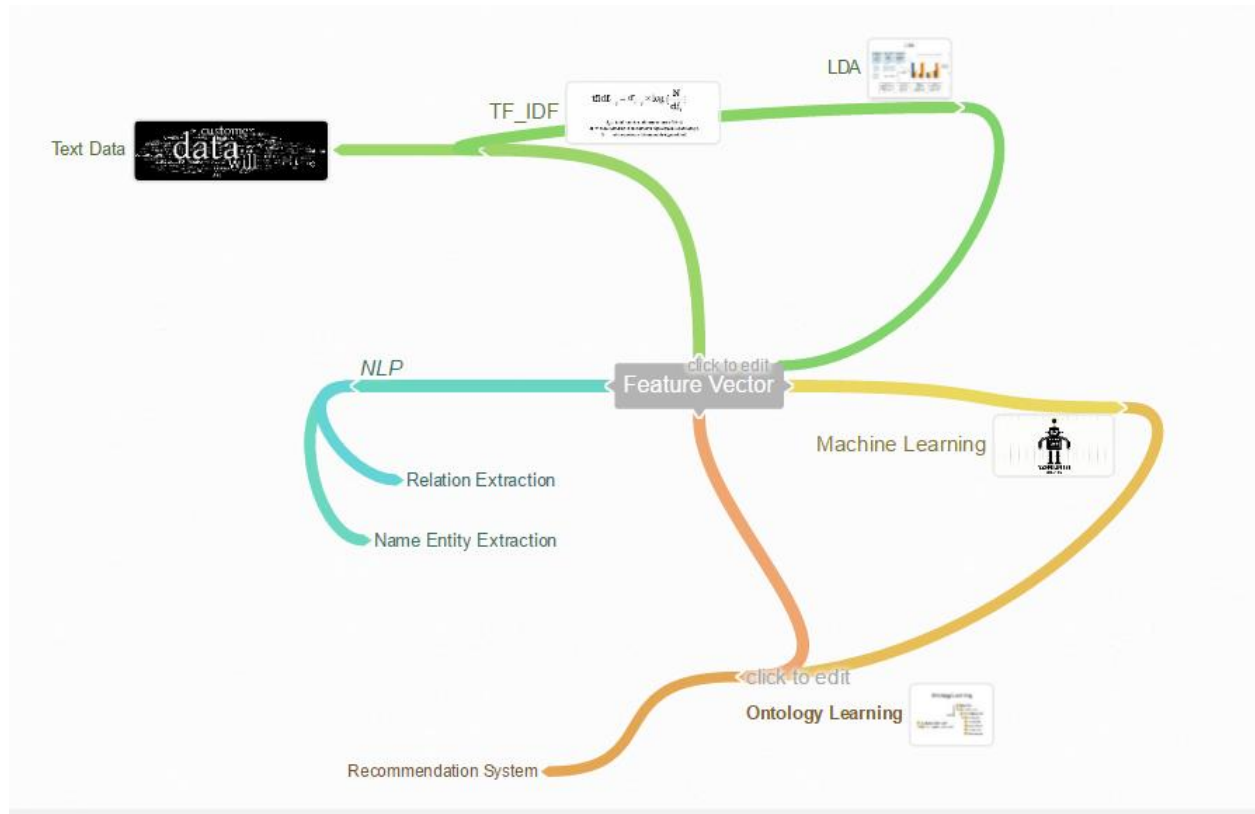
The entire sequence would then return the server side run program which would contain the list of articles which would navigate as per the machine learning algorithms.

### 3. Use Case Diagram:

The use case diagram would then show how the user is linked through the various NLP operations which will do the different API based handling where the client program will send a request to the server and focus on handling the various terms that can bring about a certain kind of linkage to draw the knowledge graph that will bind the different APIs. The final recommendation engine is built based on the machine learning algorithms that are executed.



## b. Workflow:



The user visits an article once then data from the article is extracted and send for NLP processing where different Computational Linguistics are carried out. Data is given to the NLP goes through sentence segmentation; Tokenization; Stemming; Part-of-speech tagging; parsing; Named Entity Recognition; Co-reference Resolution. Data is passed further to tf-idf class where the information retrieval process is done based on the cosine similarity and words that provide meaning to the document are retrieved and sent to the New York Article search API. The New York Article search API which takes the keywords from the tf-idf process and these keywords are given as the input to the API as a query. API writes top articles and their hits in different social networking sites.

### c. Existing Services used:

New York Times articles search API enables users to search New York times articles from 1851. Using this api users/ developers can extract data based on keyword, it also offers facet searching(multidimensional). Information retrieval contains headline, lead paragraphs, related docs, abstract, time specific details can also be retrieved etc.

The screenshot displays the Swagger UI for the 'Article Search API'. At the top, there are links for 'README', 'Documentation', and 'Console'. Below these, a dropdown menu shows 'All' and a selected endpoint 'GET /articlesearch.json' with a 'View Results' button. The main area is divided into two panels. The left panel, titled 'Sample Code', shows a JavaScript code snippet for making an AJAX request to the API endpoint, including an API key and a search query 'computer science'. Below the code is a 'Parameters' section with input fields for 'q' (containing 'computer science') and 'fq'. The right panel displays the JSON response of the API call, which includes details about the search results, such as the publication date, document type, news desk, section name, and a list of authors with their roles and names.

```
Article Search API Source: [Swagger 2.0] README Documentation Console
```

▼ All GET /articlesearch.json View Results ↻

#### Sample Code

JavaScript NodeJS PHP Ruby

```
// Built by LucyBot. www.lucybot.com
var url = "https://api.nytimes.com/svc/search/v2/articlesearch.json";
url += '?' + $.param({
  'api-key': "5ef8269c555c4fc091eb00c6bb94484a",
  'q': "computer science"
});
$.ajax({
  url: url,
  method: 'GET',
}).done(function(result) {
  console.log(result);
}).fail(function(err) {
  throw err;
});
```

#### Parameters

q  
computer science

fq

```
{
  "value": "Office of Science and Technology Policy",
},
{
  "rank": "6",
  "is_major": "N",
  "name": "organizations",
  "value": "University of Washington"
},
{
  "rank": "7",
  "is_major": "N",
  "name": "glocations",
  "value": "Seattle (Wash)"
}
],
"pub_date": "2016-05-26T00:00:00Z",
"document_type": "article",
"news_desk": "Business",
"section_name": "Technology",
"subsection_name": null,
"byline": {
  "person": [
    {
      "organization": "",
      "role": "reported",
      "firstname": "John",
      "rank": 1,
      "lastname": "MARKOFF"
    }
  ]
}
```

Field	Behavior
body	Multiple tokens
body.search	Left-edge n-grams
creative_works	Single token
creative_works.contains	Multiple tokens
day_of_week	Single token
document_type	Case-sensitive exact match
glocations	Single token
glocations.contains	Multiple tokens
headline	Multiple tokens
headline.search	Left-edge n-grams
kicker	Single token
kicker.contains	Multiple tokens
news_desk	Single token
news_desk.contains	Multiple tokens
organizations	Single token
organizations.contains	Multiple tokens
persons	Single token
persons.contains	Multiple tokens
pub_date	Timestamp (YYYY-MM-DD)
pub_year	Integer

d. **New feature implemented:**

We Intend to develop a content based recommendation system based on collaborative filtering and ranking the recommendation with scores. Our developed model can tackle cold start. Cold start is a state of the machine when there are no previous train data. We intend to get the ontology graph from the content through machine learning algorithms and recommend based on the first article the user reads. Available recommendation systems struggle during the cold start phase.

Our developed model takes user interest to train the model by asking the user to rate the rate the article before he leaves. If the user has given a negative feed back then then the particular article will not show up in other recommendation for that particular user even when the content matches his query.

6. Project Management:

a. Individual Contribution:

Name (CLASS ID)	Contribution	Number of hours
<b>SRICHARAN PAMURU (20)</b>	Word2Vec and Feature Vector extraction	20
<b>PRUDHVI RAJ MUDUNURI (22)</b>	LDA and WordNet	20
<b>SNEHAL VANTASHALA (41)</b>	NGram and WordNet	20
<b>BHARGAV KRISHNA VELAGAPUDI (42)</b>	Word2Vec and NGram	20

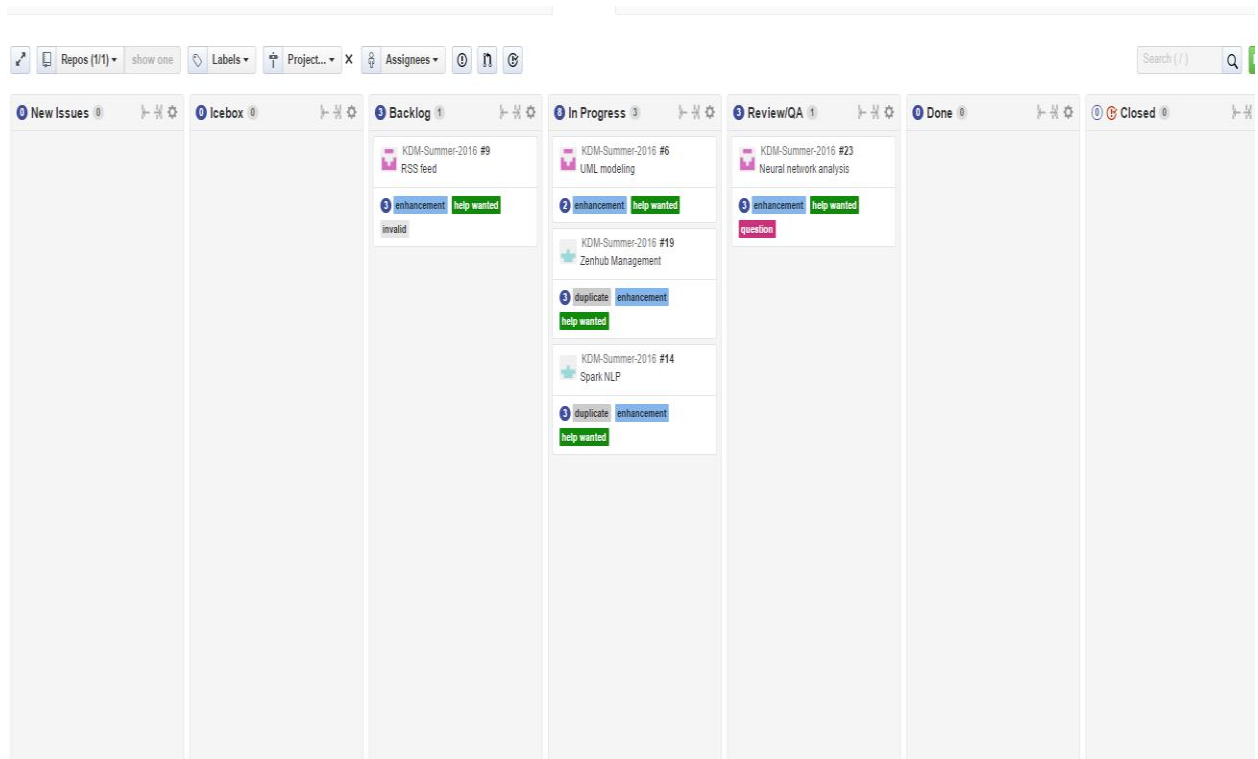
## b. Zenhub/Github Screenshots:

### i. Full timeline:

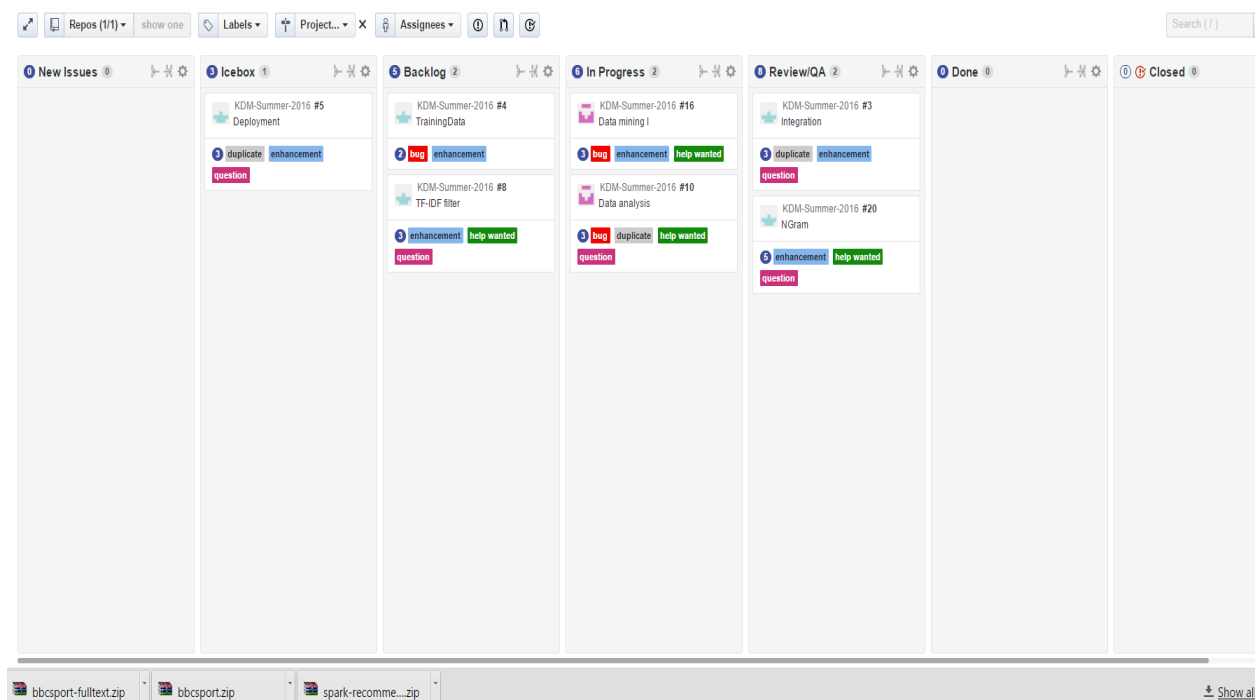
This screenshot displays a GitHub repository page for 'SricharanPamuru / KDM-Summer-2016'. The interface includes a top navigation bar with links for Pull requests, Issues, Gist, and a To Do list. Below the repository name, there are statistics for Unwatch (2), Star (0), and Fork (0). The main content area features a Kanban board with columns for New Issues (0), Icebox (4), Backlog (6), In Progress (7), Review/QA (8), Done (1), and Closed (0). Each column contains a list of issues, each with a title, a status label (e.g., bug, enhancement, question), and a 'help wanted' tag. The issues are organized into a timeline view, showing the progression of work from the Icebox to the Done column.

Column	Issue ID	Title	Status	Help Wanted
Icebox	#2	Testing	bug	help wanted
	#5	Deployment	duplicate	enhancement
	#17	Data Mining II	enhancement	help wanted
	#11	Visualization	duplicate	enhancement
Backlog	#9	RSS feed	enhancement	help wanted
	#4	TrainingData	invalid	
	#8	TF-IDF filter	enhancement	help wanted
	#13	Web model	enhancement	help wanted
	#12	Final code	bug	enhancement
	#25	Project Report	enhancement	help wanted
In Progress		UML modeling	enhancement	help wanted
	#16	Data mining I	bug	enhancement
	#19	Zenhub Management	duplicate	enhancement
	#10	Data analysis	bug	duplicate
	#14	Spark NLP	duplicate	enhancement
	#21	Information Retrieval	enhancement	help wanted
			enhancement	help wanted
Review/QA	#3	Integration	duplicate	enhancement
	#7	Collaborative filtering	bug	enhancement
	#23	Neural network analysis	enhancement	help wanted
	#20	NGram	enhancement	help wanted
	#15	Knowledge Graph	enhancement	help wanted
	#22	Deep Learning	enhancement	help wanted
Done	#1	Project Proposal	enhancement	help wanted

ii. Report 1:

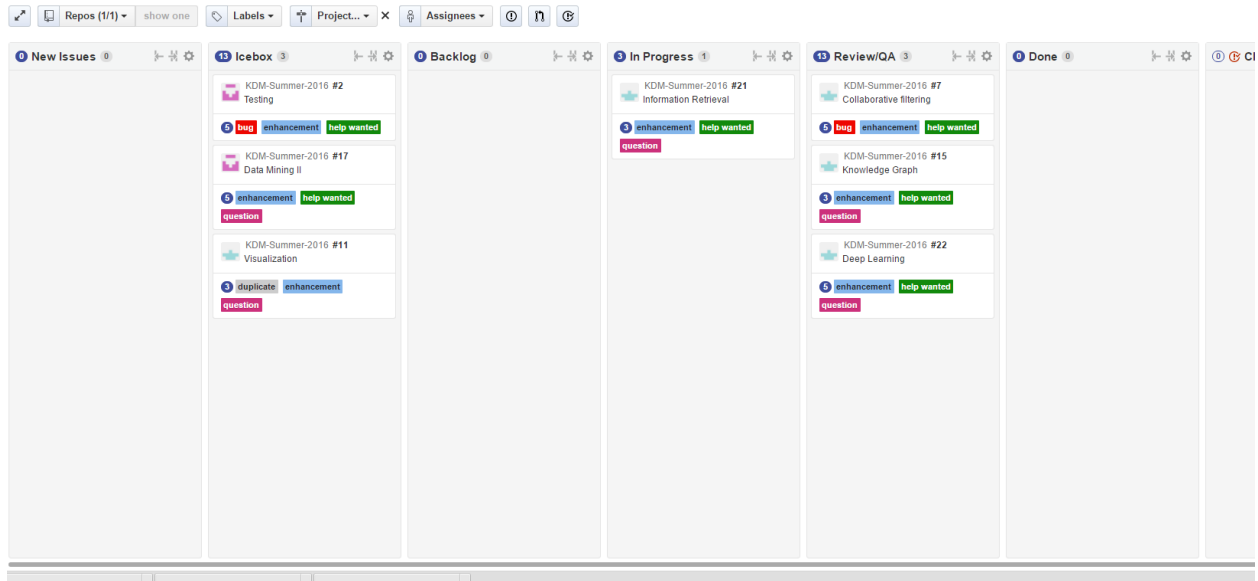


iii. Report 2:

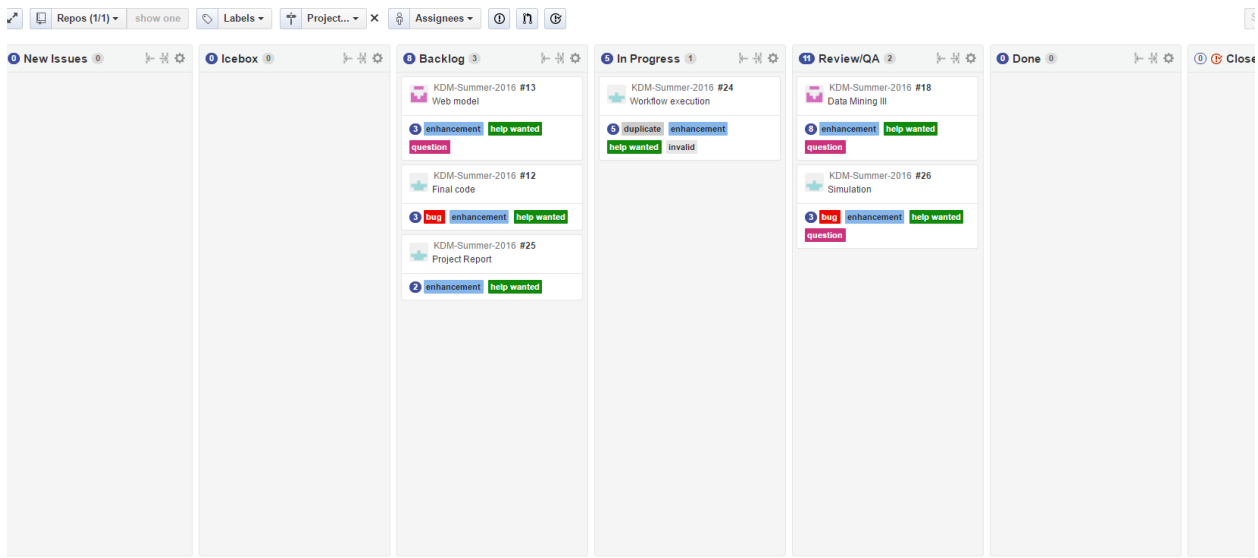




#### iv. Report 3:



#### v. Report 4:



vi. Burndown chart:



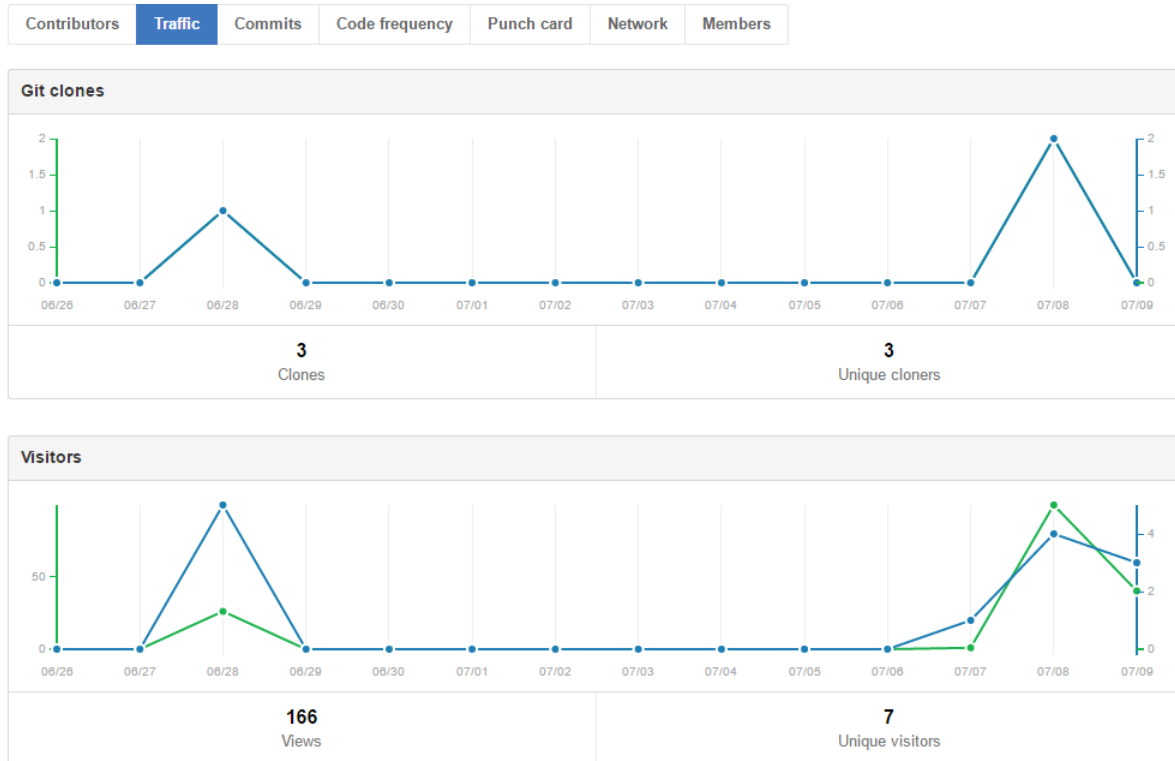
vii. Pulse:

June 17, 2016 – June 24, 2016

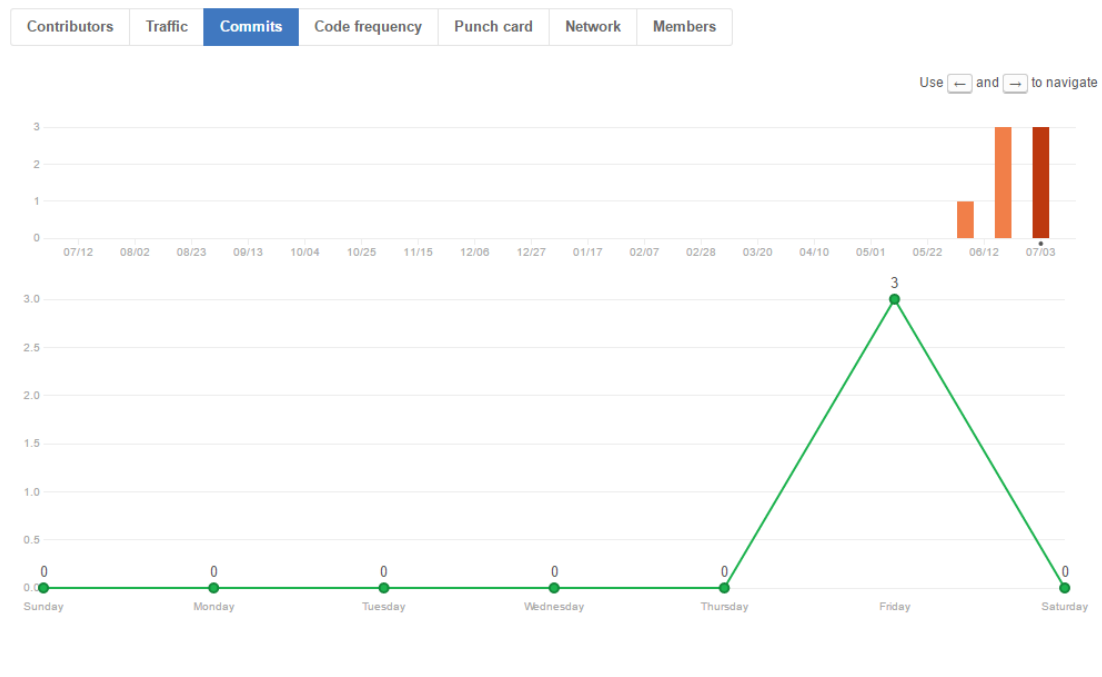
Period: 1 week

Overview			
0 Active Pull Requests		27 Active Issues	
Merged Pull Requests 0	Proposed Pull Requests 0	Closed Issue 1	New Issues 26

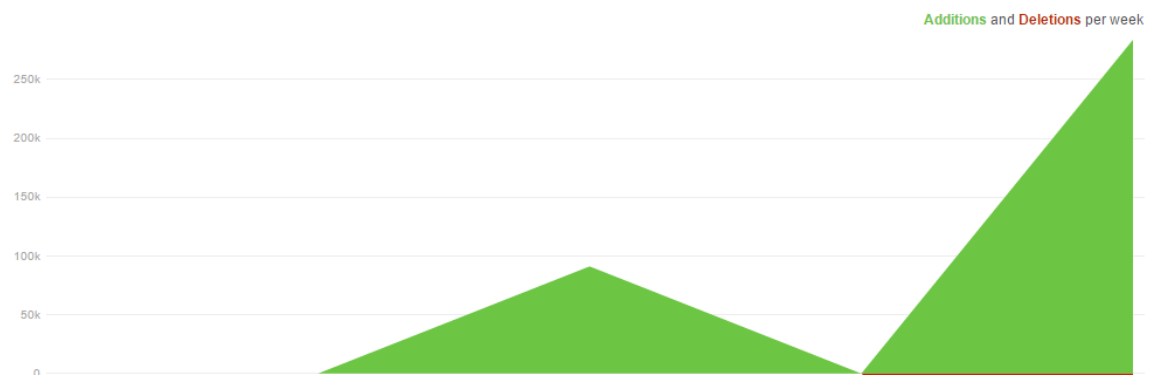
## viii. Traffic:



## ix. Commits:



x. Code Frequency:



#### c. Concerns/Issue:

Our project can face issues when it comes to deployment dynamically as we have to obtain the RSS feed which serves as the main constraint to show the update of the feed data. Moreover, the program would serve differently for various users to suit his preferences.

The machine would not be able to scale exponentially and provide incremental gains over due course. It would limit the user's self search mechanism as the recommendations are showed based on preferences.

The time taken for the training is more and hence the machine learning algorithms need to be trained first in order to run on the system. User's preferences can also be obtained through the union and intersection of various terms.

The Cold Start Problem is a major issue in recommendation systems.

#### d. Future Work:

The project can be expanded in various domains individually such as news articles, journals, scientific data and there are still certain areas of research. Efficient algorithms can be computed. Probably, the cold start problem is something that future researchers would look into so that there is no issue when it comes to certain issues existing over time.

## 7. Bibliography:

1. Spangher, Alexander. "Building the Next New York Times Recommendation Engine." Open Blog. New York Times, 11 Aug. 2015. Web. 25 June 2016.
2. Jain, Aarshay. "Quick Guide to Build a Recommendation Engine in Python." Analytics Vidhya. N.p., 02 June 2016. Web. 25 June 2016.
3. Ridwan, Mahmud. "Predicting Likes: Inside A Simple Recommendation Engine's Algorithms." Toptal Engineering Blog. N.p., 12 Mar. 2015. Web. 25 June 2016
4. Filloux, Frederic. "DeepMind Could Bring The Best News Recommendation Engine." Medium. Monday Note, 20 Mar. 2016. Web. 25 June 2016.