# Report on Search Engine(Information Retrieval)

Raj Bhatta
University of Windsor
401 Sunset Avenue
Windsor,Ontario
bhatt11y@uwindsor.ca

## ABSTRACT

This report describes the details about Search Engine that we have developed. During the development of our project, we have started from basic search engine function such as building index and query processing. To add some feature to the search engine, we have added keyword highlighting feature and metadata search. Finally, the search engine contains some advanced algorithms such as page rank, Naive Bayes classification and K-Mean clustering algorithm along with vector data generated from Python

## Keywords

Sigmod,ICSE,VLDB,JDK,Doc2Vec,K-mean,pagerank,Naive Bayes classification

## 1. INTRODUCTION

Data repositories are growing day by day. So, the need for retrieving information from those data is important.In search engine project,we have implemented full text searching feature using various algorithms. We have used Apache Lucene as API. The *Apache Lucene* is a full-text search library in java,which makes it easy to add search functionality to an application or website. In this project, we have implemented different algorithms of Information Retrieval. The search engine is created using java servlet and JSP page.

## 2. CONFIGURATION AND SERVER SETUP

We have used Apache tomcat to run our application.There are lot of versions of Apache tomcat but we have used Apache tomcat 1.7[1]. To control Tomcat Server,we have setup tomcat user and we have uploaded our project on server in the form of WAR(Web Achieve file). Architecture of our Java application file is pretty simple. It contains JSP[2],CSS and Jascript file along with index file. The index file is created with the help of console Java application and they are implemented in the dynamic web project. All our Java code is located inside SRC directory and other file such as JSP,CSS

and index file are located inside WebContent folder of our application.



## 3. DATA SOURCES

### 3.1 Citeseer Data Source

We have used different data sets for our application. The first data sets consists of 10,000 papers and that is in .txt file and .xml file.

### 3.2 Microsoft Academic Search Data

We have used ICSE,VLDB and SIGMOD Data sets released by Microsoft Academic Search.These data sets contains link in the form of citation graph which is very very useful for our search engine.

## 4. INDEXING AND SEARCHING

We have used Apache Lucene 5.2.1 for indexing and searching documents.

### 4.1 Indexing

Index organizes all the data so that it is quickly accessible. So,index consists of documents containing one or more fields.When Lucene indexes a document,it breaks it down into a number of terms. It then stores the terms in an index file where each term is associated with the documents that contain it. In our case, we have created index using java console application. First, we need to read the documents from file directory and pass that into `IndexWriter`. Finally, we have to add fields like `"title"`,`"author"`,`"contents"` and `"path"` to create index.

### 4.2 Searching

Searching is one of the important feature of the search engine and searching allow the end user to get desirable result. Apahce Lucene provides search on index so that searching is very fast and efficient. To make searching fast and efficient, various kinds of improvements is implemented in our search engine.

### 4.2.1 Highlighting keyword

Apahce Lucene provides an API to highlight the keyword in our search result. To highlight keyword, we have to store the term vector and their offsets in the index. The highlighter API will insert HTML Tag with attribute according to our desire, so that we can get result in our desiring format.

### 4.2.2 Pagination

The pagination helps to split the results into smaller chunk so that is it easier for user to get the result they are looking for. Here we have implemented pagination feature in our website. In our case every user can see only top 10 paper at once and this will reduce the server load thus accelerating the most relevant result to the user questions.Furthermore, we can control the limit of record using control given at the top of the page.

### 4.2.3 fast Search

It is slow process to search each and every time on index for each user query. Thus,we have implemented pre-data load technique in our server. The basic idea is to load data and store them in the memory with the help for servlet context initializer.

## 5. PAGE RANK

Page rank calculation process is the eminent process in search engine optimization. In our project, we have used `sigmod_subgraph.txt` and `pagerank_sigmod_id.tx` file and it contains the information required for link between the papers. The page ranking part consist of different stages like 1) calculating node 2) building adjacent matrix 3)building transition matrix 4)calculating page rank based on Markov Chain.

### 5.1 Detail Description

#### 5.1.1 Calculating Node

we are provided with `sigmod_subgraph.txt` documents and that contains information required to link between papers.So, first we have to read this file, compute total number of node and save all data into `arrayList` to create adjacent matrix.

#### 5.1.2 Building adjacent matrix

With nodes number and adjacent table we can create adjacent matrix among the papers.The adjacent matrix `AdjacentMatrix[i][j]` means, paper i have citation with paper j and the value adjacent matrix value is 1 if there is link otherwise it is set to 0. Then finally, with the help of transition matrix we can compute page rank value.

#### 5.1.3 Sorting page rank value

We can sort the page rank value using different method. In the front end, we can extract all the numeric value and sort particular column of table using function given below:
```
dataTable( "order":  [[ 0, "desc" ]])
```

or we can put our data in the map and sort the result using Java API.

### 5.2 Page Ranking of VLDB and ICSE

The page ranking algorithm is extended to calculate page rank of ICSE and VLDB documents. So, link between paper is provided from `sigmod_vldb_icse_subgraph.txt` documents. If we want to rank the ICSE documents then in that case, we have to input `icse_id.txt` and if we want to rank the VLDB documents then in that case, we have to input `vldb_id.txt`.

## 6. NAIVE BAYES CLASSIFICATION

We have used two sets of data for classification. First data set is from ICSE documents and second data set is from VLDB documents.Here we are classifying text using Multinominal Navies Bayes technique.So, Mutinominal Navies Bayes is machine Learning algorithm of supervised type that assume each and every feature is independent. Here we are classifying document,so, feature in our case is word.Thus by calculating probability of each word in the title belong to each class, and combine them together ,we can get the probability that title belong to a particular class. For this part we have used reference from group 3.

In the classification process, we need to provide train data and finally we need to test our data for classification. In our case, we have provided test data in the form of CSV file. The CSV file is generated from Java application and that is feed into application to classify it.The necessary code is attached in our web page with web link.

### 6.1 Training Data and Testing

During the training phase, we count the word in title of each document(VLDB and ICSE) and in the next step, we calculate the conditional probability.Finally, we converted our test document into CSV and we fed that into application for testing purpose.

The necessary formula for computing conditional probability is given below:

$$ConProb(word(i)) = word(i).getCountClassICSE() + 1/$$
$$\text{numberOfWordICSE} + \text{noUniqueWord} \quad (1)$$

## 7. CLUSTERING

For clustering, we have tested different kinds of algorithm available around the internet. The idea behind clustering data is that we need to convert our document to vector representation. The vector representation shows how a word or text usually is similar with other words or text.

Here in our case, we are creating cluster of ICSE and VLDB documents.We are going to use python to convert our document to vector(for this part,we have used code reference of Shane Peelar). We have tested many API[3] available around the Internet but we have successfully able to run the code in python.So, we will follow vector document generated by python to cluster the documents(VLDB,ICSE).

### 7.1 Converting Doc2Vec using Python

To convert our document to vector using Python, we have used PIP package manager.To convert document to vector

using Python, we must have different kinds of library and finally, we have setup Python programming environment and we are successfully able to run the code using Python interpreter using library listed below.

### 7.1.1 library required
1.Numpy
2.Gensim

### 7.1.2 Reading Doc File
if args.model is not None:
print("Loading " + args.model)
$model_dm = Doc2Vec.load(args.model)$
elif args.train is not None:
files = args.train
tags = [] for f in files:
print("Reading " + f)
for line in open(f):
components = line.split("") tags.append(TaggedDocument(copone [components[0]]))

### 7.1.3 building model
print("Building Model for document")
$model_dm = Doc2Vec(documents = tags, alpha = 0.025, size = 90, workers = 4, dbow_words = 1, iter = 90, min_count = 2)$
In this case, we have make vector size as 90 and minimum count number is 2.

### 7.1.4 Exporting to Vector space documents
if args.dumpdocvectors
is not None and model_dm != "":
print("Writing document vectors to
file: " + args.dumpdocvectors)
f = open(args.dumpdocvectors, "w")
for i in
$model_dm.docvecs.offset2doctag$ :
f.write(i + "")
for j in $model_dm.docvecs[i]$ :
f.write(":.17f".format(j))
f.write("")

594BCD68 vldb 0.2511026561260234 0.8209866981370544 -0.519647343254089
7EE98LD6 icse 0.2409145032061767 -0.622127950191497 -0.166981075084536
70B6412F icse 0.1001430017940150 1.829934035433959 -1.520063996315002
7DCD2E8 icse 0.2334282100200653 -2.102597951809038 0.970959722957580
7C4C6E34 icse 0.1405349641972406 -0.402694404125213 0.609229385852013
09230E61 icse -1.034776926040649 -1.158102631568908 -0.394428223371505
80C5DCE4 icse 0.6906639973121643 0.406089603900094 1.088207840919494
755907FC vldb 0.8012803792953491 -0.631868980685607 -0.990099787712097
58CF6CC2 vldb 1.0171680795089721 1.750862598419184 0.91454911231994
7BAFC3B8 icse -0.582703649997111 1.763621211051940 1.097581744194030
7FD3A8F2 icse -1.229230046272277 0.467873071326446 -2.426632221647003
7D6D6B17 icse -0.740370221930869 -1.295163273811340 0.038368068584758
5935B189 vldb -0.443999230061638 -2.246376290917396 -0.952576220035529
7F04570F vldb -0.386968314647674 -1.135639548301636 1.080207006454467
02750233 icse 0.475068569183349 -0.681521415104492 0.149166760022689
7DCCFF4B vldb -1.033794474601745 -2.384108781814576 -1.214133058697025
824697EE icse -0.198459014296531 0.524463117226501 0.772212672233581
7B826F83 sigmod -0.377513885490468 -0.956993341449228 0.553922951221460
812EDB8B vldb 0.366134673357009 -1.595465064048767 1.560071033477783
7C0052A8 icse -2.465319395065307 -0.969892203807830 -0.708529158401489
80EA2D44 vldb -0.522180795669555 -0.861458361148834 0.521474599938256
7B08C22C sigmod -0.377817362546920 -0.220548023475837 -0.433039426035688
0A5C8596 icse 0.848742961803494 -2.471519947052001 -1.265362977981567
03D80C47 sigmod 0.175767049193382 1.274924874305725 0.311370015144348
7CF41357 icse 0.257733136415491 -1.306076526641845 -0.011035036309597
7ED52CA0 sigmod -0.354995965957641 0.103924609720706 -0.117736406243171
7A762FD6 sigmod -0.065462401975555 -0.334069406063041 -0.791015923023323
82864317 vldb -0.994681656360626 -0.744118452072143 -0.010727422311902
76D0580E icse 0.032133374363183 -0.612989246845245 1.067267179491357

In this way, we will get document in the form of `documentID`, `documentType` and `vectorValue` format.Finally, we have converted this text document into csv format with the help of Java application and CSV file is cluster with the help of R programming.
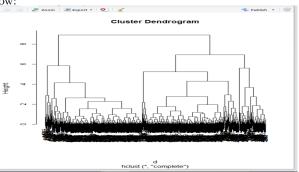
## 7.2 Clustering Vector value using Chart in R
In the R programming, we have read the CSV[] data from clusterData.csv and we have stored that in myData.

Doc2Vec document is very large about 5GB so for ease of use we have splitted it into 1200 records.



Finally we have plotted it in the Dendrogram as show below:



## 8. WEBSITE LINK
http://137.207.234.72:8080/

## 9. REFERENCES
1.https://tomcat.apache.org/tomcat-7.0-doc/building.html
2.http://www.oracle.com/technetwork/java/javaee/jsp/index.html
3.https://github.com/deeplearning4j/dl4j-examples/tree/master/dl examples/src/main/java/org/deeplearning4j/examples/nlp
4.https://deeplearning4j.org/doc2vec
5.https://rare-technologies.com/doc2vec-tutorial/
6.https://radimrehurek.com/gensim/models/doc2vec.html