

# Lab Report 1

## Zipfs Law

Irum Hussain  
Data Science  
IIIT Vadodara  
Gandhinagar, Gujarat  
202162005

Karre Kamal  
Artificial Intelligence  
IIIT Vadodara  
Gandhinagar, Gujarat  
202161004

Prashant Kumar  
Data Science  
IIIT Vadodara  
Gandhinagar, Gujarat  
202162007

**Abstract**—Zipf's law defines how often words appear in a large corpus.

### I. INTRODUCTION

Most of the electronic text content may be encoded into format supported by various text processing modules. A raw text should be generated from documents in binary formats must - a set of characters as a primary step prior to tokenization. Frequency distribution of words is very skewed. There are a few words that have a very high frequencies and many words that have low frequencies. This frequency distribution of words described by Zipf's law.

### II. APPROACH

In natural language, few terms are very frequent and many are very rare terms. According to Zipf's law the  $n$ th most frequent term has frequency proportional to  $1/n$ .

$$f(r) \propto 1/r \quad (1)$$

$$f(r) \propto 1/r = K/r \quad (2)$$

where  $K$  is a normalizing constant. In this equation,  $r$  is rank of a word, and  $f(r)$  is frequency of the word having rank  $r$ .

#### A. Tokenization and Porter Stemming

```
with open ('IR_1.txt', encoding="utf8") as fin:
    data = re.sub('[^a-zA-Z-]', ' ', fin.read())
    tokens = word_tokenize(data.lower())
    corpus = [ps.stem(token) for token in tokens]
```

#### B. Frequency and Rank relation

If the most frequent term (the) occurs  $f(r)$  times

- then the second most frequent term (of) occurs  $f(r)/2$  times
- the third most frequent term (and) occurs  $f(r)/3$  times ...

Equivalent:  $f(r) = K/r$  where  $K$  is a normalizing factor.

- $\log f(r) = \log K - \log r$
- Linear relationship between  $\log f(r)$  and  $\log r$

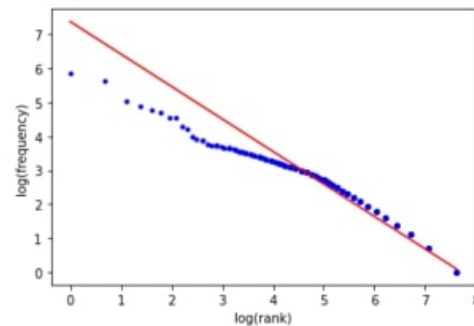


Fig. 1.

#### C. Word Cloud

The size of each word represents its frequency or relevance in a **Word cloud**, which is a data visualisation tool for visualising text data. A word cloud can be used to highlight key textual data points. Data from social networking websites is frequently analysed using word clouds.

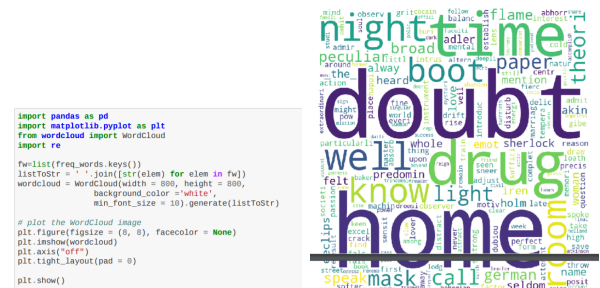


Fig. 2.

#### D. Hadoop Map Reduce

Hadoop map reduce is the Apache Hadoop processing component. It processes data parallelly in distributed environment. Without removing stopwords and stemming, Map-Reduce on the Large dataset reveals that the same word with symbols attached has separate counts. We discovered the true counts of the words in the dataset after stemming and deleting symbols.



# Lab Report 2

## BSBI and TF-IDF

Irum Hussain  
Data Science  
IIIT Vadodra  
Gandhinagar, Gujarat  
202162005

Karre Kamal  
Artificial Intelligence  
IIIT Vadodara  
Gandhinagar, Gujarat  
202161004

Prashant Kumar  
Data Science  
IIIT Vadodara  
Gandhinagar, Gujarat  
202162007

**Abstract**—BSBI divide the large corpus into a systematic matrix where terms and documents are rows and columns. TF-IDF is a method for calculating the number of words in a corpus. We usually assign each word a score to indicate its significance in the document and corpus. This method is mainly used in information retrieval and text processing.

## I. INTRODUCTION

Indexing a huge corpus with limited mainframe resources is a major problem. We can't index the entire corpus because of memory constraints. Block Sort Based Indexing (BSBI) comes to tackle this problem. It creates an index of a large corpus using small blocks. We divide the corpus into equal-sized blocks in BSBI. Then we can do TF-IDF to how significant a word is to a document. We do inverse document frequency which measures the informativeness of term.

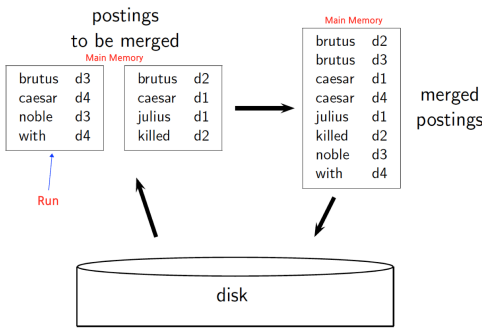


Fig. 1.

### A. Porter Stemming and BSBI

We do Porter Stemming to reduce a word to its word stem. We are dividing the corpus into small blocks where each block contains word, document id and its corresponding term frequency. We are initially fixing a block size as a threshold to increase the main memory and cache optimization. If the block size exceeds the so far created posting list of block size will be updated in a new file and stored locally. In the same way we will be getting some files based on the corpus size. At last we will merge all the posting list files to a single file.

## Pseudocode

```
def bsbi():
    current_block=0
    freq_dict = defaultdict(dict)
    for doc in Corpus:
        id = doc["id"]
        text = doc["body"]
        for word in text:
            word = porter.stem(word)
            if word not in freq_dict:
                current_block += 1
            if not freq_dict[word].__contains__(id):
                freq_dict[word].update({id:word_count})
                current_block += 1
    if current_block >= BLOCK_SIZE:
        Create Posting Lists using freq_dict
        word->(doc_id, tf)->(doc_id, tf)->
```

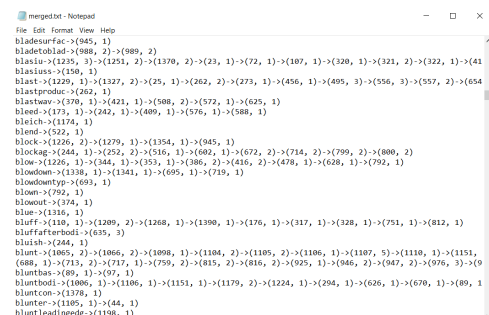


Fig. 2.

## Metrics

- Precision tells us how many number of retrieved docs are relevant
- Recall tells us how many number of relevant docs are retrieved
- F1 score tells us weighted average of Precision and Recall

### B. TF(Term Frequency)-IDF(Inverse Document Frequency)

The TF-IDF for a word in a document is calculated by multiplying two metrics:

- The term frequency of a word in a document (tf)
- The inverse document frequency of the word (idf)

So, if the word is very common and presents in many documents, this number will tends 0. Otherwise, it will tends 1.

TF-IDF score of a word in a document will be given by  $tf * idf$ . The higher the score, the more significant that term becomes in that particular document.

$$tfidf(t, d) = tf(t, d) * idf(t) \quad (1)$$

$$idf(t) = \log[n/df(t)] + 1 \quad (2)$$

where  $n$  = no. of documents in the corpus.

## II. CONCLUSION

For large corpus BSBI is a good solution to do indexing. BSBI is efficient in memory wise. In BSBI we do sorting 2 times. First we store Posting lists and then we do merging.

## REFERENCES

- [1] <https://towardsdatascience.com/tf-idf-for-document-ranking-from-scratch-in-python-on-real-world-dataset-796d339a4089>

# Lab Report 3

## Boolean Retrieval

Irum Hussain  
Data Science  
IIIT Vadodra  
Gandhinagar, Gujarat  
202162005

Karre Kamal  
Artificial Intelligence  
IIT Vadodra  
Gandhinagar, Gujarat  
202161004

Prashant Kumar  
Data Science  
IIT Vadodra  
Gandhinagar, Gujarat  
202162007

**Abstract**—Explore Boolean Retrieval model and how to process Boolean queries and do Weighted Zone Scoring.

## I. INTRODUCTION

Any boolean expression can be answered using the Boolean Retrieval paradigm. Query is usually specified using operators from Boolean logic ( AND, OR, NOT ).It views each document as a set of terms. For any query evaluation, there are two possible results ( TRUE and FALSE ).. Boolean Retrieval is also known as exact-match retrieval.

- Assumption : All documents in the retrieved set are equivalent in terms of relevance.

## II. APPROACH

### A. Porter Stemming and BSBI

We do Porter Stemming to reduce a word to its word stem. We are dividing the corpus into small blocks where each block contains word, document id and its corresponding term frequency. We are initially fixing a block size as a threshold to increase the main memory and cache optimization. If the block size exceeds the so far created posting list of block size will be updated in a new file and stored locally. In the same way we will be getting some files based on the corpus size. At last we will merge all the posting list files to a single file.

```
File Edit Format View Help
bladesurfuc<(945, 1)
bladedcolad<(120, 2)>(989, 2)
bladedcolad<(1229, 2)>(1237, 2)>(1370, 2)>(23, 1)>(72, 1)>(107, 1)>(320, 1)>(321, 2)>(322, 1)>(41
bladesun<(150, 1)
blast<(1229, 1)>(1227, 2)>(25, 1)>(262, 2)>(273, 1)>(456, 1)>(409, 3)>(556, 3)>(557, 2)>(654
blastw<(102, 1)
blastwuv<(370, 1)>(421, 1)>(508, 2)>(572, 1)>(625, 1)
blend<(37, 2)>(362, 1)>(409, 1)>(576, 1)>(568, 1)
bleich<(3174, 1)
blend<(522, 1)
blend<(1226, 2)>(1279, 1)>(1354, 1)>(945, 1)
block<(244, 1)>(252, 2)>(516, 1)>(602, 1)>(672, 2)>(714, 2)>(799, 2)>(800, 2)
blue<(1226, 1)>(1344, 1)>(1353, 1)>(1389, 2)>(416, 2)>(478, 1)>(626, 1)>(792, 1)
bluendom<(1310, 1)>(1314, 1)>(1065, 1)>(719, 1)
bluendomtyp<(693, 1)
blue<(792, 1)
bluend<(374, 1)
blue<(1316, 1)
blue<(1140, 1)>(1209, 2)>(1268, 1)>(1390, 1)>(1376, 1)>(1317, 1)>(1328, 1)>(751, 1)>(812, 1)
bluesherid<(102, 1)
blush<(244, 1)
blush<(1805, 2)>(1806, 2)>(1809, 1)>(1804, 2)>(1815, 2)>(1816, 1)>(1817, 5)>(1810, 1)>(1815,
blush<(1805, 2)>(717, 1)>(759, 2)>(815, 2)>(816, 2)>(925, 1)>(946, 2)>(947, 2)>(976, 3)>(9
blunthas<(80, 1)>(97, 1)
blunthas<(102, 1)>(1106, 1)>(1151, 1)>(1179, 2)>(1224, 1)>(284, 1)>(626, 1)>(670, 1)>(89, 1)
blunth<(1378, 1)
blunter<(1805, 1)>(444, 1)
blueLauLindem<(102, 1)
```

Fig. 1.

### B. Boolean retrieval model and Zone Scoring

Usually boolean queries use AND, OR and NOT to join query terms

- Document matches condition or not. This tells us the Precision.

Perhaps the most basic model for constructing an IR system on the primary commercial retrieval tool for the past three decades.

Search systems that uses Boolean model till date:

- Emails, Library catalogs, Mac Spotlight

```

34 return doc titles
35
36 d = query(include=['flow', 'field'], dont_include=['transient'])
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1
```

Fig. 2.

## Zone Scoring

Given a Boolean query and a document, Weighted Zone Scoring computes a linear combination of zone scores, where each document’s zone adds a Boolean value, and awards a score in the interval  $[0, 1]$  to the pair  $(q_1, q_2)$ .

```

1  ZONESCORE( $q_1, q_2$ )
2    1 float scores[ $N$ ] = [0]
3    2 constant  $\hat{g}[i]$ 
4    3  $p_1 \leftarrow \text{postings}(q_1)$ 
5    4  $p_2 \leftarrow \text{postings}(q_2)$ 
6    5 // scores[] is an array with a score entry for each document, initialized to zero.
7    6 //  $p_1$  and  $p_2$  are initialized to point to the beginning of their respective postings.
8    7 // Assume  $\hat{g}[i]$  is initialized to the respective zone weights.
9    8 while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$ 
10   9   do if  $\text{docID}(p_1) = \text{docID}(p_2)$ 
11      10   then  $\text{scores}[\text{docID}(p_1)] \leftarrow \text{WEIGHTEDZONE}(p_1, p_2, \hat{g})$ 
12      11    $p_1 \leftarrow \text{next}(p_1)$ 
13      12    $p_2 \leftarrow \text{next}(p_2)$ 
14   13   else if  $\text{docID}(p_1) < \text{docID}(p_2)$ 
15      14   then  $p_1 \leftarrow \text{next}(p_1)$ 
16      15   else  $p_2 \leftarrow \text{next}(p_2)$ 
17   16 return scores

```

Fig. 3.

Let  $s_i$  be the Boolean score denoting a match (or absence thereof) between  $q$  and the  $i$  th zone. Let  $g_1, g_2, \dots, g_n \in [0, 1]$  such that  $\sum_{i=1}^n g_i = 1$ . Then the weighted zone score is defined to be

$$\sum_{i=1}^n g_i * s_i \in [0, 1] \quad (1)$$

### III. CONCLUSION

Results are very predictable and easy to explain to users. Boolean query operands can be of any document attribute. Documents can be quickly discarded from consideration from consideration in the scoring process, so more efficient.

### REFERENCES

- [1] <https://nlp.stanford.edu/IR-book/html/htmledition/weighted-zone-scoring-1.html>

# Lab Report 4

## Vector Space Retrieval

Irum Hussain  
202162005

Karre Kamal  
202161004

Prashant Kumar  
202162007

**Abstract**—The aim here is to install and experiment with Elasticsearch and Kibana.

### I. INTRODUCTION

Elasticsearch is an open source ,distributed, JSON based search and analytics engine which is capable of addressing a large number of use cases.Elasticsearch is based on Apache Lucene search engine and allows you to perform and combine different types of searches as needed, including structured, unstructured, geographic, and metric.Elasticsearch uses aggregation to zoom out to see trends and patterns in your data. It can quickly store, search, and analyze large amounts of data in near real time and get answers in milliseconds. It uses document-based structures instead of tables and schemas and has extensive REST APIs for storing and retrieving data. Basically, Elasticsearch can be thought of as a server that can process JSON requests and return JSON data.Kibana is a visualization layer built on Elasticsearch that allows users to analyze and visualize their data.This allows you to merge Elasticsearch data into multiple indexes and combine it with other SQL / NoSQL / REST API data sources to create visualizations with a user-friendly user interface.We run queries and using dev tools in kibana try to understand how search results are returned based on our query. We also do vector space retrieval where set of documents are represented as vectors and is fundamental to a large number of information retrieval ranging from scoring documents on a query, document classification and document clustering.

#### A. Similarity model in Elasticsearch

Elasticsearch calls the Lucene Practical Scoring Function. This function generates a relevance score that Elasticsearch uses to sort documents when data is requested.It is based on vector space retrieval model which uses tf-idf(term frequency-inverse document frequency).

##### 1) A. Vector Space Model:

Document Frequency ( $d_f$ ):The number of documents in the collection that contains a particular term .

Term Frequency( $t_f$ ): Counting the number of occurrences of the term in the document.

Inverse document frequency ( $idf$ ): The logarithm of number of documents in the collection (or index) divided by the number of documents that contain the word. It is highest when a term

occurs many times within a small collection of documents but is lowest when term occurs in all documents i.e. zero.

$$tf-idf_{t,d} = tf_{t,d} \times idf_t.$$

Overlap Score measure: the score of a document is the sum over all query terms of the number of times each of the query terms occurs in a document. Now to quantify the

$$Score(q, d) = \sum_{t \in q} tf-idf_{t,d}.$$

similarity between two documents we use cosine similarity which takes care of the length of the document.So cosine similarity between two documents d1 and d2 is shown below.

$$sim(d_1, d_2) = \frac{\vec{V}(d_1) \cdot \vec{V}(d_2)}{|\vec{V}(d_1)| |\vec{V}(d_2)|},$$

Query boost:At search time users can specify boosts to each query, sub-query, and each query term, hence the contribution of a query term to the score of a document is multiplied by the boost of that query term.

coord-factor(q,d):A document may match a multi term query without containing all the terms of that query (this is correct for some of the queries), and users can further reward documents matching more query terms through a coordination factor, which is usually larger when more terms are matched. Query normalization (queryNorm): Total of the squared weights of the query.

normalization (norm): The inverse square root of the number of terms in the field is used to calculate the field length normalisation (norm).

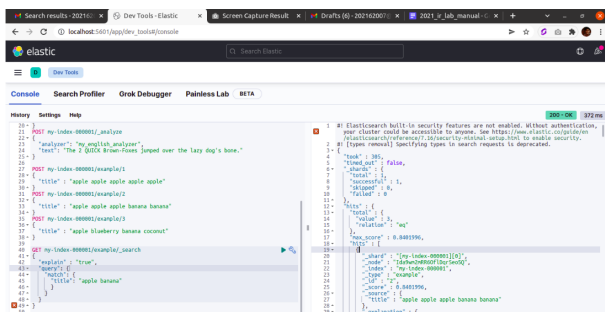
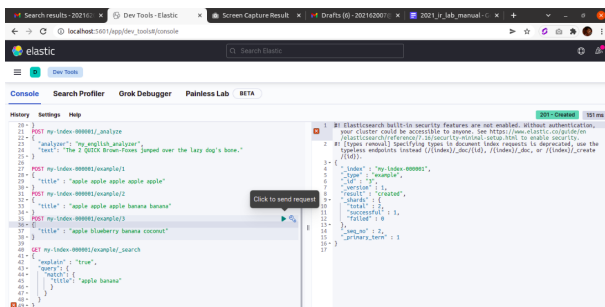
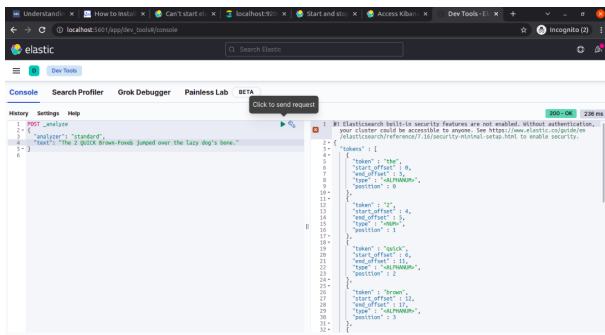
Lucene's Conceptual scoring formula: For query q and document d is defined as follows:

$$score(q,d) = queryNorm(q) \quad coord(q,d) \quad SUM(t_f(ind), idf_f(t), t.getBoost(), norm(t,d))(t \text{ in } q)$$

### II. RESULTS

We found that we could change the priority or ranking by specifying fields such as indexBoost and queryBoost in

Elasticsearch. Kibana helps in visualisation by providing easy-to-understand graphic representation of our dataset.



## CONCLUSION

The combination of Elasticsearch and Kibana makes it a very powerful tool for information retrieval. They can be together used to analyse data quickly and efficiently.

## REFERENCES

1. <https://www.elastic.co/blog/found-similarity-in-elasticsearch>
2. <https://lucene.apache.org/core/451/core/org/apache/lucene/search/similarities/TFIDFSimilarity.html>
3. Manning, C. D., Raghavan, P., Schu tze, H. (2008). Introduction to information retrieval. New York: Cambridge University Press.



# Lab Report 5

## Text Classification

Irum Hussain  
202162005

Karre Kamal  
202161004

Prashant Kumar  
202162007

**Abstract**—The large volume of unstructured documents and text are difficult to deal with. Obtaining specific information from this large set of data takes a lot of time. One of the techniques is text classification. Text classification task is to assign a document to one or more category. Text classifiers can be used to organize, structure, and categorize any data.

### I. INTRODUCTION

The growth of the Internet has had a major impact on data generation. Most of the world's data is in text format. There is need to access and use this data efficiently and easily hence, text classification is widely studied problem in research community. Text classification is the process of classifying a text document into a fixed number of predefined classes. Text classification applications include spam filtering, email routing, sentiment analysis, voice recognition and more. It is a supervised learning approach in which we predict the class label of incoming document based on the training model built from a training set of documents labelled with classes. The most important step of the text classification pipeline is choosing the best classifier. So we are going to understand Naive Bayes classifier, Rocchio classifier, K Nearest Neighbor classifier by training and testing on 20 Newsgroup dataset.

#### A. Classification Algorithms

##### 1) A. Naive Bayes Classification:

Naive Bayes classification is probabilistic machine learning model based on the Bayes theorem. The probability of a document  $d$  being in class  $c$  is  $P(t_k|c)$  is the conditional

$$P(c|d) \propto P(c) \prod_{1 \leq k \leq n_d} P(t_k|c)$$

probability of term  $t_k$  occurring in a document of class  $c$ . We interpret  $P(t_k|c)$  as a measure of how much evidence  $t_k$  contributes that  $c$  is the correct class.  $P(c)$  is the prior probability of a document occurring in class  $c$ .

2) B. K Nearest Neighbour: KNN is a supervised learning algorithm which determines the decision boundary locally. It is based on the hypothesis that it expects a test document  $d$  to have the same label as the training documents located in the local region surrounding  $d$ . Here  $K$  is the number of nearest neighbours to be considered. For a new data point in vector space, kNN looks for  $k$  points which are nearest to that

```
TRAINMULTINOMIALNB(C,D)
1  V ← EXTRACTVOCABULARY(D)
2  N ← COUNTDOCS(D)
3  for each c ∈ C
4  do Nc ← COUNTDOCSINCLASS(D,c)
5  prior[c] ← Nc/N
6  textc ← CONCATENATETEXTOFALLDOCSINCLASS(D,c)
7  for each t ∈ V
8  do Ttc ← COUNTTOKENSOFTERM(textc,t)
9  for each t ∈ V
10 do condprob[t][c] ←  $\frac{T_{tc}}{\sum_{c' \in C} T_{tc'}}$ 
11 return V, prior, condprob

APPLYMULTINOMIALNB(C,V,prior,condprob,d)
1  W ← EXTRACTTOKENSFROMDOC(V,d)
2  for each c ∈ C
3  do score[c] ← log prior[c]
4  for each t ∈ W
5  do score[c] += log condprob[t][c]
6  return arg maxc ∈ C score[c]
```

data by selecting distance measure such as euclidean distance, manhattan distance and minkowski distance. These  $K$  points then becomes the nearest neighbours of that new data point and is assigned to the class which represents the most points among those  $k$  neighbours. For finding most optimum  $k$  value, first differentiate the training and validation dataset and then plot the validation error curve.

```
TRAIN-KNN(C,D)
1  D' ← PREPROCESS(D)
2  k ← SELECT-K(C,D')
3  return D', k

APPLY-KNN(C,D',k,d)
1  Sk ← COMPUTENEARESTNEIGHBORS(D',k,d)
2  for each cj ∈ C
3  do pj ← |Sk ∩ cj|/k
4  return arg maxj pj
```

3) C. Rocchio Classifier: One of the main task in vector space is to create good class boundaries between the classes and rocchio classification is one of the ways. It uses centroids to define the boundaries. It is a form of relevance feedback where the average of the relevant documents, corresponds to the most important component of the Rocchio vector in relevance feedback, is the centroid of relevant documents.

```
TRAINROCCHIO(C,D)
1  for each cj ∈ C
2  do Dj ← {d : (d, cj) ∈ D}
3   $\vec{\mu}_j \leftarrow \frac{1}{|D_j|} \sum_{d \in D_j} \vec{v}(d)$ 
4  return { $\vec{\mu}_1, \dots, \vec{\mu}_J$ }

APPLYROCCHIO({ $\vec{\mu}_1, \dots, \vec{\mu}_J$ }, d)
1  return arg minj | $\vec{\mu}_j - \vec{v}(d)$ |
```

### II. APPROACH

Dataset: The 20 newsgroups dataset comprises of around 18000 newsgroups posts on 20 topics split into two subsets. One is for training and the other one is for testing. The

split between the train and test data set is based upon a messages posted before and after a specific date. Each record in the corpus is actually a text file. Then we create tokens using CountVectorizer() which counts the number of times a word appears in each file giving us a count matrix. Then we transform our matrix using TfidfTransformer() which transforms the count matrix into Tf-idf representation. The goal of using tf-idf instead of the count matrix is to scale down the impact of tokens that occur very frequently in a given corpus and that are hence less informative than features that occur in a small portion of the training corpus. This comprises of our training dataset.

**Naive Bayes:** We use MultinomialNB (the multinomial Naive Bayes classifier) present in sklearn.naive\_bayes which is suitable for discrete classification. Scikit-learn has a class called Pipeline, which allows us to add the functions that we want to use on our input data into the pipeline for the classifier. Then we take the training dataset and feed it in the classifier. We then get a matrix from it for which we calculate the accuracy, precision and recall using the test set.

**Rocchio Classifier:** We use NearestCentroid present in sklearn.neighbors.nearestCentroid where each class is represented by its centroid, with test samples classified to the class with nearest centroid. Scikit-learn has a class called Pipeline, which allows us to add the functions that we want to use on our input data into the pipeline for the classifier. Then we take the training dataset and feed it in the classifier. We then get a matrix from it for which we calculate the accuracy, precision and recall using the test set.

**K-Nearest Neighbour Classifier:** We use KNeighborsClassifier present in sklearn.neighbors with k and weights as parameters. Scikit-learn has a class called Pipeline, which allows us to add the functions that we want to use on our input data into the pipeline for the classifier. Then we take the training dataset and feed it in the classifier. We then get a matrix from it for which we calculate the accuracy, precision and recall using the test set.

### III. RESULTS

The precision, recall and f1 scores for few categories of 20 newsgroup dataset for naive bayes, rocchio and K nearest neighbour classifier have been computed.

Accuracy 0.9357565511411665	precision	recall	f1-score	support
sci.electronics	0.95	0.90	0.92	393
sci.med	0.94	0.94	0.94	396
sci.space	0.92	0.97	0.94	394
accuracy			0.94	1183
macro avg	0.94	0.94	0.94	1183
weighted avg	0.94	0.94	0.94	1183
(None, array([[353, 14, 26], [ 17, 371, 8], [ 3, 8, 383]]))				

Fig. 1. Naive Bayes Classifier

Accuracy 0.8571428571428571	precision	recall	f1-score	support
sci.electronics	0.84	0.80	0.82	393
sci.med	0.85	0.84	0.85	396
sci.space	0.88	0.93	0.90	394
accuracy			0.86	1183
macro avg	0.86	0.86	0.86	1183
weighted avg	0.86	0.86	0.86	1183
array([[313, 47, 33], [ 44, 334, 18], [ 14, 13, 367]])				

Fig. 2. K Nearest Neighbour Classifier

Accuracy 0.801352493660186	precision	recall	f1-score	support
sci.electronics	0.66	0.91	0.76	393
sci.med	0.88	0.63	0.74	396
sci.space	0.96	0.86	0.91	394
accuracy			0.80	1183
macro avg	0.83	0.80	0.80	1183
weighted avg	0.83	0.80	0.80	1183
array([[359, 24, 10], [142, 250, 4], [ 45, 10, 339]])				

Fig. 3. Rocchio classifier

### CONCLUSION

Classification tasks are one of the most essential issues in machine learning. As text And the document dataset grows year by year, having a better document categorization system for this growing information requires learning and discerning these text classification algorithms. The classification of the news articles could be expanded to multi-label text articles, as there are many news articles that do not belong to a single category and are instead a mixture of various categories.

### REFERNECES

1. Manning, C. D., Raghavan, P., Schütze, H. (2008). Introduction to information retrieval. New York: Cambridge University Press.
2. <https://archive.ics.uci.edu/ml/datasets/Twenty+Newsgroups>
3. Ali, N., Neagu, D., Trundle, P. Evaluation of k-nearest neighbour classifier performance for heterogeneous data sets. SN Appl. Sci. 1, 1559 (2019). <https://doi.org/10.1007/s42452-019-1356-9>
4. [https://www.researchgate.net/publication/228961211\\_Explaining\\_naiveBayes\\_classifications](https://www.researchgate.net/publication/228961211_Explaining_naiveBayes_classifications)

# Lab Report 6

## Latent Semantic Indexing

Irum Hussain  
202162005

Karre Kamal  
202161004

Prashant Kumar  
202162007

**Abstract**—To help analyse relationships between terms in the documents or between the documents in the corpus we aim to experiment with latent semantic indexing based on singular value decomposition.

### I. INTRODUCTION

Latent Semantic Indexing (LSI) is a retrieval and indexing mechanism which uses a mathematical technique called Singular Value Decomposition (SVD) to identify patterns of relationships between terms and concepts in unstructured collections of text data. Singular Value Decomposition is a factorization of a real or complex matrix into simpler matrices where each matrix is a column vector times row vector. The SVD chooses rank pieces in order of importance.

#### A. Singular Value Decomposition

Singular value decomposition decomposes a matrix into three components where two of these components are representations of word and documents in the approximated vector space while the third components provides weights to these vectors. The decomposition of a matrix is useful when the matrix is not of full rank. That is, the rows or columns of the matrix are linearly dependent.

#### B. Latent Semantic Indexing

The latent semantic indexing tries to overcome the problems of polysemy (same word but multiple meanings) and synonymy (two words but same meaning). It takes the low rank term document matrix from SVD which gives a new representation of each document in the corpus. We then cast queries into this low rank representation as well and compute similarities between document and query. If two terms are used in similar contexts or have similar meaning then they will have similar vectors in the reduced dimension representation. It is thus a method for dimensionality reduction where we take objects from high dimensional space and cast them in low dimensional space. The projection into the latent semantic space is such that when measured by the sum of the squares of the differences the representations in the original space are changed as little as possible.

### II. APPROACH

A large sparse term-document matrix  $A$  is created from the corpus which is not a full rank matrix. Let  $n$  be the rank of  $t \times d$  matrix  $A$  where  $t$  is the number of terms and  $d$  is the number of documents. SVD takes this matrix  $A$  and in lower dimension  $A'$  "distance" between the two matrices as measured by the 2-norm is minimized. Then we apply

$$\Delta = \|A - \hat{A}\|_2$$

Singular value decomposition as  $T$  and  $D$  have orthonormal

$$A_{t \times d} = T_{t \times n} S_{n \times n} (D_{d \times n})^T$$

columns where  $(TT^T) = (D^T D) = I$  and  $\text{rank}(A_k) = r$  where  $r \leq n$ .

$$\hat{A}_{t \times k} = T_{t \times k} S_{k \times k} (D_{d \times k})^T$$

After constructing a low-rank approximation  $A'$  to the term-document matrix, for a value of  $r$  that is far smaller than the original rank of  $A$ . We use this new  $r$ -dimensional representation to compute similarities and map the query vector  $q$  to  $q'$ .  $q'$  is a simple vector in space of terms. The query vector

$$\hat{q} = q^T T_{t \times k} S^{-1}_{k \times k}$$

is then compared to all document vectors, and the documents ranked by their similarity (nearness) to the query. One measure of similarity is the cosine similarity between the query vector and document vector. Now if a new document arrives then we either recompute the SVD of a new term-document matrix or fold the new terms and documents when they arrive. Folding new documents and terms require less time and space so new documents can be folded into this LSI representation

thus we keep on incrementally adding new documents to LSI representation as follows.

$$A = TSD^T$$

$$T^T A = T^T TSD^T$$

$$T^T A = SD^T$$

### III. RESULTS

```

1 #finding the most similar document for the query vector by taking the dot product between query and
2 #each column in the vt matrix and normalize them
3 similarities = []
4 for i in range(len(vt)):
5     similarities.append(list(np.dot(q,vt[:,i])/(np.linalg.norm(q)*np.linalg.norm(vt[:,i])))[0])
6
7 print('Most similar document:')
8 idx = similarities.index(max(similarities))
9 print(idx)
10 corpus[idx]

```

Most similar document:  
1030  
'note on creep buckling of columns .'

### CONCLUSION

In a nutshell, Latent semantic indexing is able to overcome the problem of synonymy and polysemy. Also since its a mathematical approach therefore it is independent of the languages. It is also tolerant to noise data which is misspelled words,unreadable characters etc. Though it has its own disadvantages such as a large sparse matrix , time consuming but it still has a wide scope in future.

### REFERENCES

- 1.Manning, C. D., Raghavan, P., Schu tze, H. (2008). Introduction to information retrieval. New York: Cambridge University Press.
- 2.Strang, G. (1976). Linear algebra and its applications. New York: Academic Press
- 3.G.W. Furnas et al., Information retrieval using a singular value decomposition model of latent semantic structure, in Proceedings of SIGIR, 1988.
- 4.Deerwester, Dumais, Furnas, Lanouauer, and Harshman, Indexing by latent semantic analysis, Journal of the American Society for Information Science, 41 (1990), pp. 391-407.

# Lab Report 8

## Support Vector Machine

Irum Hussain  
202162005

Karre Kamal  
202161004

Prashant Kumar  
202162007

**Abstract**—The aim of this lab is to classify two classes which are non linearly separable via support vector machines. Also learn and experiment with kernel SVM

### I. INTRODUCTION

Support Vector Machine is a supervised learning used for regression and classification. Support Vector Machine (SVM) tries to maximize predictive accuracy while automatically avoiding over fitting to the data. SVM is based on vector space model where we try to find a decision boundary between classes that is maximally far away from any of the point of the training dataset. A classifier with a large margin ensures high certainty classification decisions. Non-linear data set are difficult to classify using a linear hyperplane. Thus we want to apply a function that projects / transforms the data in such a manner that the data becomes linearly separable. The idea of the kernel SVM is to enable operations to be performed in the input space rather than the high dimensional feature space. We want the function to perform mapping of the attributes of the input space to the feature space.

#### A. Theory

##### 1) A. Support Vector Machine:

For a two-class, linearly separable training data sets, there are lots of possible linear separators but only one of them achieves maximum separation. To choose among all the hyperplanes (linear separators), we specify the intercept term  $b$ . Because the hyperplane is perpendicular to the normal vector  $w$ , all points  $x'$  on the hyperplane satisfy the equation

$$\bar{w}^T \bar{x} = -b.$$

Now for a set of training data points  $D(\text{point } x_i, \text{class label } y_i)$  the two data classes are  $+1$  and  $-1$  and the intercept term is  $b$ . The linear classifier is then represented as :

$$f(\bar{x}) = \text{sign}(\bar{w}^T \bar{x} + b)$$

Each point's distance from the hyperplane is

$$r_i = y_i(\bar{w}^T \bar{x}_i + b) / |\bar{w}|.$$

$$\rho = 2 / |\bar{w}|$$

The geometric margin is Now we get a quadratic optimization problem and need to solve it for  $w$  and  $b$ . Thus we optimize the quadratic function with linear constraints and construct a dual problem where a Lagrange's multiplier  $i$  is associated. SVM as a minimisation problem is to find  $w$  and  $b$  such that:

- $\frac{1}{2} \bar{w}^T \bar{w}$  is minimized, and
- for all  $\{(\bar{x}_i, y_i)\}, y_i(\bar{w}^T \bar{x}_i + b) \geq 1$

2) *B. kernel SVM*: Now for data that cannot be classified by a linear classifier, the idea is to map the original feature space to some higher-dimensional feature space where the training set is separable while preserving the relations between the data points. Thus a kernel function  $K$  is a function that corresponds to a dot product in expanded feature space. A kernel function  $K$  is continuous, symmetric, and a positive definite gram matrix.

$$K(\bar{x}_i, \bar{x}_j) = \phi(\bar{x}_i)^T \phi(\bar{x}_j).$$

#### B. Approach

For the first part, we randomly generate dataset which follows a quadratic distribution. Then we split this dataset into (20:80) test and training datasets. Then we use polynomial kernel svm that represents the similarity of vectors in a feature space over polynomials of the original variables enabling learning of non-linear models.

Dataset: The 20 newsgroups dataset comprises of around 18000 newsgroups posts on 20 topics split into two subsets. One is for training and the other one is for testing. The split between the train and test data set is based upon a messages posted before and after a specific date. Each record in the corpus is actually a text file. Then we create tokens using CountVectorizer() which counts the number of times a word appears in each file giving us a count matrix. Then we transform our matrix using TfidfTransformer() which transforms the count matrix into Tf-idf representation. The goal of using tf-idf instead of the count matrix is to scale down the impact of tokens that occur very frequently in a given corpus and that are hence less informative than features that occur in a small portion of the training corpus. This comprises of our training dataset.

Then we use SGD classifier from the `sklearn.linear_model` which implements regularized linear models with stochastic gradient descent (SGD) learning. The gradient of the loss is estimated each sample at a time and the model is updated with a decreasing strength schedule (aka learning rate). SGD allows minibatch learning via the `partial_fit` method.

We then get a matrix from it for which we calculate the accuracy, precision and recall using the test set.

## II. RESULTS

For part 1 the results are as shown below:

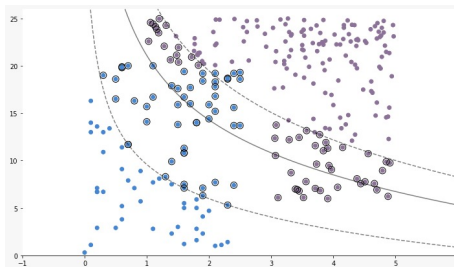


Fig. 1. SVM plot

	precision	recall	f1-score	support
0	0.66	0.91	0.76	393
1	0.88	0.63	0.74	396
2	0.96	0.86	0.91	394
accuracy			0.80	1183
macro avg	0.83	0.80	0.80	1183
weighted avg	0.83	0.80	0.80	1183

Fig. 2. precision ,recall and f1 score

For part 2 the precision, recall and f1 score for the same categories used in naive bayes have been done for SVM as well

1 SVM(categories)				
*** SVM Model ***				
Newsgroup Categories : ['sci.med', 'sci.space', 'sci.electronics']				
Accuracy : 93.99830938292477%				
	precision	recall	f1-score	support
sci.electronics	0.89	0.97	0.93	393
sci.med	0.95	0.91	0.93	396
sci.space	0.98	0.94	0.96	394
accuracy			0.94	1183
macro avg	0.94	0.94	0.94	1183
weighted avg	0.94	0.94	0.94	1183
Confusion Matrix :				
[[381 9 3]				
[ 32 360 4]				
[ 13 10 371]]				

Fig. 3. Newsgroup SVM accuracy

## CONCLUSION

Support Vector Machines acts is one of the best approach to data modeling and classification. The kernel mapping provides a common base for most of the commonly employed model architectures based on different datasets for comparisons to be performed and is able to classify non linear datasets. The SVM is dependent on very few support vectors leading to less memory and concise models . It is affected only by the points on the margin and works well with high dimensional

data. Kernel SVM are also very versatile and able to classify large no of databases of many types.

## REFERENCES

1. Manning, C. D., Raghavan, P., Schütze, H. (2008). Introduction to information retrieval. New York: Cambridge University Press.
2. Strang, G. (1976). Linear algebra and its applications. New York: Academic Press
3. [https://en.wikipedia.org/wiki/Polynomial\\_kernel](https://en.wikipedia.org/wiki/Polynomial_kernel)
4. <https://archive.ics.uci.edu/ml/datasets/Twenty+Newsgroups>
5. Tutorial slides by Andrew Moore. [Http://www.cs.cmu.edu/~awm](http://www.cs.cmu.edu/~awm)

# LAB CODE LINK:

<https://github.com/Kamal-prog-code/IR-Lab>