

Document Retrieval using TF-IDF & Ranking using Matching Score and Cosine Similarity



Submitted to:
Dr. Arun Kumar Yadav

Submitted by: Group 1
Khushwant Kaswan - 195030
Monu Kumar - 195083
Mohd Salmaan - 195088

Content

- TF - IDF
- Workflow
- Data Set Used
- Preprocessing
- TF - IDF for body and title
- Matching Score Calculation
- Cosine Similarity
- Comparison and Conclusion

TF - IDF

- “Term Frequency - Inverse Document Frequency”
- Term Frequency: frequency of a word in a document.
- Document Frequency: frequency of documents in which the word is present.
- $\text{TF-IDF} = \text{Term Frequency (TF)} * \text{Inverse Document Frequency (IDF)}$
- $\text{tf-idf}(t, d) = \text{tf}(t, d) * \log(N/(\text{df}))$
- The goal of tf-idf is to quantify how uniquely important a word of interest is to a given document within a collection of documents

Workflow

- Data scraping from local fs
- Collecting the file paths, titles and creating dataset
- Data Preprocessing
- Calculating DF for all word
- Calculating TF-IDF for body text
- Calculating TF-IDF for Title text
- Merging the TF-IDF of title and body
- Ranking
 - Matching Score
 - Cosine Similarity

Dataset Used

The dataset we have used are collection of files stored in different directories.

This means we have documents in different formats i.e. html, txt, hac, hum, key....

Dataset : <http://archives.textfiles.com/>

Files are present in different folders. There's an index.html in each folder (including the root), which contains all the document names and their titles.

Total files in dataset : 467

After Preprocessing, Total Vocab Size : 32350

Dataset at a glance

Filename	Size	Description of the Textfile
<u>FARNON</u>		The Storles of Tristan Farnon
<u>SRE</u>		The Solar Realms Elite, by Josh Renaud
<u>100west.txt</u>	20839	Going 100 West by 53 North by Jim Prentice (1990)
<u>13chil.txt</u>	8457	The Story of the Sly Fox
<u>14.lws</u>	5261	A Smart Bomb with a Language Parser
<u>16.lws</u>	15294	Two Guys in a Garage, by M. Peshota
<u>17.lws</u>	10853	The Early Days of a High-Tech Start-up are Magic (November 18, 1991) by M. Peshota
<u>18.lws</u>	26624	The Couch, the File Cabinet, and the Calendar, by M. Peshota (December 9, 1991)
<u>19.lws</u>	17902	Engineering the Future of American Technology by M. Peshota (January 5, 1992)
<u>20.lws</u>	13588	What Research and Development Was Always Meant to Be, by M. Peshota
<u>3gables.txt</u>	33985	The Adventure of the Three Gables
<u>3lpigs.txt</u>	5403	The Story of the 3 Little Pigs
<u>3sonnets.vrs</u>	2433	A Collection of Sonnets by Staeorra Rokraven (March 9, 1989)

Filename and Description are in the [index.html](#), we need to extract those names and titles.

Preprocessing

- `convert_lower_case(data)` : The One With Unagi => the one with unagi
- `remove_punctuation(data)` : `!\"#$%&()*+-./:;<=>?@[\\]^_`{|}~\n` => `<space>`
- `remove_apostrophe(data)` : ``` => `<empty>`
- `remove_stop_words(data)` : the one with unagi => one unagi
- `convert_numbers(data)` : 100 => hundred
- `stemming(data)`
 - Porter Stemmer

Calculating DF for all words

Measures the frequency of documents in which the word is present.

We will use a dictionary and we can use the word as the key and a set of documents as the value.

But for DF we don't actually need the list of docs, we just need the count. so we are going to replace the list with its count.

```
{'sharewar': 1,  
'trial': 1,  
'project': 4,  
'freewar': 1,  
'need': 6,  
'support': 2,  
'continu': 4,  
'one': 10,  
'hundr': 8,
```



```
1  DF = {}  
2  
3  for i in range(N):  
4      tokens = processed_text[i]  
5      for w in tokens:  
6          try:  
7              DF[w].add(i)  
8          except:  
9              DF[w] = {i}  
10  
11     tokens = processed_title[i]  
12     for w in tokens:  
13         try:  
14             DF[w].add(i)  
15         except:  
16             DF[w] = {i}  
17 for i in DF:  
18     DF[i] = len(DF[i])
```


TF - IDF for body and title

We need to iterate over all the documents, we can calculate the frequency of the tokens, tf-idf and finally store as a (doc, token) pair in tf_idf.



```
1  tf_idf = {}
2
3  for i in range(N):
4      tokens = processed_text[i]
5      counter = Counter(tokens + processed_title[i])
6      words_count = len(tokens + processed_title[i])
7
8      for token in np.unique(tokens):
9          tf = counter[token]/words_count
10         df = doc_freq(token)
11         idf = np.log((N+1)/(df+1))
12         tf_idf[doc, token] = tf*idf
13     doc += 1
```

```
{(0, 'fifti'): 0.005434960598980563,
 (0, 'go'): 0.0002906893990853149,
 (0, 'hundr'): 0.002570392381970895,
 (0, 'jim'): 0.005269857144642146,
 (0, 'nine'): 0.0008420698058556812,
 (0, 'nineti'): 0.001312716278834434,
 (0, 'north'): 0.021919902239379185,
 (0, 'one'): 0.0003992734536048051,
 (0, 'prentic'): 0.008085184948722846,
 (0, 'thousand'): 0.0008961984314476824,
 (0, 'three'): 0.0015785688576535318,
 (0, 'west'): 0.0033256596840258424,
 (1, 'fox'): 0.11198195635330804,
 (1, 'sli'): 0.11239056533822733,
 (1, 'stori'): 0.0007682063585522353,
```

Merging the TF-IDF of title and body

document = body + title

$\text{TF-IDF}(\text{document}) = \text{TF-IDF}(\text{body}) * (\alpha) + \text{TF-IDF}(\text{title}) * (1-\alpha)$

$\alpha = 0.3$

Flow :

- Calculate DF
- Calculate TF-IDF for Body for all docs
- Calculate TF-IDF for title for all docs
- Merging the TF-IDF of title and body

So, finally, we have a dictionary `tf_idf` which has the values as a (doc, token) pair

Ranking: Matching Score

We need to check in every document if these query_tokens exist and if the query_tokens exists, then the tf_idf value is added to the matching score of that particular doc_id. Format of tf_idf is (doc, token).



```
1 query_weights = {}
2 for key in tf_idf:
3     if key[1] in query_tokens:
4         try:
5             query_weights[key[0]] += tf_idf[key]
6         except:
7             query_weights[key[0]] = tf_idf[key]
8 query_weights = sorted(query_weights.items(), key=lambda x: x[1], reverse=True)
```

Query: One day, he noticed that one of the peacocks had dropped a feather. When the

Tokens: 36

Matching Score

429	c:\Users\kk910\OneDrive\Desktop\IR\data\vaincrow.txt	Matching Score 0.4
186	c:\Users\kk910\OneDrive\Desktop\IR\data/foxncrow.txt	Matching Score 0.21
405	c:\Users\kk910\OneDrive\Desktop\IR\data/tailbear.txt	Matching Score 0.1
339	c:\Users\kk910\OneDrive\Desktop\IR\data/quarter.c11	Matching Score 0.08
351	c:\Users\kk910\OneDrive\Desktop\IR\data/quarter.c5	Matching Score 0.05
416	c:\Users\kk910\OneDrive\Desktop\IR\data/the-tree.txt	Matching Score 0.03
291	c:\Users\kk910\OneDrive\Desktop\IR\data/musibrem.txt	Matching Score 0.02
176	c:\Users\kk910\OneDrive\Desktop\IR\data/fish.txt	Matching Score 0.02
294	c:\Users\kk910\OneDrive\Desktop\IR\data/narciss.txt	Matching Score 0.02
264	c:\Users\kk910\OneDrive\Desktop\IR\data/lionbird	Matching Score 0.01

Vectorising Documents

Used `total_vocab` (list of unique tokens) to generate an index for each token, and we will use numpy array of size `(docs, total_vocab)` to store the document vectors.



```
1 D = np.zeros((N, len(total_vocab)))
2 for i in tf_idf:
3     try:
4         ind = total_vocab.index(i[1])
5         D[i[0]][ind] = tf_idf[i]
6     except:
7         pass
8
9 row, col = D.shape
10 print(row, col)
```

Ranking: Cosine Similarity

Measure of similarity between two non-zero vectors in a multi-dimensional space. It measures the cosine of the angle between the two vectors, which ranges from -1 to 1.



```
1 d_cosines = []
2
3 query_vector = gen_vector(query_tokens)
4
5 for d in D:
6     d_cosines.append(cosine_sim(query_vector, d))
7
8 d_cosines_sorted=sorted(d_cosines)[-k:][::-1]
```



```
1 def gen_vector(tokens):
2
3     Q = np.zeros((len(total_vocab)))
4
5     counter = Counter(tokens)
6     words_count = len(tokens)
7
8     query_weights = {}
9
10    for token in np.unique(tokens):
11
12        tf = counter[token]/words_count
13        df = doc_freq(token)
14        idf = math.log((N+1)/(df+1))
15
16        try:
17            ind = total_vocab.index(token)
18            Q[ind] = tf*idf
19        except:
20            pass
21    return Q
```

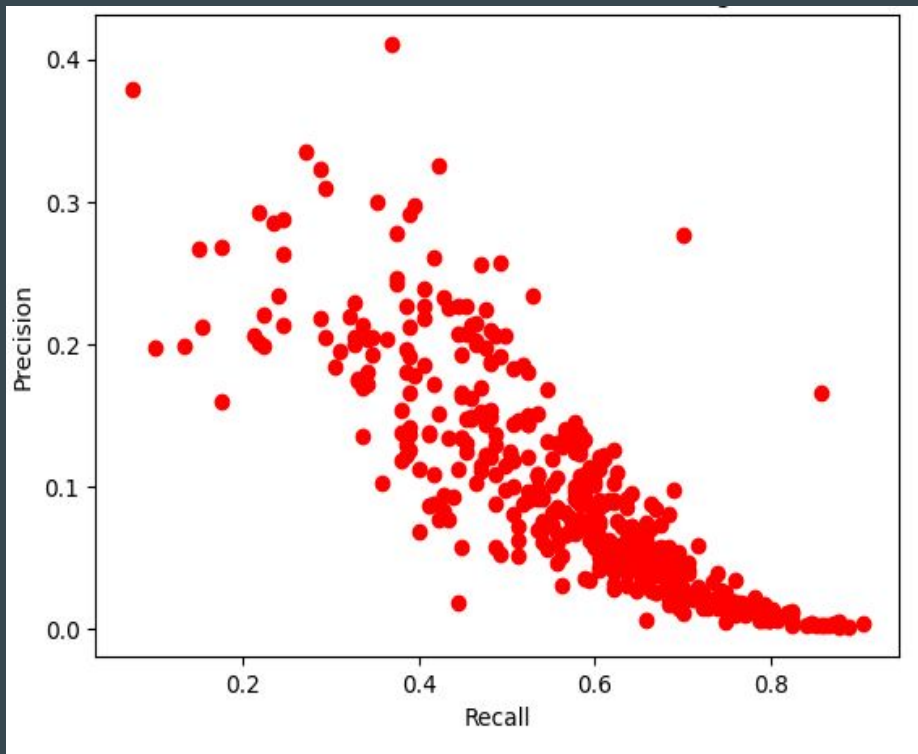
Cosine Similarity

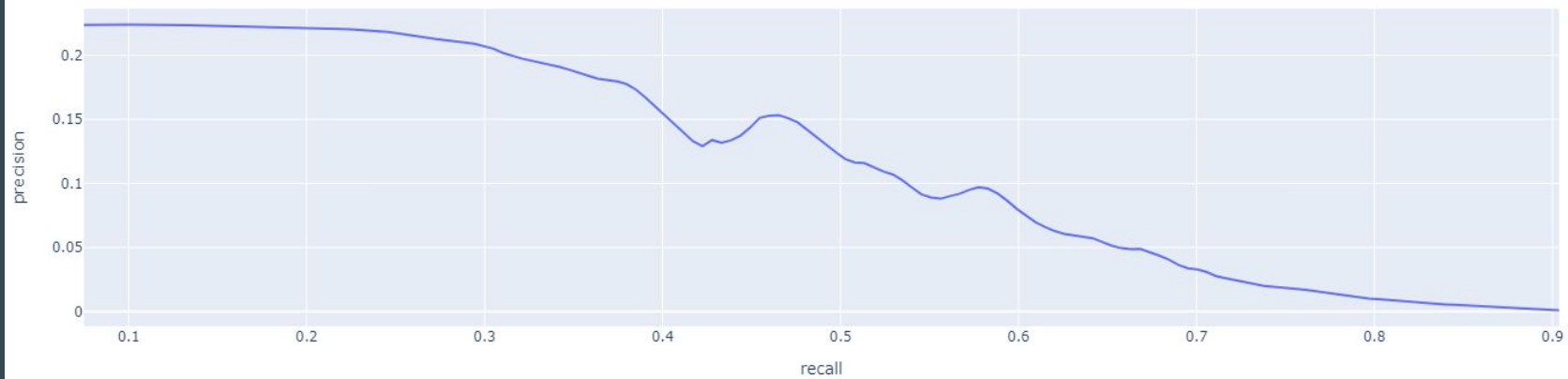
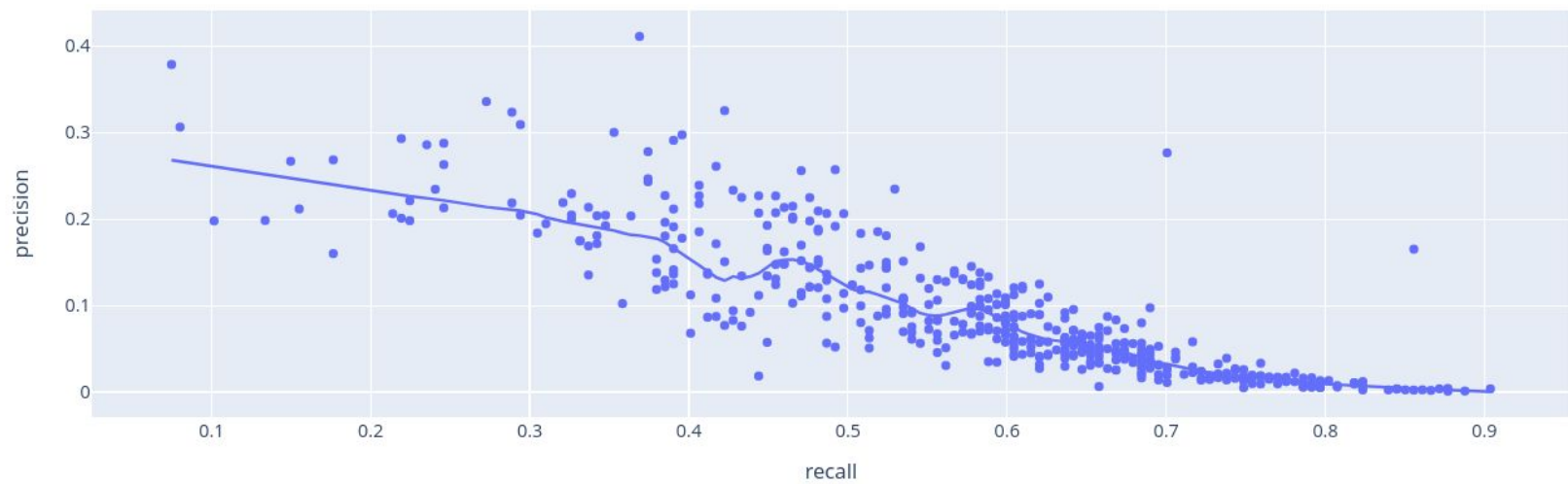
Query: One day, he noticed that one of the peacocks had dropped a feather. When the sun

Tokens: 44

429	c:\Users\kk910\OneDrive\Desktop\IR\data\vaincrow.txt	Cosine Similarity 0.49
186	c:\Users\kk910\OneDrive\Desktop\IR\data/foxncrow.txt	Cosine Similarity 0.13
27	c:\Users\kk910\OneDrive\Desktop\IR\data/aesopa10.txt	Cosine Similarity 0.1
229	c:\Users\kk910\OneDrive\Desktop\IR\data/home.fil	Cosine Similarity 0.07
405	c:\Users\kk910\OneDrive\Desktop\IR\data/tailbear.txt	Cosine Similarity 0.06
416	c:\Users\kk910\OneDrive\Desktop\IR\data/the-tree.txt	Cosine Similarity 0.06
339	c:\Users\kk910\OneDrive\Desktop\IR\data/quarter.c11	Cosine Similarity 0.04
194	c:\Users\kk910\OneDrive\Desktop\IR\data/frogp.txt	Cosine Similarity 0.04
26	c:\Users\kk910\OneDrive\Desktop\IR\data/aesop11.txt	Cosine Similarity 0.04
450	c:\Users\kk910\OneDrive\Desktop\IR\data/yukon.txt	Cosine Similarity 0.03

Comparison: Mixed Query from doc(429) + doc(439)





Matching Score

Query: One day, he noticed that one of the peacocks had dropped a feather. When the sun went down,
Tokens: 98

429	c:\Users\kk910\OneDrive\Desktop\IR\data\vaincrow.txt	Matching Score 0.44
186	c:\Users\kk910\OneDrive\Desktop\IR\data/foxncrow.txt	Matching Score 0.22
439	c:\Users\kk910\OneDrive\Desktop\IR\data/weeprncs.txt	Matching Score 0.18
400	c:\Users\kk910\OneDrive\Desktop\IR\data/sucker.txt	Matching Score 0.11
405	c:\Users\kk910\OneDrive\Desktop\IR\data/tailbear.txt	Matching Score 0.11
309	c:\Users\kk910\OneDrive\Desktop\IR\data/omarsheh.txt	Matching Score 0.1
161	c:\Users\kk910\OneDrive\Desktop\IR\data/fantas.hum	Matching Score 0.08
339	c:\Users\kk910\OneDrive\Desktop\IR\data/quarter.c11	Matching Score 0.08
382	c:\Users\kk910\OneDrive\Desktop\IR\data/sleprncs.txt	Matching Score 0.07
293	c:\Users\kk910\OneDrive\Desktop\IR\data/myeyes	Matching Score 0.05

Cosine Similarity

Query: One day, he noticed that one of the peacocks had dropped a feather. When the sun w
Tokens: 98

429	c:\Users\kk910\OneDrive\Desktop\IR\data\vaincrow.txt	Cosine Similarity 0.37
439	c:\Users\kk910\OneDrive\Desktop\IR\data/weeprncs.txt	Cosine Similarity 0.24
382	c:\Users\kk910\OneDrive\Desktop\IR\data/sleprncs.txt	Cosine Similarity 0.1
186	c:\Users\kk910\OneDrive\Desktop\IR\data/foxncrow.txt	Cosine Similarity 0.1
27	c:\Users\kk910\OneDrive\Desktop\IR\data/aesopa10.txt	Cosine Similarity 0.08
161	c:\Users\kk910\OneDrive\Desktop\IR\data/fantas.hum	Cosine Similarity 0.06
229	c:\Users\kk910\OneDrive\Desktop\IR\data/home.fil	Cosine Similarity 0.06
405	c:\Users\kk910\OneDrive\Desktop\IR\data/tailbear.txt	Cosine Similarity 0.05
450	c:\Users\kk910\OneDrive\Desktop\IR\data/yukon.txt	Cosine Similarity 0.05
416	c:\Users\kk910\OneDrive\Desktop\IR\data/the-tree.txt	Cosine Similarity 0.04

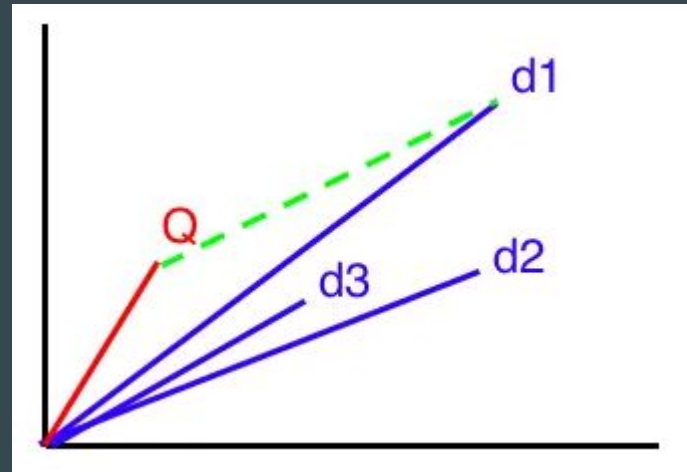
Conclusion

- Mixed Query : doc(429) + doc(439)
 - Matching Score shows 429 , 186 , 439 , 400 , ..
 - Cosine Similarity shows 429, 439, 382, 186 , ..
- Matching Score gives relevant documents but it quite fails when we give long queries, it will not be able to rank them properly.
- Vector POV : Matching Score computes Manhattan distance (straight line from tips) while Cosine score considers the angle of the vectors.

Conclusion (contd..)

- Matching Score will return document d3 but that is not very closely related.
- Cosine Similarity will return document d1

Cosine similarity learns the context more. In every type of query (small , medium , large) cosine similarity ranking is better than matching score ranking.



Thank you