

IRWA Project Part II: Indexing and Evaluation

Statement

Indexing

1. **Build inverted index:** After having pre-processed the data, you can then create the inverted index.

HINT - you may use the vocabulary data structure, like the one seen during the Practical Labs:

```
{  
Term_id_1: [document_1, document_2, document_4],  
Term_id_2: [document_1, document_3, document_5, document_6],  
etc...  
}
```

Documents information: Since we are dealing with conjunctive queries (AND), each of the returned documents should contain all the words in the query.

2. **Propose test queries:** Define five queries that will be used to evaluate your search engine (e.g., "covid pandemic", "covid vaccine")

HINT: How to choose the queries? The selection of the queries is up to you but it's suggested to select terms based on the popularity (keywords ranked by term frequencies or by TF-IDF, etc...).

3. **Rank your results:** Implement the TF-IDF algorithm and provide ranking based results.

Evaluation

- There will be 2 main evaluation components:
 1. A baseline with 3 queries and the ground truth files for each query will be given to you, using a subset of documents from the dataset.
 - a. Query 1: Landfall in South Carolina
 - b. Query 2: Help and recovery during the hurricane disaster
 - c. Query 3: Floodings in South Carolina
 2. You will be the expert judges, so you will be setting the ground truth for each document and query in a binary way for the test queries that you defined in step 2 at the indexing stage.
- For the prior evaluation components you must evaluate your algorithm by using different evaluation techniques and only for the second component (your queries) comment in each of them how they differ, and which information gives each of them:

- **Precision@K (P@K)**
- **Recall@K (R@K)**
- **Average Precision@K (P@K)**
- **F1-Score**
- **Mean Average Precision (MAP)**
- **Mean Reciprocal Rank (MRR)**
- **Normalized Discounted Cumulative Gain (NDCG)**

• Choose one vector representation, TF-IDF or word2vec, and represent the tweets in a two-dimensional scatter plot through the T-SNE (T-distributed Stochastic Neighbor Embedding) algorithm. To do so, you may need first to represent the word as a vector, and then the tweet, i.e., resulted as the average value over the words involved. Any other option rather than T-SNE may be used, but needs to be justified.

HINT: You don't have to know all the theoretical details used in T-SNE, just use the proper library and generate the output and play with it.

Also, you can choose to perform an alternative method to generate a 2-dimensional representation for the word embeddings (like PCA).

Here some T-SNE examples which may be good guidelines for the task:

1. <https://towardsdatascience.com/google-news-and-leo-tolstoy-visualizing-word2vec-word-embeddings-with-t-sne-11558d8bd4d>
2. <https://towardsdatascience.com/visualizing-word-embedding-with-pca-and-t-sne-961a692509f5>
3. <https://stackoverflow.com/questions/40581010/how-to-run-tsne-on-word2vec-created-from-gensim>

GitHub Repository

All the code and resources for the project will be submitted to the following repository:
<https://github.com/homexiang3/IRWA-2022-u172769-u172801>

And the repository TAG is :

IRWA-2022-u172769-u172801-part-2

Code development

Previous work considerations

This is the second part of the IRWA 2022 project so we skip some previous code documentation provided in the first report such as Google Drive connection, how to load JSON data and how we preprocess the data.

For further details visit:

<https://github.com/homexiang3/IRWA-2022-u172769-u172801/blob/main/Project/P1/IRWA-2022-u172769-u172801-part-1.pdf>

Adjustments

For this lab we added some packets used for mathematical calculations. The list of all the packets used in this project is placed below:

```
import nltk
nltk.download('stopwords')
nltk.download('punkt')
from collections import defaultdict
from array import array
from nltk.stem import PorterStemmer
from nltk.corpus import stopwords
import collections
import json
import re
from tabulate import tabulate
stemmer = nltk.stem.SnowballStemmer('english')
stopwords = set(stopwords.words('english'))
# Packets needed for lab 2
import math
import numpy as np
import collections
import pandas as pd
from numpy import linalg as la
```

We implemented the map provided in file “*tweet_document_ids_map.csv*”, using panda dataframe to load the data in the file.

```
id_map= pd.read_csv("/content/drive/Shared drives/IRWA/PROJECT/data/tweet_document_ids_map.csv", sep='\t',
                    engine='python', names = ["doc_id", "tweet_id"])
id_map.head()
```

	doc_id	tweet_id
0	doc_1	1575918182698979328
1	doc_2	1575918151862304768
2	doc_3	1575918140839673873
3	doc_4	1575918135009738752
4	doc_5	1575918119251419136

Then, in our Tweet class we added a new attribute name `doc_id`, that stores the `doc_id` value from our dataframe based on the map between the `tweet_id` from our JSON and the dataframe.

```
doc_id = id_map.loc[id_map['tweet_id'] == lines[i]['id'], 'doc_id'].iloc[0]
```

Also, we decided to modify the preprocess function to return the list of terms preprocessed instead of the whole sentence, since we just need the terms for this part of the project. Left image is the code used in part 1 and right image is the code used in part 2.

```
# Preprocess text
def preprocess(text):
    text = text.replace('\n', '')
    text = remove_emojis(text)
    text = remove_punctuation(text)
    text = remove_numbers(text)
    text = remove_white_space(text)
    words = nltk.tokenize.word_tokenize(text)
    words = [stemmer.stem(word) for word in words]
    words = remove_stopwords(words)
    words = remove_https(words)
    text = " ".join(words)
    return text
```

```
# Preprocess text
def preprocess(text):
    text = text.replace('\n', '')
    text = remove_emojis(text)
    text = remove_punctuation(text)
    text = remove_numbers(text)
    text = remove_white_space(text)
    words = nltk.tokenize.word_tokenize(text)
    words = [stemmer.stem(word) for word in words]
    words = remove_stopwords(words)
    words = remove_https(words)
    return words
```

Finally, we deleted some debug / checking parts of code since we already proved the correctness of our code in the first part and we want to maintain a clear code.

Indexing

Create Index

```
# Create index function
def create_simple_index(tweets):
    index = defaultdict(list)
    tweet_index = {} # dictionary to map tweet id with index in tweets list
    counter = 0 # keep track of index inside tweets
    for t in tweets: # For each tweet

        tweet_id = t.id
        terms = preprocess(t.tweet) #preprocess tweet and return list of terms
        tweet_index[tweet_id] = counter # Save original tweets position with tweet id to recover all the information
        counter = counter + 1 # Move to next tweets position

    current_page_index = {}
    for position, term in enumerate(terms): # Loop over all terms
        try:
            # if the term is already in the index for the current page (current_page_index)
            # append the position to the corresponding list
            current_page_index[term][1].append(position)
        except:
            # Add the new term as dict key and initialize the array of positions and add the position
            current_page_index[term] = [tweet_id, array('I', [position])] # 'I' indicates unsigned int (int in Python)

    # merge the current page index with the main index
    for term_page, posting_page in current_page_index.items():
        index[term_page].append(posting_page)
    return index, tweet_index
```

```
# Apply index function for all the tweets
index, tweet_index = create_simple_index(tweets)
# Print first 10 results of word 'hurrican' (stemmed word)
print("First 10 Index out of",len(index['hurrican']),"results for the term 'hurrican': {}".format(index['hurrican']))
# Print first 10 results of word 'Hurricane'(not stemmed word)
print("First 10 Index out of",len(index['Hurricane']),"results for the term 'Hurricane': {}".format(index['Hurricane']))

First 10 Index out of 796 results for the term 'hurrican': [[1575918105854984192, array('I', [7, 11]), [1575918088473788429, array('I', [7])],
First 10 Index out of 0 results for the term 'Hurricane': []
```

Search Query

```
def search(query, index):
    query = preprocess(query) #create list of query terms (each term is preprocessed to match terms in index)
    docs = set()
    for term in query:
        try:
            # store in term_docs the ids of the docs that contain "term"
            term_docs = [posting[0] for posting in index[term]]
            # docs = docs Union term_docs
            docs |= set(term_docs)
        except:
            #term is not in index
            pass
    docs = list(docs)
    return docs
```

Display Query

We create a loop of 5 queries, for each of them we execute the function search to retrieve the documents matching the terms and we prepare our visualization table for our top 5 tweets.

```
# Define 5 queries to visualize - display top 5 tweets (without any rank or order)
# example used in our report: 1. covid pandemic (15) 2. hurricane ian (1087) 3. south carolina (354) 4. god bless (45) 5. help victims (417)
for i in range(5):
    print("Insert your query (i.e.: 'covid pandemic'):\n")
    query = input()
    docs = search(query, index)
    top = 5
    visualization_tweets = []
    #create table headers
    headers = ['DOC_ID', 'ID', 'TWEET', 'USERNAME', 'DATE', 'HASHTAGS', 'LIKES', 'RETWEETS', 'URL']
    print("\n=====Sample of {} results out of {} for the searched query:\n".format(top, len(docs)))
    #create table of tweets for each match
    for d_id in docs[:top]:
        t = tweet_index[d_id]
        visualization_tweets.append(tweets[t])
    #print table
    print(tabulate(visualization_tweets, headers=headers, tablefmt='grid'))
```

Results

We use the set of 5 queries “covid pandemic”, “hurricane ian”, “south carolina”, “god bless” and “help victims”. The results in this case seems to be low relevant for the queries, given that we are not applying any weighting schema (first matches of any of the terms are retrieved first).

```
covid pandemic
=====
Sample of 3 results out of 15 for the searched query:
```

DOC_ID	ID	TWEET
doc_917	1575907538192674816	Naziesque mutilation of children PAYS WELL in Liberal Hospitals https://t.co/haccgu7KQe #news #BTC,#God #bb24 #USA #NYC #Bitcoin #COVID19 #Russia,#MAGAt,#gay,#NFTs,Alito #Metaverse,Masks,#Nuclear,#Biden,#audioleak
doc_652	1575910184794304513	Mut ila t ion of children PAYS WELL in Liberal Hos pital s https://t.co/haccgu7KQe #news #BTC,#God #bb24 #USA #NYC #Bitcoin #COVID19 #Russia,#MAGAt,#gay,#NFTs,Alito #Metaverse,Masks,#Nuclear,#Biden,#audioleak
doc_1115	1575904328891518978	Can U live w/o ur phone? Nuclear war will show u how. US NATO Ukraine Russia https://t.co/AkpF680rLD #news #BTC,#God #bb24 #USA #NYC #Bitcoin #COVID19 #Russia,#MAGAt,#gay,#NFTs,#Metaverse Masks #Nuclear,#Biden,#audioleak,#NATO

Sergi Solís Valdivieso u172769 - Hong-ming Xiang Vico u172801

hurricane ian

Sample of 3 results out of 1087 for the searched query:

DOC_ID	ID	TWEET
doc_6	1575918105854984192	Ace Handyman Services hopes everyone was safe during the Hurricane. Any damages caused by the hurricane is our first priority. #HurricaneIan #AHS #BringingHelpfulToYourHome https://t.co/BfoQ97tJE9
doc_18	1575917983062380545	#BREAKING Hurricane Ian has made landfall in #SouthCarolina. The storm reformed into a hurricane over the Atlantic
doc_24	1575917833573179392	Just heard from my husband. He and his unit are in North Port assisting the Fire Department and rescuing people. TF

south carolina

Sample of 3 results out of 354 for the searched query:

DOC_ID	ID	TWEET
doc_10	1575918057037303808	How pissed is GOD to send #HurricaneIan to Florida and South Carolina!?
		The #MAGA cult has angered GOD and are paying for their sins.
		@RonDeSantisFL @scgovernorpress #Florida #SouthCarolina #MyrtleBeach
doc_30	1575917717600681984	#HurricaneIan Ian makes 3rd landfall S of Georgetown, South Carolina as Cat. 1 hurricane; widespread damage, e #wotwitter #scwx WPDE: https://t.co/EUV6TAMq8i https://t.co/hH0KK4uaUM
doc_36	1575917617281658880	@itsbethbooker Hi Beth. So happy & relieved to see that you finally got to hug your mom! Lots of love &

god bless

Sample of 3 results out of 45 for the searched query:

DOC_ID	ID	TWEET
doc_10	1575918057037303808	How pissed is GOD to send #HurricaneIan to Florida and South Carolina!?
		The #MAGA cult has angered GOD and are paying for their sins.
		@RonDeSantisFL @scgovernorpress #Florida #SouthCarolina #MyrtleBeach
doc_246	1575915002573205505	God reacts to "Don't Say Gay Bill" and the treatment of humans.
		#HurricaneIan
doc_567	1575911109998964736	Releasing merchandise early from @impressink to help #HurricaneIan victims. Half the profits

help victims

Sample of 3 results out of 417 for the searched query:

DOC_ID	ID	TWEET
doc_20	1575917943426535424	@xfinitysupport is busy doing nothing to help those of us who have been affected by #HurricaneIan. Our
doc_27	1575917773376540672	In the aftermath of a Disaster some just can't resist taking advantage of the vulnerability.
		Disaster related scams can happen to anyone, awareness helps reduce the chances it's you.
		#HurricaneIan #PuertoRico #Florida #Georgia #SouthCarolina https://t.co/4ncfVvmu6F
doc_42	1575917468996239360	If you need to be out on the road, #OnStar Advisors are here to help with routing assistance. Just pus

TF-IDF Implementation

```

# Create index function
def create_tfidf_index(tweets, num_tweets):
    index = defaultdict(list)
    tweet_index = {} # dictionary to map tweet id with index in tweets list
    counter = 0 # keep track of index inside tweets
    tf = defaultdict(list) # term frequencies of terms in documents (documents in the same order as in the main index)
    df = defaultdict(int) # document frequencies of terms in the corpus
    idf = defaultdict(float) # inverse document frequency for each term
    for t in tweets: # for all tweets

        tweet_id = t.id
        terms = preprocess(t.tweet) #preprocess tweet and return list of terms
        tweet_index[tweet_id] = counter # Save original tweets position with tweet id to recover all the information
        counter = counter + 1 # Move to next tweets position
        current_page_index = {}

        for position, term in enumerate(terms):
            try:
                # if the term is already in the dict append the position to the corresponding list
                current_page_index[term][1].append(position)
            except:
                # Add the new term as dict key and initialize the array of positions and add the position
                current_page_index[term] = [tweet_id, array('I', [position])] # 'I' indicates unsigned int (int in Python)

        # normalize term frequencies
        # Compute the denominator to normalize term frequencies (formula 2 above)
        # norm is the same for all terms of a document.
        norm = 0
        for term, posting in current_page_index.items():
            # posting will contain the list of positions for current term in current document.
            # posting ==> [current_doc, [list of positions]]
            # you can use it to infer the frequency of current term.
            norm += len(posting[1]) ** 2
        norm = math.sqrt(norm)

        #calculate the tf(dividing the term frequency by the above computed norm) and df weights
        for term, posting in current_page_index.items():
            # append the tf for current term (tf = term frequency in current doc/norm)
            tf[term].append(np.round(len(posting[1]) / norm, 4)) ## SEE formula (1) above
            #increment the document frequency of current term (number of documents containing the current term)
            df[term] += 1 # increment DF for current term

        # Compute IDF
        for term in df:
            idf[term] = np.round(np.log(float(num_tweets / df[term])), 4)

        #merge the current page index with the main index
        for term_page, posting_page in current_page_index.items():
            index[term_page].append(posting_page)

    return index, tf, df, idf, tweet_index

```



```

def rank_documents(terms, docs, index, idf, tf, title_index):

    # I'm interested only on the element of the docvector corresponding to the query terms
    # The remaining elements would become 0 when multiplied to the query_vector
    doc_vectors = defaultdict(lambda: [0] * len(terms)) # I call doc_vectors[k] for a nonexistent key
    query_vector = [0] * len(terms)

    # compute the norm for the query tf
    query_terms_count = collections.Counter(terms) # get the frequency of each term in the query.

    query_norm = la.norm(list(query_terms_count.values()))

    for termIndex, term in enumerate(terms): #termIndex is the index of the term in the query
        if term not in index:
            continue
        # query_vector[termIndex]=idf[term] # original
        ## Compute tf*idf(normalize TF as done with documents)
        query_vector[termIndex] = query_terms_count[term] / query_norm * idf[term]

        # Generate doc_vectors for matching docs
        for doc_index, (doc, postings) in enumerate(index[term]):
            if doc in docs:
                doc_vectors[doc][termIndex] = tf[term][doc_index] * idf[term]

    # Calculate the score of each doc
    # compute the cosine similarity between queryVector and each docVector:

    doc_scores = [[np.dot(curDocVec, query_vector), doc] for doc, curDocVec in doc_vectors.items()]
    doc_scores.sort(reverse=True)
    result_docs = [x[1] for x in doc_scores]
    result_rank = [x[0] for x in doc_scores] #get rank value
    #print document titles instead of document id's
    #result_docs=[ title_index[x] for x in result_docs ]
    if len(result_docs) == 0:
        print("No results found, try again")
        query = input()
        docs = search_tfidf(query, index)
    return result_docs, result_rank

def search_tfidf(query, index):
    query = preprocess(query)#create list of query terms (each term is preprocessed to match terms in index)
    docs = set()
    for term in query:
        try:
            # store in term_docs the ids of the docs that contain "term"
            term_docs = [posting[0] for posting in index[term]]

            # docs = docs Union term_docs
            docs |= set(term_docs)
        except:
            #term is not in index
            pass
    docs = list(docs)
    ranked_docs, ranked_score = rank_documents(query, docs, index, idf, tf, tweet_index)#rank docs
    return ranked_docs, ranked_score

```

Similarly as we did in the simple case, we prepare a loop of 5 queries but now we search for the top ranked docs. Notice that we also take the ranked score to print it on screen and check if the result is correct.


```
# Define 5 queries to visualize - top 3 ranked tweets displayed
# example used in our report: 1. covid pandemic (15) 2. hurricane ian (1087) 3. south carolina (354) 4. god bless (45) 5. help victims (417)
for i in range(5):
    print("Insert your query (i.e.: 'covid pandemic'):\n")
    query = input()
    ranked_docs, ranked_score = search_tfidf(query, index)
    top = 3
    visualization_tweets = []
    #create table headers
    headers = ['DOC_ID', 'ID', 'TWEET', 'USERNAME', 'DATE', 'HASHTAGS', 'LIKES', 'RETWEETS', 'URL']
    print("\n=====Sample of {} results out of {} for the searched query:\n".format(top, len(ranked_docs)))
    #create table of tweets for each match
    for d_id in ranked_docs[:top]:
        t = tweet_index[d_id]
        visualization_tweets.append(tweets[t])
    #print ranked score
    print("Ranked Scores:", ranked_score[:top])
    #print table
    print(tabulate(visualization_tweets, headers=headers, tablefmt='grid'))
```

Results

We run again the set of 5 queries “covid pandemic”, “hurricane ian”, “south carolina”, “god bless” and “help victims”. The results in this case are much more accurate of the queries since we are taking into account TF-IDF weights, specially, we can notice how short tweets with query words appear in the top 3 most relevant documents. We can also confirm that ranked scores are in descending order.

covid pandemic

=====

Sample of 3 results out of 15 for the searched query:

Ranked Scores: [12.160545944094801, 6.044344086303483, 6.044344086303483]

DOC_ID	ID	TWEET
doc_3088	1575868406603350017	This is so heartbreaking. Fl was one of my last trips to the US in 2019, b
		#Florida #HurricaneIan https://t.co/pr5npit3ce
doc_1414	1575899131154821120	#Hurricanes, #COVID19, #disability, #ableism, government neglect, #climate
doc_2796	1575871626540818432	#MentalHealth matters more than ever, especially given #COVID and #Hurrica

hurricane ian

=====

Sample of 3 results out of 1087 for the searched query:

Ranked Scores: [2.1383822503267504, 1.8517338503002754, 1.8517338503002754]

DOC_ID	ID	TWEET	USERNAME
doc_634	1575910361298968576	Hurricane IAN #Ian #HurricaneIan #HurricaneIan #Huracan #HuracanIan #Hurricane https://t.co/Hb1l04Q3v8	cesarharamillo
doc_640	1575910304159977472	Hurricane Ian before and after #HurricaneIan https://t.co/XZstkI2oH2	kadenfields8
doc_1217	1575902689040666626	Hurricane Ian on tour 🌪️ #HurricaneIan	besmarterpeople

south carolina

=====

Sample of 3 results out of 354 for the searched query:

Ranked Scores: [5.750153366279658, 4.979349988118859, 4.8059638892671845]

DOC_ID	ID	TWEET	USERNAME
doc_254	1575914929898782720	South Carolina #HurricaneIan https://t.co/YA4dFUC2v	webgyr12
doc_174	1575915969913839616	South Carolina #HurricaneIan here we go	TheAstuteGaloot
doc_493	1575912058163408896	Just south of Myrtle Beach in South Carolina. #HurricaneIan #Ian #Sciix https://t.co/ErHrSXSc00	Damian_vix

god bless

=====

Sample of 3 results out of 45 for the searched query:

Ranked Scores: [11.008039858203505, 11.008039858203505, 10.557137146207896]

DOC_ID	ID	TWEET
doc_3782	1575859119042416640	Good morning Patriotsus If you have the ability to help, please join me and support the https://t.c
doc_3800	1575858935281881088	My thoughts and prayers go out special to those affected by the hurricane Ian.. God Bless. 🙏 #Florida #HurricaneIan https://t.co/8dcUI9MRA7
doc_2372	1575877434712334336	@KellyClarksonTV I just wanted to let @kellyclarkson know that we are praying for those who are aff

help victims

=====

Sample of 3 results out of 417 for the searched query:

Ranked Scores: [8.030304493960108, 6.862258668688771, 6.349543088247527]

DOC_ID	ID	TWEET
doc_321	1575914189071138818	A list of ways you can help the victims of #HurricaneIan https://t.co/D77WkvKhI
doc_632	1575910393767133184	@DonaldJTrumpjr #MAGA doesn't believe in #climatechange so how can you politicize these folks? Trump's would just
doc_3057	1575868689768857601	#HurricaneIan -> How to help victims of Hurricane Ian - CBS News https://t.co/3D7qU9A79

Evaluation

Precision@K (P@K)

```
def precision_at_k(doc_score, y_score, k=10): #binary relevance, predicted relevance, k for a given query
    """
    Parameters
    -----
    doc_score: Ground truth (true relevance labels).
    y_score: Predicted scores.
    k : number of doc to consider.

    Returns
    -----
    precision @k : float

    """
    order = np.argsort(y_score)[::-1] #we get the ranking of the documents according to the predicted score/ use np
    doc_score = np.take(doc_score, order[:k]) # align the binary relevance to the corresponding document / use the i
    relevant = sum(doc_score == 1) #get number of relevant documents
    return float(relevant) / k #calculate precision at k, which is the number of relevant documents retrieved at k
```

Recall@K (R@K)

```
def recall_at_k(doc_score, y_score, k=10): #binary relevance, predicted relevance, k for a given query
    """
    Parameters
    -----
    doc_score: Ground truth (true relevance labels).
    y_score: Predicted scores.
    k : number of doc to consider.

    Returns
    -----
    recall @k : float
    """
    order = np.argsort(y_score)[::-1] #we get the ranking of the documents according to the predicted score/ use np.argsort and [::-1]
    doc_score = np.take(doc_score, order[:k]) # align the binary relevance to the corresponding document / use the indexes of point 1 to k
    all_relevants = sum(np.take(doc_score) == 1) # take all relevant docs
    relevant = sum(doc_score == 1) #get number of relevant documents
    return float(relevant) / all_relevants #calculate recall at k, which is the number of relevant documents among all relevant docs
```

Average Precision@K (P@K)

```
def avg_precision_at_k(doc_score, y_score, k=10): #binary relevance, predicted relevance, k for a given query
    """
    Parameters
    -----
    doc_score: Ground truth (true relevance labels).
    y_score: Predicted scores.
    k : number of doc to consider.

    Returns
    -----
    average precision @k : float
    """
    gtp = np.sum(doc_score == 1) #Total number of gt positives
    order = np.argsort(y_score)[::-1] #same as for precision
    doc_score = np.take(doc_score, order[:k]) #same as for precision
    ## if all documents are not relevant
    if gtp == 0:
        return 0
    n_relevant_at_i = 0
    prec_at_i = 0
    for i in range(len(doc_score)):
        if doc_score[i] == 1: #only add the P@k when the doc is relevant
            n_relevant_at_i += 1
            prec_at_i += n_relevant_at_i / (i + 1) #calculate P@K (#docs relevant at k/k)
    return prec_at_i / gtp #return ap
```

F1-Score

```
def f1_score(precision, recall):
    return 2*(precision*recall)/(precision+recall)

f1 = f1_score(0.5, 1)
print(f1)
```

Mean Average Precision (MAP)

```
def map_at_k(search_res, k=10): #receives all the search results dataframe containing all the queries and the results and relevances
    """
    Parameters
    -----
    search_res: search results dataset containing:
        query_id: query id.
        doc_id: document id.
        predicted_relevance: relevance predicted through LightGBM.
        doc_score: actual score of the document for the query (ground truth).

    Returns
    -----
    mean average precision @ k : float
    """
    avp = []
    for q in search_res["query_id"].unique(): # loop over all query ids
        curr_data = search_res[search_res["query_id"] == q] # select data for current query (get a slice of the dataframe keeping only the data for the current query)
        avp.append(avg_precision_at_k(np.array(curr_data["is_relevant"]),
                                           np.array(curr_data["predicted_relevance"]), k)) #append average precision for current query
    return np.sum(avp) / len(avp), avp # return mean average precision
```

Mean Reciprocal Rank (MRR)

```
def rr_at_k(doc_score, y_score, k=10):
    """
    Parameters
    -----
    doc_score: Ground truth (true relevance labels).
    y_score: Predicted scores.
    k : number of doc to consider.

    Returns
    -----
    Reciprocal Rank for current query
    """

    order = np.argsort(y_score)[::-1] # get the list of indexes of the predicted scores in descending order
    doc_score = np.take(doc_score, order[:k]) # sort the actual relevance label of the documents according to the predicted scores
    if np.sum(doc_score) == 0: # if there are not relevant documents return 0
        return 0
    return 1 / (np.argmax(doc_score == 1) + 1) # hint: to get the position of the first relevant document
```

Normalized Discounted Cumulative Gain (NDCG)

```
def dcg_at_k(doc_score, y_score, k=10): #doc_score are the labels (ground truth)
    order = np.argsort(y_score)[::-1] # get the list of indexes of the predicted
    doc_score = np.take(doc_score, order[:k]) # sort the actual relevance labels
    gain = 2 ** doc_score - 1 # First we calculate the upper part of the formula
    discounts = np.log2(np.arange(len(doc_score)) + 2) # Compute denominator (n)
    return np.sum(gain / discounts) #return dcg@k

def ndcg_at_k(doc_score, y_score, k=10):
    dcg_max = dcg_at_k(doc_score, doc_score, k) #ideal dcg
    #print(dcg_max)
    if not dcg_max:
        return 0
    return np.round(dcg_at_k(doc_score, y_score, k) / dcg_max, 4)
```

Vector Representation