

Information Retrieval of a *Harry Potter* Dataset

Lino Virgen
Texas Tech University

Patrice Fote
Texas Tech University

Abstract—Using a dataset of summaries of Harry Potter chapters, this report describes in detail the processes required to perform queries. The dataset will go through a series of stages such as tokenization, filtering, lemmatization and stemming to prepare it for processing. Afterward, each term of an input query was ranked according to its frequency and inverse document frequency. Finally, a tf-idf score was calculated for each document. A cosine similarity matrix was also produced to illustrate the similarities between the documents.

Keywords: *information retrieval, tf-idf, data analysis, Harry Potter, query*

I. INTRODUCTION

Information retrieval is done every single second all around the globe thanks to search engines such as Google. These search engines look up words and sentences through millions and millions of websites. Then, they present the results to the user according to a unique match score. This report will present a number of stages that are intended to retrieve and analyze a dataset of Harry Potter summary chapters utilizing Python and a number of imported modules.

Section II will describe the dataset that was used along the project and the reasons of its choosing.

Section III will describe the steps took toward ranking documents given an input query.

II. DATASET

The dataset used to complete the task set for our project is a collection of summaries to the chapters of the books in the Harry Potter Series By author JK Rowling's with a concentration on the 1st 3 books in the series of 7 namely: Harry Potter: The Philosopher's Stone, Harry Potter: The Chamber of Secrets, Harry Potter: The Prisoner of Azkaban. It consists of 57 documents, each containing on average 1000 to 1376 words which is the reason we chose to work with it. Most chapters were obtained from SparkNotes diligently split into individual *txt* files in order to facilitate processing.

III. TF-IDF RANKING

The importance increases proportionally to the number of times a word appears in the document but is offset by the frequency of the word in the corpus. Variations of the tf-idf weighting scheme are often used by search engines as a central tool in scoring and ranking a document's relevance given a user query. One of the simplest ranking functions is

computed by summing the tf-idf for each query term; many more sophisticated ranking functions are variants of this simple model.

A. Pre-Processing

Data pre-processing means preparing the data for analysis. Most unstructured datasets consist of long texts without a formal structure. In this case, all the Harry Potter chapter summaries are long paragraphs that are easy to read for humans, but not for computers. Therefore, each document must be “cleaned” so we can extract the desired information. To achieve it, each document passed to a series of four processes. An imported module, *nltk* or natural language toolkit, was used in every preprocessing stage.

First, each document was tokenized. Tokenization is the process of segmenting strings in small units. Hence, each document was broken down into words and these words were put into an array. At this point, every word can be worked with individually for the following processes.

Second, each document was filtered. As most texts, punctuation and stop words (i.e. and, for, to, of, etc.) are abundant and comprise a big percentage the total word count. Moreover, they are non-relevant to the performed search. Therefore, they were all filtered from each document array.

Third, each document was lemmatized. This is process consists of changing every word to its base form. For example, “better” would be changed to “good” and “playing” would be changed to “play”. This step is necessary so that a word can be found regarding of its form.

Forth, each document was stemmed. Stemming consists of changing a word to a root form. Similar to lemmatization, stemming tries to get rid of derived words and produce terms that are easier to work with. For our project, the Porter method was used to stem every word in the array.

B. Query

In order to test the tf-idf ranking project, a query of 11 terms was came up with. The query was “magic owl wand study train school learn class fear teacher friend”. In this manner, the ranker would check every term in the query and determine the term frequency and rarity. These processes are defined in the following subsection.

C. Term Frequency and Inverse Document Frequency

According to Jin [1], the term frequency or tf “of term t in document d is the number of times that t occurs in d .” This property shows quantitatively how abundant a term is in each document. Although a high term frequency denotes relevance, it does not mean a document is more informative to the user. Hence the inverse document frequency is also calculated. Document frequency or df is “the number of documents N that contain t .” The larger df is, the less important that term is. This is due to the commonality of the term and the irrelevance it has compared to other rare and key terms.

The term frequency of each term in the query was calculated by simply counting how many times a term appeared in each document. As shown in Figure 1, the most frequent term among all the documents was *magic* while the least frequent was *study*.

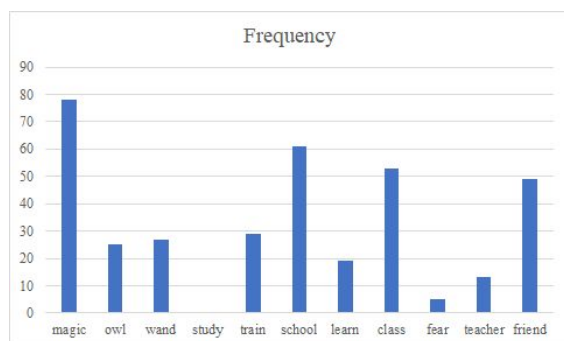


Figure 1. Total term frequency of each term in the query.

The inverse document frequency of each document was calculated as follows. First, for each term, the number of documents where it appeared was summed. Then the idf of each term was calculated utilizing the formula:

$$idf = \log_{10}(N/df)$$

The log function’s function is to “dampen” the magnitude of idf since the rarity of a term does not grow proportionally with the number of documents matched. Figure 2 is a visualization of the idf scores. Thus, it can be noted that the term *fear* had a higher “rarity” score compared to the rest.

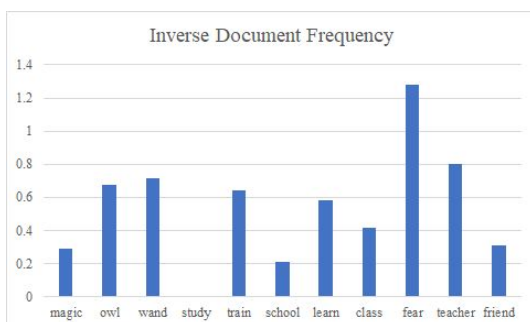


Figure 2. Inverse document frequency of each term in the query.

D. TF-IDF

Also called the term frequency-inverse document frequency, it can be seen as a statistical measure used to evaluate how important a word is to a document in a collection or corpus. In our example, we can observe with the rarity and frequency of the word *fear* in the Harry Potter Chapter summary dataset. In order to calculate the tf - idf weight, we first compute the normalized Term Frequency (TF) (the number of times a word appears in a document) divided by the total number of words in that document. With an average number of words per chapter being 1000 words, we measure how frequent a term occurs in that document. As seen in Figure 1, the rarity of words like *study* or *fear* is low while the frequency of *magic* is low. Since every document is different in length, it is possible that a term would appear much more times in long documents like chapters 5 and 57 would display a higher probability for the words *fear* or *teacher*, for instance than shorter ones as seen in short chapters like chapters 3 and 51. Thus, the term frequency is divided by the document length (the total number of terms in the document) as a way of normalization; the second part is the Inverse Document Frequency (IDF), computed as the logarithm of the number of the documents in the corpus divided by the number of documents where the specific term appears like the words *magic*, *class*, *friend* or *school* seen in Figure 1 to have the highest occurrence in our dataset which measures how important a term is. Whereas in computing TF, all terms are considered equally important in finding the inverse document frequency it is known that certain terms, such as “is”, “of”, and “that”, may appear a lot of times but have little importance, therefore, our setup will need to weigh down the frequent terms while scaling up the rare ones like *fear*, *study*, and *teacher*, by computing the following: $IDF(t) = \log_{10}(\text{Total number of documents} / \text{Number of documents with term } t \text{ in it})$

It can be noted in Figure 3 that chapter 5 has the greatest matching score. Chapters 6, 32, and 57 are close behind. On the other hand, chapters 15 and 29 show no relevance to the query. For example, consider a chapter in our dataset document 57 containing 1000 words wherein the word *study* appears 0 times. The term frequency (i.e., tf) for *study* is then $(0 / 1000) = 0$. Now, assume we have 56 more documents and the word *study* appears in 1 of these. Then, the inverse document frequency (i.e., idf) is calculated as $\log_{10}(57 / 1) = 1.75$. Thus, the Tf - idf weight is the product of these quantities: $0 * 1.75 = 0$ as seen in Figure 1 and Figure 2.

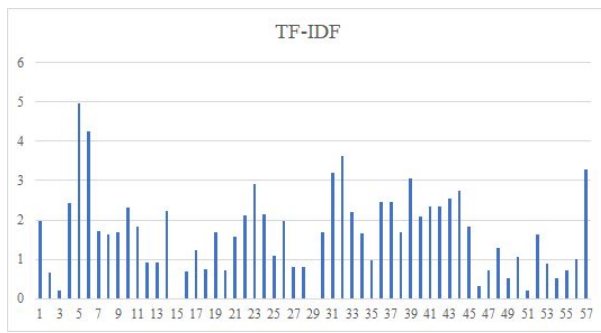


Figure 3. TF-IDF score of each document in the collection.

E. Vector Space

We model each document in vector space in order to facilitate tf-idf transformations. For that we calculate the tf-idf vector for the query, and compute the score of each document relative to this query, we multiply the tf scores by the idf values of each term, obtaining the matrix in *Figure 4*: (All the terms appeared only once in each document in our small collection, so the maximum value for normalization is 1.22 on the word fear” as seen in *Figure 2*. When computing the tf-idf values for the query terms we divide the frequency by the maximum frequency calculated and multiply with the idf values. TF-IDF values for each word of the document in our example and the N is the dimension of the vectors:

$$\vec{a} \cdot \vec{b} = \|\vec{a}\| \|\vec{b}\| \cos \theta$$

$$\cos \theta = \frac{\vec{a} \cdot \vec{b}}{\|\vec{a}\| \|\vec{b}\|}$$

As observed, the definition of the dot product is a simple multiplication of each component from both vectors added together

F. Cosine Similarity

The cosine similarity between two vectors (or two documents on the Vector Space) is a measure that calculates the cosine of the angle between them. This metric is a measurement of orientation and not magnitude, it can be seen as a comparison between documents on a normalized space because we’re not taking into the consideration only the magnitude of each word count (tf-idf) of each document, but the angle between the documents. What we have to do to build the cosine similarity equation is to solve the equation of the dot product for all vectors involved.

Figure 4 illustrates the similarity of all 57 chapters with each other. It can be observed that all chapters in the middle have a high or medium similarity while chapters at the beginning and at the end show very low similarity.

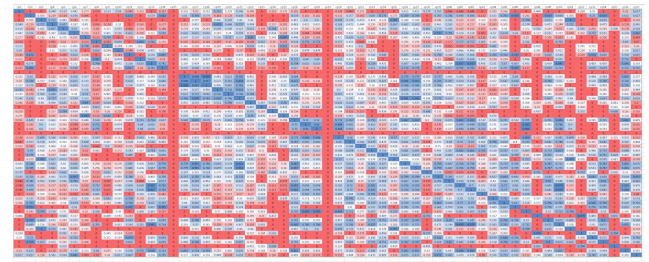


Figure 4. Heatmap showing the cosine similarity of the 57 documents against each other. Blue denotes high similarity, white denotes medium similarity and red denotes low similarity.

REFERENCES

- [1] F. Jin, “TF-IDF and Vector Space.” CS 5364 - Information Retrieval. PowerPoint. Computer Science Department of Texas Tech University. Fall 2018