

Introduction to TF-IDF: To Represent Importance of Keyword within whole Dataset

Dipti D. Mehare¹, Prof. A. V. Deorankar²

¹, P. G. Scholar, Department of Computer Science, Government college of Engineering, Amravati, Maharashtra, India

²Associate Professor, Department of Computer Science, Government college of Engineering, Amravati, Maharashtra, India

Abstract: In this paper, we examine the results of applying Term Frequency Inverse Document Frequency (TF-IDF) to the whole dataset for representation of importance of a keyword. 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. TF-IDF calculates values for each word in a document through an inverse proportion of the frequency of the word in a particular document to the percentage of documents the word appears in. In this process term weight algorithm plays an important role. It greatly interferes the precision and recall results of the natural language processing (NLP) systems. Currently, TF-IDF term weight algorithm is widely applied into language models to build NLP Systems. NLP is a natural language parser that works out the grammatical structure of sentences.

Keywords: TF, IDF, Term Weight, Natural Language Processing.

I. INTRODUCTION

Nowadays, TF-IDF is the most widely used term weight algorithm. In the process of document formalization, documents are represented by document vectors which are expected to indicate as much information of the documents as possible. To make the representation accurate, term weight algorithm plays an important role in the process.

TF-IDF is one of the most commonly used term weighting algorithms in today's information retrieval systems. Two parts of the weighting were proposed by Gerard Salton and Karen Spärck Jones respectively. TF, the term frequency, also called Local Term Weight, is defined as the number of times a term in question occurs in a document. IDF, the inverse document frequency, also called Global Term Weight, is based on counting the number of documents in the collection being searched that are indexed by the term. The product of TF and IDF, known as TF-IDF, is used as an indicator of the importance of a term in representing a document.

A. How to Compute TF-IDF

The TF-IDF weight is composed by two terms: the first computes 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. The second term 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.

- 1) **TF:** Term Frequency, which measures how frequently a term occurs in a document. Since every document is different in length, it is possible that a term would appear much more times in long documents than shorter ones. Thus, the term frequency is often divided by the document length. $TF(t) = (\text{Number of times term } t \text{ appears in a document}) / (\text{Total number of terms in the document})$
- 2) **IDF:** Inverse Document Frequency, which measures how important a term is. While computing TF, all terms are considered equally important. However it is known that certain terms, such as "is", "of", and "that", may appear a lot of times but have little importance. Thus we need to weigh down the frequent terms while scale up the rare ones, by computing the following: $IDF(t) = \log_e(\text{Total number of documents} / \text{Number of documents with term } t \text{ in it})$

II. AN OVERVIEW OF TF-IDF

TF is defined as the number of times term t occurs in the document d , it is given as $t_{f,d}$ and Inverse Document Frequency Estimate the rarity of a term in the whole document collection. If the term occurs in all the document then IDF will be zero. IDF is given as,

$$idf_i = \log \frac{|D|}{|\{j : t_i \in d_j\}|}$$

Where, D is the total number of documents and $\{|j : t_i \in d_j\}$: number of documents where the term t_i appears. Now we can combine the term frequency and inverse document frequency, to produce a composite weight for each term in each document. TF-IDF is defined as the product of its tf weight and its idf weight. And it is given as $tf-idf_{t,d} = tf_{t,d} \times idf_t$. In other words, $tf-idf_{t,d}$ assigns to term t a weight in document d that is highest when t occurs many times within a small number of documents (thus lending high discriminating power to those documents); lower when the term occurs fewer times in a document, or occurs in many documents (thus offering a less pronounced relevance signal); 3. lowest when the term occurs in virtually all documents. Now we can view each document as a vector with one component corresponding to each term in the dictionary, together with a weight for each component. If the term do not occur in dictionary then weight of that term is zero and this vector form will prove to be crucial to scoring and ranking. Scoring in the document d is the sum of all query terms, of the number of times each of the query terms occurs in document d . We can refine this idea so that we add up not the number of occurrences of each query term t in d , but instead the $tf-idf$ weight of each term in d .

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

The well-known equation to compute the relevance score as there are numerous variations of the TF-IDF weighting scheme is given below,

$$Score(w, F_d) = \frac{1}{|F_d|} (1 + \ln f_{d,w}) \ln(1 + \frac{N}{f_w}),$$

where w denotes the given keyword, $|F_d|$ is the length of file F_d , $f_{d,w}$ denotes the TF of w in file, f_w denotes the number of files containing w , and N denotes the total number of files in the collection.

A. Drawbacks of TF

TF-IDF term weight algorithm is widely applied into language models to build NLP Systems. Since TF-IDF only takes term frequency into consideration, therefore it has the following drawbacks.

- 1) TF algorithm calculates term weight based on their frequency. That is, term weight is positive correlated to their frequency. Term with higher frequency may be only intensively distributed in a part of the document.
- 2) The intuitive meaning of IDF algorithm is that terms which rarely occur in a collection of documents which are valuable and The importance of each term is assumed to be inversely proportional to the number of documents that the term occurs. However, obviously, the term which occurs widely in the document collection but intensively appears in a few documents much probably represents the topic of a document category and is significant for text classifying.
- 3) Empty terms and function terms, including conjunctive, preposition, some adverbs, auxiliary term, modal particles, are usually existed with high frequency. Which leads to inaccurate weight assignments to such terms.

B. TF-IDF improvements

The improvement in TF-IDF concentrated on two topics, which are given below.

- 1) By taking additional term statistical information into consideration
- 2) Introducing additional techniques into this field.

III. EXPERIMENT

A. Data Formatting and Collecting

We can test our TF-IDF implementation on a collection 1400 of documents. These documents were gathered from a larger collection of documents from the database. The documents were encoded with the SGML text format, so we decided to leave in the formatting tags to account for noisy data and to test the robustness of TF-IDF. We simulated more noise by enforcing case-sensitivity. Due to certain constraints, we had to limit the number of queries used to perform information retrieval on to 86. We calculate TF-IDF weights for these queries according to equation, and then return the first 100 documents that maximize equation. The returned documents were returned in descending order, with documents with higher weight sums appearing first. To compare our results, we also performed in parallel the brute force method of performing query retrieval based only on the term f_w, d . Since

this approach would simply return documents where non relevant words appear most. We will provide evidence that TF-IDF, though relatively simple, is a big improvement over this.

B. Experimental Results

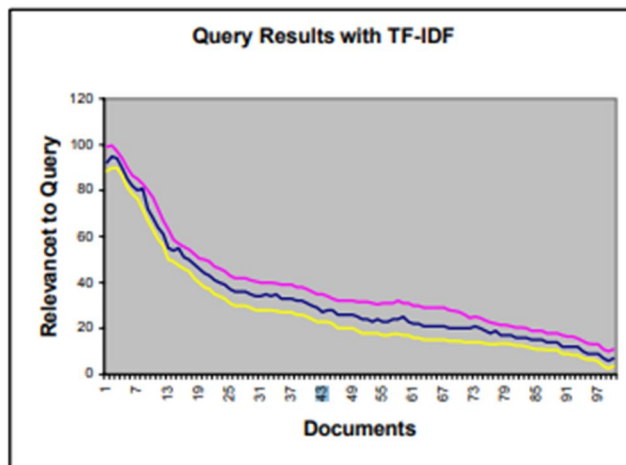


Figure 1. Results of retrieval with TF-IDF on our data.

IV. CONCLUSIONS AND FUTURE RESEARCH

In this paper, we introduce TF-IDF i.e Term Frequency and Inverse Document Frequency. Also examine the result of applying TF-IDF to the whole dataset for representation of importance of a keyword. But the author felt the results were not considered significant, the paper shows that there is still interest in enhancing the simple TF-IDF scheme. This paper gives some improvement in TF and IDF algorithm respectively through introducing term distribution data into term weighting research. By examining our data, the easiest way for us to enhance TF-IDF would be case-sensitivity and equate words. Future research might also include TF-IDF to performing searches in documents written in a different language than the query. Enhancing the already powerful TF-IDF algorithm would increase the success of query retrieval systems, which have quickly risen to become a key element of present global information exchange.

REFERENCES

- [1]. Juan Ramos," Using TF-IDF to Determine Word Relevance in Document Queries", Department of Computer Science, Rutgers University, 23515 BPO Way, Piscataway, NJ, 08855.
- [2]. Mingyong Liu and Jiangang Yang," An improvement of TFIDF weighting in text categorization", 2012 International Conference on Computer Technology and Science (ICCTS 2012) IPCSIT vol. 47 (2012) © (2012) IACSIT Press, Singapore.
- [3]. Tian Xia, Yanmei Chai," An Improvement to TF-IDF: Term Distribution based Term Weight Algorithm", JOURNAL OF SOFTWARE, VOL. 6, NO. 3, MARCH 2011.
- [4]. Zhangjie Fu, Member, IEEE, Xinle Wu, Qian Wang," Enabling Central Keyword-based Semantic Extension Search over Encrypted Outsourced Data", IEEE TRANSACTIONS ON INFORMATION FORENSICS AND SECURITY.
- [5]. Handout of the paper "Exploring the Similarity Space" by Zobel and Moffat
- [6]. Scoring, term weighting and the vector space model, Online edition (c) 2009 Cambridge UP