

Extractive Summarization of Text Using TF-IDF

Kanchan Lalchandani¹, Rekha Jain^{2,*}

Abstract

Natural language getting ready a subfield of Artificial intelligence (AI) which facilities round how machines could realize and system human language. Content summarization is one of the big elements of Natural Language preparation (NLP). The objective of the synopsis is to create a shorter form of a unique book by safeguarding the significance and the key substance of the first report. This research work is set a calculation of Extractive Text Summarization, that is, TF-IDF. TF-IDF (term recurrence reverse archive recurrence) is examined in detail utilizing calculation and model. A significant spotlight is on how we can apply the calculation on various reports. An elegantly composed outline can fundamentally lessen the amount of work expected to process a lot of content.

Keywords: Text summarization, TF-IDF, Metrix, Natural Language preparation (NLP), Artificial intelligence (AI)

INTRODUCTION

A synopsis is a book that is extricated from at least one archive, that passes on large statistics of the file and it' is a precise model of the primary record. Rundown is valuable for us in this day and age as we have the main part of information accessible on the web and it gets hard to snatch and comprehend the best possible importance of the data that we need. A bridged content should comprise one of a kind sentences. Content outline is, where information is given as a content report and it restores the rundown of a book. The synopsis ought to be short and exact [1]. These days, content outline is crucial and extremely helpful gratitude to its different sorts of utilizations simply like the synopsis of books, digest (synopsis of stories), the stock trade, news, highlights-(meeting, event, sport), abstract of scientific papers, newspaper articles, magazines, etc. There are many techniques or algorithms which can be used to summarize the data but this study is concentrated on TF-IDF. TF-IDF may be a technique that quantifies a word in documents, we compute a frequency of every word which signifies the importance of the word within the document. It provides those keywords which are used to identify some specific documents or categories [2].

TEXT SUMMARIZATION

Text summarization is that the technique for generating a specific summary of voluminous data that

*Author for Correspondence

Rekha Jain
E-mail: rekha.jain@poornima.org

¹Student, Department of Computer Science and Engineering, Poornima Institute of Engineering and Technology, Jaipur, Rajasthan, India

²Professor, Department of Computer Science and Engineering, Poornima Institute of Engineering and Technology, Jaipur, Rajasthan, India

Received Date: August 22, 2021

Accepted Date: October 24, 2021

Published Date: December 24, 2021

Citation: Kanchan Lalchandani, Rekha Jain. Extractive Summarization of Text Using TF-IDF. Journal of Artificial Intelligence Research & Advances. 2021; 8(3): 21–25p.

concentrate on sections that contain useful information, and without converting the general which suggests of the text. Text summarizers are able to extract beneficial data that leaves out unimportant statistics from the document. Summarization can beautify the clarity and understandability of documents, and additionally reduces the time spent in finding out information [3]. Most of this statistics is redundant, insignificant, and wish to now not bring the meant meaning. For instance, if you are trying to find specific information from a web article, you will need to search tons through its content and spend plenty of some time removing unnecessary data before getting the knowledge you would like. It

aims to convert lengthy documents into shortened versions, something which could become far more difficult and dear if done manually. In today's advancing world, volumes of data are increasing and it is very difficult to read and understand the required data in less time. It is a task to collect the required information and then convert it into a summarized form [4]. It saves time and helps in avoiding retrieving massive text. Therefore, text summarization came into demand.

Types of Text Summarization

Extractive Text Summarization

Extractive outlines are made by reusing words, sentences, and so on of the primary content archive. The framework removes content from the entire assortment, without changing the content archive. A large portion of the content rundown frameworks uses extraction strategies to supply a synopsis. We use sentence extraction strategies to supply extraction outlines [5].

Extractive summarization method is often partitioned into stages:

1. Pre-Processing: Pre-processing is an organized portrayal of the first content.
2. Preparing: Processing-highlights affecting the significance of sentences are chosen and determined and when loads are doled out to those highlights. The last score of each sentence is chosen. Top positioned sentences are chosen for definite synopsis.

Abstractive Text Summarization

Abstractive synopsis calls for profound information and questioning over the content. It gives its own synopsis over info content without a similar word or sentence in the information content [6]. The principle point of abstractive content synopsis is to deliver an important abbreviated adaptation of the first record. The adjective abstractive is utilized because it denotes that the generated summary is not a combination or selection of some repeated sentences. Abstractive summarization is a totally hard hassle aside from machine translation [7]. Abstractive summarization is better and provides quality than extractive summarization as it takes data from multiple documents and then generates a precise summary in own words/sentences.

Abstractive summarization is again achieved in two ways. They are:

- a. Structure-based approach, and
- b. Semantic-based approach.

EXTRACTIVE SUMMARIZATION METHODS

- a. Term Frequency-Inverse Document Frequency (TF-IDF) method.
- b. Cluster-based method.
- c. Graph-theoretic approach.
- d. Machine Learning approach.
- e. Text summarization with neural networks.
- f. Automatic text summarization based on fuzzy logic.

Term Frequency-Inverse Document Frequency (TF-IDF)

TF-IDF stands for "*Term Frequency-Inverse Document Frequency*". The TF-IDF weight is applied in facts healing and content material mining. This weight might be a measure that is utilized to check how significant a word is to a record during an assortment or corpus [8]. The significance of a word expands relatively to the number of times that word shows up inside the archive, however is balanced by the recurrence of the word inside the corpus.

TF-IDF = Term Frequency (TF) * Inverse Document Frequency (IDF)

Wording:

t: term (word),

d: record (set of words),

N: check of the corpus,
corpus: the whole archive set.

Term Frequency: This estimates the recurrence of a phrase at some stage in an archive. We will define TF as follows.

$$tf(t, d) = \text{include of } t \text{ in } d / \text{number of words in } d \quad (1)$$

Document Frequency: DF is the wide variety of news at some stage in which the phrase is available.

$$df(t) = \text{event of } t \text{ in records} \quad (2)$$

Inverse Document Frequency: Inverse record recurrence (IDF) is the manner thrilling or unusual a phrase is.

$$idf(t) = \log(N / (df + 1)) \quad (3)$$

At last, through taking a multiplicative estimation of TF and IDF, we get the TF-IDF rating as:

$$tf-idf(t, d) = tf(t, d) \times \log(N / (df + 1)) \quad (4)$$

Now we will see the implementation and result of this algorithm, in which an input text document is given and its summary is calculated with this method.

THE 9-STEPS IMPLEMENTATION

1. Tokenize the sentences:

```
sentences = sent_tokenize(text)
```

```
total_documents = len(sentences)
```

2. Create the Frequency matrix:

```
freq_matrix = _create_frequency_matrix(sentences)
```

3. Calculate TF and generate matrix:

```
tf_matrix = _create_tf_matrix(freq_matrix)
```

4. Create a table for documents per words:

```
count_doc_per_words = _create_documents_per_words(freq_matrix)
```

5. Calculate IDF and generate matrix:

```
idf_matrix = _create_idf_matrix(freq_matrix, count_doc_per_words, total_documents)
```

6. Calculate TF-IDF and generate matrix:

```
tf_idf_matrix = _create_tf_idf_matrix(tf_matrix, idf_matrix)
```

7. Score the sentences:

```
sentence_scores = _score_sentences(tf_idf_matrix)
```

8. Find the threshold:

```
threshold = _find_average_score(sentence_scores)
```

9. Generate summary:

```
summary = _generate_summary(sentences, sentence_scores, 1.3 * threshold)
```

```
return summary
```

INPUT AND OUTPUT

Input Document

The people who are resilient stay in the game longer.

"On the mountains of reality you may in no way circulate futile: possibly you will arrive at some degree above today, else you will be preparing your forces all together that you will be prepared to move higher tomorrow." —Friedrich Nietzsche. Challenges and mishaps are not intended to

overcome you, yet advance you. Be that as it may, I understand following a couple of long stretches of annihilations, it can pulverize your soul and it is simpler to present than hazard further difficulties and dissatisfactions [9]. Have you at any point encountered this previously? To be absolutely forthright, I do not have the right responses. I cannot mention to you what the best possible strategy is; just you will know. Notwithstanding, it is essential to not be disheartened by disappointment while seeking an objective or a fantasy, since disappointment itself implies different things to various individuals [10]. To a person with a firm mindset, disappointment might be a hit to their confidence, yet to a person with a growth mind-set, it is an opportunity to upgrade and find better approaches to beat their hindrances. Same disappointment, yet various reactions [11, 12]. Who is valid and who is not right? Not one or the other. Everybody includes an alternate mentality that chooses their result. Those that are strong remain inside the game longer and draw on their inward way to succeed.

Output Summary

Have you at any point encountered this previously? Who is valid and who is not right? Not one or the other.

RESULT

It is clearly visible that this summary is an extractive summary as it has taken the sentences from the input document only and has not made any changes in it. The sentences are presented exactly the way they are present in the original input document. Hence this extractive summarization is straightforward to enforce and really beneficial to generate a précis for any quantity of data.

Advantages

- The straightforward calculation for coordinating words inside the inquiry to record applicable to address.
- It is very efficient.
- From the research, many scholars have proven the TF-IDF document which is very relevant to question.
- Easy to compute the problem.
- Have some essential measurements to remove the graphic report.
- Effectively parent the similitude among the reviews utilizing it.

Disadvantages

- TF-IDF is completely founded on BOW (pack of words) and does not catch the situation in semantics, co-events in the report, and so forth.
- The TF-IDF is hired due to the lexical stage feature.
- They cannot capture the semantic.

Limitations

- Used to figure the similitude straightforwardly in word-check which might be delayed for vocabularies.
- They expect the tally of various words that give free proof of the comparability.
- It utilizes similarities between the words.

CONCLUSION

The TF-IDF calculation is clear to execute and is amazingly incredible. This calculation is normally looked at and gives precise outcomes since it is unmistakably demonstrated during this study. There is a few impediments yet we will say it is one among the least difficult calculations for Extractive Text Summarization. In this day and age of gigantic information, we require some new procedures for handling and creating precise synopses. Numerous analysts have proposed an improved kind of TF-IDF calculation alluded to as Adaptive TF-IDF. In future research, the planet goes to observe some new strategies which will beat the limitations of TF-IDF and give the best and exact outcomes varying.

REFERENCES

1. Kanitha DK, Muhammad Noorul Mubarak D, Shanavas SA. Survey on Text Summarization Methods. *International Journal of Computer Engineering & Technology (IJCET)*. 2018 Jan–Feb; 9(1): 26–36. Article ID: IJCET_09_01_004.
2. Samrat Babar, et al. Text Summarization: An Overview. 2013. Available from https://www.researchgate.net/publication/257947528_Text_SummarizationAn_Overview
3. Shahzad Qaiser, Ramsha Ali. Text Mining: Use of TF-IDF to Examine the Relevance of Words to Documents. *International Journal of Computer Applications*. 2018;181(1):25-29.
4. Moratanch N, Chitrakala S. A Survey on Extractive Text Summarization. *IEEE International Conference on Computer, Communication, and Signal Processing (ICCCSP-2017)*. 2017; 1–6.
5. Ahmed Elrefaiy, Ahmed Rafat Abas, Ibrahim Elhenawy. Review of Recent Techniques for Extractive Text Summarization. *J Theor Appl Inf Technol*. 2018 Dec 15; 96(23): 7739–7759.
6. Towards Data Science. William Scott (15 Feb, 2019). TF-IDF from scratch in python on a real-world dataset [Online]. Available from <https://towardsdatascience.com/tf-idf-for-document-ranking-from-scratch-in-python-on-real-world-dataset-796d339a4089>
7. Towards Data Science. Akash Panchal (10 Jun, 2019). NLP — Text Summarization using NLTK: TF-IDF Algorithm [Online]. Available from <https://towardsdatascience.com/text-summarization-using-tf-idf-e64a0644ace3>
8. Chen F, Han K, Chen G. An approach to sentence-selection-based text summarization. In *TENCON'02 Proceedings; 2002 IEEE Region 10 Conference on Computers, Communications, Control, and Power Engineering*, IEEE. 2002; 1: 489–493.
9. Sankarasubramaniam Y, Ramanathan K, Ghosh S. Text summarization using Wikipedia. *Inf Process Manag*. 2014; 50(3): 443–461.
10. Mashechkin, Petrovskiy M, Popov D, Tsarev DV. Automatic text summarization using latent semantic analysis. *Program Comput Softw*. 2011; 37(6): 299–305.
11. Gupta V, Lehal GS. A survey of text summarization extractive techniques. *J Emerg Technol Web Intell*. 2010; 2(3): 258–268.
12. Hingu D, Shah D, Udmale SS. Automatic text summarization of Wikipedia articles. In: *IEEE International Conference on Communication, Information & Computing Technology (ICCICT)*, 2015. 2015; 1–4.