

ff

by Surekha Kumari

Submission date: 30-Apr-2023 11:56PM (UTC-0500)

Submission ID: 1851114041

File name: big_data.pdf (3.07M)

Word count: 1654

Character count: 8957

Document Summarization Using NLP

Sai Abhishek Thota

⁴ Computer Engineering
University of Texas at Dallas
Dallas, Texas

saiabhishek.thota@utdallas.edu

Saicharan Reddy Ragannagari

⁴ Computer Engineering
University of Texas at Dallas
Dallas, Texas

saicharanreddy.ragannagari@utdallas.edu

Sai Shantanu Nagelli

⁵ Computer Engineering
University of Texas at Dallas
Dallas, Texas

saishantan.nagelli@utdallas.edu

Surekha Kumari

Computer Engineering
University of Texas at Dallas
Dallas, Texas

surekha.kumari@utdallas.edu

Abstract—In this paper we discuss the implementation of a method to summarize documents irrespective of its size. With abundance of data and information, it is essential that there be a method to automatically summarize documents in order to save people the effort and time in selecting relevant information. Automating summarizing is crucial to the amount of information that can be processed when compared with manually summarizing and this method is also wary of human biases. This paper further discusses the algorithm used and implementation steps.

Index Terms—implementation, summarize, algorithm [1]

I. INTRODUCTION

We have worked on the Python implementation of the text summarization algorithm, which seeks to extract the most important sentence from a given paragraph and generate summaries. The algorithm uses (NLP) techniques such as word frequency, sentence relevance, and graph-based algorithms to identify the most salient information in the text. That extraction approach is the most widely used method of summarising text, given that it provides greater clarity and interpretation compared with abstract methods, which generate summaries based on paraphrasing and rewording. However, the whole meaning of the original text cannot be taken into account by extractive methods. we achieve this by using libraries such as spaCy and NLP library NLTK, we built a code around it to paraphrase long paragraphs into meaningful sentences to form a summary. This paper elaborates on the libraries used and algorithms implemented to achieve the desired result.

Index Terms—NLP, NLTK, spaCy

II. BACKGROUND WORK

A. NLP

NLP stands for Natural Language Processing, which is a subfield of Artificial Intelligence used to train machines to understand Natural languages used by humans. It helps machines understand human language to perform redundant tasks post automation.

B. NLTK

NLTK stands for Natural Language ToolKit which is a python library used to aid NLP.

C. SpaCy

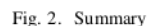
spaCy is a highly advanced, open source library used to perform NLP. It uses the best algorithms and is faster than NLTK when it comes to word tokenizing.

D. Stopwords

Stopwords are commonly used words which are irrelevant to the machine when constituting meaning from the text. It generally consists of articles, prepositions, conjunctions, pronouns, etc which are present in a text to aid human understanding. As machines do not need typical grammar, stopwords are filtered out for NLP.

III. THEORETICAL AND CONCEPTUAL STUDY

Text summarization, which can be used for a variety of applications in different areas, such as news summaries, documents, and chatbot responses, is an extensively discussed problem within the field of natural language processing NLP. The summarization of texts is made up of two main approaches: extractive and abstract. The best sentences or phrases in the original texts are selected by an extract summarization algorithm that uses them to produce a summary. These algorithms typically use statistical methods such as word frequency, sentence relevance, and graph-based algorithms to identify the most salient information. Generally, extraction methods are quicker to interpret than abstract methods, but they can't describe the whole meaning of a text. On the other hand, paraphrasing and rewording an existing text is also commonly done in abstract summarization algorithms that use deep learning models like neural networks to produce summaries. Abstractive methods require more computational resources and may introduce errors or biases but are able to provide better clarity and consistency of summaries. Before,



Line score is the sum of normalized frequency scores in a sentence. We are now going to get sentences with the top-line score. We wanted our algorithm to get 30% as a summary of the original document, so we selected the 30. We have used the nlargest function in the heapq library. nlargest will return the list of sentences with the amount of selected percentage we want.

```
1 summarizedata = nlargest(30, length, line_scores, key = line_scores.get)
2 summa = '\n'.join(summarizedata)

In[18]: Amarendra is accepted as guard at the royal palace of Kuntala while Bhallaladeva perceives about Amarendra's acts and upon chasing Devasen's parents, later f
or her and receives a promise from Kattappa, who ensures Bhallaladeva's marriage with Devasen's daughter. In the fact that Amarendra has already fallen in love with h
Lord Vijjaladeva, Amarendra's brother and the next in line for the throne, is denied the position due to his scheming nature, and as such Bhallaladeva's wife, Lady Siv
agami, assumes power with the intention of raising both her son Bhallaladeva and the orphaned Baahubali in an equal manner to select the next heir to the throne.
In Kuntala, as attack by Pandur's daughter on the royal palace ensues, Amarendra and Kattappa's bravery and they nullify the attack with assistance from Devasen's mother
at court's honor. Devasen and their father's love's story shall be threatened by Bhallaladeva's motivations and manipulates him into entering the palace at a point of sig
to for assassinating Bhallaladeva to do good.
In the present day, Sivudu's adoptive parents, impressed by Kattappa's story, wish to meet Bhallaladeva.
In the process, he is recognized by a woman, who then leads the other workers in chasing "Bhallaladeva," to Sivudu's confusion and Bhallaladeva's displeasure.
With Kattappa's and Devasen's assistance, the army lays siege to Mahishmati.
While Bhallaladeva kills Devasen, Bhallaladeva's valor and concern for the people of his kingdom convince Sivudu to make Bhallaladeva his apparent.
After vanquishing the Kalakeyas, Amarendra is declared as the heir apparent to the throne while Bhallaladeva is associated as the kingdom's future commander-in-chief.
He also enters his son to join the head of Bhallaladeva's statue out of the palace walls, where it is kept in the great waterfall.
The wife of the tribe's chieftain, Sanga, decides to adopt the boy and names him Sivudu.
Bhallaladeva sends his son, Bhadra, and the royal family's loyal slave Kattappa to rescue Devasen.
He falls into Amarendra's trap, protecting him to be "Bhallaladeva".
```

Fig. 3. line frequency calculation

The summary of amount of words in the actual paragraph and the amount of words in the summarized paragraph is shown below

VII. CONCLUSION AND FUTURE WORK

Our code is a result of frequency based approach to extract text summaries. Algorithm used identifies the most relevant sentences in a text and then extract those sentences to combine and form a summary of the whole input provided. In conclusion, though we do have limitations, by using extractive text summarization, this approach is quick in extracting most important information from a large body of text, we now know the significance/effectiveness of the NLP in producing exact summaries of huge paragraphs. This algorithm selects the most important sentence in the paragraphs and adds all those into one summary. Our current implementation is limited to summarization only has room for improvement in areas where more precise information is to be extracted. The project could be extended to create presentation from extracted summaries, post automated online reviews.

ACKNOWLEDGMENT

It is my sincere thanks to Professor Anurag Nagar for his invaluable assistance, valuable observations, and constant support in the course of this project. He was instrumental in shaping our understanding of the subject and helping us achieve our objectives through his expertise in Natural Language Processing. Our heartfelt thanks go out to our teacher's assistant, Truong Quang Pham, who has been doing everything he can to let us know our questions and help us keep track of the project.

We would like to thank our team for their contribution, SaiCharan Reddy, Sai Abhishek Thota, Sai Shantan Nagelli, and Surekha Kumari, who have played a major role in this project. This project's successful achievement was largely made possible through their dedication, hard work, and creativity. Lastly, we'd like to say our deepest thank you to each and everyone who has supported us in this endeavor directly or

```
1 print("Input Document")
2 print(in_file)

Input Document
Near the ancient Indian kingdom of Mahishmati, an injured woman named Sivagami ex
fore drowning, she holds the baby aloft and prays to Lord Shiva, explaining that
Lord Shiva. The wife of the tribe's chieftain, Sanga, decides to adopt the boy and
Sivudu grows up to be an ambitious and mischievous child, obsessed with ascending
Siva and places it at the foot of the mountain. A mask then falls from the cliffs,
oman named Avantika fighting Mahishmati soldiers. He discovers that she is a mem
n with Avantika and secretly follows her, even managing to draw a tattoo on her h
feelings and later they have sex.
After she quietly faints Sivudu and leaves, she gets attacked by more soldiers. At
e city on Bhallaladeva's birthday and assists in erecting a gigantic statue of the
rates the royal palace disguised as a soldier and distracts Bhallaladeva and his
s Bhadra as both the Amburi tribe and resistance warriors arrive. Kattappa lunges
The next morning, Kattappa reveals to Sivudu that Sivudu is actually Mahendra Baal
ther died giving birth to him. Lord Bijjaladeva, Vikramadeva's brother and the ne
llaladeva and the orphaned Baahubali in an equal manner to select the next heir t
comes beloved by the kingdom.
Command took 0.09 seconds -- by ssn210083@utdallas.edu at 4/30/2023, 8:15:12 PM on Project
Cwd 25
1 print("Summarized text of Input Document")
2 print(summarizedData)

Summarized text of Input Document
Amarendra is accepted as guard at the royal palace of Kuntala while Bhallaladeva g
the fact that Amarendra has already fallen in love with her. Lord Bijjaladeva, Vik
n of raising both her son Bhallaladeva and the orphaned Baahubali in an equal manne
ce from Devasena's maternal cousin Kumara Varma who overcomes his cowardice. that
hallaladeva to do good. In the present day, Sivudu's adoptive parents, impressed t
a's displeasure. With Kattappa's and Avantika's assistance, the army lays siege to
After vanquishing the Kalakeyas, Amarendra is declared as the heir apparent to th
s swept to the great waterfall. The wife of the tribe's chieftain, Sanga, decides
Bhallaladeva sends his son, Bhadra, and the royal family's loyal slave Kattappa t
Kattappa is revealed to be the side companion of Mahendra's father.
It breaks as it falls and crashes against the cliff's walls and lands near the li
llaladeva's treachery and exposes it to Shivagami, who regrets while Devasena deliv
river after being hit by an arrow shot by Bhallaladeva which leads her to a passag
he tries many times to scale the cliffs but always falls. Devasena, during her vi
```

Fig. 4. final output

```
Command took 0.09 seconds -- by ssn210083@utdallas.edu at 4/30/2023, 8:15:12 PM on Project
Cwd 26
1 len(in_file)

Out[22]: 10097
Command took 0.09 seconds -- by ssn210083@utdallas.edu at 4/30/2023, 8:15:12 PM on Project
Cwd 27
1 len(summarizedData)

Out[23]: 3227
Command took 0.10 seconds -- by ssn210083@utdallas.edu at 4/30/2023, 8:15:12 PM on Project
```

Fig. 5. length of summary

indirectly. Encouragement and feedback have been the source of motivation for us, and we appreciate your support.

3

REFERENCES

- [1] Dan W. Patterson, Introduction to Artificial Intelligence and Expert Systems , PHI, 2001, Chapter 12
- [2] <https://en.m.wikipedia.org/wiki/India>
- [3] <https://en.m.wikipedia.org/wiki/Dallas>
- [5] https://en.m.wikipedia.org/wiki/Spider-Man:_Far_From_Home
- [5] https://en.m.wikipedia.org/wiki/Baahubali:_The_Beginning
- [6] [https://en.m.wikipedia.org/wiki/Titanic_\(1997_film\)](https://en.m.wikipedia.org/wiki/Titanic_(1997_film))

ORIGINALITY REPORT

7%

SIMILARITY INDEX

3%

INTERNET SOURCES

4%

PUBLICATIONS

4%

STUDENT PAPERS

PRIMARY SOURCES

1	Larice Toko Lumanda. "Reflexivity in English, French and Kinshasa Lingala: Similarities and Differences", Communication and Linguistics Studies, 2019 Publication	1 %
2	Submitted to UT, Dallas Student Paper	1 %
3	krchowdhary.com Internet Source	1 %
4	Delva Culp. "Developmental apraxia and augmentative or alternative communication—a case example", Augmentative and Alternative Communication, 2009 Publication	1 %
5	Submitted to Mentone Grammar Student Paper	1 %
6	Submitted to Kingston University Student Paper	1 %
7	www.analyticsvidhya.com Internet Source	1 %

Exclude quotes Off

Exclude matches

< 3 words

Exclude bibliography Off