

AI and ML Techniques for Cyber Security

Assignment 1: Text Summarization using Python & NLTK: TF-IDF Algorithm

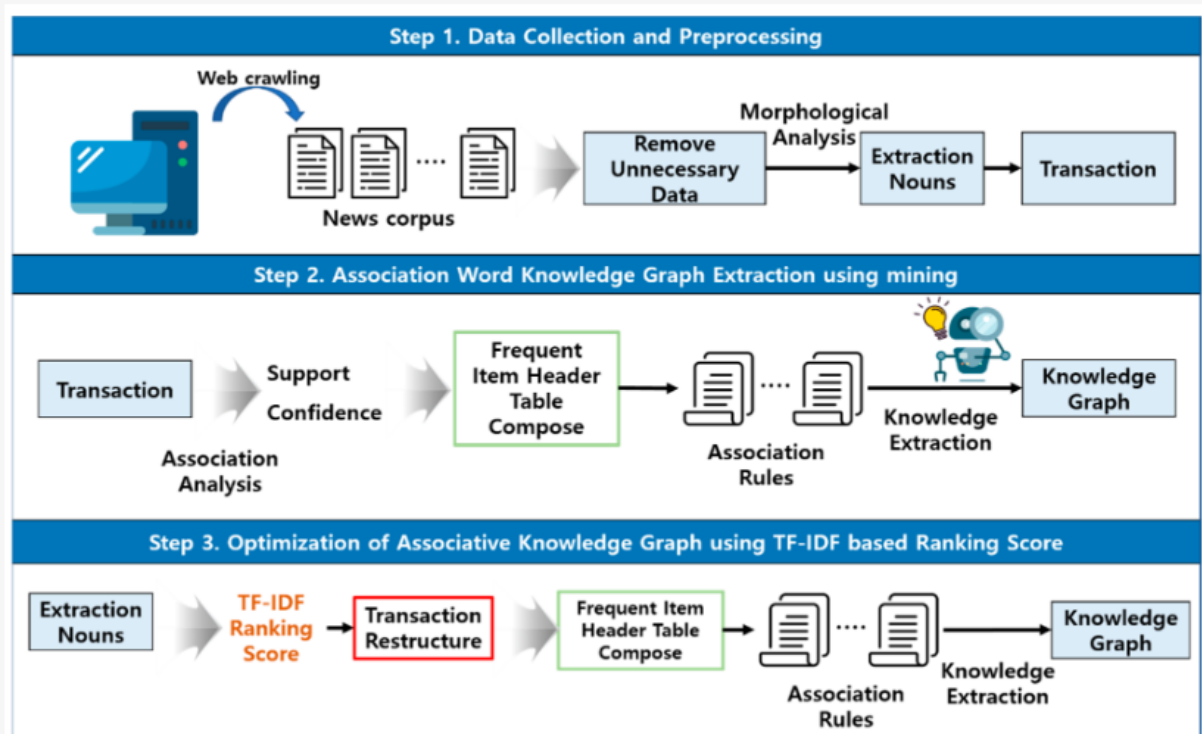
Contents

1.	Overall process description & solution approach	3
1.1	The text summarization solution described in the provided file involves several key steps:	3
1.2	Process diagram:	4
2.	Tool used and reasons to use this specific tool	5
2.1	Tools used	5
2.2	Reasons to use above tools	5
3.	Source code snippets	6
3.1	Imports	7
3.2	Read the pdf file.	7
3.3	Tokenize the sentences.....	8
3.4	Create the Frequency matrix of the words in each sentence	8
3.5	Calculate TermFrequency and generate a matrix	8
3.6	Creating a table for documents per words.....	8
3.7	Calculate IDF and generate a matrix	8
3.8	Calculate TF-IDF and generate a matrix.....	9
3.9	Score the sentences.....	9
3.10	Find the threshold	9
3.11	Generate the summary	10
3.12	Execution of the text summerization	10
4.	Final output results and analysis of results	11
4.1	Final Output Results:.....	11
4.2	Analysis of Results:	11

Objective:

A TF-IDF-based text summarization model can automatically extract and highlight the most relevant sentences from a given document. The model aims to leverage Term Frequency-Inverse Document Frequency (TF-IDF) to identify key sentences that best represent the core content of the text. By scoring sentences based on the significance of their words, the model will generate concise summaries that effectively capture the main ideas while minimizing redundancy. This approach is intended to provide a quick and efficient means of distilling important information from large volumes of text.

Figure 4. Optimizing process of association graph using TF-IDF ranking score.



1. Overall process description & solution approach

1.1 The text summarization solution described in the provided file involves several key steps:

1. **Sentence Tokenization:** The text is split into individual sentences using the ``sent_tokenize`` function from the NLTK library. This forms the basis for analysing each sentence separately.
2. **Frequency Matrix Creation:** For each sentence, a frequency matrix is created that tracks the frequency of each word. The words are first converted to lowercase and stemmed to their root form using the PorterStemmer. Stop words (common words that don't contribute much to the meaning, like "and" or "the") are removed.
3. **Term Frequency (TF) Calculation:** The term frequency for each word in a sentence is calculated by dividing the frequency of the word by the total number of words in the sentence. This results in a matrix where each word has a term frequency score.

$TF(t) = (\text{Number of times term } t \text{ appears in a document}) / (\text{Total number of terms in the document})$

4. **Documents per Words Table Creation:** A table is generated that keeps track of how many sentences (documents) each word appears in. This is useful for calculating the IDF (Inverse Document Frequency).

5. **Inverse Document Frequency (IDF) Calculation:** The IDF for each word is calculated using the formula:

$$IDF(t) = \log_{10}(\text{Total number of documents} / \text{Number of documents with term})$$

This measures how unique or rare a word is across all sentences.

6. **TF-IDF Matrix Calculation:** The TF and IDF values are multiplied to generate a TF-IDF matrix. This matrix helps in identifying the importance of words in the context of the entire text.

7. **Sentence Scoring:** Each sentence is scored by summing up the TF-IDF values of the words it contains. The score reflects the overall importance of the sentence within the text.

8. Threshold Calculation: An average score is calculated from all sentence scores. This threshold helps to determine which sentences are important enough to be included in the summary.

9. Summary Generation: Sentences that have a score above a certain threshold (1.3 times the average score) are selected and concatenated to form the final summary.

This approach leverages the TF-IDF technique to identify and extract the most significant sentences, creating a concise summary of the original text.

1.2 Process diagram:

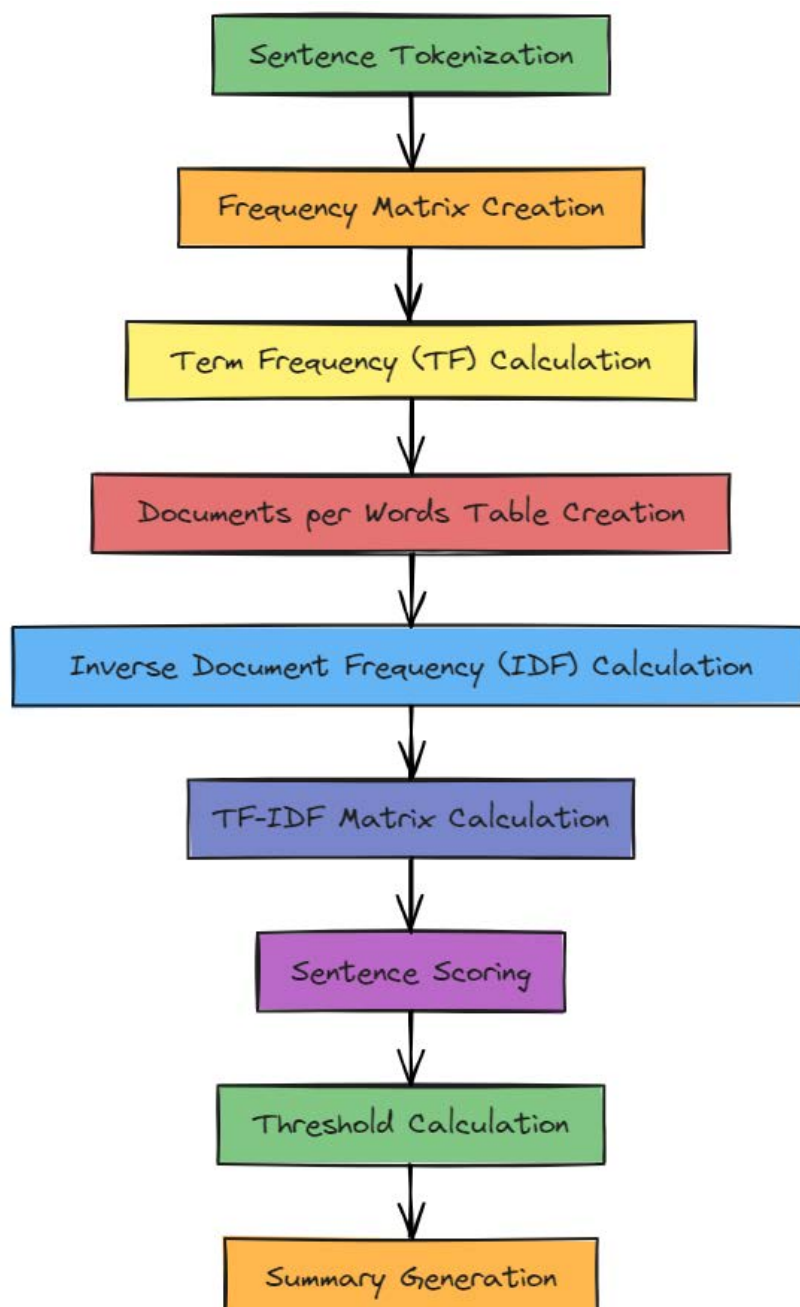


Figure 1: process diagram

2. Tool used and reasons to use this specific tool

2.1 Tools used

Name	Description
Python	The programming language to implement the algorithm.
NLTK (Natural Language Toolkit)	For sentence and word tokenization, stopwords removal, and stemming.

2.2 Reasons to use above tools

Name	Reason
Python	Ease of Use Python's syntax is clear and readable. Extensive Libraries There are many libraries available for natural language processing and machine learning.
NLTK (Natural Language Toolkit)	Comprehensive Toolkit It provides easy-to-use interfaces to over 50 corpora and lexical resources. Preprocessing Capabilities Includes tools for text preprocessing, such as tokenization, stemming, lemmatization, and stopword removal.

3. Source code snippets

3.1 Imports

3.2 Read the pdf file

3.3 Tokenize the sentences

3.4 Create the Frequency matrix of the words in each sentence

3.5 Calculate TermFrequency and generate a matrix

3.6 Creating a table for documents per words

3.7 Calculate IDF and generate a matrix

3.8 Calculate TF-IDF and generate a matrix

3.9 Score the sentences

3.10 Find the threshold

3.11 Generate the summary

3.12 Execution of the text summerization

Imports

```

1 import math
2
3 import nltk
4 nltk.download('punkt')
5 nltk.download('stopwords')
6
7 from nltk import sent_tokenize, word_tokenize, PorterStemmer
8 from nltk.corpus import stopwords

```

```

[ nltk_data ] Downloading package punkt to /root/nltk_data...
[ nltk_data ] Package punkt is already up-to-date!
[ nltk_data ] Downloading package stopwords to /root/nltk_data...
[ nltk_data ] Package stopwords is already up-to-date!

```

```

1 pip install pdfplumber

```

```

Requirement already satisfied: pdfplumber in /usr/local/lib/python3.10/dist-packages (0.11.4)
Requirement already satisfied: pdfminer.six==20231228 in /usr/local/lib/python3.10/dist-packages (from pdfplumber) (20231228)
Requirement already satisfied: Pillow>=9.1 in /usr/local/lib/python3.10/dist-packages (from pdfplumber) (9.4.0)
Requirement already satisfied: pypdfium2>=4.18.0 in /usr/local/lib/python3.10/dist-packages (from pdfplumber) (4.30.0)
Requirement already satisfied: charset-normalizer>=2.0.0 in /usr/local/lib/python3.10/dist-packages (from pdfminer.six==20231228->pdfpl)
Requirement already satisfied: cryptography>=36.0.0 in /usr/local/lib/python3.10/dist-packages (from pdfminer.six==20231228->pdfplumber)
Requirement already satisfied: cffi>=1.12 in /usr/local/lib/python3.10/dist-packages (from cryptography>=36.0.0->pdfminer.six==20231228)
Requirement already satisfied: pycparser in /usr/local/lib/python3.10/dist-packages (from cffi>=1.12->cryptography>=36.0.0->pdfminer.six==20231228)

```

Read the pdf file

```

1 #text = "Those Who Are Resilient Stay In The Game Longer. On the mountains of truth you can never climb in vain: either you will reach a
2
3 import pdfplumber
4
5 # Path to the PDF file
6 pdf_file_path = '/content/Paytm secures finance ministry approval for investment in payment services business - Times of India.pdf'
7
8 # Open the PDF file
9 with pdfplumber.open(pdf_file_path) as pdf:
10     # Initialize a variable to store the extracted text
11     text = ''
12
13     # Iterate through the pages and extract text
14     for page in pdf.pages:
15         text += page.extract_text()
16
17 # Display the extracted text
18 print(text)
19

```

```

Printed from
Paytm secures finance ministry approval for
investment in payment services business
TIMESOFINDIA.COM | Aug 28, 2024, 08.46 PM IST
Paytm, a prominent fintech company, has secured approval from
the finance ministry to invest in its payment services business, the
company announced on Wednesday.
The company, officially known as One 97 Communications, has
been under the watchful eye of banking regulator and financial
crime-fighting agency following the central bank's directive in
January to shut down its payments bank.
With the recent approval, Paytm plans to submit a new
application to the ministry to regain the license for its payments
services business. The company said that Paytm Payment
Services will continue to provide online payment aggregation
services to existing partners.
Although Paytm did not disclose the specifics of the approved investment, Reuters reported in July, based on
information from a senior finance ministry official, that the company had obtained approval for a 500 million
rupee (approximately $6 million) investment in its payments division.
Paytm Payment Services, a significant component of the fintech firm's operations, contributed a quarter of its
consolidated revenue in the fiscal year that ended in March 2023.
In July, Vivek Joshi, financial services secretary, said that the company could approach India's central bank to apply for a payment agg
On the day of the announcement, Paytm's shares closed 1.3% lower. Since the central bank's order to wind down
the payment bank in January, the company's shares have declined by more than 29%.

```

1. Tokenize the sentences

```
1 sentences = sent_tokenize(text) # NLTK function
2 total_documents = len(sentences)
```

2. Create the Frequency matrix of the words in each sentence.

```
1 def _create_frequency_matrix(sentences):
2     frequency_matrix = {}
3     stopWords = set(stopwords.words("english"))
4     ps = PorterStemmer()
5
6     for sent in sentences:
7         freq_table = {}
8         words = word_tokenize(sent)
9         for word in words:
10            word = word.lower()
11            word = ps.stem(word)
12            if word in stopWords:
13                continue
14
15            if word in freq_table:
16                freq_table[word] += 1
17            else:
18                freq_table[word] = 1
19
20        frequency_matrix[sent[:15]] = freq_table
21
22    return frequency_matrix
```

3. Calculate TermFrequency and generate a matrix

$TF(t) = (\text{Number of times term } t \text{ appears in a document}) / (\text{Total number of terms in the document})$

```
1 def _create_tf_matrix(freq_matrix):
2     tf_matrix = {}
3
4     for sent, f_table in freq_matrix.items():
5         tf_table = {}
6
7         count_words_in_sentence = len(f_table)
8         for word, count in f_table.items():
9             tf_table[word] = count / count_words_in_sentence
10
11        tf_matrix[sent] = tf_table
12
13    return tf_matrix
```

4. Creating a table for documents per words

```
1 def _create_documents_per_words(freq_matrix):
2     word_per_doc_table = {}
3
4     for sent, f_table in freq_matrix.items():
5         for word, count in f_table.items():
6             if word in word_per_doc_table:
7                 word_per_doc_table[word] += 1
8             else:
9                 word_per_doc_table[word] = 1
10
11    return word_per_doc_table
```

5. Calculate IDF and generate a matrix

$IDF(t) = \log_e(\text{Total number of documents} / \text{Number of documents with term } t \text{ in it})$


```

1 def _create_idf_matrix(freq_matrix, count_doc_per_words, total_documents):
2     idf_matrix = {}
3
4     for sent, f_table in freq_matrix.items():
5         idf_table = {}
6
7         for word in f_table.keys():
8             idf_table[word] = math.log10(total_documents / float(count_doc_per_words[word]))
9
10    idf_matrix[sent] = idf_table
11

```

6. Calculate TF-IDF and generate a matrix

TF-IDF algorithm is made of 2 algorithms multiplied together.

```

1 def _create_tf_idf_matrix(tf_matrix, idf_matrix):
2     tf_idf_matrix = {}
3
4     for (sent1, f_table1), (sent2, f_table2) in zip(tf_matrix.items(), idf_matrix.items()):
5
6         tf_idf_table = {}
7
8         for (word1, value1), (word2, value2) in zip(f_table1.items(),
9                                                     f_table2.items()): # here, keys are the same in both the table
10            tf_idf_table[word1] = float(value1 * value2)
11
12    tf_idf_matrix[sent1] = tf_idf_table
13
14    return tf_idf_matrix

```

7. Score the sentences

```

1 def _score_sentences(tf_idf_matrix) -> dict:
2     """
3     score a sentence by its word's TF
4     Basic algorithm: adding the TF frequency of every non-stop word in a sentence divided by total no of words in a sentence.
5     :rtype: dict
6     """
7
8     sentenceValue = {}
9
10    for sent, f_table in tf_idf_matrix.items():
11        total_score_per_sentence = 0
12
13        count_words_in_sentence = len(f_table)
14        for word, score in f_table.items():
15            total_score_per_sentence += score
16
17        sentenceValue[sent] = total_score_per_sentence / count_words_in_sentence
18
19    return sentenceValue

```

8. Find the threshold

```

1 def _find_average_score(sentenceValue) -> int:
2     """
3     Find the average score from the sentence value dictionary
4     :rtype: int
5     """
6
7     sumValues = 0
8     for entry in sentenceValue:
9         sumValues += sentenceValue[entry]
10
11    # Average value of a sentence from original summary_text
12    average = (sumValues / len(sentenceValue))
13
14    return average

```

9. Generate the summary

```

1 def _generate_summary(sentences, sentenceValue, threshold):
2     sentence_count = 0
3     summary = ''
4
5     for sentence in sentences:
6         if sentence[:15] in sentenceValue and sentenceValue[sentence[:15]] >= (threshold):
7             summary += " " + sentence
8             sentence_count += 1
9
10    return summary

```

Execution of the text summerization

```

1 import math
2
3 from nltk import sent_tokenize, word_tokenize, PorterStemmer
4 from nltk.corpus import stopwords
5
6 '''
7 We already have a sentence tokenizer, so we just need
8 to run the sent_tokenize() method to create the array of sentences.
9 '''
10 # 1 Sentence Tokenize
11 sentences = sent_tokenize(text)
12 total_documents = len(sentences)
13 #print(sentences)
14
15 # 2 Create the Frequency matrix of the words in each sentence.
16 freq_matrix = _create_frequency_matrix(sentences)
17 #print(freq_matrix)
18
19 '''
20 Term frequency (TF) is how often a word appears in a document, divided by how many words are there in a document.
21 '''
22 # 3 Calculate TermFrequency and generate a matrix
23 tf_matrix = _create_tf_matrix(freq_matrix)
24 #print(tf_matrix)
25
26 # 4 creating table for documents per words
27 count_doc_per_words = _create_documents_per_words(freq_matrix)
28 #print(count_doc_per_words)
29
30 '''
31 Inverse document frequency (IDF) is how unique or rare a word is.
32 '''
33 # 5 Calculate IDF and generate a matrix
34 idf_matrix = _create_idf_matrix(freq_matrix, count_doc_per_words, total_documents)
35 #print(idf_matrix)
36
37 # 6 Calculate TF-IDF and generate a matrix
38 tf_idf_matrix = _create_tf_idf_matrix(tf_matrix, idf_matrix)
39 #print(tf_idf_matrix)
40
41 # 7 Important Algorithm: score the sentences
42 sentence_scores = _score_sentences(tf_idf_matrix)
43 #print(sentence_scores)
44
45 # 8 Find the threshold
46 threshold = _find_average_score(sentence_scores)
47 #print(threshold)
48
49 # 9 Important Algorithm: Generate the summary
50 summary = _generate_summary(sentences, sentence_scores, 1.3 * threshold)
51 print(summary)

```



The company said that Paytm Payment Services will continue to provide online payment aggregation services to existing partners. On the day of the announcement, Paytm's shares closed 1.3% lower.

4. Final output results and analysis of results

4.1 Final Output Results:

The final output of the text summarization process is a concise summary of the input text. This summary is generated by selecting the most significant sentences based on their TF-IDF scores. The sentences that were scored above a certain threshold were included in the summary.

4.2 Analysis of Results:

- **Efficiency:** The summarization process efficiently reduces the text size while retaining the most important information. The use of TF-IDF helps to prioritize sentences that contain key concepts.
- **Quality:** The quality of the summary is dependent on the accuracy of the TF-IDF calculations. The method effectively identifies and retains sentences with high information density.
- **Limitations:** The summarization is purely extractive, meaning it only selects and compiles existing sentences. It does not rephrase or condense ideas, which could sometimes lead to a summary that is less fluid or coherent.

Overall, the approach provides a straightforward and effective way to generate summaries, especially useful for large texts where quick comprehension of key points is needed.