

We use summarizations in many places in our lives. While reading about a book, the blurb at the back provides us with a general idea about it. News apps often have small captions that talk about the article and movie reviews are another place where summaries prove to be extremely valuable.

With the advent of technology such as natural language processing (NLP) and machine learning, why not use them to generate summaries without human intervention?

Before you decide to summarize your textbook for the upcoming exam, let us learn some important concepts.

Summarization of a text can be of two types — extractive and abstractive.

Extractive Summarization

It can be defined as hand-picking any important sentence and adding it to the summary as it is.

1. Read the text and split it into sentences.
2. Parse through each sentence, identify which ones are more important than others based on certain parameters and assign them a score.
3. Pick out the sentences that rank above others.

4. Add these sentences to the summary and TA DA!
You're done.

Abstractive Summarization

It is more like how humans summarize a piece of text, understanding and analyzing the article.

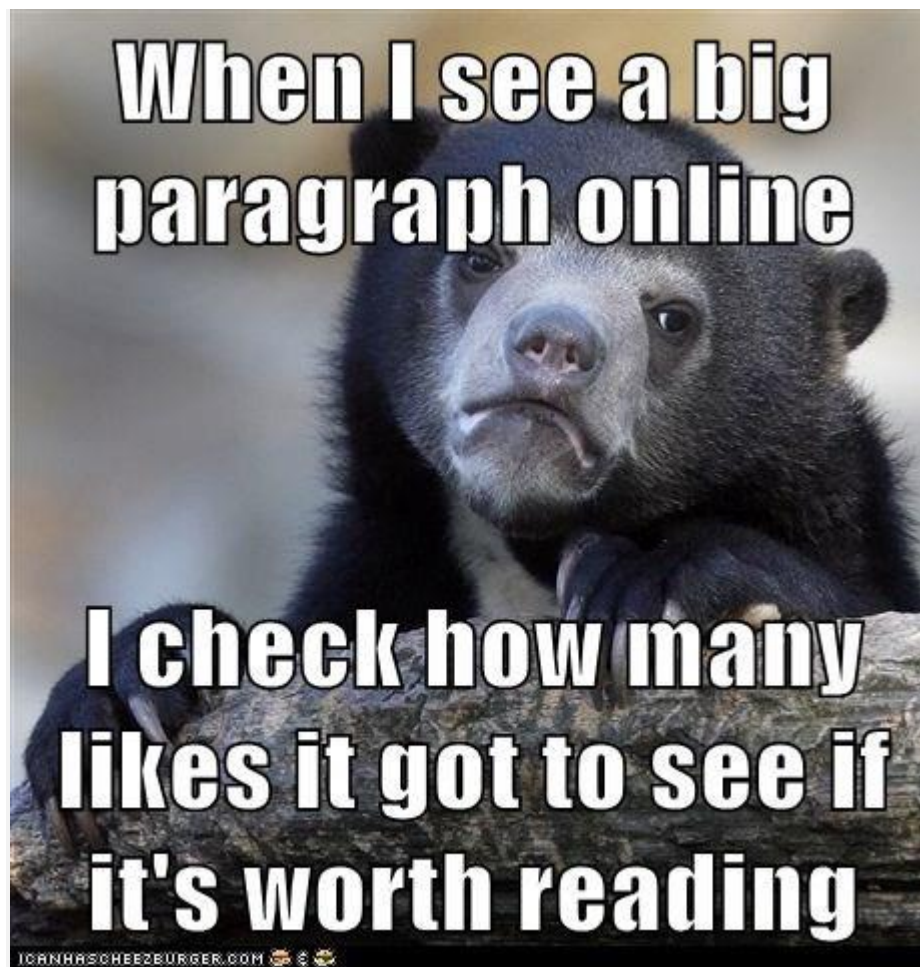
1. Read the text.
2. Analyze the underlying meaning of the text and the sentences.
3. Pick out the important topics and create new sentences (may or may not use vocabulary from the article).
4. Add these sentences to the summary and voila!

As you can see, extractive summarization is relatively easier when compared to abstractive summarization since it doesn't have to deal with semantics or vocabulary.

In this article, we'll discuss extractive summarization and then proceed to create our very own text summarizer from scratch. Excited?

ONE OF THE
MOST
IMPORTANT

THE
CLOCK TOWER
WAS A PART OF
A GREAT
CITY



[credits](#)

Term Frequency — Inverse Document Frequency (TF-IDF)

Given a sentence, we can easily identify what it's talking about and what it's trying to imply. However, a computer can't do it. So, what do we do? As discussed above, we rank the sentences and pick the highest-ranked ones. Who decides the rank of the sentences? TF-IDF.

Term frequency-inverse document frequency is what we'll be using as the basis for selecting sentences that'll make it to the final summary. The TF-IDF value is calculated as given below:

$$\text{TF-IDF} = \text{Term Frequency} \times \text{Inverse Document Frequency}$$

Term Frequency (TF)

This is a measure of the number of times a word appears in our document. Taking only the term frequency will give a highly inaccurate result because a longer document will have a greater TF for a word 'w' when compared to a shorter document. However, this can be normalized if we divide the frequency by the total number of words. This makes TF a more valuable quantity for us.

$$\text{TF}(w, d) = \frac{\text{occurrences of word 'w' in the document 'd'}}{\text{total number of words in the document 'd'}}$$

This is it right? Why do we need anything else? Well, if you think about it, all the words that are extremely common such as "is, am, was" etc. will have a very high TF value. The term TF on

its own will once again prove to be redundant, so behold IDF, our saviour!

Inverse Document Frequency (IDF)

Before we jump to IDF, let's discuss document frequency. DF is the number of documents that contain the word 'w' out of the total N documents. 'D' is the set of all documents.

$DF(w, D)$ = number of documents containing word 'w'

However, DF is not of concern to us, it's the inverse which is relevant. It tells us how much information each term carries. As a result of inversion, common words will have a low IDF value while terms with a lesser occurrence gain higher precedence. We take the logarithm of this value because, for large documents, the value of this term would sky-rocket.

$$IDF(w, D) = \log(N / (DF(w, D) + 1))$$

We do division by (DF + 1) to avoid division by zero. It can prove to be a serious prick for some occasions.

$$TF-IDF(w, d) = TF(w, d) * \log(N / (DF + 1))$$

Once all this is calculated, we rank each document (here, sentences) based on their TF-IDF scores and a threshold value, we finally get a summary of our text. Pretty cool right? Now, let's get coding!

Code Segment

Importing Libraries

```
1 import re
2 import nltk
3 import math
4 from nltk.tokenize import word_tokenize
```

Cleaning the Text

```
6 def clean_text(file_name):
7     file = open(file_name, "r")
8     filedata = file.readlines()
9     article = filedata[0].split(". ")
10    sentences = []
11    # removing special characters and extra whitespaces
12    for sentence in article:
13        sentence = re.sub('[^a-zA-Z]', ' ', str(sentence))
14        sentence = re.sub('[\s+]', ' ', sentence)
15        sentences.append(sentence)
16    sentences.pop()
17    display = " ".join(sentences)
18    print('Initial Text: ')
19    print(display)
20    print('\n')
21    return sentences
```

Converting the text file into individual sentences which will be our document

Number of Words in Each Sentence

```

24 def cnt_words(sent):
25     cnt = 0
26     words = word_tokenize(sent)
27     for word in words:
28         cnt = cnt + 1
29     return cnt
30
31 # getting data about each sentence (frequency of words)
32 def cnt_in_sent(sentences):
33     txt_data = []
34     i = 0
35     for sent in sentences:
36         i = i + 1
37         cnt = cnt_words(sent)
38         temp = {'id' : i, 'word_cnt' : cnt}
39         txt_data.append(temp)
40     return txt_data

```

Number of words in each sentence

Creating List of Frequencies for each Word in all Documents

```

43 def freq_dict(sentences):
44     i = 0
45     freq_list = []
46     for sent in sentences:
47         i = i + 1
48         freq_dict = {}
49         words = word_tokenize(sent)
50         for word in words:
51             word = word.lower()
52             if word in freq_dict:
53                 freq_dict[word] = freq_dict[word] + 1
54             else:
55                 freq_dict[word] = 1
56             temp = {'id' : i, 'freq_dict' : freq_dict}
57             freq_list.append(temp)
58     return freq_list

```

Here, we create a list which stores the frequency of each word of the document

Calculating TF and IDF Values


```

61 def calc_TF(text_data, freq_list):
62     tf_scores = []
63     for item in freq_list:
64         ID = item['id']
65         for k in item['freq_dict']:
66             temp = {
67                 'id': item['id'],
68                 'tf_score': item['freq_dict'][k]/text_data[ID-1]['word_cnt'],
69                 'key': k
70             }
71             tf_scores.append(temp)
72     return tf_scores
73
74 #calculating the inverse document frequency
75 def calc_IDF(text_data, freq_list):
76     idf_scores = []
77     cnt = 0
78     for item in freq_list:
79         cnt = cnt + 1
80         for k in item['freq_dict']:
81             val = sum([k in it['freq_dict'] for it in freq_list])
82             temp = {
83                 'id': cnt,
84                 'idf_score': math.log(len(text_data)/(val+1)),
85                 'key': k}
86             idf_scores.append(temp)
87     return idf_scores

```

The first function calculates the term frequency of words in each document (here, sentences). The second function calculates the inverse document frequency for each word in the sentences.

Calculating TF-IDF Values

```

90 def calc_TFIDF(tf_scores, idf_scores):
91     tfidf_scores = []
92     for j in idf_scores:
93         for i in tf_scores:
94             if j['key'] == i['key'] and j['id'] == i['id']:
95                 temp = {
96                     'id': j['id'],
97                     'tfidf_score': j['idf_score'] * i['tf_score'],
98                     'key': j['key']
99                 }
100                 tfidf_scores.append(temp)
101     return tfidf_scores

```

Computing the TF-IDF scores for each term

Ranking all the Documents

```

104 def sent_scores(tfidf_scores, sentences, text_data):
105     sent_data = []
106     for txt in text_data:
107         score = 0
108         for i in range(0, len(tfidf_scores)):
109             t_dict = tfidf_scores[i]
110             if txt['id'] == t_dict['id']:
111                 score = score + t_dict['tfidf_score']
112         temp = {
113             'id': txt['id'],
114             'score': score,
115             'sentence': sentences[txt['id']-1]}
116         sent_data.append(temp)
117     return sent_data

```

We calculate the score for each sentence based on the TF-IDF value for each word in the sentence

Generating the Summary

```

120 def summary(sent_data):
121     cnt = 0
122     summary = []
123     for t_dict in sent_data:
124         cnt = cnt + t_dict['score']
125     avg = cnt / len(sent_data)
126     for sent in sent_data:
127         if sent['score'] >= (avg * 0.9):
128             summary.append(sent['sentence'])
129     summary = ". ".join(summary)
130     return summary

```

The threshold is calculated as a linear function of the average of the TF-IDF scores

(Note: the threshold for average can be modified by multiplying a scalar to reduce or increase the size of the summary)

And...Time to call all the functions!

```

136 sentences = clean_text('text.txt')
137 text_data = cnt_in_sent(sentences)
138
139 freq_list = freq_dict(sentences)
140 tf_scores = calc_TF(text_data, freq_list)
141 idf_scores = calc_IDF(text_data, freq_list)
142
143 tfidf_scores = calc_TFIDF(tf_scores, idf_scores)
144
145 sent_data = sent_scores(tfidf_scores, sentences, text_data)
146 result = summary(sent_data)
147
148 print('Final Summary: ')
149 print(result)

```

Function calls

(Note: Ensure that your textual matter is in a file (*.txt) and is in the same directory as your Python script)

Testing...testing...1, 2, 3!

```

In [13]: runfile('/Users/anjaneyatripathi/Desktop/nlp_1/prac.py', wdir='/Users/anjaneyatripathi/Desktop/nlp_1')
Initial Text:
America has changed dramatically during recent years Not only has the number of graduates in traditional engineering disciplines such as
mechanical civil electrical chemical and aeronautical engineering declined but in most of the premier American universities
engineering curricula now concentrate on and encourage largely the study of engineering science As a result there are declining
offerings in engineering subjects dealing with infrastructure the environment and related issues and greater concentration on high
technology subjects largely supporting increasingly complex scientific developments While the latter is important it should not be at
the expense of more traditional engineering Rapidly developing economies such as China and India as well as other industrial countries
in Europe and Asia continue to encourage and advance the teaching of engineering Both China and India respectively graduate six and
eight times as many traditional engineers as does the United States

Final Summary:
America has changed dramatically during recent years. As a result there are declining offerings in engineering subjects dealing with
infrastructure the environment and related issues and greater concentration on high technology subjects largely supporting
increasingly complex scientific developments. While the latter is important it should not be at the expense of more traditional
engineering. Both China and India respectively graduate six and eight times as many traditional engineers as does the United States

```

Oooh, La La!

Conclusion — The End of a Wonderful Journey

Well, I hope you had fun coding along with me and learning about text summarization while we were at it. However, this isn't the only way to summarize a text, there are many more techniques that can be implemented to accomplish the task. Try checking them out too!

Hope you enjoyed it and see you next time!