

- Text Summarization
 - Research
 - Overview
 - Approach
 - Assumptions
 - Bibliography
 - Output
 - Python
 - Example
 - Sample Input
 - Sample Output
 - Code

Text Summarization

Extractive summary of input article.

Author: Nishit Jain

Email: nishitjain1997@gmail.com

Source: [Github](#)

Research

Overview

Extractive text summarization is a type of Automated Text Summarization (summarization using computers) in which the summary is made up of complete sentences picked from the original text. To achieve this, each sentence is given a relevance ranking and the top most relevant sentences are picked for the summary.

Approach

This is an unsupervised summary algorithm which is an amalgamation of the TextRank algorithm given by Mihalcea et al. and the Feature Term based method given by Suneetha Manne et al. The scores given by TextRank algorithm are enhanced by extracting additional word-level and sentence-level features to incorporate both semantic and syntactic meaning in the scores.

1. TextRank algorithm

- Each sentence is converted into a **Sentence Vector** by taking the average of the Glove embeddings of each word in that sentence.
- Pairwise Cosine Similarity of Sentence Vectors is used to create a similarity matrix of sentences.

The similarity matrix is used to build a directed graph, which is then fed into Google's PageRank algorithm to give a score to each sentence. The sentence that is most similar to all sentences in the text supposedly contains ideas from them all and should be considered in the summary.

2. Feature Term based method. The features extracted include:

- Resultant Term Weight: This denotes the amount of information conveyed by each word in the sentence. It is a product of the Term Weight (normalized frequency of word) and Inverse Sentence Frequency (log of ratio of number of sentences and number of appearances of the term).
- Sentence Weight: Parts of Speech tagging to ensure include syntactic meaning of words. Ratio of number of noun and verb terms in a sentence to the total number of terms in all sentences.
- Sentence Position: Sentences coming at the beginning contain more generalized ideas while those coming in the middle contain more specific ideas.
- Sentence Length: A normalization term to prevent longer sentences from dominating shorter sentences. Gives score per unit length of sentence

$$\text{Feature Rank} = \frac{(\sigma(\text{Resultant Term Weight}) + \text{Sentence Weight} + \text{Sentence Position})}{\text{Sentence Length}}$$

$$\text{Final Score} = 0.8 * \text{Text Rank} + 0.2 * \text{Feature Rank}$$

Assumptions

1. The extractive summary has been limited to 3 lines.
2. This is a single-document summary algorithm.

Bibliography

1. Mihalcea R., & Tarau P. (2004). TextRank: Bringing Order into Texts. (W04-3252). Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing
2. Manne S., & Fatima S. (2012). A Feature Terms based Method for Improving Text Summarization with Supervised POS Tagging. (10.5120/7494-0541). International Journal of Computer Applications (0975 – 8887).
3. Jagadeesh J., & Pingali P., & Varma V. (2005). Sentence Extraction Based Single Document Summarization. (IIIT/TR/2008/97). Workshop on Document Summarization, 19th and 20th March, 2005, IIIT Allahabad.

Output

The summary for the article will be stored in 'output.txt' file in root directory.

Python

This script has been developed in and tested for Python3. It has not been tested for Python2.

Example

Sample Input

S.No. **News Article**

S.No. News Article

1 The U.N. Security Council approved a resolution Monday to send 4,200 peacekeepers to Abyei, Sudan, as part of a recent agreement between Sudan and Southern Sudan.

2 The resolution will establish, for six months, the United Nations Interim Security Force for Abyei (UNISFA), comprising "a maximum of 4,200 military personnel, 50 police personnel, and appropriate civilian support," the resolution states.

3 It passed the council unanimously, 15-0.

4 In a statement released by the State Department, Secretary Hiliary Clinton commended the swift passage of the resolution.

5 "Abyei has been a source of regional tension for many years," the statement said.

6 "We urge the parties to reach an immediate cease-fire and to provide aid workers with the unfettered access required to deliver humanitarian assistance to innocent civilians affected by the conflict."

7 A week ago, the Sudanese government and the Sudan People's Liberation Movement signed an agreement to allow peacekeepers in Abyei, aimed at ending strife that has ravaged much of the country.

8 The two sides agreed in principle on the need for a third party to monitor the ill-defined border between north and south before the scheduled July 9 independence for the south.

9 The U.N. peacekeepers will "monitor and verify the redeployment of any Sudan Armed Forces, Sudan People's Liberation Army or its successor" from the Abyei area, among other tasks, the Security Council resolution states.

Sample Output

The U.N. Security Council approved a resolution Monday to send 4,200 peacekeepers to Abyei, Sudan, as part of a recent agreement between Sudan and Southern Sudan. A week ago, the Sudanese government and the Sudan People's Liberation Movement signed an agreement to allow peacekeepers in Abyei, aimed at ending strife that has ravaged much of the country. The U.N. peacekeepers will "monitor and verify the redeployment of any Sudan Armed Forces, Sudan People's Liberation Army or its successor" from the Abyei area, among other tasks, the Security Council resolution states.

```

1 import nltk
2 import os
3 import re
4 import spacy
5 import sys
6 import unicodedata
7
8 import networkx as nx
9 import numpy as np
10 import pandas as pd
11
12 nltk.download('punkt')
13 nltk.download('stopwords')
14
15 from nltk.corpus import stopwords
16 from nltk.tokenize import sent_tokenize
17 from sklearn.metrics.pairwise import cosine_similarity
18
19 # Root directory
20 root = '.'
21 # Loading Spacy for Parts-of-Speech tagging.
22 nlp_spacy = spacy.load('en', parse=True, tag=True)
23 # Loading list of stopwords
24 stop_words = stopwords.words('english')
25
26
27
28
29 """
30     For TextRank Algorithm.
31 """
32 def remove_extraneous_text(sentence:str)->str:
33     """
34         Input: String
35         Output: String
36         Takes a news article as input and removes extra spaces and reporting
37         location from it.
38     """
39     # Remove multiple spaces
40     sentence = re.sub(" +", " ", sentence)
41
42     # Remove reporting location
43     if ") --" in sentence:
44         sentence = sentence.split(") --")[-1]
45
46     # Remove media name from article
47     if "(CNN)" in sentence:
48         sentence = sentence.split("(CNN)")[-1]
49
50     return sentence
51
52 def remove_stopwords(sentence:str)->str:
53     """
54         Input: String
55         Output: String
56         Takes a sentence as input and returns the sentence after removing all
57         stopwords.
58     """

```

```

57     sentence = " ".join([word for word in sentence.split() if word not in
stop_words])
58
59     return sentence
60
61 def lemmatize_text(sentence:str)->str:
62     """
63     Input: String
64     Output: String
65     Takes a sentence as input and uses Spacy to convert each word into it's
lemma.
66     """
67     sentence = nlp_spacy(sentence)
68     sentence = ' '.join([word.lemma_ if word.lemma_ != "-PRON-" else word.text
for word in sentence])
69     return sentence
70
71 def clean_text(sentence:str)->str:
72     """
73     Input: String
74     Output: String
75     Takes a sentence and cleans by:
76     - Converting to lowercase
77     - Remove non alphabetic characters
78     - Removing extraneous characters
79     - Removing stopwords
80     - Lemmatizing words
81     """
82     sentence = sentence.lower()
83     sentence = re.sub("[^a-zA-Z]", " ", sentence)
84     sentence = remove_extraneous_text(sentence)
85     sentence = remove_stopwords(sentence)
86     sentence = lemmatize_text(sentence)
87     return sentence
88
89
90
91
92     """
93     For Feature Term enhancements.
94     """
95 def get_total_terms(cleaned_sentences:list)->int:
96     """
97     Input: List
98     Output: Int
99     Takes in a list of sentences and returns total number of tokens in those
sentences.
100    """
101    total_terms = 0
102
103    for sentence in cleaned_sentences:
104        total_terms += len(sentence.split())
105
106    return total_terms
107
108 def get_term_frequencies(cleaned_sentences:list)->dict:
109    """
110    Input: List
111    Output: Dict

```

```

112     Takes in a list of sentences and returns a dictionary containing Tokens
as keys and their frequencies as values.
113     """
114     freq_dict = {}
115
116     for sentence in cleaned_sentences:
117         for word in sentence.split():
118             freq_dict[word] = freq_dict.get(word, 0) + 1
119
120     return freq_dict
121
122 def get_term_weights(cleaned_sentences:list)->dict:
123     """
124     Input: List
125     Output: Dict
126     Takes in a list of sentences and returns a dictionary containing Tokens
as keys and their weightage as values.
127     The weight is calculated using formula:
128          $TW(t_i) = (TF(t_i) * 1000) / (N_t)$ 
129     where  $t_i$  is each token,  $TW$  is term weight,  $TF$  is term frequency and  $N_t$  is
total number of terms
130     """
131     total_terms = get_total_terms(cleaned_sentences)
132     term_freq_dict = get_term_frequencies(cleaned_sentences)
133     term_weights = dict()
134
135     for key, value in term_freq_dict.items():
136         term_weights[key] = (value * 1000) / total_terms
137
138     return term_weights
139
140 def inverse_sentence_frequency(cleaned_sentences:list)->dict:
141     """
142     Input: List
143     Output: Dict
144     Takes in a list of sentences and returns a dictionary containing Tokens
as keys and their inverse sentence frequency as values.
145     The inverse sentence frequency is calculated as:
146          $ISF(t_i) = \log((N_s) / N_{t_i})$ 
147     where  $t_i$  is each token,  $ISF$  is inverse sentence frequency,  $N_s$  is total
number of sentences in paragraph and  $N_{t_i}$  are the total number of
148     sentences in which  $t_i$  appeared in that paragraph.
149     """
150     vocabulary = set()
151
152     for sentence in cleaned_sentences:
153         vocabulary = vocabulary.union(set(sentence.split()))
154
155     isf = dict()
156     number_of_sentences = len(cleaned_sentences)
157
158     for word in vocabulary:
159         number_of_appearances = 0
160
161         for sentence in cleaned_sentences:
162             if word in sentence:
163                 number_of_appearances += 1
164
165         isf[word] = np.log(number_of_sentences / number_of_appearances)
166

```

```

167     return isf
168
169 def word_weights(cleaned_sentences:str)->dict:
170     """
171     Input: List
172     Output: Dict
173     Takes in a list of sentences and returns a dictionary containing Tokens
as keys and their resultant weightage as values.
174     The weightage is calculated as:
175          $RW(t_i) = ISF(t_i) * TW(t_i)$ 
176     where  $t_i$  is each token,  $RW$  is resultant weightage,  $ISF$  is inverse
sentence frequency and  $TW$  is term weightage.
177     """
178
179     term_weights = get_term_weights(cleaned_sentences)
180     inverse_sentence_freq = inverse_sentence_frequency(cleaned_sentences)
181
182     resultant_weights = dict()
183
184     for word in term_weights.keys():
185         resultant_weights[word] = term_weights[word] *
inverse_sentence_freq[word]
186
187     return resultant_weights
188
189 def pos_tagging(cleaned_sentences:list)->list:
190     """
191     Input: List
192     Output: List
193     Takes in a list of sentences and returns a list of lists, where each
Token is represented as a tuple of the form (Token, POS tag).
194     """
195     tagged_sentences = []
196
197     for sentence in cleaned_sentences:
198         sentence_nlp = nlp_spacy(sentence)
199
200         tagged_sentence = []
201
202         for word in sentence_nlp:
203             tagged_sentence.append((word, word.pos_))
204
205         tagged_sentences.append(tagged_sentence)
206
207     return tagged_sentences
208
209 def sentence_weights(tagged_sentences:list, total_terms:int)->list:
210     """
211     Input: List, Int
212     Output: List
213     Takes in a list of POS tagged sentences and total number of terms.
Returns a list containing the sentence weight of each sentence.
214     The sentence weight is calculated as:
215          $SW(s_i) = \text{Number of nouns and verbs in sentence} / \text{total number of}$ 
terms in paragraph.
216     """
217     sent_weights = []
218
219     for sentence in tagged_sentences:
220         relevance_count = 0

```

```

221
222     for word, tag in sentence:
223         if tag == 'NOUN' or tag == 'VERB':
224             relevance_count += 1
225
226     sent_weights.append(relevance_count / total_terms)
227
228     return sent_weights
229
230 def sentence_position(cleaned_sentences:list)->list:
231     """
232     Input: List
233     Output: List
234     Takes in a list of sentences and returns weight for each sentence based
on it's position.
235     """
236     sent_position = []
237     number_of_sentences = len(cleaned_sentences)
238
239     weights = [0, 0.25, 0.23, 0.14, 0.08, 0.05, 0.04, 0.06, 0.04, 0.04, 0.15]
240
241     for i in range(1, len(cleaned_sentences)+1):
242         sent_position.append(weights[int(np.ceil(10 * (i /
number_of_sentences)))]))
243     return sent_position
244
245 def sentence_length(cleaned_sentences:list)->list:
246     """
247     Input: List
248     Output: List
249     Takes in a list of sentences and returns a list containing length of each
sentence.
250     """
251     sent_len = []
252
253     for sentence in cleaned_sentences:
254         sent_len.append(len(sentence.split()))
255
256     return sent_len
257
258
259
260
261     """
262     Functions to rank sentences.
263     """
264 def text_rank(sentences:list, word_embeddings:dict)->dict:
265     """
266     Input: List, Dict
267     Output: Dict
268     Takes a list of sentences and Glove word embeddings as input and returns
a dictionary containing sentences index as key and rank as value.
269     The ranking is done based on the PageRank algorithm
270     """
271     # Clean sentences for PageRank algorithm.
272     clean_sentences = pd.Series(sentences).str.replace("[^a-zA-Z]", " ")
273     clean_sentences = [s.lower() for s in clean_sentences]
274     clean_sentences = [remove_stopwords(r) for r in clean_sentences]
275

```

```

276 # Replace each word with Glove embeddings. The Sentence vector is the
average of the sum of embeddings of all words in that
277 # sentence.
278 sentence_vectors = []
279 for i in clean_sentences:
280     if len(i) != 0:
281         v = sum([word_embeddings.get(w, np.zeros((100, ))) for w in i.split()])
/ (len(i.split()) + 0.001)
282     else:
283         v = np.zeros((100, ))
284     sentence_vectors.append(v)
285
286 # Initialize a similarity matrix for pair of sentences
287 sim_mat = np.zeros([len(sentences), len(sentences)])
288
289 # Calculate cosine similarity for each pair of sentences
290 for i in range(len(sentences)):
291     for j in range(len(sentences)):
292         if i != j:
293             sim_mat[i][j] = cosine_similarity(sentence_vectors[i].reshape(1,
100), sentence_vectors[j].reshape(1, 100))[0, 0]
294
295 # Create a PageRank graph using similarity matrix
296 nx_graph = nx.from_numpy_array(sim_mat)
297 scores = nx.pagerank(nx_graph)
298
299 return scores
300
301 def feature_rank(sentences:list)->dict:
302     """
303     Input: List
304     Output: Dict
305     Takes a list of sentences as input and returns a dict containig ranking
of each sentence.
306     The ranking is calculated using word and sentence level features.
307     """
308     cleaned_sentences = [clean_text(sentence) for sentence in sentences]
309
310     term_weights = word_weights(cleaned_sentences)
311     tagged_sentences = pos_tagging(cleaned_sentences)
312     total_terms = get_total_terms(cleaned_sentences)
313     sent_weights = sentence_weights(tagged_sentences, total_terms)
314     sent_position = sentence_position(cleaned_sentences)
315     sent_len = sentence_length(cleaned_sentences)
316
317     sentence_scores = []
318
319     for index, sentence in enumerate(cleaned_sentences):
320         score = 0
321
322         for word in sentence.split():
323             score += term_weights[word]
324
325         score *= sent_weights[index]
326         score += sent_position[index]
327
328         if sent_len[index] != 0:
329             score /= sent_len[index]
330         else:
331             score = 0

```

```

332     sentence_scores.append(score)
333
334 sentence_scores = sentence_scores / np.sum(sentence_scores)
335
336 final_scores = dict()
337
338 for i in range(len(sentence_scores)):
339     final_scores[i] = sentence_scores[i]
340
341 return final_scores
342
343
344
345
346 def main()->None:
347     """
348     The driver function.
349     """
350
351     # Path to input file
352     input_filepath = os.path.join(root, sys.argv[-1])
353
354     if not os.path.exists(input_filepath):
355         # Check if input file does not exist.
356         print("Could not find input file at location '%s'" % (input_filepath))
357         return
358
359     input_text = ""
360
361     with open(input_filepath, 'r') as f:
362         input_text = f.read()
363
364
365
366     # Location of Glove word embeddings.
367     glove_location = os.path.join(root, 'embeddings', 'glove.6B.100d.txt')
368
369     if not os.path.exists(glove_location):
370         # Check if word embeddings do not exist.
371         print("Could not find Glove Word Embeddings. Kindly download from
372         'https://drive.google.com/open?id=1cQBYwoLHZzHk4w8zdgCSPFmOP5Xq-x0z' \
373         and save in './embeddings' location.")
374         return
375
376     print("Loading Glove Word embeddings.")
377
378     # Dictionary to store embeddings
379     word_embeddings = {}
380
381     # Open file and load embeddings in memory
382     f = open(glove_location, encoding='utf-8')
383     for line in f:
384         values = line.split()
385         word = values[0]
386         coefs = np.asarray(values[1:], dtype='float32')
387         word_embeddings[word] = coefs
388     f.close()
389
390     print("Embeddings loaded.")

```

```
391
392 print("Creating summary.")
393
394 sentences = sent_tokenize(input_text)
395 text_rank_scores = text_rank(sentences, word_embeddings)
396 feature_rank_scores = feature_rank(sentences)
397
398 final_scores = dict()
399 for i in range(len(text_rank_scores.keys())):
400     final_scores[i] = 0.8 * text_rank_scores[i] + 0.2 *
feature_rank_scores[i]
401
402 ranked_sentences = sorted(((final_scores[i], s, i) for i, s in
enumerate(sentences)), reverse=True)[:3]
403 ranked_sentences = sorted(ranked_sentences, key=lambda x: x[2])
404
405 output_text = ""
406 for i in range(len(ranked_sentences)):
407     output_text += ranked_sentences[i][1] + ' '
408
409 with open('output.txt', 'w') as f:
410     f.write(output_text.strip())
411 print("Summary stored in 'output.txt'.")
412
413
414 if __name__ == '__main__':
415     if len(sys.argv) != 2:
416         print("The syntax to run this program is: 'python run.py file_name.txt'")
417     else:
418         main()
```