

News Retrieval and Multimedia Exploration System

Smart Systems

Faculty of Computer and Data Sciences,
Alexandria University.

Authored by: Ahmed Ehab Abdelhalem
ID:2203147

Authored by: Youssef Mohamed Helmy
ID:2203163



Libraries

```
1 import feedparser
2 import re
3 from translate import Translator
4 from gtts import gTTS
5 import time
6 from pytube import YouTube
7 from pytube import Search
8 import pygame
9 from sklearn.feature_extraction.text import TfidfVectorizer
10 from sklearn.metrics.pairwise import linear_kernel
11
```

1. **Feedparser:** Used for parsing RSS and Atom feeds
 - We used it to parse RSS from sky news
2. **Re:** Stands for regular expressions, a powerful tool for text manipulation
 - We used it to search in feeds of sky news
3. **Translate:** Used for translating texts
 - We used it to translate news in sky news feeds
4. **Gtts (Google Text-to-Speech):** Converts text into spoken words using Google Text-to-Speech API.
 - We used it to convert news into spoken words
5. **Time:** A standard Python module for working with time-related functions.
 - We used it to set time for spoken words
6. **Pytube:** A library for downloading YouTube videos
 - We used it to download videos about our news
7. **Pygame:** Used for multimedia applications.
 - We used it to load and play voice records of gtts
8. **Sklearn (TfidfVectorizer):** Used for applying TF-IDF method in natural language processing
 - it helped us in recommendation system because it finds unique words and turns them into vectors
9. **Sklearn (linear_kernel):** Used for applying linear kernel method which calculates dot product
 - it helped us in recommendation system because it takes word vectors and finds cosine similarity of these vectors

Functions

```
12 def get_news_feeds(feedURL):
13     feeds = feedparser.parse(feedURL)
14     news_list = []
15     for line in feeds.entries:
16         news_list.append({
17             'title': line.title,
18             'description': line.description
19         })
20     return news_list
21
```

This function takes string (feedURL) we use that URL to get feeds using parse(), then create empty list to save news in it, then make for loop to check for titles and descriptions in RSS and get it and append it in our list, then we get this list

```
22 def recommend_news(user_input, news_list):
23     news_corpus = [f"{item['title']} {item['description']}" for item in news_list]
24
25     tfidf_vectorizer = TfidfVectorizer(stop_words='english')
26
27     tfidf_matrix = tfidf_vectorizer.fit_transform(news_corpus)
28
29     user_tfidf = tfidf_vectorizer.transform([user_input])
30
31     cosine_similarities = linear_kernel(user_tfidf, tfidf_matrix).flatten()
32
33     top_indices = cosine_similarities.argsort()[:-6:-1]
34
35     recommendations = [news_list[i] for i in top_indices]
36     return recommendations
37
```

This function takes input from user and our news list, then it initializes our corpus using combinations of titles and descriptions in our list, then we set our TF-IDF function to use English as a language then we will turn our unique words and user input to vectors using TF-IDF method

Note: there is simple example explain how TF-IDF method works

Example

- Sen 1 : good boy
- Sen 2 : good girl
- Sen 3 : boy girl good
- First we calculate TF $\rightarrow TF = \frac{\text{No.of rep in a sen}}{\text{No.of words in a sen}}$
- $IDF = \log\left(\frac{\text{no. of sentences}}{\text{No.of sentences containing words}}\right)$

	sen1	sen2	sen3
Good	1/2	1/2	1/3
Boy	1/2	0	1/3
girl	0	1/2	1/3

X

words	IDF
Good	$\log(3/3)=0$
Boy	$\log(3/2)$
girl	$\log(3/2)$

	good	boy	girl	output
Sen 1	0	$\frac{1}{2} * \log(3/2)$	0	
Sen 2	0	0	$\frac{1}{2} * \log(3/2)$	
Sen 3	0	$\frac{1}{3} * \log(3/2)$	$\frac{1}{3} * \log(3/2)$	

Sen 1, Sen2, Sen3 are our corpus

First: we apply TF to every word in every sen and insert outputs in matrix, outputs with higher value are more important to specific sen

Second: we apply IDF to every unique word in corpus and insert outputs in matrix, outputs with higher value are more unique in corpus

Finally: we multiply the two matrices and the new matrix is a matrix of word vectors

Returning to our main function the best way to check similarity between user input vector and news words vectors is cosine similarity so from machine learning methods we used linear kernel which is essentially the same as the cosine similarity, it's a simple dot product of the vectors that calculates the dot product between the TF-IDF vector of the user's input and each news item in the TF-IDF matrix, resulting in a measure of similarity. The higher the value, the more similar the user input is to a particular news item, then flatten() to make results in one dimensional array

Then, we sort our top 5 recommendations news from with highest values and update our news list

```

38 def download_video(youtube_url, output_path):
39     try:
40         youtube_object = YouTube(youtube_url)
41         youtube_stream = youtube_object.streams.get_highest_resolution()
42         if youtube_stream:
43             youtube_stream.download(output_path)
44             print("Video downloaded successfully.")
45         else:
46             print("No suitable streams found for the video.")
47     except Exception as e:
48         print(f"Error downloading video: {e}")
49

```

This function takes YouTube URL and path to download video in it, it's main idea of it to download selected video from YouTube so it search for YouTube video and get highest resolution from it and if that done it download video in selected path and if there was an unexpected error it print an error message

```

50 def NewsDownloader(lang, feedURL, WordToSearch):
51     pygame.mixer.init()
52     translator = Translator(to_lang=lang)
53
54     feeds = feedparser.parse(feedURL)
55
56     for line in feeds.entries:
57         if re.search(WordToSearch, line.description, re.IGNORECASE) or re.search(WordToSearch, line.title, re.IGNORECASE):
58
59             translated_title = translator.translate(line.title)
60             title_tts = gTTS(translated_title, lang=lang)
61             title_tts_file = f"{line.title[:10]}_title.mp3"
62             title_tts.save(title_tts_file)
63             print("Newspaper title:", line.title)
64             pygame.mixer.music.load(title_tts_file)
65             pygame.mixer.music.play()
66             time.sleep(10)
67
68             translated_desc = translator.translate(line.description)
69             desc_tts = gTTS(translated_desc, lang=lang)
70             desc_tts_file = f"{line.title[:10]}_description.mp3"
71             desc_tts.save(desc_tts_file)
72             print("Newspaper description:\n", line.description)
73             pygame.mixer.music.load(desc_tts_file)
74             pygame.mixer.music.play()
75             time.sleep(10)
76

```

This is the function that search and display our news from sky news

First: it take feed URL and language name and input word then, we initialize our player and translator using selected language

Second: we search in feeds about the input word in descriptions or titles in feeds.

- Translate the title using translator
- Read title with spoken words using (Google Text-to-Speech)
- Save voice recode to mp3 file
- Print feed title
- Load and play our mp3 record of reading title
- Set 10 sec to talk
- Do the same thing to description of feed

```

77 search_query = line.title
78 search_results = Search(search_query)
79
80 if search_results.results:
81     first_video = search_results.results[0]
82     print("Title:", first_video.title)
83     print("URL:", first_video.watch_url)
84
85     # Download video
86     output_path = f"{line.title[:10]}_video.mp4"
87     download_video(first_video.watch_url, output_path)
88
89 else:
90     print("No videos found for the search query.")
91
92 print("\n" * 60)

```

Third: take title line and search in YouTube for video using it, if search is done we will get first result and print video title and URL, then download and save video in the input path if there is no results print message

```

94 if __name__ == "__main__":
95
96     feed_urls = {
97         '1': "https://feeds.skynews.com/feeds/rss/home.xml",
98         '2': "https://feeds.skynews.com/feeds/rss/uk.xml",
99         '3': "https://feeds.skynews.com/feeds/rss/world.xml",
100        '4': "https://feeds.skynews.com/feeds/rss/us.xml",
101        '5': "https://feeds.skynews.com/feeds/rss/business.xml",
102        '6': "https://feeds.skynews.com/feeds/rss/politics.xml",
103        '7': "https://feeds.skynews.com/feeds/rss/technology.xml",
104        '8': "https://feeds.skynews.com/feeds/rss/entertainment.xml",
105        '9': "https://feeds.skynews.com/feeds/rss/strange.xml",
106    }
107
108     while True:
109         print("Categories:")
110         print("1. Home")
111         print("2. UK")
112         print("3. World")
113         print("4. US")
114         print("5. Business")
115         print("6. Politics")
116         print("7. Technology")
117         print("8. Entertainment")
118         print("9. Strange")
119         print("0. Exit")
120

```

This is feeds links that we will use and we will make menu to select feed category

```

121     option = input("Enter category to search about using feeds.skynews.com: ")
122
123     if option == '0':
124         break # Exit the loop
125
126     feed_url = feed_urls.get(option)
127
128     if feed_url:
129         user_input = input("Enter keywords to search for: ")
130
131         news_list = get_news_feeds(feed_url)
132         recommendations = recommend_news(user_input, news_list)
133
134         print("\nTop 5 Recommended News:")
135         for i, news in enumerate(recommendations, 1):
136             print(f"{i}. Title: {news['title']}")
137
138         choose = input("Enter a keyword from titles to search for: ")
139         NewsDownloader("ar", feed_url, choose)
140     else:
141         print("Enter a valid number, please try again")

```

Our input will get one of feeds URL and we will use our previous functions to get feed URL and search for recommendations and get news in spoken words and videos



THANK YOU!