

```

1 import pandas as pd
2 import requests
3 import regex as re

--NORMAL--

1 import os
2 os.environ['NotebookApp.iopub_data_rate_limit']

```

```
1 !pip install transformers
```

```

Looking in indexes: https://pypi.org/simple,
Collecting transformers
  Downloading transformers-4.28.1-py3-none-any
Requirement already satisfied: requests in /u
Collecting huggingface-hub<1.0,>=0.11.0
  Downloading huggingface-hub-0.14.0-py3-none
Requirement already satisfied: filelock in /u
Requirement already satisfied: numpy>=1.17 in
Requirement already satisfied: pyyaml>=5.1 in
Requirement already satisfied: packaging>=20.0
Collecting tokenizers!=0.11.3,<0.14,>=0.11.1
  Downloading tokenizers-0.13.3-cp39-cp39-man
Requirement already satisfied: tqdm>=4.27 in
Requirement already satisfied: regex!=2019.12
Requirement already satisfied: fsspec in /usr
Requirement already satisfied: typing-extendi
Requirement already satisfied: idna<4,>=2.5 in
Requirement already satisfied: certifi>=2017.
Requirement already satisfied: urllib3<1.27,>
Requirement already satisfied: charset-normal
Installing collected packages: tokenizers, hu
Successfully installed huggingface-hub-0.14.0

```

xt	1931.txt	1930.txt	...
1	27.64137828350067101		

Downloading	687/687
(...)lve/main/config.json:	[00:00<00:00,
100%	25.9kB/s]
Downloading	1.42G/1.42G
pytorch_model.bin:	[00:14<00:00,
100%	150MB/s]
Downloading	256/256

```
1 !pip install symspellpy
```

Page 2 of 13

```

1 from symspellpy import SymSpell
2 import pkg_resources

1 symspell = SymSpell()
2 dict_path = pkg_resources.resource_filename("symspell", "dict.txt")
3 symspell.load_dictionary(dict_path, 0, 1) # Load dictionary
4 eng_dictionary = symspell.words # Get keyed list of words

1 from bs4 import BeautifulSoup

1 # read in the CSV file
2 df = pd.read_csv('/content/4_20_CA_newspapers - 4_20_CA_newspapers.csv')
3
4 print(df.columns)
5 df.columns = [c.strip() for c in df.columns]
6
7 # convert the 'First Issue Date' and 'Last Issue Date' to datetime
8 df['First Issue Date'] = pd.to_datetime(df['First Issue Date'])
9 df['Last Issue Date'] = pd.to_datetime(df['Last Issue Date'])
10
11 # extract the year from the 'First Issue Date' and 'Last Issue Date'
12 df['First Issue Year'] = df['First Issue Date'].dt.year
13 df['Last Issue Year'] = df['Last Issue Date'].dt.year
14
15 # filter the dataframe to only keep entries with 'First Issue Year' < 1940
16 df_filtered = df[(df['First Issue Year'] < 1940)]
17
18 # print out the filtered dataframe
19 print(df_filtered)

Index(['Persistent Link', 'State', 'Title', 'No. of Issues', 'First Issue Date',
      'Last Issue Date', 'First Issue Year', 'Last Issue Year'],
      dtype='object')

```

	Persistent Link	State	Title	No. of Issues	First Issue Date	Last Issue Date	First Issue Year	Last Issue Year
2	<a href="https://chroniclingamerica.loc.gov/lccn/sn83-03427-1/1940-01-01">https://chroniclingamerica.loc.gov/lccn/sn83-03427-1/1940-01-01</a>	CA	The San Francisco Chronicle	1	1940-01-01	1940-01-01	1940	1940
22	<a href="https://chroniclingamerica.loc.gov/lccn/sn83-03427-1/1940-01-01">https://chroniclingamerica.loc.gov/lccn/sn83-03427-1/1940-01-01</a>	CA	The San Francisco Chronicle	1	1940-01-01	1940-01-01	1940	1940
26	<a href="https://chroniclingamerica.loc.gov/lccn/sn83-03427-1/1940-01-01">https://chroniclingamerica.loc.gov/lccn/sn83-03427-1/1940-01-01</a>	CA	The San Francisco Chronicle	1	1940-01-01	1940-01-01	1940	1940
27	<a href="https://chroniclingamerica.loc.gov/lccn/sn83-03427-1/1940-01-01">https://chroniclingamerica.loc.gov/lccn/sn83-03427-1/1940-01-01</a>	CA	The San Francisco Chronicle	1	1940-01-01	1940-01-01	1940	1940
34	<a href="https://chroniclingamerica.loc.gov/lccn/sn83-03427-1/1940-01-01">https://chroniclingamerica.loc.gov/lccn/sn83-03427-1/1940-01-01</a>	CA	The San Francisco Chronicle	1	1940-01-01	1940-01-01	1940	1940
...								
3992	<a href="https://chroniclingamerica.loc.gov/lccn/sn83-03427-1/1940-01-01">https://chroniclingamerica.loc.gov/lccn/sn83-03427-1/1940-01-01</a>	CA	The San Francisco Chronicle	1	1940-01-01	1940-01-01	1940	1940
3999	<a href="https://chroniclingamerica.loc.gov/lccn/sn83-03427-1/1940-01-01">https://chroniclingamerica.loc.gov/lccn/sn83-03427-1/1940-01-01</a>	CA	The San Francisco Chronicle	1	1940-01-01	1940-01-01	1940	1940
4000	<a href="https://chroniclingamerica.loc.gov/lccn/sn83-03427-1/1940-01-01">https://chroniclingamerica.loc.gov/lccn/sn83-03427-1/1940-01-01</a>	CA	The San Francisco Chronicle	1	1940-01-01	1940-01-01	1940	1940
4001	<a href="https://chroniclingamerica.loc.gov/lccn/sn83-03427-1/1940-01-01">https://chroniclingamerica.loc.gov/lccn/sn83-03427-1/1940-01-01</a>	CA	The San Francisco Chronicle	1	1940-01-01	1940-01-01	1940	1940

4004 <https://chroniclingamerica.loc.gov/lccn/2011-01-01>

2 The Birmingham age-herald. [volume] (Birmingham, Ala.)  
 22 The "J" bird. [volume] (Juneau, Alaska)  
 26 The Alaska daily empire. [volume] (Juneau, Alaska)  
 27 The Alaska fisherman. [volume] (Juneau, Alaska)  
 34 The Alaskan. (Petersburg, Alaska)  
 ...  
 3992 The Cody enterprise. (Cody, Wyo.)  
 3999 The Northern Wyoming herald. (Cody, Wyo.)  
 4000 Park County herald. (Cody, Wyo.)  
 4001 Rawlins Republican. (Rawlins, Wyo.)  
 4004 The Saratoga sun. (Saratoga, Carbon Co., Wyo.)

	OCLC	ISSN	No. of Issues
2	12607279.0	2692-6318	8237.0
22	767526793.0	NaN	88.0
26	3039521.0	2576-9227	4340.0
27	29593944.0	NaN	97.0
34	31214772.0	NaN	320.0
...	...	...	...
3992	25224927.0	NaN	184.0
3999	25224904.0	NaN	406.0
4000	25224920.0	NaN	117.0
4001	27184809.0	NaN	1617.0
4004	9965363.0	0740-4948	1748.0

2 <https://chroniclingamerica.loc.gov/lccn/2011-01-01>

22

26 <https://chroniclingamerica.loc.gov/lccn/2011-01-01>

27

34

...

3992 <https://chroniclingamerica.loc.gov/lccn/2011-01-01>

3999 <https://chroniclingamerica.loc.gov/lccn/2011-01-01>

4000 <https://chroniclingamerica.loc.gov/lccn/2011-01-01>

4001 <https://chroniclingamerica.loc.gov/lccn/2011-01-01>

4004 <https://chroniclingamerica.loc.gov/lccn/2011-01-01>

Last Issue Year

2 1924.0

22 1936.0

26 1926.0

Loop through all the sn# & years

```
1 def get_url(lccn, year):
2     # Example url:
3     url = f'https://chroniclingamerica.loc.gov/'
4     # url = 'https://chroniclingamerica.loc.gov/'
5     lccn_label = lccn.strip()
6
7     session = requests.Session()
8     response = session.get(url).text
9     # print(response)
10
11     # Parse out issue dates from the response using re
12     # pattern = rf'<a href="/lccn/{lccn}/{year}'
13     # pattern = rf'<a href="/lccn/{lccn}/{year}'
14
15     pattern = r'<a href="/lccn/' + lccn.strip() +
16     matches = re.findall(pattern, response)
17
18     # Return a dictionary indexed by lccn and year
19     return {'lccn': lccn, 'year': year, 'matches': matches}
```

```

1 import os
2
3 def save_matches_to_file(lccn, year, matches):
4     # Helper function to save matches to file
5     directory = f'/content/urls/{lccn.strip()}/'
6     print(f'{directory}{year}.txt')
7     if not os.path.exists(directory):
8         os.makedirs(directory)
9     with open(f'{directory}/{year}.txt', 'w') as outfile:
10         outfile.writelines(matches)
11
12 def scrape_urls(df_filtered):
13     # Main function to scrape URLs and save matches
14     # for index, row in df_filtered.iterrows():
15     #     print(row)
16     for index, row in df_filtered.iterrows():
17         lccn = row['LCCN']
18         # first_year = row['First Issue Year']
19         # last_year = row['Last Issue Year']
20         first_year = int(row['First Issue Year'])
21         last_year = int(row['Last Issue Year'])
22         for year in range(first_year, last_year+1):
23             # result = get_url(lccn, year)
24             # matches = result['matches']
25             if int(year) >= 1924 and int(year) <= 1940:
26                 print("LCCN",lccn,"year:",year)
27                 texts_score, chunk_number = get_text(lccn, year)
28                 if (chunk_number > 0):
29                     score = texts_score/chunk_number
30                     print("LCCN",lccn,"year:",year, 'score:',score)
31                     # save_matches_to_file(lccn, year, matches)
32                     save_matches_to_file(lccn, year, score)
33             else:
34                 continue

```

```

1 import re
2
3 # html_text = '<tr><td class="single sun"><a href=
4 html_text = ' <a href="/lccn/sn86072192/1897-08
5 # Use regular expression to find all matches
6 # pattern = r'<a href="/lccn/sn86072192/1897-.*?
7 lccn = ' sn86072192 '
8 year = 1897
9 pattern = r'<a href="/lccn/' + lccn.strip() + '/'
10 matches = re.findall(pattern, html_text)
11
12 # Print all matches
13 for match in matches:
14     print(match)

```

```

    <a href="/lccn/sn86072192/1897-08-01/ed-1/seq

```

```

1 def get_word_rate(text, dictionary): # Function
2     num_words = 0
3     for word in text.split():
4         if word.lower() in dictionary.keys():
5             num_words += 1
6
7     return num_words / len(text.split())

```

```

1 def chunk_text(text, chunk_size):
2     chunks = []
3     words = text.split()
4     i = 0
5     while i + chunk_size < len(words):
6         chunks.append(' '.join(words[i: i + chunk_s
7         i += chunk_size
8     chunks.append(' '.join(words[i:]))
9     return chunks

```

```

1 chunk_size = 100
2 WORD_RATE_THRESHOLD = .6 # Set this to whateve

```

```

1 def get_text(lccn, year):
2     # Fetch the URLs of all issues in the given
3     url_data = get_url(lccn, year)

```

```

4     issue_urls = url_data['matches']
5     head = 'https://chroniclingamerica.loc.gov/'
6     tail = 'ocr/'
7     text = ''
8     score = 0
9     chunk_number = 0
10    old_url = head
11    # Loop through all issue URLs and extract the
12    for issue_url in issue_urls:
13        # Extract the URL from the issues_url variable
14        url = issue_url.strip('<a href="').strip('"')
15
16        # Combine the URL with the head variable
17        full_url = head + url
18        response = requests.get(full_url)
19        html_content = response.text
20        print(full_url)
21        # Use regular expression to find all matches
22        pattern = r'<a href="/lccn/' + lccn.strip('"') + '/'
23        text_matches = re.findall(pattern, html_content)
24
25        for text_url in text_matches:
26            t_url = text_url.strip('<a href="').strip('"')
27            ocr_url = head + t_url + tail
28            print(ocr_url)
29            if (ocr_url == old_url):
30                continue
31            else:
32                response = requests.get(ocr_url)
33                ocr_content = response.text
34
35                # Parse the HTML content using BeautifulSoup
36                soup = BeautifulSoup(ocr_content, 'html.parser')
37
38                # Extract the text content of the page
39                for elem in soup.find_all('div'):
40                    if elem.find('p') is not None:
41                        text = elem.find('p').get_text()
42                        chunks = chunk_text(text)
43                        for chunk in chunks:
44                            rate = get_word_rate(chunk)
45                            if (rate > WORD_RATE):
46                                chunk_number += 1
47                                if chunk_number == 100:
48                                    break

```



```

49         # score += senti
50         data =sentiment_
51         print(data)
52
53         if data[0]['label'] == 'positive':
54             score += 1
55         elif data[0]['label'] == 'negative':
56             score -= 1
57         if chunk_number > 100:
58             break
59         print('current score', score)
60         old_url = ocr_url
61         text = ''
62
63         if chunk_number > 100:
64             break
65
66         if chunk_number > 100:
67             break
68
69     return str(score), str(chunk_number)
70

```

```

1 def get_text(lccn, year):
2     # Fetch the URLs of all issues in the given year
3     url_data = get_url(lccn, year)
4     issue_urls = url_data['matches']
5     head = 'https://chroniclingamerica.loc.gov/'
6     tail = 'ocr/'
7     text = ''
8     score = 0
9     chunk_number = 0
10    old_url = head
11    # Loop through all issue URLs and extract the text
12    for issue_url in issue_urls:
13        # Extract the URL from the issues_url variable
14        url = issue_url.strip('<a href="').strip('"')
15
16        # Combine the URL with the head variable
17        full_url = head + url
18        response = requests.get(full_url)
19        html_content = response.text
20        print(full_url)
21        # Use regular expression to find all matches

```

```

22     # pattern = r'<a href="/lccn/'+ lccn.sti
23     pattern = r'<a href="/lccn/'+ lccn.stri
24     # r'<a href="/lccn/sn86072192/1897-08-22
25     text_matches = re.findall(pattern, html_
26
27     for text_url in text_matches:
28         t_url = text_url.strip('<a href="').
29         ocr_url = head+ t_url + tail
30         print(ocr_url)
31         if (ocr_url == old_url):
32             continue
33         else:
34             # Fetch the HTML content of the page wit
35             # ocr_url = issue_url.replace('/ed-', '/')
36             response = requests.get(ocr_url)
37             ocr_content = response.text
38
39             # Parse the HTML content using B
40             soup = BeautifulSoup(ocr_content
41
42             # Extract the text content of th
43
44             for elem in soup.find_all('div')
45                 if elem.find('p') is not Nor
46                     # text += elem.find('p')
47                     text = elem.find('p').ge
48                     # print(len(text))
49                     chunks = chunk_text(text
50                     for chunk in chunks:
51                         rate = get_word_rate(c
52                         # print(rate)
53                         if (rate > WORD_RATE_T
54                             # print('I am here
55                             # data =sentiment_
56                             # # print(data)
57                             chunk_number += 1
58                             # if data[0]['labe
59                             #     score += -1
60                             # elif data[0]['la
61                             #     score += dat
62                     print('current score', s
63                     # print('ocr url:', ocr_
64                     old_url = ocr_url
65                     text = ''
66             if chunk_number > 100:

```


```
1 scrape_urls(df_filtered)
```

[illegible]



[Colab paid products](#) - [Cancel contracts here](#)

--VISUAL--

 1h 33m 49s    completed at 11:06 PM