# Web Scraper

## Gurpreet Singh

## March 6, 2025

# 1 Introduction

The web scraper, an integral part of the project, as it allows to fetch article data for the implementation of the models. The web scraper was initially designed to scrape publicly accessible data from various websites in order to minimise bias and increase credibility of data. However, along with other changes, the scraper was made to be specific to one website, the BBC news.

BBC news, a leading news agency in the UK. (**Add more context here**).

# 2 Design

The design of the scraper has been modified several times over the course of the project to ensure that the most optimal and efficient design is chosen to strictly adhere to time guidelines.

## 2.1 Initial Design

### 2.1.1 Methodology

Many methodologies were explored such as traditional and hybrid models. It was decided that the hybrid models would be best hence the spiral model was chosen as it allowed continuous improvement until the sought for product was made.

**add figure showing visual representation of model and description using references**

### 2.1.2 Structure

The use of different python scripts was firstly discussed to ensure that each component would work well independently and therefore not become an unsolvable issue where at the last step of the implementation all components would be imported into one file and that would be the executable.

The design revolved around each script being able to do only one component hence the python scripts were first designed as expressed in table 1.

| Script Name | Description |
|---|---|
| scraper | loading the API or bs4. |
| fetching | accessing different data from the website. |
| dataFormat | import and export of data to and from the python data types. |
| main | executing all scripts in order to run the web scraper. |

Table 1: Python scripts

### 2.1.3 Python Modules

It was important to understand that to gain access to data, a secure path was needed to the website, in light of this there were two potential pathways that were explored:

1. API;

2. Selenium;

3. BeautifulSoup4(bs4).

All three of the options had their advantages and limitation, and predominantly all deemed as not viable options. API's were either specific to the different projects worked in, outdated or simply missing components. Whereas, selenium and bs4, due to lack of prior knowledge, became difficult to comprehend, establish a successfully connection and essentially scrape the website.

To transfer data to files the csv module was initially selected as the most expertise were with that file type, however this was deemed inefficient as the data was not stored to be easily accessible.

Request module was also needed in coherence to the above outlined options 1-3. This module was used to send requests which allowed to gain access to specific data from the HTML.

## 2.2 Current design model

The current design model was completed after careful consideration of deadlines, limitations and through trial and error of software.

### 2.2.1 Methodology

The methodology remained the same, the spiral method. The hybrid model allowed constant critique of implemented code whilst allowing the betterment of the scraper. Each process would be put through the process individually to ensure that before the scraper work as a whole they work individually. The process are than collected and executed together and the method is implemented on the code again.

### 2.2.2 Structure

The processes were not to be implemented the same way as a decision was made before implementing to use the OOP structure to ensure that a scraper object can be called and used as easily on future projects as well. The differentiation is that for this to be successfully be operational, all components were needed to be further broken down so that each component had subroutines which could be called individually. The subroutines which are listed in Table 2 are carefully outlined so that the implementation would be completed with great regards to time efficiency.

| Subroutine name | Function |
|---|---|
| init | set self variables |
| fetchdata | load and manipulate data from json files |
| loadHtml | get HTML of page |
| getarticledata | execution of all subroutines to get article data |
| getlinks | getting all article links from results page |
| main | executing all scripts in order to run the web scraper. |

Table 2: WebScraper class subroutines

# 3 Implementation

The implementation was carried out over several months and completed to satisfy basic requirements of the scraper. An OOP model was created to ensure smooth modification at later stages. The scraper collects data into two json files, an article links file and and article data file. (**add figures and make metadata**)

The webscraper class is the only class in the script defined on line 16, in Appendix A.1. The webscraper class has one parameter, an initial url needs to be inputted for it to work. This is done on lines 194, 208, 228, 246 in appendix A.1.

# A   Appendix- Web Scraper code

```python
###
#WebScraper
#OOP designed software to scrape data off given website and collect publicly available d
#Gurpreet singh
#21131818
###

from requests_html import HTMLSession
from time import sleep
from datetime import datetime
from dateutil import parser
import json
import os
#import all required modules

class Webscraper():
    def __init__(self, url) -> None:
        #initial constructor creation
        self.session = HTMLSession()
        self.base_url = "https://www.bbc.co.uk/"
        self.url = url

    def load_html(self):
        #load the html for the webpage
        self.reader = self.session.get(self.url)
        self.reader.html.render(sleep=1, scrolldown=0)

    def close(self):
        #used to complete session and free all resources
        self.session.close()

    def fetch_aldata(self, div_finder):
        #this selects the link and title from the div by finding correct data in the htm
        divs = self.reader.html.find(div_finder)
        data = []
        for div in divs:
            # Goes through specific objects to fetch data
            link_element = div.find("a", first=True)
            if link_element:
                link = link_element.attrs.get('href', '')
                title = link_element.text
                data.append({'title':title, 'link': link})
        return data

    def fetch_adata(self):

        article_data = []
        msg = "No"
        article = ''
        title = (self.reader.html.find('title', first=True) or self.reader.html.find('h1

        try:
            paragraphs = self.reader.html.find('article')
            if paragraphs:
                for p in paragraphs:
                    article += p.text.strip()
```

3

```
57              else:
58                    article = msg
59          except Exception as e:
60              print(f"An error occurred: {e} article")
61
62          date = None
63          try:
64              time_elements = self.reader.html.find('time')
65              if time_elements:
66
67                  date_published = time_elements[0].text.strip() or time_elements[0].dateti
68                  print (date_published)
69                  try:
70                      date_published = parser.parse(date_published, default=datetime(dateti
71                      date_published = date_published.strftime('%d-%m-%Y')
72                  except:
73                      date_published = msg
74
75
76              else:
77                  date_published = msg
78          except Exception as e:
79              print(f"An error occurred: {e} time")
80
81
82
83          try:
84              author = self.reader.html.find('.ssrcss-68pt20-Text-TextContributorName', fi
85              author = author.text
86
87              if not author:
88                  author = msg
89          except Exception as e:
90              print(f"An error occurred: {e} author")
91
92          article_data.append({'title': title, 'article': article, 'Publishdate': date_pub
93
94          return article_data
95
96      def pagination(self):
97          #get links for all successive pages of the results
98          self.reader = self.session.get(self.url)
99          npg_links = []
100         for link in self.reader.html:
101             if link.search("page="):
102                 print("yes")
103                 #sleep(0.5)
104                 link = str(link)
105                 link = link.split("'")[1]
106                 npg_links.append(link)
107                 sleep(1)
108
109             else:
110                 break
111         return npg_links
112
```

```python
113         def link_format(self):
114             #check url is correct format
115             search_ext = self.url
116             if search_ext and not search_ext.startswith("https://"):
117                     self.url = self.base_url + search_ext
118             else:
119                 self.url = search_ext
120                 return self.url
121             pass
122
123         def file_checker(self, found, filename):
124
125             path = os.path.dirname(os.path.realpath(__file__))
126             destination = ""
127             if found == False:
128
129                 for root, dirs, files in os.walk(path):
130                     for file in files:
131                         if file.endswith('.json') and file.startswith(filename):
132                             print ('yes')
133                             destination =  (root+'/'+str(file))
134                             print (destination)
135                             found = True
136
137             return found, destination
138
139         def datatype_conversion(self, found, data, destination):
140
141             if found  == False:
142
143                 json_obj = json.dumps(data, indent=4)
144                 with open(destination, "w") as crfile:
145                     crfile.write(json_obj)
146                     found = True
147                     return found, 'File-created-and-hyperlinks-stored-successfully'
148
149             elif found == True:
150                     with open(destination, 'r+')  as file:
151                         json_obj = json.load(file)
152                     for item in data:
153                         json_obj.append(item)
154
155                     print(json_obj)
156
157                     with  open(destination, 'w') as file:
158                         json.dump(json_obj, file,
159                                         indent=4,
160                                         separators=(',',':-'))
161
162                     return found, 'Data-stored-successfully'
163
164
165
166         def fetch_adata_links(self,filename):
167             with open(filename, 'r+') as file:
168                 link_data = json.load(file)
```

```
169                return link_data
170
171        def div_select(self):
172            #div selector for specific data needed from html
173            if self.url and not self.url.endswith("NEWS_PS") or self.url.find("page="):
174                div_finder = ".ssrcss-tq7xfh-PromoContent > *"
175
176            else:
177                div_finder= "h1.ssrcss-1j5vay3-Heading.e1hq9lx0"
178                #should be for individual article page - find specific div that correllates
179                #might need more thn one check page
180
181            return div_finder
182
183 def get_article_link_data ():
184     #runs all the function and evokes the object
185     has_run = False
186     div_finder = ".ssrcss-tq7xfh-PromoContent > *"
187     found  = False
188     filename = 'article_links'
189
190
191
192     while has_run  == False:
193         url = "search?q=s%26p+500&seqId=e1005640-2774-11ef-b757-6398eaf17df6&d=NEWS_PS"
194         scraper = Webscraper(url)
195         scraper.link_format()
196         npg_links = scraper.pagination()
197         has_run = True
198         found, destination = scraper.file_checker(found, filename)
199         print("—")
200         print(found)
201         scraper.close()
202
203
204     try:
205
206         for link in npg_links:
207             url = link
208             scraper = Webscraper(url)
209             print(scraper.url)
210             div_finder = scraper.div_select()
211             scraper.load_html()
212             data = scraper.fetch_aldata(div_finder)
213             found, message = scraper.datatype_conversion(found, data, destination)
214             print (message)
215             for item in data:
216                 print(item)
217             sleep(1)
218             print("_____")
219             scraper.close()
220             print(found)
221
222     except Exception as e:
223         print(f"An error occurred: {e}")
224
```

```python
225
226  def get_article_data():
227      url = "search?q=s%26p+500&seqId=e1005640-2774-11ef-b757-6398eaf17df6&d=NEWS_PS"
228      scraper = Webscraper(url)
229      al_file = 'article_links.json'
230      a_file = 'article_data.json'
231      url = scraper.fetch_adata_links(al_file)
232      found = False
233      found, filename = scraper.file_checker(found, al_file)
234      link_data = scraper.fetch_adata_links(filename)
235      found = False
236      found, destination = scraper.file_checker(found, a_file)
237      if found == False:
238          destination = a_file
239
240
241      try:
242
243          for link in link_data:
244              url = (link['link'])
245              print(url)
246              scraper = Webscraper(url)
247              scraper.load_html()
248              data = scraper.fetch_adata()
249              found, message = scraper.datatype_conversion(found, data, destination)
250              print(message)
251              print('----------------------------')
252
253              scraper.close()
254
255
256      except Exception as e:
257          print(f"An error occurred: {e}")
258
259  def main():
260      get_article_data()
261      #get_article_link_data()
262
263
264
265
266  if __name__ == "__main__":
267      main()
```