

# 使用 Selenium 爬取 ChinaDaily 搜索结果

搜索网页：

CDMPC 2022-2023：

<https://newssearch.chinadaily.com.cn/en/search?cond=%7B%22publishedDateFrom%22%3A%222022-08-24%22%2C%22publishedDateTo%22%3A%222023-08-23%22%2C%22fullMust%22%3A%22marine+pollution%22%2C%22sort%22%3A%22dp%22%2C%22duplication%22%3A%22on%22%7D&language=en>

CDMPC      2023-2024      :      <https://newssearch.chinadaily.com.cn/en/search?cond=%7B%22publishedDateFrom%22%3A%222023-08-24%22%2C%22publishedDateTo%22%3A%222024-08-23%22%2C%22fullMust%22%3A%22marine+pollution%22%2C%22sort%22%3A%22dp%22%2C%22duplication%22%3A%22on%22%7D&language=en>

## 1. 任务流程

自行进行详细搜索后复制网址并进行爬取，为了获取全文内容，需要从左侧列表里获取标题，并通过超链接进入网络页面，然后再爬取文章内容。

## 2. 爬取文章标题列表

## 2.1. 获取当前静态页面的文章列表

在总览页面下，文章标题列表被保存在一个 `body > div.main > div.main_l > div.cs_result > div.lft_art` 中，每个文章通过一个 `<div>` 类进行封装，里面包括图像（`<a class="art_pic">`）、文章信息（`<span class="intro">`），链接包括在里面的 `<a>` 类中，通过调用 `href` 属性获得。

```
1 <div class="art_detail">
2     <a class="art_pic"
  href="https://www.chinadaily.com.cn/a/202308/23/WS64e55928a31035260b81d
  b2f.html" target="_blank">... </a>
3     <span class="intro">
4         <h4>
5             <a
  href="https://www.chinadaily.com.cn/a/202308/23/WS64e55928a31035260b81d
  b2f.html" target="_blank">S. Koreans condemn Japan's nuke wastewater
  dumping decision
6             </a>
7         </h4>
8         <b> (Xinhua) 2023-08-23 09:28</b>
9         <b style="margin-bottom:2px"> London Convention on the
  Prevention of
10             <em>Marine</em>
11             <em>Pollution</em> by Dumping of Wastes and Other Matter.
  The
12         </b>
13     </span>
14 </div>
```

为了方便处理，通过 `css selector` 获取链接，保存到一个字典列表，每个元素是一个 {标题: 链接} 的字典。

```
1 def get_art_urls(driver, urls):
2     articles = driver.find_elements("css selector", "span.intro")
3     # print(len(results))
4     # print(len(results))
5     for r in articles:
6         # 每个 span 类
```

```
7         art = r.find_element("css selector", 'a')
8         # 标题
9         print(art.text)
10        print(art.get_attribute("href"))
11        article = {
12            'title': art.text,
13            'href': art.get_attribute('href')
14        }
15        # art.text = art.get_attribute("href")
16        urls.append(article)
17    return urls
```

## 2.2. 翻页

在CHINADAILY 的搜索页面提供了翻页按键实现翻页，封装在左边的列表框的右上方，如下图所示。在开头和结尾的时候都会在相应的端口去掉 << PREVIOUS 和 NEXT >> 按钮，但是在当前任务的翻页过程中，只需要 NEXT >> 按钮。

Page: <<PREV6 7 8 9 10

Page:1 2 3 4 5 NEXT >>

```

1  <div class="page rt">
2      <span>Page:
3          <b>1</b>&nbsp;
4          <a href="javascript:void(0);"
onclick="$SearchController.pagging(1)" title="page 2">2</a>&nbsp;
5          <a href="javascript:void(0);"
onclick="$SearchController.pagging(2)" title="page 2">3</a>&nbsp;
6          <a href="javascript:void(0);"
onclick="$SearchController.pagging(3)" title="page 2">4</a>&nbsp;
7          <a href="javascript:void(0);"
onclick="$SearchController.pagging(4)" title="page 2">5</a>&nbsp;
8          <a href="javascript:void(0);"
onclick="$SearchController.pagging(1)" title="next">NEXT &gt;&gt;</a>
9      </span>
10 </div>

```

将尾页鉴定和翻页操作分开到两个函数中，以便处理当前页面的信息。在翻页的时候，如果页面栏中不存在 NEXT >> 之后，返回一个空值一表示到尾页。翻页后给定时间以便处理。

```

1  ##### Part 2: Page down
2  def current_pages(driver):
3      # 上下都有一个换页栏，选择一个
4      page_block = driver.find_element("css selector", "div.page.rt")
5      pages = page_block.find_element("css selector", "span")
6
7      # 翻页，如果没有 NEXT 就不翻页了
8      if "NEXT" not in pages.text:
9          return None
10     else:
11         return pages
12
13  def next_page(driver, pages):
14      page_buttons = pages.find_elements("css selector", "a")
15      for a in page_buttons:
16          # print(a.text)
17          if "NEXT" in a.text:
18              print("Go to the next page.")
19              driver.get(a.get_attribute("href"))

```

```
20         # actions.click(a)
21         a.click()
22         driver.implicitly_wait(10)
23         import time
24         time.sleep(5)
25         break
```

## 2.3. 从搜索结果中获取对应的所有文章

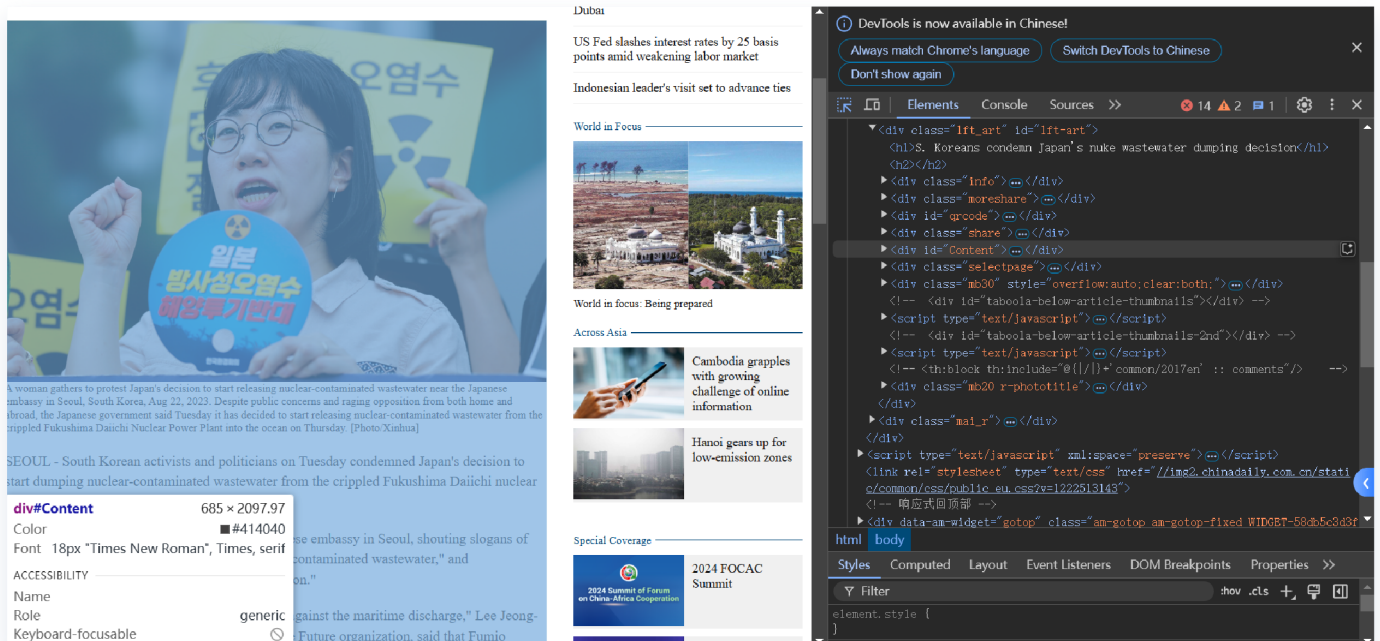
```
1  def get_articles_from(driver, url):
2      driver.get(url)
3      driver.implicitly_wait(10)
4      # store urls
5      urls = []
6      # 获取页面信息
7      pages = current_pages(ff_driver)
8      while True:
9          # 获取当前页面的所有文章链接
10         urls = get_art_urls(ff_driver, urls)
11         # 提取完最后一页需要跳出
12         if not pages:
13             break
14         # 翻页并重新获取页面信息
15         next_page(ff_driver, pages)
16         pages = current_pages(ff_driver)
17     return urls
```

之后在主函数中将文章链接存在对应的 csv 文件中。

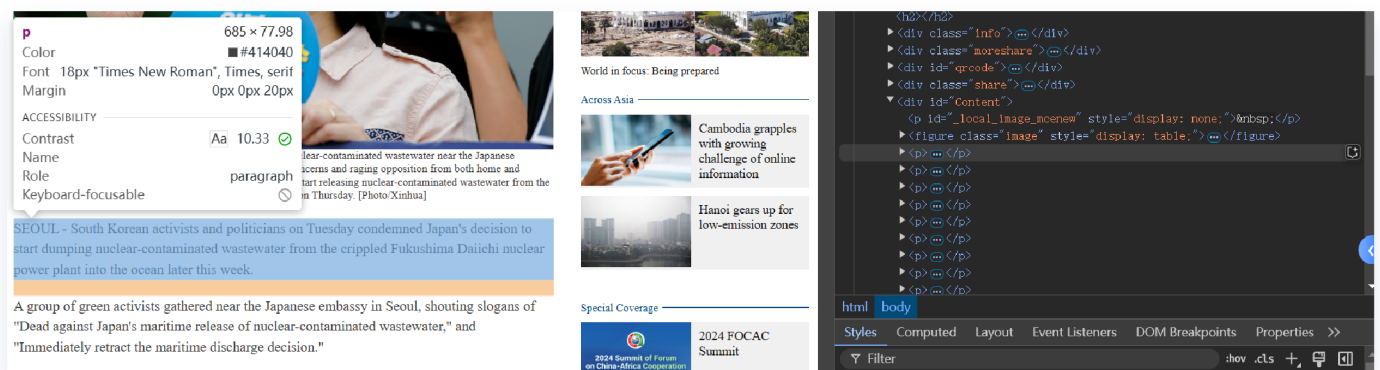
## 3. 获取文章具体内容

### 3.1. 获取单一网页下的文章内容

以第一篇文章为例，如果我们直接输出 `<div class="Content">` 的文本，则会将图像标题包括在里面。



为了解决这个问题，使用 `xpath` 选择这个板块下的 `<p>` 类，这是每个自然段分开的。



之后再将所有内容拼接起来即可。

```
1 def article_content(driver, url, title, folder, idx):
2     driver.get(url)
3     driver.implicitly_wait(20)
4     # 去除一些“声明”
5     def is_normal(t):
6         if "necessarily reflect those of China Daily." in t \
7             or "send us your writings" in t:
8             print(f"Announcement in ", url)
```

```

9         return False
10    else:
11        return True
12    # 获取段落, 在 content.p 里面
13    texts = driver.find_elements("xpath", "//*[@id='Content']/p")
14    # 拼接段落并输出
15    text = "\n".join([t.text for t in texts if is_normal(t.text)])
16    # 添加标题
17    text = "\n".join([f"Title: {title}", text])
18    # 替换冒号为连接号
19    title = title.replace(":", "-")
20    output_file = os.path.join(folder, str(idx + 1) + "-" + title +
11    ".txt")
21    with open(output_file, mode='w', encoding='utf-8') as txtfile:
22        txtfile.write(text)
23    return text

```

## 3.2. 通过遍历获取全年的相关文献

首先从上一章中保存的文件读取每个年份的列表，逐个超链接打开，然后再调用 `article_content` 函数即可。这里为了防止加载不成功还设置了一个循环用于打不开时刷新页面。

```

1  def extract_by_year(driver, df, tgt_path, together_path):
2      cnt = 1
3      for idx, art in df.iterrows():
4          with open(together_path, mode='a', encoding='utf-8') as
12  txtfile:
13
14      for t in range(retry_times):
15          try:
16              content = article_content(ff_driver, art["href"],
17  art["title"], tgt_path, idx)
18              break
19          except:
20              print(f"Time out, retrying time {t+1}.")
21              time.sleep(1)

```



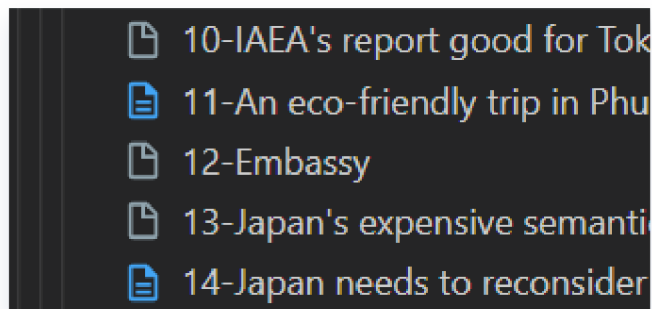
```

13         txtfile.write(str(idx + 1).center(4, " ").center(40, "#") +
14         "\n" + content + '\n')
14         time.sleep(2)

```

### 3.3. 具体问题解决

发现内容截断（得到空文件），主要是因为部分标题中会有 `:` 存在，在代码中将其替换成 `-` 解决。



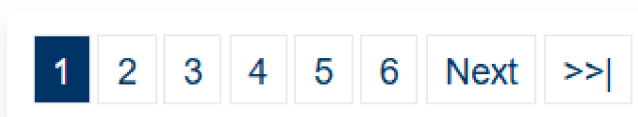
对部分“声明”进行了筛选。

```

1     def is_normal(t):
2         if "necessarily reflect those of China Daily." in t \
3             or "send us your writings" in t:
4             print(f"Announcement in ", url)
5             return False
6         else:
7             return True

```

还有一个问题就是部分的文章是分页式的结构，这个时候就需要获取这个分页代码然后分页读取。在这里是在文章中设置了一个分页按钮，长这样。



实际上是因为属于习总书记的讲话、活动一类要闻都进行了分页。根据标题把“Xi”提取出来之后进行又出现报错，修改代码逻辑，如果遇到报错则返回用单页的处理方式。具体代码如下：

[illegible]

```

33         # 确保页面已经完全加载, 可以在这里使用
driver.get() 强制加载
34         driver.get(current_url)
35
36         # 等待新页面加载完成
37         WebDriverWait(driver, 10).until(
38
39             EC.presence_of_element_located((By.XPATH, "//*[ @id='div_currpage']"))
40         )
41
42         # 重新获取分页块元素
43         page_block =
driver.find_element(By.XPATH, "//*[ @id='div_currpage']")
44         page_buttons =
page_block.find_elements(By.CSS_SELECTOR, "a")
45         break
46
47         except Exception as err:
48             print(f"Error: {err}")
49             driver.quit()
50             break
51
52         else:
53             break
54
55         except Exception as err:
56             print("Error: when extracting in a multi-page manner,
exception occured:")
57             print(err)
58             print("Try to solve it in single-page manner.")
59             # 获取段落, 在 content.p 里面
60             texts = driver.find_elements("xpath", "//*[
[id='Content']/p")
61             # 拼接段落并输出
62             text = "\n".join([t.text for t in texts if
is_normal(t.text)])
63             # 添加标题
64             text = "\n".join([f>Title: {title}", text])
65
66         # .

```

