

Cheating Bot Detector

Fall Term Retrospective

Grein, Cameron Gareth

Department of Electrical Engineering
and Computer Science,
Oregon State University
Corvallis, Oregon, United States
greinc@oregonstate.edu

Yidong Lin

Department of Electrical Engineering
and Computer Science,
Oregon State University
Corvallis, Oregon, United States
linyid@oregonstate.edu

Section 1

Due to the continuous improvement of the Internet, learning aid websites such as Chegg and Coursehero have appeared. They claim to be a social learning network. Users can access old test questions, assignments, answers, notes, etc. through these sites. Of course, these are not free, users need to pay a certain amount to continue to access the content. However, you can still browse the content by publishing the materials, and the 40 documents will provide users with 1 month of free browsing. This has led to the leakage of teaching content in many schools. The reason why our product is beneficial that not only the content of Oregon State University's courses has been compromised. On these websites you can find almost all university courses, tests, and answers. Our current work is only for Oregon State University, if the product can achieve the desired effect. The beneficiaries of the product can be more schools. In relation to the scope of the assignment, we should be able to make a basic scraper that can identify websites that contain the course information. Scraping things such as scrape should be quite achievable in the time span that we are provided. Doing things such as scraping images might be a little too ambitious for the time that we have. There also is the question of if we want this to be used as a detector to manually inform that the content needs to be taken down, or are we planning on implementing an automatic system. An automatic system may be out of the scope of the time that we have to complete this assignment.

Section 2

We have made some progress on the implementation of our project. We have completed a python script to search the web for our project. We have not created the user interface yet. The scope of our project seems doable in the timeframe that we have. We should be able to implement all of the features we listed, and what we have also looks good. In terms of scope the current scope of our project seems fine, maybe if we can implement all that we stated in our sprint 1 we could add some more features. Where we are heading in terms of the next sprint, is to finish the basic functionality, and start working on the user interface. We will be working on the project during winter break, hopefully to get more work done. Overall, I think we are doing a good job progression wise.

Section 3

The current problem is that Google will limit the frequency of access. The solution is to set the access delay.

Section 4

```
spider.py  => x
1  from requests_html import HTMLSession
2  from parsel import Selector
3  import pandas as pd
4  import time
5
6
7  class Crawler():
8      def __init__(self, delay):
9          self.session = HTMLSession()
10         self.keyword_list = self.get_keywords()
11         self.base_url = 'https://www.google.com.hk/'
12         self.data_list = []
13         self.delay = delay
14
15     def get_keywords(self):
16         with open('keywords.txt', 'r', encoding='utf-8-sig') as f:
17             keywords = [i.strip() for i in f.readlines() if i.strip()]
18         return keywords
19
20     def search(self, url, keyword, page):
21         time.sleep(self.delay)
22         print(f'>>>>>Processing keyword <{keyword}> {page}: {url}')
23         r = self.session.get(url)
24         # with open('test.html', 'w', encoding='utf-8-sig') as f:
25         #     f.write(r.text)
26         for i in Selector(r.text).xpath('//*[@id="rso"]/div[@class="g"]'):
27             item = {}
28             item['keyword'] = keyword
29             item['page'] = page
30             item['title'] = i.xpath('./div/div/a/h3/span/text()').get('').strip()
31             item['desc'] = i.xpath('string(./div/div[2]/div/span)').get('').replace('\n', '').replace('\t', '').replace(
32                 '\r', '').replace(' ', '').replace('\xa0', '').strip()
33             item['url'] = i.xpath('./div/div[2]/div/div/div/a[1]/@href').get('').strip()
34             if not item['url']:
35                 item['url'] = i.xpath('./div/div[1]/a/@href').get('').strip()
36             print(item)
37             self.data_list.append(item)
```

```

38
39
40     next_page_url = r.html.xpath('//a[@id="pnnext"]/@href')
41     if next_page_url:
42         next_page_url = self.base_url + next_page_url[0]
43         self.search(next_page_url, keyword, page + 1)
44
45     def save(self):
46         df = pd.DataFrame(self.data_list)
47         path = f'{str(int(time.time()))}.xlsx'
48         with pd.ExcelWriter(path, engine='xlsxwriter', options={'strings_to_urls': False}) as writer:
49             df.to_excel(writer, index=False, encoding='utf-8-sig')
50
51     def main(self):
52         for keyword in self.keyword_list:
53             url = self.base_url + 'search?q=' + keyword
54             page = 1
55             self.search(url, keyword, page)
56             self.save()
57
58 if __name__ == '__main__':
59     delay = input('>>>>>Please enter the request delay interval (seconds):\n').strip
60     try:
61         delay = int(delay)
62     except:
63         delay = 1
64     c = Crawler(delay)
65     c.main()
66

```

- Open the command prompt and run pip install -r requirements.txt
- Confirm the search term: keywords.txt
- Execute python spider.py

Section 5

Regarding some feedback we received when we displayed our project. People were confused at the goals of the project. Since then we have tried to pinpoint the goals of our project more.

There was also some feedback on what language we should use for which parts of our project. Though we feel that the languages and technologies we are using now will suffice.

The last feedback that was given to us was an efficient user interface, and what faculty may want in terms of features for our web crawler.