

《Python 程序设计高阶》

(2023-2024 学年第 1 学期)

实验报告二

学号:2021329600006 姓名: 陈昊天 班级: 计算机科学与技术 21(4) 班

学号:2021329621213 姓名: 陈佳伟 班级: 计算机科学与技术 21(3) 班

学号:2021329621257 姓名: 冯佳钧 班级: 计算机科学与技术 21(4) 班

§1 实验题 1

修改课件中的下载安全标准信息的爬虫，使之能够正常工作，将下载的信息保存在“安全.csv”文件中；参考的结果，看课件中安全.csv 这一节的输出结果；

§1.1 解答：

解答要对你实验结果截图，加上你的代码。

修改了 SSL 证书问题 [1]，在 requests.get 加入 verify=False 参数。

1.1.1 代码

```
1 import requests # 关键库，用于发送Http请求 # pip install requests
2 from lxml import etree # xml解析包，ElementTree # pip install lxml
3 import time
4 import random
5 import re
6 import csv
7 import ssl
8
9 # 忽略 SSL 证书警告
10 requests.packages.urllib3.disable_warnings(
11     requests.packages.urllib3.exceptions.InsecureRequestWarning
12 )
13
14
15 class Spider:
16     def __init__(self, name):
17         self.name = name
18         self.url = (
19             "https://openstd.samr.gov.cn/bzgk/gb/std_list?p.p1=0&p.p2="
20             + name
```

```

21         + "&p.p90=circulation_date&p.p91=desc"
22     )
23     print(self.url)
24     self.headers = [
25         "Mozilla/5.0 (Windows NT 6.3; WOW64) AppleWebKit/537.36 (KHTML,
           like Gecko) Chrome/39.0.2171.95 Safari/537.36",
26         "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_2)
           AppleWebKit/537.36 (KHTML, like Gecko) Chrome/35.0.1916.153
           Safari/537.36",
27         "Mozilla/5.0 (Windows NT 6.1; WOW64; rv:30.0) Gecko/20100101
           Firefox/30.0",
28         "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_2)
           AppleWebKit/537.75.14 (KHTML, like Gecko) Version/7.0.3
           Safari/537.75.14",
29         "Mozilla/5.0 (compatible; MSIE 10.0; Windows NT 6.2; Win64;
           x64; Trident/6.0)",
30         "Mozilla/5.0 (Windows; U; Windows NT 5.1; it; rv:1.8.1.11)
           Gecko/20071127 Firefox/2.0.0.11",
31         "Opera/9.25 (Windows NT 5.1; U; en)",
32         "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1; SV1; .NET
           CLR 1.1.4322; .NET CLR 2.0.50727)",
33         "Mozilla/5.0 (compatible; Konqueror/3.5; Linux) KHTML/3.5.5
           (like Gecko) (Kubuntu)",
34         "Mozilla/5.0 (X11; U; Linux i686; en-US; rv:1.8.0.12)
           Gecko/20070731 Ubuntu/dapper-security Firefox/1.5.0.12",
35         "Lynx/2.8.5rel.1 libwww-FM/2.14 SSL-MM/1.4.1 GNUTLS/1.2.9",
36         "Mozilla/5.0 (X11; Linux i686) AppleWebKit/535.7 (KHTML, like
           Gecko) Ubuntu/11.04 Chromium/16.0.912.77 Chrome/16.0.912.77
           Safari/535.7",
37         "Mozilla/5.0 (X11; Ubuntu; Linux i686; rv:10.0) Gecko/20100101
           Firefox/10.0 ",
38     ]
39
40     def get_each_link(self):
41         """获得首页链接，返回列表"""
42         index = 1 # 翻页
43         link_list = [] # 储存每一个链接的列表
44         while True:
45             try:
46                 headers = {
47                     "User-Agent": random.choice(self.headers) #
48                                     随机选择请求头，模拟不同的浏览器
49                 }
49                 session = requests.session()

```

```

50 session.verify = False
51 r = session.get(
52     self.url, headers=headers, verify=False
53 ) # 请求数据,返回给Response对象
54 # r = requests.get(self.url, headers=headers) #
    请求数据,返回给Response对象
55 r.encoding = r.apparent_encoding # 设置编码
56 # print(r.status_code)
57 html = r.text
58 selector = etree.HTML(html) # 将网页源码转换为 XPath
    可以解析的格式
59 Ids = selector.xpath(
60     '//tr/td[@style="text-align: left;"]/a/@onclick'
61 ) # 选择结点
62 print(Ids)
63
64 # 判断是否获取结束
65 if len(Ids) == 0:
66     break
67
68 for each_id in Ids:
69     each_id = each_id.replace("showInfo('",
    "").replace("'", "");
70     # link = 'http://www.gb688.cn/bzgk/gb/newGbInfo?hcno='
    + each_id
71     link = (
72         "https://openstd.samr.gov.cn/bzgk/gb/newGbInfo?hcno="
    + each_id
73     )
74
75     print(link)
76     link_list.append(link)
77 index += 1
78 if index > 1: # 只请求1页
79     break
80 # self.url =
    'http://www.gb688.cn/bzgk/gb/std_list?r=0.40987171406906286&page='
    + str(index) + '&pageSize=10&p.p1=0&p.p2=' + self.name +
    '&p.p90=circulation_date&p.p91=desc'
81 self.url = (
82     "https://openstd.samr.gov.cn/bzgk/gb/std_list?r=0.4098717140690628
83     + str(index)
84     + "&pageSize=10&p.p1=0&p.p2="
85     + self.name

```

```

86         + "&p.p90=circulation_date&p.p91=desc"
87     )
88     print(self.url)
89     time.sleep(random.random() + random.randint(0, 1))
90 except Exception as e:
91     print(e)
92     break
93 return link_list
94
95 def get_detial_from_linkList(self, link_list):
96     """从列表的链接下载文件内容"""
97     path = self.name + ".csv"
98     i = 0 # 序号
99     # 写入表头
100    info = [
101        "序号",
102        "标准编号",
103        "标准名称",
104        "英文名称",
105        "状态",
106        "中国标准分类号 (CCS) ",
107        "国际标准分类号 (ICS) ",
108        "发布日期",
109        "实施日期",
110        "主管部门",
111        "归口单位",
112        "发布单位",
113        "备注",
114    ]
115    with open(path, "a") as f:
116        csv_writer = csv.writer(f)
117        csv_writer.writerow(info)
118    # with open(path, "a", newline="", encoding="utf-8-sig") as f:
119    #     csv_writer = csv.writer(f)
120    #     csv_writer.writerow(info)
121    for link in link_list:
122        try:
123            headers = {"User-Agent": random.choice(self.headers)}
124            r = requests.get(link, headers=headers, verify=False)
125            r.encoding = r.apparent_encoding
126            html = r.text
127            selector = etree.HTML(html)
128            print(selector.xpath("//td/h1/text()"))
129            number =

```

```

        selector.xpath("//td/h1/text()")[0].strip().replace("标准号: ",
        "")
130 name = selector.xpath("//td/b/text()")[0].strip()
131 E_name = re.findall("<td>英文标准名称: (.+)</td>", html)[
132     0
133 ] # 这里使用正则表达式, 因为xpath莫名其妙使用不了
134 status =
        selector.xpath('//td[@align="left"]/span[@class]/text()')[
135     0
136 ].strip()
137 ccs = selector.xpath(
138     '//div[@class="col-xs-12 col-md-4 content"]/text()'
139 )[
140     0
141 ].strip() # 中国标准分类号 (CCS)
142 ics = selector.xpath(
143     '//div[@class="col-xs-12 col-md-4 content"]/text()'
144 )[
145     1
146 ].strip() # 国际标准分类号 (ICS)
147 release_date = selector.xpath(
148     '//div[@class="col-xs-12 col-md-4 content"]/text()'
149 )[
150     2
151 ].strip() # 发布日期
152 implement_date = selector.xpath(
153     '//div[@class="col-xs-12 col-md-4 content"]/text()'
154 )[
155     3
156 ].strip() # 实施日期
157 department = selector.xpath(
158     '//div[@class="col-xs-12 col-md-4 content"]/text()'
159 )[
160     4
161 ].strip() # 主管部门
162 unit = selector.xpath(
163     '//div[@class="col-xs-12 col-md-4 content"]/text()'
164 )[
165     5
166 ].strip() # 归口单位
167 f_unit = selector.xpath(
168     '//div[@class="col-xs-12 col-md-10 content"]/text()'
169 )[
170     0

```

```

171     ].strip() # 发布单位
172     other = selector.xpath(
173         '//div[@class="col-xs-12 col-md-10 content"]/text()'
174     )[
175         1
176     ].strip() # 备注
177
178     i += 1
179     info = []
180     info.append(i)
181     info.append(number)
182     info.append(name)
183     info.append(E_name)
184     info.append(status)
185     info.append(ccs)
186     info.append(ics)
187     info.append(release_date)
188     info.append(implement_date)
189     info.append(department)
190     info.append(unit)
191     info.append(f_unit)
192     info.append(other)
193
194     # # 用None替换不存在的项目
195     # for j in range(len(info)):
196     #     if len(info) == 0:
197     #         info = "None"
198     # 正确处理 info 列表为空的情况
199     for j in range(len(info)):
200         if info[j] == "":
201             info[j] = "None"
202
203     with open(path, "a") as f:
204         csv_writer = csv.writer(f)
205         csv_writer.writerow(info)
206     # with open(path, "a", newline="", encoding="utf-8-sig") as
207     #     f:
208     #         csv_writer = csv.writer(f)
209     #         csv_writer.writerow(info)
210
211 except Exception as e:
212     print("Downloa Error:", link, e)
213

```

```

214 def main():
215     print("开始下载有关 '安全' 的国家标准")
216     scrapy = Spider("安全")
217
218     print("正在获取链接...")
219     link_list = scrapy.get_each_link()
220     print("获取完毕共%d个" % len(link_list))
221     print("正在下载...")
222
223     # link_list =
224         ['http://www.gb688.cn/bzgk/gb/newGbInfo?hcno=E6CD13ACF8E3B3BEFAA16852499676D0']
225     scrapy.get_detial_from_linkList(link_list)
226
227     print("有关 '安全' 的国家标准下载完毕")
228
229 if __name__ == "__main__":
230     main()

```

1.1.2 截图

- (base) nanmener@Haotians-MacBook-Pro 10. Spider % /Applications/Anaconda3/anaconda3/bin/python "/Users/nanmener/Github/zstu-study/Python程序设计-高阶/Python programming@ZSTU by Laitude/10. Spider/1.py"
 开始下载有关 '安全' 的国家标准
 https://openstd.samr.gov.cn/bzgk/gb/std_list?p.p1=0&p.p2=安全&p.p90=circulation_date&p.p91=desc
 正在获取链接...
 ["showInfo('B3460D59D0BF065AADE72E239964574F');", "showInfo('61473F823BB6A9F954DEDF020D85111');", "showInfo('17D551A108AA6EDB56633768FCD9310F');", "showInfo('37D0509D06E0F6243DE88A54351426DE');", "showInfo('C4AB705DA0063542548212CD76070D2D');", "showInfo('1271E1730AF7C01110E54AE5F77A18D0');", "showInfo('CCC8109C92E1B03B92A18F8B87B14C5D');", "showInfo('31ECB92BAE10BBA109E2FA227D68B81D');", "showInfo('559E267F161240CE953889C802A5518');", "showInfo('AF72BB94A5536336467DB7D49D5210B3');"]
 https://openstd.samr.gov.cn/bzgk/gb/newGbInfo?hcno=B3460D59D0BF065AADE72E239964574F
 https://openstd.samr.gov.cn/bzgk/gb/newGbInfo?hcno=61473F823BB6A9F954DEDF020D85111
 https://openstd.samr.gov.cn/bzgk/gb/newGbInfo?hcno=17D551A108AA6EDB56633768FCD9310F
 https://openstd.samr.gov.cn/bzgk/gb/newGbInfo?hcno=37D0509D06E0F6243DE88A54351426DE
 https://openstd.samr.gov.cn/bzgk/gb/newGbInfo?hcno=C4AB705DA0063542548212CD76070D2D
 https://openstd.samr.gov.cn/bzgk/gb/newGbInfo?hcno=1271E1730AF7C01110E54AE5F77A18D0
 https://openstd.samr.gov.cn/bzgk/gb/newGbInfo?hcno=CCC8109C92E1B03B92A18F8B87B14C5D
 https://openstd.samr.gov.cn/bzgk/gb/newGbInfo?hcno=31ECB92BAE10BBA109E2FA227D68B81D
 https://openstd.samr.gov.cn/bzgk/gb/newGbInfo?hcno=559E267F161240CE953889C802A5518
 https://openstd.samr.gov.cn/bzgk/gb/newGbInfo?hcno=AF72BB94A5536336467DB7D49D5210B3
 获取完毕共10个
 正在下载...
 ['标准号: GB 16798-2023 ']
 ['标准号: GB 16994.4-2023 ']
 ['标准号: GB 29753-2023 ']
 ['标准号: GB 43069-2023 ']
 ['标准号: GB 43068-2023 ']
 ['标准号: GB 43066-2023 ']
 ['标准号: GB 43067-2023 ']
 ['标准号: GB 7957-2023 ']
 ['标准号: GB 43203-2023 ']
 ['标准号: GB 43049-2023 ']
 有关 '安全' 的国家标准下载完毕

```
(base) nanmener@Haotians-MacBook-Pro 10. Spider % cat 安全.csv
序号,标准编号,标准名称,英文名称,状态,中国标准分类号 (CCS),国际标准分类号 (ICS),发布日期,实施日期,主管部门,归口单位,发
布单位,备注
1,GB 16798-2023,食品机械安全要求,Food machinery safety requirements,即将实施,X90,67.260,2023-09-08,2024-04-01,工业和信
息化部,工业和信息化部,国家市场监督管理总局、国家标准化管理委员会,None
2,GB 16994.4-2023,港口作业安全要求 第4部分: 普通货物集装箱,Safety requirements for port operation-Part 4: General carg
o container,即将实施,R43,03.220.40,2023-09-08,2024-04-01,交通运输部,交通运输部,国家市场监督管理总局、国家标准化管理委员
会,None
3,GB 29753-2023,道路运输 易腐食品与生物制品 冷藏车安全要求及试验方法,Road transportation—Perishable foodstuffs and bio
logical products—Safety requirement and test methods of refrigerated vehicle,即将实施,T54,43.160,2023-09-08,2024-01-01
,工业和信息化部,工业和信息化部,国家市场监督管理总局、国家标准化管理委员会,None
4,GB 43069-2023,矿用电缆安全技术要求,Safety requirements for mine cable,即将实施,K13,29.060.20,2023-09-08,2025-04-01,
国家矿山安全监察局,国家矿山安全监察局,国家市场监督管理总局、国家标准化管理委员会,None
5,GB 43068-2023,煤矿用跑车防护装置安全技术要求,Safety technical requirements for coal mine-used car catcher,即将实施,D1
8,73.100.40,2023-09-08,2024-04-01,国家矿山安全监察局,国家矿山安全监察局,国家市场监督管理总局、国家标准化管理委员会,None
6,GB 43066-2023,煤矿用被动式隔爆设施安全技术要求,Safety technique requirements for passive flameproof facilities of c
oal mine,即将实施,D09,13.100,2023-09-08,2024-04-01,国家矿山安全监察局,国家矿山安全监察局,国家市场监督管理总局、国家标
准化管理委员会,None
7,GB 43067-2023,煤矿用仪器仪表安全技术要求,Technical requirements for safety of environmental monitoring instruments fo
r coalmines,即将实施,N21,13.100,2023-09-08,2024-04-01,国家矿山安全监察局,国家矿山安全监察局,国家市场监督管理总局、国家
标准化管理委员会,None
8,GB 7957-2023,煤矿用矿灯安全技术要求,Safety technical requirement for cap lamps in coal mine,即将实施,K35,29.260.20,20
23-09-08,2025-04-01,国家矿山安全监察局,国家矿山安全监察局,国家市场监督管理总局、国家标准化管理委员会,None
9,GB 43203-2023,选煤厂安全规程,Code of practice for coal preparation plant safety,即将实施,D16,73.120,2023-09-08,2024-0
7-01,国家矿山安全监察局,国家矿山安全监察局,国家市场监督管理总局、国家标准化管理委员会,None
10,GB 43049-2023,连铸机安全技术条件,Safety technical conditions for continuous casting machines,即将实施,H93,77.180,202
3-09-08,2024-10-01,工业和信息化部,工业和信息化部,国家市场监督管理总局、国家标准化管理委员会,None
```

§2 实验题 2

选择一个在线视频网站，使用 selenium 编写一个爬虫，进行视频下载。

§2.1 解答

2.1.1 代码

```
1 import requests
2 from selenium import webdriver
3 from selenium.webdriver.chrome.service import Service
4 from selenium.webdriver.common.by import By
5 from win32api import Sleep
6
7 chromedriver_path = r"C:\Program
    Files\Google\Chrome\Application\chromedriver.exe"
8 service = Service(chromedriver_path)
9 driver = webdriver.Chrome(service=service)
10
11 driver.get('https://v.jstv.com/xw360/')
12 Sleep(2000)
13 driver.find_element(By.XPATH,
    "//*[ @id=\"J_bodyBd\"]/div/div/ul/li[1]/div[1]/a/img").click()
14 Sleep(2000)
15 driver.find_element(By.XPATH,
    "//*[ @id=\"J_bodyBd\"]/div/div/ul/li[2]/div[1]/a/img").click()
16 Sleep(3000)
17 driver.find_element(By.XPATH,
    "//*[ @id=\"J_bodyBd\"]/div/div/ul/li[3]/div[1]/a/img").click()
18 Sleep(3000)
19 driver.find_element(By.XPATH,
    "//*[ @id=\"J_bodyBd\"]/div/div/ul/li[4]/div[1]/a/img").click()
```



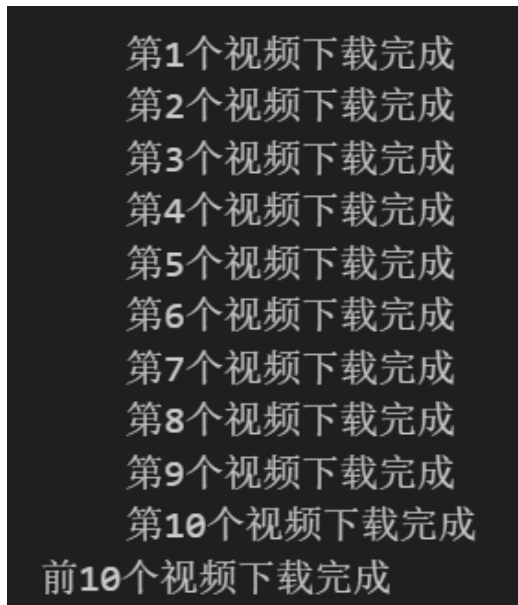
```

20 Sleep(3000)
21 driver.find_element(By.XPATH,
    "//*[@id=\"J_bodyBd\"]/div/div/ul/li[5]/div[1]/a/img").click()
22 Sleep(3000)
23 driver.find_element(By.XPATH,
    "//*[@id=\"J_bodyBd\"]/div/div/ul/li[6]/div[1]/a/img").click()
24 Sleep(3000)
25 driver.find_element(By.XPATH,
    "//*[@id=\"J_bodyBd\"]/div/div/ul/li[7]/div[1]/a/img").click()
26 Sleep(3000)
27 driver.find_element(By.XPATH,
    "//*[@id=\"J_bodyBd\"]/div/div/ul/li[8]/div[1]/a/img").click()
28 Sleep(3000)
29 driver.find_element(By.XPATH,
    "//*[@id=\"J_bodyBd\"]/div/div/ul/li[9]/div[1]/a/img").click()
30 Sleep(3000)
31 driver.find_element(By.XPATH,
    "//*[@id=\"J_bodyBd\"]/div/div/ul/li[10]/div[1]/a/img").click()
32
33 Sleep(3000)
34 windows = driver.window_handles # 获取当前所有页面句柄
35 for i in range(1, len(windows)):
36     driver.switch_to.window(windows[i]) # 切换当前新页面
37     Sleep(3000)
38
39     # print(driver.title)
40
41     driver.find_element(By.XPATH, "//*[@id=\"videoxg\"]/xg-start").click()
42     Sleep(3000)
43     video_element = driver.find_element(By.TAG_NAME, 'video')
44     video_url = video_element.get_attribute('src')
45     # 使用Requests库下载视频
46     response = requests.get(video_url)
47
48     # 指定视频文件的保存路径
49     save_path = str(driver.title)+'.mp4'
50
51     # 将视频内容写入文件
52     with open(save_path, 'wb') as file:
53         file.write(response.content)
54     # print(driver.title)
55     driver.close()
56     print(f"    第{i}个视频下载完成")
57     driver.switch_to.window(windows[0])

```

```
58     Sleep(3000)
59 driver.quit()
60 print("前10个视频下载完成")
```

2.1.2 截图



§3 总结

§3.1 分工情况

以下说明小组的分工情况，每个人完成了什么功能，完成的情况如何

陈昊天

负责报告的整合撰写；

负责实验题 1 的研究与解答；

陈佳伟

负责实验题 2 的研究与解答；

冯佳钧

解读实验题 1 的课件代码；

§3.2 解决方案评价

以下说明小组对解决方案的评价，是否完成了题目的要求，怎么证明完成了题目的要求，好的地方在哪里，哪里有不足。

实验题 1

完成了题目的要求，修改课件中的下载安全标准信息的爬虫，使之能够正常工作。

优点：忽略 SSL 证书验证，使 TLS 握手成功。

缺点：对于下载失败的情况没有做重试处理。

实验题 2

完成了题目的要求，能够使用 selenium 从一个在线视频网站下载视频。

优点：使用 selenium 模拟真人操作浏览器，并使用 Sleep() 函数模拟操作间隔，最终下载到视频文件，使网站服务方无法察觉爬虫的存在。

缺点：只下载了网站的前 10 个视频，没有对网页列出的所有视频进行下载。

§3.3 个人心得

收获与不足的个人总结

陈昊天

收获：学习了 TLS/SSL 相关知识，了解跳过证书验证的利弊

不足：在处理下载失败的情况时，缺乏相应的异常处理机制和重试策略。

陈佳伟

收获：实践 selenium 功能和代码实现，初步使用 chromedriver，熟悉了自动化测试工具的实际应用。

不足：未能全面覆盖网站上所有视频内容的下载。

冯佳钧

收获：学习浏览器 UA、request 请求相关知识，理解课件代码

不足：对网络爬虫工作原理的理解还不够深刻。

参考文献

- [1] Wikipedia contributors. Transport layer security — Wikipedia, the free encyclopedia, 2023. [Online; accessed 28-November-2023].