# 数据分析及实践实验二实验报告

PB18061443 江昊霖
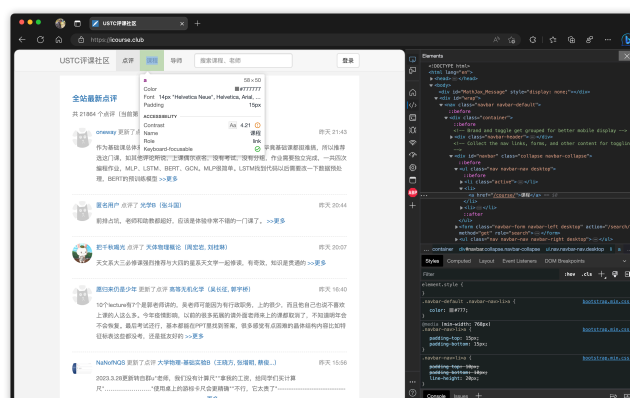
2023 年 3 月 29 日

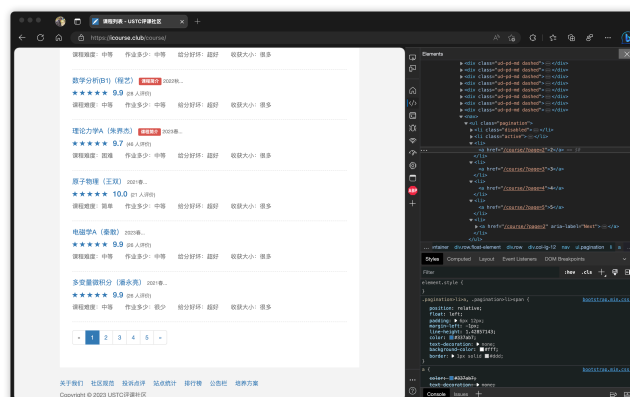## 1 实验目的

给定评课社区网站，需要设计一个网站遍历策略，爬取至少 200 个课程的详细信息，记录于 json 格式的文件中.
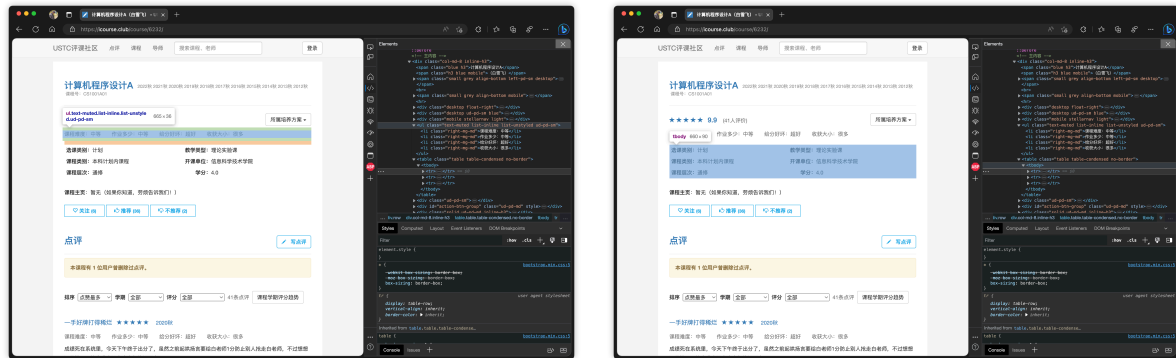
## 2 实验步骤

由给定的 url 进入课程列表网站爬取课程信息。



通过观察，课程列表的页面规律为：`https://icourse.club/course/?page={pagenumber}`，每页为 10 个课程。



递归进入每个课程页面后爬取信息.

# 3 代码

由于在进行请求时需要请求头，使用 `fake_useragent` 生成 user-agent.

## 3.1 获取课程列表网页链接

```python
def courseUrl(url: str) -> str:
    """get course list page suffix from given page

    Args:
        url (str): given url

    Returns:
        str: course list page
    """

    headers = {'User-Agent': UserAgent().random}
    ret = Request(url, headers=headers)
    res = urlopen(ret)
    data = res.read().decode('utf-8')
    soup = bf(data, features="lxml")

    suffix = soup.find('a', string=re.compile('.*?课程。*?'))['href']
    return url + suffix
```

## 3.2

```python
def getCoursesID(url: str) -> tuple[list, list]:
    """get courses' links and courses' names

    Args:
        url (str): courses list page

    Returns:
        tuple[list, list]: courses' links and names
    """

    headers = {'User-Agent': UserAgent().random}
```

```python
12    ret = Request(url, headers=headers)
13    res = urlopen(ret)
14    data = res.read().decode('utf-8')
15    soup = bf(data, features="lxml")
16
17    links = soup.find_all(class_="px16")
18    courses = []
19
20    for index, link in enumerate(links):
21        courses.append(link.text)
22        links[index] = link['href']
23
24    return links, courses
```

```python
1  def courseInfo(courseID: str, courseName: str, serial: int) -> list:
2      """get courses info
3
4      Args:
5          courseID (str): the courses id to course page
6          courseName (str): course's name and teacher
7          serial (int): serial of course
8
9      Returns:
10          info(list): info of the course
11      """
12
13      url = URL+courseID
14
15      headers = {'User-Agent': UserAgent().random}
16      ret = Request(url, headers=headers)
17      res = urlopen(ret)
18      data = res.read().decode('utf-8')
19      soup = bf(data, features="lxml")
20
21      info = []
22      info.append(str(serial))
23      info.append(courseName)
24
25      block1 = soup.find_all(class_="right-mg-md", limit=5)
26      del (block1[0])
27      for name in block1:
28          info.append(name.string.split(': ')[-1])
29
30      links = soup.find_all('strong', limit=6)
31
32      for link in links:
33          info.append(link.nextSibling)
34
35      return info
```