# STAT243 Lecture 2.4 Working with Information on the Web

> 🔥 **Logic** ⌄
>
> - 以编程方式与 Web 交互：以下载为主，也可上传
> - 常见格式：HTML, XML, JSON, YAML。
> - 核心工具：`requests` 、 `BeautifulSoup` / `lxml` 、 `json` 、 `yaml` 、 `pandas.read_html` / `read_json` 。

## 1 Reading HTML（抓取 HTML 表格）

- 关键思路：获取 HTML 源码 → 解析 DOM → 抽取 `<table>` 或节点文本。
- 浏览器中可先查看 **View Source / Developer Tools** 了解结构与是否存在大量 Javascript。

## 1.1 使用 `find_all` 通过 HTML 标签或属性搜索

```python
import io, requests, pandas as pd
from bs4 import BeautifulSoup as bs

URL = "https://en.wikipedia.org/wiki/List_of_countries_and_dependencies_by_population"
user_agent = "stat243_educational_bot/0.1 (paciorek@berkeley.edu)"
headers = {'User-Agent': user_agent}
response = requests.get(URL, headers=headers)
html = response.content

soup = bs(html, 'html.parser')
html_tables = soup.find_all('table')

# pandas.read_html 需要字符串/文件句柄，不直接接受 Tag 对象
pd_tables = [pd.read_html(io.StringIO(str(tbl)))[0] for tbl in html_tables]
[x.shape for x in pd_tables]
# 输出:
# [(242, 6), (13, 2), (1, 2)]

pd_tables[0].head()
# 输出:
#     Location  ...  Notes
# 0      World  ...    NaN
# 1      India  ...    [b]
# ...
```

- `BeautifulSoup` 构建树后可按标签或属性搜索，再交给 `pandas.read_html` 解析。

## 1.2 Extracting hyperlinks（提取超链接）

```python
import requests
from bs4 import BeautifulSoup as bs

URL = "http://www1.ncdc.noaa.gov/pub/data/ghcn/daily/by_year"
response = requests.get(URL)
soup = bs(response.content, 'html.parser')

# 方法 1: 所有 <a> 标签
```

```
 9   a_elements = soup.find_all('a')
10   links1 = [x.get('href') for x in a_elements]
11
12   # 方法 2: 具有 href 属性的 <a>
13   href_elements = soup.find_all('a', href=True)
14   links2 = [x.get('href') for x in href_elements]
15
16   links2[:9]
17   # 输出:
18   # ['?C=N;O=D', '?C=M;O=A', '?C=S;O=A', '?C=D;O=A', '/pub/data/ghcn/daily/',
19   #  '1750.csv.gz', '1763.csv.gz', '1764.csv.gz', '1765.csv.gz']
```

## 1.3 CSS selectors（用选择器抽取）

**Python**

```python
1   # 所有 <tr> 内部的 <th>
2   soup.select("tr th")
3   # 输出:
4   # [<th>...</th>, <th>...</th>, ...]
5
6   # 所有父元素为 <th> 的 <a>
7   soup.select("th > a")
8   # 输出:
9   # [<a href="?C=N;O=D">Name</a>, <a href="?C=M;O=A">Last modified</a>, ...]
```

## 1.4 XPath（用 lxml 进行 XPath 查询）

**Python**

```python
1    import lxml.html
2
3    # 将 BeautifulSoup object 转为 lxml object
4    lxml_doc = lxml.html.fromstring(str(soup))
5
6    # 所有带 href 的 <a>
7    a_elements = lxml_doc.xpath('//a[@href]')
8    links = [x.get('href') for x in a_elements]
9    links[:9]
10   # 输出:
11   # ['?C=N;O=D', '?C=M;O=A', '?C=S;O=A', '?C=D;O=A', '/pub/data/ghcn/daily/',
12   #  '1750.csv.gz', '1763.csv.gz', '1764.csv.gz', '1765.csv.gz']
```

## 2 XML, JSON, and YAML

- 三者均支持键值、数组与层级结构
- 读取需用对应库解析为 Python 结构（dict/list 等）。

## 2.1 XML（结构化自描述）

- XML 是一种以自描述格式储存数据的 markup language, 通常有 hierarchical structure, 不需要 metadata
- XML 文档具有树状结构, 由元素（节点）组成
- 常见存档/办公文档/空间信息（如 KML）。

**XML**

```xml
1   <?xml version="1.0"?>
2   <catalog>
3      <book id="bk101">
```

```xml
      <author>Gambardella, Matthew</author>
      <title>XML Developer's Guide</title>
      <genre>Computer</genre>
      <price>44.95</price>
      <publish_date>2000-10-01</publish_date>
      <description>An in-depth look at creating applications with XML.</description>
   </book>
   <book id="bk102">
      <author>Ralls, Kim</author>
      <title>Midnight Rain</title>
      <genre>Fantasy</genre>
      <price>5.95</price>
      <publish_date>2000-12-16</publish_date>
      <description>A former architect battles corporate zombies, an evil sorceress, and her
own childhood to become queen of the world.</description>
   </book>
</catalog>
```

**≣ Example ∨**

示例：Kiva 最新贷款数据（注意：在线接口可能返回 403，需本地保存演示）。我们采用两种方法:

1. 暴力解法 (将数据视作列表而不是树结构)
2. 使用 XPath 来遍历树结构

**Python**

```python
import xmltodict
# 假设 'newest.xml' 已手动下载（或用浏览器另存）
with open('newest.xml', 'r') as file:
    content = file.read()

# 一些 XML 含有裸 '&'，需先替换以避免解析错误
content = content.replace("&", "and")
data = xmltodict.parse(content)

data.keys()
# 输出:
# dict_keys(['response'])

data['response'].keys()
# 输出:
# dict_keys(['paging', 'loans'])

len(data['response']['loans']['loan'])
# 输出:
# 20

type(data['response']['loans']['loan'][2])
# 输出:
# dictionary

data['response']['loans']['loan'][2]['activity']
# 输出:
# 'Retail'
```

**Python**

```python
# 同样可用 lxml + XPath 提取指定字段
```

```
2    from lxml import etree
3
4    doc = etree.fromstring(content)
5    loans = doc.xpath("//loan")
6    [loan.xpath("activity/text()") for loan in loans][:3]
7    # 输出:
8    # [['Poultry'], ['Retail'], ['Retail']]
9
10   ## 假设只想要 country locations of the loans (using XPath)
11   [loan.xpath("location/country/text()") for loan in loans]
12   # 输出:
13   # [['Uganda'], ['Ecuador'], ['Ecuador'] ...]
14
15   ## or extract the geographic coordinates
16   [loan.xpath("location/geo/pairs/text()") for loan in loans]
17   # 输出:
18   # [['-0.352537 31.552699'], ['-1.054723 -80.45249'] ...]
```

## 2.2 JSON（层级键值，较 XML 简洁）

- JSON 文件以 **"attribute-value" pairs**（也称为"键-值"对）结构化，通常具有层次结构
- 可以使用 `json` 包将 JSON 读取到 Python 中
- JSON 的最外层可以是 **对象（object，用花括号 {}）** 或 **数组（array，用方括号 []）**。如果最外层是数组，那么就没有键（key）去命名它的元素。

**{ }    JSON**

```json
1    {
2        "firstName": "John",
3        "lastName": "Smith",
4        "isAlive": true,
5        "age": 25,
6        "address": {
7            "streetAddress": "21 2nd Street",
8            "city": "New York",
9            "state": "NY",
10           "postalCode": "10021-3100"
11       },
12       "phoneNumbers": [
13           { "type": "home", "number": "212 555-1234" },
14           { "type": "office", "number": "646 555-4567" }
15       ],
16       "children": [],
17       "spouse": null
18   }
```

**🐍    Python**

```python
1    import json
2    # 假设 'newest.json' 已手动下载
3    with open('newest.json', 'r') as file:
4        content = file.read()
5
6    data = json.loads(content)
7    list(data.keys())
8    # 输出:
9    # ['loans']
10
11   type(data['loans']), data['loans'][0]['location']['country']
12   # 输出:
```

```
13  # (<class 'list'>, 'Uganda')
14
15  [c['location']['country'] for c in data['loans']][:5]
16  # 输出:
17  # ['Uganda', 'Ecuador', 'Ecuador', 'Tajikistan', 'Mali']
```

> ⚠ **Remark** ⌄
>
> 注意：JSON 不原生支持缺失、无穷大等特殊值。

## 2.3 YAML（常用于配置）

- 以缩进表达层级，人类可读；缩进易出错；部分关键字（如 `on`）在某些实现中会被当作布尔。

**YAML**
```yaml
1   name: deploy-book
2
3   # Only run this when the master branch changes
4   on:
5     push:
6       branches:
7       - main
8
9   # This job installs dependencies, build the book, and pushes it to `gh-pages`
10  jobs:
11    deploy-book:
12      runs-on: ubuntu-latest
13      steps:
14      - uses: actions/checkout@v2
15
16      # Install dependencies
17      - name: Set up Python 3.9
18        uses: actions/setup-python@v1
19        with:
20          python-version: 3.9
21
22      - name: Install dependencies
23        run: |
24          pip install -r book-requirements.txt
```

**Python**
```python
1   import yaml
2   with open("book.yml") as stream:
3       config = yaml.safe_load(stream)
4   print(config)
5   # 输出:
6   # {'name': 'deploy-book', True: {'push': {'branches': ['main']}}, 'jobs': {'deploy-book':
    ...}}
7
8   print(config.get('name'))
9   # 输出:
10  # deploy-book
11
12  len(config['jobs']['deploy-book']['steps'])
13  # 输出:
14  # 3
```

> ⚠️ **Remark** ⌄
>
> 注意 `on` 会被视作 boolean value

## 3 Web APIs and webscraping

- 目标：获取 Web 数据时，优先使用正式 API；当 API 不可用时，再考虑 webscraping，并遵守网站条款与伦理。
- 常用库：`requests`，`json`，`pandas`，`BeautifulSoup`，`lxml`，`yaml`。
- 核心操作路径：理解 HTTP → REST 风格 API → 参数拼接与分页 → 处理压缩与归档 → POST 与认证 → 第三方封装 → 动态页面。

## 3.1 What is HTTP?

- **请求方法**：常用 GET、POST、PUT、DELETE，实际数据提取以 GET 为主。
- **状态码**：200 成功，4xx 客户端错误（如 403/404），5xx 服务器错误。
- **URL 查询字符串**：`?` 之后为参数，`&` 分隔键值对，空格常编码为 `+` 或 `%20`。

  > ☰ **Example** ⌄
  >
  > 1. `www.somewebsite.com?param1=arg1&param2=arg2`
  > 2. `https://www.yelp.com/search?find_desc=plumbers&find_loc=Berkeley+CA&ns=1`

- **响应内容**: 通常包含文本形式的内容（例如，HTML、XML、JSON）或原始字节

```Python
import requests

url = "https://httpbin.org/status/200"
r = requests.get(url)
print(r.status_code)
# 输出:
# 200
```

## 3.2 APIs: REST-based web services

> 🔥 **Logic** ⌄
>
> - 理想情况下，一个网络服务会用提供其 API (Application Programming Interface) 文档，该接口用于提供数据或允许其他交互
> - REST 是一种流行的 API 标准/风格

- **资源导向**：以 URL (也叫 endpoint) 作为资源，通过 query string 过滤，通过 `GET` 实现 request, 常返回 JSON。
- **分页与每页条数**：注意 `page`，`per_page`，`limit`，`offset` 等字段。
- **两种构造 request 的方式**：
  1. 直接拼接查询字符串。
  2. 使用 `params` 以 dict 传参。

```Python
import json, requests

# 方法一：直接拼接查询字符串（World Bank 示例）
url = "https://api.worldbank.org/V2/country?incomeLevel=MIC&format=json"
resp = requests.get(url)
data = json.loads(resp.content)
```

```
7
8    ## 注意 data truncation/pagination
9    if False:
10       url = "https://api.worldbank.org/V2/country?incomeLevel=MIC&format=json&per_page=1000"
11       response = requests.get(url)
12       data = json.loads(response.content)
13
14   # 方法二: Programmatic control
15   baseURL = "https://api.worldbank.org/V2/country"
16   group = 'MIC'
17   format = 'json'
18   args = {'incomeLevel': group, 'format': format, 'per_page': 1000}
19   url = baseURL + '?' +
20       '&'.join(['='.join([key, str(args[key])])
21               for key in args])
22   response = requests.get(url)
23   data = json.loads(response.content)
24
25   # 方法三: params 传参（更稳妥）
26   baseURL = "https://api.worldbank.org/V2/country"
27   params = {"incomeLevel": "MIC", "format": "json", "per_page": 1000}
28   resp = requests.get(baseURL, params=params)
29   data = resp.json()
30
31   print(type(data), len(data))
32   # 输出:
33   # <class 'list'> 2
34
35   print(len(data[1]), isinstance(data[1][5], dict), data[1][5]['name'])
36   # 输出:
37   # 104 True Benin
```

**Python**

```
1    # 简单的分页遍历（若 API 需要翻页）
2    all_rows = []
3    page = 1
4    while True:
5        params = {"incomeLevel": "MIC", "format": "json", "per_page": 100, "page": page}
6        resp = requests.get(baseURL, params=params)
7        data = resp.json()
8        rows = data[1]
9        if not rows:
10           break
11       all_rows.extend(rows)
12       page += 1
13
14   print(len(all_rows) >= 100)
15   # 输出:
16   # True
```

## 3.3 HTTP requests by deconstructing an (undocumented) API

- **思路**：用浏览器 DevTools 的 Network 面板观察实际请求（URL、headers、查询参数、cookies），在代码中复现。
- **典型场景**：下载链接返回 zip 压缩，需内存解压后再读取 CSV。

**Python**

```
1    import io, zipfile, requests, pandas as pd
2
```

```python
3    itemCode = 526
4    baseURL = "https://data.un.org/Handlers/DownloadHandler.ashx"
5    yrs = ",".join(str(yr) for yr in range(2012, 2018))
6    filter_ = f"?DataFilter=itemCode:{itemCode};year:{yrs}"
7    args1 = "&DataMartId=FAO&Format=csv&c=2,3,4,5,6,7&"
8    args2 = "s=countryName:asc,elementCode:asc,year:desc"
9    url = baseURL + filter_ + args1 + args2
10
11   resp = requests.get(url)
12   # 把 zip 文件放在内存中，而不是保存成 .zip 文件
13   with io.BytesIO(resp.content) as stream:
14       # 在内存中打开这个 zip 文件
15       with zipfile.ZipFile(stream, "r") as archive:
16           # 从压缩包中读取第一个文件，并用 pandas 加载
17           name = archive.filelist[0].filename
18           with archive.open(name, "r") as f:
19               dat = pd.read_csv(f)
20
21   print(dat.head(2))
22   # 输出:
23   #   Country or Area  Element Code  ...   Value  Value Footnotes
24   # 0      Afghanistan           432  ...  202.19              NaN
25   # 1      Afghanistan           432  ...   27.45              NaN
```

## 3.4 Webscraping ethics and best practices

- **是否应该抓**：优先使用公开下载文件或正式 API；抓取是下策。
- **是否允许抓**：遵守网站条款与 `robots.txt`；尊重速率限制与版权、隐私。
- **实践建议**：请求加 `User-Agent`，缓存响应避免重复请求；对高频请求使用 `time.sleep`；谨慎处理认证信息。

```python
1    import time, requests
2
3    headers = {"User-Agent": "stat243_educational_bot/0.1 (contact@example.com)"}
4    for page in range(1, 4):
5        r = requests.get("https://httpbin.org/get", params={"page": page}, headers=headers,
     timeout=10)
6        print(r.status_code, r.json()["args"]["page"])
7        # 输出:
8        # 200 1
9        # 200 2
10       # 200 3
11       time.sleep(1)  # 友好限速
```

## 3.5 More details on HTTP requests

- **结构化参数**：`params` 适合 GET，避免手写拼接错误。
- **复杂下载**：大文件/二进制内容可用 `stream=True` 分块下载。
- **错误处理**：`response.raise_for_status()`，或根据 `status_code` 分支处理。

```python
1    # 分块下载（示意）
2    import requests
3
4    url = "https://speed.hetzner.de/100MB.bin"
5    with requests.get(url, stream=True) as r:
6        r.raise_for_status()
```

```
 7        total = 0
 8        for chunk in r.iter_content(chunk_size=8192):
 9            if chunk:
10                total += len(chunk)
11    print(total > 0)
12    # 输出:
13    # True
```

## 3.6 POST example（创建 GitHub issue）

- **说明**：需个人 access token；最小权限原则；切勿把 token 提交到仓库。
- **两种方式**：裸 `requests.post` 与封装库 `PyGitHub` 。

**Python**

```python
import requests

with open(".github-access-token.txt", "r") as file:
    ghtoken = file.read().strip()

owner, repo = "paciorek", "test"
url = f"https://api.github.com/repos/{owner}/{repo}/issues"
issue = {
    "title": "This is an example issue",
    "body": "This is the body of the issue created via API."
}
headers = {
    "Authorization": f"token {ghtoken}",
    "Accept": "application/vnd.github+json"
}
resp = requests.post(url, json=issue, headers=headers)
print(resp.status_code in (200, 201))
# 输出:
# True
```

**Python**

```python
# 使用 PyGitHub 封装
from github import Github

with open(".github-access-token.txt", "r") as file:
    ghtoken = file.read().strip()

g = Github(ghtoken)
repo = g.get_repo("paciorek/test")
issue = repo.create_issue(
    title="Test Issue Created Programmatically",
    body="This is an issue filed programmatically using PyGitHub."
)
print(f"#{issue.number}", issue.html_url)
# 输出:
# #18 https://github.com/paciorek/test/issues/18
g.close()
```

## 3.7 Accessing dynamic pages

- **适用场景**：内容需 Javascript 渲染或需要模拟用户交互。
- **方案**：`selenium` 驱动浏览器；或 `scrapy` 框架配合 `splash` 渲染。

- **注意**：渲染成本高、速率慢，更应重视限速、重试与缓存；若站点有公开接口，优先 API。