Where are data from?

# 网络数据获取入门

Zhang Li
Nanjing University

# 网络数据爬取

# 用Python获取网络数据



API获取数据

## 网络数据如何获取（爬取）？

### 抓取网页，解析网页内容

- 抓取
  - Requests第三方库
  - Scrapy框架
- 解析
  - Beautiful Soup库
  - re模块

# html页面

```
<html>

<head>
<title>我是页面标题</title>
</head>

<body>
<p>我是新段落</p>
数据，我就是数据
</body>

</html>
```
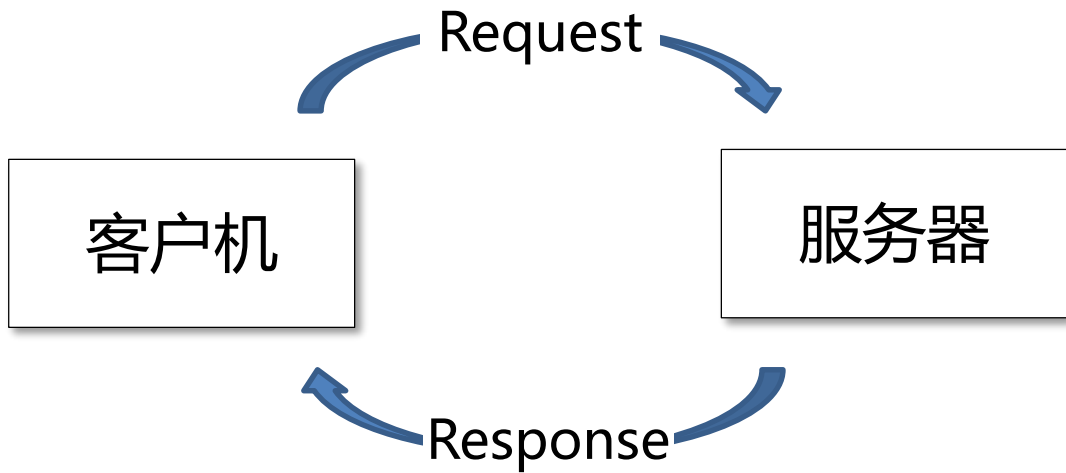
浏览器显示：
我是页面标题

我是新段落

数据，我就是数据

参考：
http://www.w3school.com.cn/html/html_backgrounds.asp

# 网页抓取

# Requests库

- Requests库是更简单、方便和人性化的Python HTTP第三方库

- Requests官网：
[http://www.python-requests.org/](http://www.python-requests.org/)

  $ pip install requests

  （Anaconda中预装）

豆瓣读书
　《小王子》短评

# Requests库

| requests.get() | 请求获取指定URL位置的资源，对应HTTP协议的GET方法，返回一个Response对象 |

**S**ource

```
>>> import requests
>>> r = requests.get('https://www.nju.edu.cn')
>>> r.status_code
200
>>> print(r.text)
>>> r = requests.get('https://book.douban.com/subject/1084336/comments/')
>>> r.status_code
418
```

# Requests

**S**ource

```
>>> headers = {'User-Agent': 'Mozilla/5.0 (Macintosh: Intel Mac
OS X 10_14_0) AppleWebKit/537.36 (KHTML, like Gecko)
Chrome/71.0.3578.98 Safari/537.36'}
>>> r = requests.get(url, headers = headers)  #定制请求头
>>> payload = {'key1': 'value1', 'key2': 'value2'}
>>> r = requests.get('http://httpbin.org/get', params = payload)
>>> r = requests.post(url, data = payload)
```

# **Requests库**

\>\>\> r.encoding      *# 根据HTTP头部自动推测*

'UTF-8'

\>\>\> r.encoding = 'gb2312'

\>\>\> r.encoding = r.apparent_encoding

\>\>\> r.content      *# 以字节方式访问Response对象*

\>\>\> r.json()

# Requests库

r.content　　　　# 以字节方式访问Response对象

**F**ile

```
r = requests.get('http://...sample.jpg')
with open('pic.jpg', 'wb') as f:
    f.write(r.content)
```

# JSON格式

- **JSON格式**
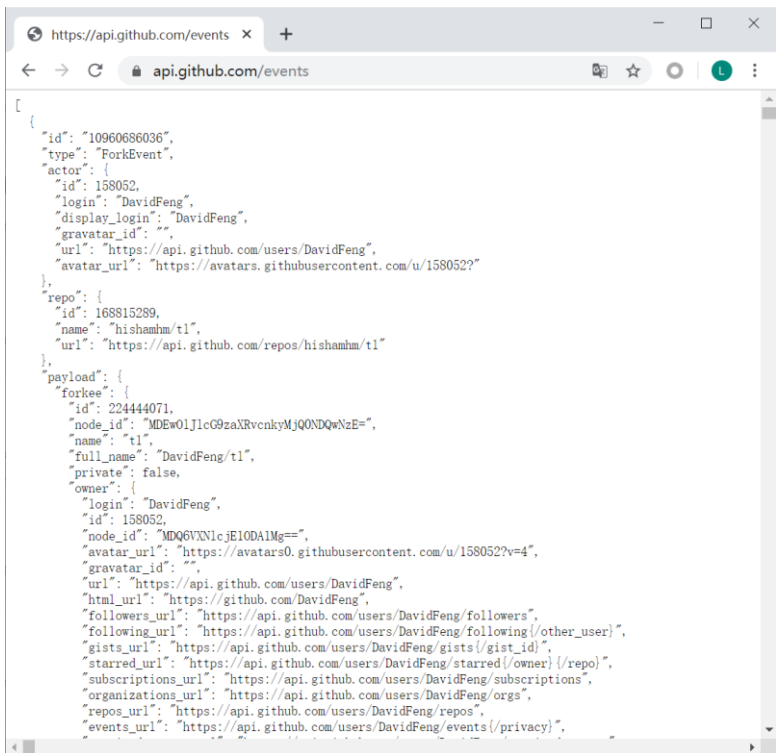  - JavaScript Object Notation，JS对象标记）
  - 一种轻量级的数据交换格式

r.json()

'{"name":"Niuyun", "address":{"city": "Beijing", "street": "Chaoyang Road"}}'

解析后
```
>>> x = {"name":"Niuyun",
         "address":{"city":"Beijing","street":"Chaoyang Road"}
        }
>>> x['address']['street']
'Chaoyang Road'
```

\>>> r = requests.get('https://api.github.com/events')

\>>> data = r.json()

\>>> data[0]['id']

# Robots协议

- Robots协议也称为爬虫协议，全称为爬虫排除协议（The Robots Exclusion Protocol）

- 检查站点根目录下是否存在robots.txt

# Robots协议-豆瓣网

User-agent: *
Disallow: /subject_search
Disallow: /amazon_search
Disallow: /search
…
Disallow: /link2/
Disallow: /recommend/
Disallow: /doubanapp/card
Disallow: /update/topic/
Allow: /ads.txt
Sitemap: https://www.douban.com/sitemap_index.xml
Sitemap: https://www.douban.com/sitemap_updated_index.xml
# Crawl-delay: 5

User-agent: Wandoujia Spider
Disallow: /
…

# 网页数据解析

# 网页数据解析

- **Beautiful Soup**是一个可以从HTML或XML文件中提取数据的Python库

- 官方网站：

https://www.crummy.com/software/BeautifulSoup/bs4/doc/

# Beautiful Soup

- $ pip install beautifulsoup4 （Anaconda中预装）

\>>> import requests

\>>> from bs4 import BeautifulSoup

\>>> r = requests.get('https://book.douban.com/subject/1084336/comments/', headers = headers)

\>>> soup = BeautifulSoup(r.text, 'lxml')

lxml: HTML解析器
$ pip install lxml

Python内置的HTML解析器
BeautifulSoup(markup, 'html.parser')

# Beautiful Soup

- BeautifulSoup对象
  - Tag
  - NavigableString
  - BeautifulSoup
  - Comment

```
>>> markup = '<p class="title"><b>The Little Prince</b></p>'
>>> soup = BeautifulSoup(markup, 'lxml')
>>> soup.b
<b>The Little Prince</b>
```

标签内容访问方式
BeautifulSoup对象.Tag

```
>>> markup = '<p class="title"><b>The Little Prince</b></p>'
>>> soup = BeautifulSoup(markup, 'lxml')
>>> tag = soup.p
>>> tag
<p class="title"><b>The Little Prince</b></p>
>>> tag.name
'p'
>>> tag['class']
['title']
>>> tag.attrs
{'class': ['title']}
>>> tag.string
'The Little Prince'
>>> soup.find_all('b')
[<b>The Little Prince</b>]
```

# 网页数据解析

```
<span class="short">不知道第几次重读。每过一
段时间再读，都有新的收获。心变得很柔软，脑里的迷
雾被驱散。更多的关注他人，关心这个世界，自私是多
么无趣的事情啊。我想，写一本能温暖人心，帮助困难
的人们的书，比世界上很多事情都有意义。</span>
```

```python
pattern = soup.find_all('span', {'class':'short'})
for item in pattern:
    print(item.string)
```

# 网页数据解析

r = requests.get('https://book.douban.com/subject/1084336/comments/', headers=headers)

```
soup = BeautifulSoup(r.text, 'lxml')

pattern = soup.find_all('span', {'class':'short'})

for item in pattern:

    print(item.string)
```

# 网页数据解析



豆瓣读书
　《小王子》推荐星级

# 网页数据解析

- 正则表达式是对字符串（包括普通字符和特殊字符）操作的一种逻辑公式

- **re**正则表达式模块进行各类正则表达式处理

- 参考网站：
  https://docs.python.org/3.5/library/re.html

| 元字符 | 描述 |
|---|---|
| . | 匹配除换行符外的任意字符 |
| * | 重复前面的子表达式0次或多次 |
| + | 重复前面的子表达式1次或更多次 |
| ? | 重复前面的子表达式0次或1次 |
| ^ | 匹配字符串的开始 |
| $ | 匹配字符串的结束 |
| {n} | 重复n次 |
| {n,} | 重复n次或更多次 |
| {n,m} | 重复n到m次 |
| \b | 匹配单词的开始或结尾即单词边界，"\B"匹配非单词边界 |
| \d | 匹配数字，"\D"匹配任意非数字字符 |
| \s | 匹配任意空白符，"\S"匹配任意非空白符 |
| \w | 匹配任意字母、数字或下划线的标识符字符，"\W"匹配任意非标识符字符 |
| [a-z] | 匹配指定范围内的任意字符 |
| [^a-z] | 匹配任何不在指定范围内的任意字符 |

# 网页数据解析



```
<span class="user-stars
allstar50 rating" title="力荐
"></span>
```

'<span class="user-stars allstar(.*?) rating"'

pattern = re.compile('<span class="user-stars allstar(.*?) rating"')
p = re.findall(pattern, r.text)

# 网页数据解析

```
r = requests.get('https://book.douban.com/subject/1084336/comments', headers = headers)
soup = BeautifulSoup(r.text, 'lxml')
pattern = soup.find_all('span', {'class':'short'})
for item in pattern:
    print(item.string)

pattern_s = re.compile('<span class="user-stars allstar(.*?) rating"')
p = re.findall(pattern_s, r.text)
```
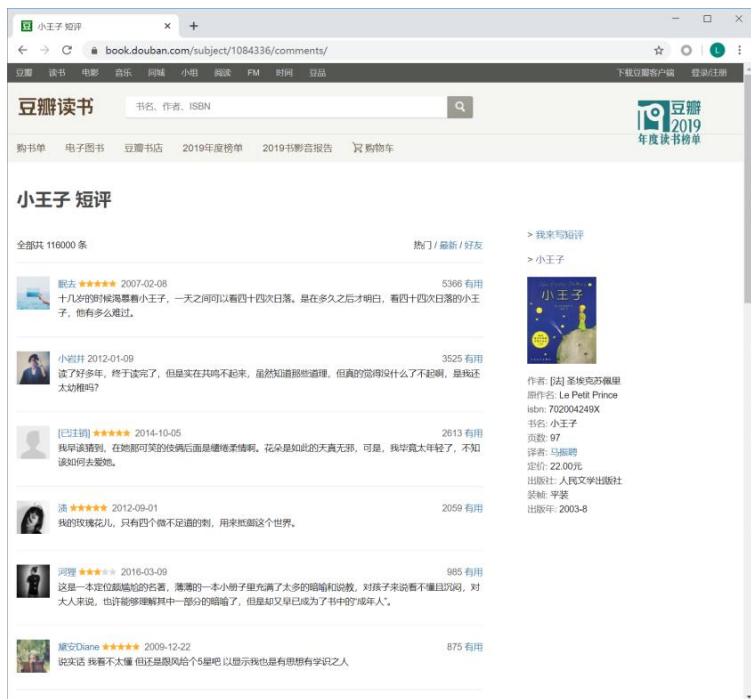
# 思考：抓取图书短评前5页



**https://book.douban.com/subject/1084336/**

r =
requests.get('https://book.douban.com/subject/1084336/comments/hot?p=' + str(i+1))

# 抓取例1

http://money.cnn.com/data/dow30/
抓取道指成分股数据并将30家公司的代码、公司名称和最近一次成交价放到一个列表中输出