# Web Scraping using BeautifulSoup

**Objectives:**

- Using `requests` to download server-side rendered HTML code
- Using `BeautifulSoup` to parse HTML code

## Scrape for Chinese Novels

We will demonstrate following website to demonstrate how to download and save articles.

- [http://www.shuku.net:8082/novels/jinyong/yi/yi.html](http://www.shuku.net:8082/novels/jinyong/yi/yi.html)
  [(http://www.shuku.net:8082/novels/jinyong/yi/yi.html)](http://www.shuku.net:8082/novels/jinyong/yi/yi.html)

For example, first article is in following website.

- [http://www.shuku.net:8082/novels/jinyong/yi/yitian01.html](http://www.shuku.net:8082/novels/jinyong/yi/yitian01.html)
  [(http://www.shuku.net:8082/novels/jinyong/yi/yitian01.html)](http://www.shuku.net:8082/novels/jinyong/yi/yitian01.html)

In [1]:
```python
!pip install beautifulsoup4
```

```
Requirement already satisfied: beautifulsoup4 in c:\users\isszq\anaconda3\l
ib\site-packages (4.8.2)
Requirement already satisfied: soupsieve>=1.2 in c:\users\isszq\anaconda3\l
ib\site-packages (from beautifulsoup4) (1.9.5)
```

## Extract List of Titles

Import libraries.

In [4]:
```python
from bs4 import BeautifulSoup
import requests
```

Use `requests` to send GET request to server and download HTML.

- Use status code to make sure request is successful.

In [14]:
```python
base_url = 'http://www.shuku.net:8082/novels/jinyong/yi/'

resp = requests.get(base_url + 'yi.html')
print(resp.status_code)
```

```
200
```

Make a soup from HTML code, which is in `resp.text`.

```
In [10]:    1  soup = BeautifulSoup(resp.text)
```

Examine the article items in the HTML code. Each title are contained in a `<td>` element with attribute `align='center'`.

```
In [12]:    1  items = soup.find_all('td', {'align':'center'})
```

Extract title and link for one item.

```
In [23]:    1  item = items[0]
            2  print(item.a)
            3  print(item.a.text)
            4  link = item.find('a').attrs['href']
            5  print(link)
```

```
<a href="yitian01.html">天涯思君不可忘</a>
天涯思君不可忘
yitian01.html
```

Use for-loop to extract title and link pairs into a list.

```
In [25]:    1  result = []
            2  for item in items:
            3      if item.find('a'):
            4          title = item.a.text
            5          link = item.a.attrs['href']
            6          result.append((title, link))
```

Examine the final `result` list.

```
In [1]:    1  print(len(result))
           2  result[:5]
```

```
---------------------------------------------------------------------------
NameError                                 Traceback (most recent call last)
<ipython-input-1-930a380903a4> in <module>
----> 1 print(len(result))
      2 result[:5]

NameError: name 'result' is not defined
```

## Extract One Article

Try to extract first article.

- Use relative link to form the complete link to the article
- Use `requests` to download HTML and make a soup.

In [47]: ▶|

```
1  title, link = result[0]
2
3  url = base_url + link
4  html = requests.get(url)
5  soup = BeautifulSoup(html.text)
```

The article is contained in `<pre>` tag.

Extract text from the tag and examine it.

- There are some extra text in the article. This is because of some redundant tags within the `<pre>` tag, e.g. `<script>` and `<a>` tags.

In [48]: ▶|

```
1  item = soup.find('pre')
2  print(item.text)
```

第一回　天涯思君不可忘

　　"春游浩荡，是年年寒食，梨花时节。白锦无纹香烂漫，玉树琼苞堆雪。静夜沉沉，浮光霭霭，冷浸溶溶月。人间天上，烂银霞照通彻。浑似姑射真人，天姿灵秀，意气殊高洁。万蕊参差谁信道，不与群芳同列。浩气清英，仙才卓荦，下土难分别。瑶台归去，洞天方看清绝。"

　　作这一首《无俗念》词的，乃南宋末年一位武学名家，有道之士。此人姓丘，名处机，道号长春子，名列全真七子之一，是全真教中出类拔萃的人物。《词品》评论此词道："长春，世之所谓仙人也，而词之清拔如此"。这首词诵的似是梨花，其实词中真意却是赞誉一位身穿白衣的美貌少女，说她"浑似姑射真人，天姿灵秀，意气殊高洁"，又说她"浩气清英，仙才卓荦"，"不与群芳同列"。词中所颂这美女，乃古墓派传人小龙女。她一生爱穿白衣，当真如风拂玉树，雪裹琼苞，兼之生性清冷，实当得起"冷浸溶溶月"的形容，以"无俗念"三字赠之，可说十分贴切。长春子丘处机和她在终南山上比邻而居，当年一见，便写下这首词来。

　　这时丘处机逝世已久，小龙女也已嫁与神雕大侠杨过为妻。在河南少室山山道

Remove those redundant tags from existing soup.

In [49]: ▶|

```
1  for s in soup.select('script'):
2      s.extract()
3
4  for s in soup.select('a'):
5      s.extract()
```

Extract the text again. The text is what we need to save to a text file.

```python
In [50]:
1  item = soup.find('pre')
2  print(item.text)
```

第一回　　天涯思君不可忘

"春游浩荡，是年年寒食，梨花时节。白锦无纹香烂漫，玉树琼苞堆雪。静夜沉沉，浮光霭霭，冷浸溶溶月。人间天上，烂银霞照通彻。浑似姑射真人，天姿灵秀，意气殊高洁。万蕊参差谁信道，不与群芳同列。浩气清英，仙才卓荦，下土难分别。瑶台归去，洞天方看清绝。"

作这一首《无俗念》词的，乃南宋末年一位武学名家，有道之士。此人姓丘，名处机，道号长春子，名列全真七子之一，是全真教中出类拔萃的人物。《词品》评论此词道："长春，世之所谓仙人也，而词之清拔如此"。这首词诵的似是梨花，其实词中真意却是赞誉一位身穿白衣的美貌少女，说她"浑似姑射真人，天姿灵秀，意气殊高洁"，又说她"浩气清英，仙才卓荦"，"不与群芳同列"。词中所颂这美女，乃古墓派传人小龙女。她一生爱穿白衣，当真如风拂玉树，雪裹琼苞，兼之生性清冷，实当得起"冷浸溶溶月"的形容，以"无俗念"三字赠之，可说十分贴切。长春子丘处机和她在终南山上比邻而居，当年一见，便写下这首词来。

这时丘处机逝世已久，小龙女也已嫁与神雕大侠杨过为妻。在河南少室山山道

Save text to file using `<title>.text` as file name.

```python
In [51]:
1  with open(title+'.txt', 'w', encoding='utf-8') as f:
2      f.write(item.text)
```

## Complete Solution

In [ ]: ▶|

```python
 1  import requests
 2  from bs4 import BeautifulSoup
 3
 4  base_url = 'http://www.shuku.net:8082/novels/jinyong/yi'
 5
 6  # Step 1
 7  def get_title_list():
 8      resp = requests.get(base_url + '/yi.html')
 9      html = resp.text
10      soup = BeautifulSoup(html)
11      items = soup.find_all('td', {'align':'center'})
12      result = []
13
14      for item in items:
15          t = item.find('a')
16          title = t.text
17          link = t.attrs['href']
18          result.append((title, link))
19      return result
20
21  result = get_title_list()
22  print(result)
23
24  # Step 2
25
26  for item in result:
27      title, link = item
28      url = base_url + '/' + link
29      resp = requests.get(url)
30      html = resp.text
31      soup = BeautifulSoup(html)
32      item = soup.find('pre')
33      article = item.text
34      print('writing... ', title)
35      with open(title+'.txt', 'w', encoding='utf-8') as f:
36          f.write(article)
```