

تمرین‌های استخراج اطلاعات از وب

انجام تمرین‌های زیر

تمرین: هر کتابخانه‌ای که فکر می‌کنید برای استخراج اطلاعات از یک وبسایت نیاز دارید، وارد کنید.

In [28]:

```
1 # CODE HERE
```

In [1]:

```
1 import requests, bs4
```

تمرین: از کتابخانه requests و BeautifulSoup برای اتصال به <http://quotes.toscrape.com/> استفاده کنید و متن HTML صفحه اصلی را دریافت کنید.

In [30]:

```
1 # CODE HERE
```

In [2]:

```
1 res = requests.get("http://quotes.toscrape.com/")
```

In [3]:

```
1 res.text
```

```
<title>Quotes to Scrape</title>\n      <link rel="stylesheet" href="/static/bootstrap.min.css">\n      <link rel="stylesheet" href="/static/main.css">\n</head>\n<body>\n      <div class="container">\n        <div class="row header-box">\n          <div class="col-md-8">\n            <h1>\n              <a href="/" style="text-decoration: none">Quotes to Scrape</a>\n            </h1>\n          </div>\n          <div class="col-md-4">\n            <p>\n              <a href="/login">Login</a>\n            </p>\n          </div>\n        </div>\n        <div class="row">\n          <div class="col-md-8">\n            <div class="quote" itemscope itemtype="http://schema.org/CreativeWork">\n              <span class="text" itemprop="text">“The world as we have created it is a process of our thinking. It cannot be changed without changing our thinking.”</span>\n              <span>by <small class="author" itemprop="author">Albert Einstein</small>\n              <a href="/author/Albert-Einstein">(about)</a>\n            </span>\n          </div>\n          <div class="tags">\n            Tags:\n            <meta class="keywords" itemprop="keywords" content="change,deep-thoughts,thinking,world" />\n            <a class="tag" href="/tag/change/page/1/">change</a>\n            <a class="tag" href="/tag/deep-thoughts/page/1/">deep-thoughts</a>\n            <a class="tag" href="/tag/thinking/page/1/">thinking</a>\n          </div>\n        </div>\n      </div>\n    </body>\n  </html>
```

تمرین: نام نویسندگان موجود در صفحه اول را بدست آورید.

In [4]:

```
1 # CODE HERE
2 soup = bs4.BeautifulSoup(res.text, 'lxml')
```

In [5]:

```
1 soup.select('.author')
```

Out[5]:

```
[<small class="author" itemprop="author">Albert Einstein</small>,
 <small class="author" itemprop="author">J.K. Rowling</small>,
 <small class="author" itemprop="author">Albert Einstein</small>,
 <small class="author" itemprop="author">Jane Austen</small>,
 <small class="author" itemprop="author">Marilyn Monroe</small>,
 <small class="author" itemprop="author">Albert Einstein</small>,
 <small class="author" itemprop="author">André Gide</small>,
 <small class="author" itemprop="author">Thomas A. Edison</small>,
 <small class="author" itemprop="author">Eleanor Roosevelt</small>,
 <small class="author" itemprop="author">Steve Martin</small>]
```

In [6]:

```
1 authors = set()
2
3 for name in soup.select(".author"):
4     authors.add(name.text)
```

In [7]:

```
1 authors
```

Out[7]:

```
{'Albert Einstein',
 'André Gide',
 'Eleanor Roosevelt',
 'J.K. Rowling',
 'Jane Austen',
 'Marilyn Monroe',
 'Steve Martin',
 'Thomas A. Edison'}
```

تمرین: یک لیست از تمامی نقل قول ها در صفحه اول ایجاد کنید.

In [13]:

```
1 #CODE HERE
```

In [8]:

```
1 quotes = []
2
3 for quote in soup.select(".text"):
4     quotes.append(quote.text)
```

In [9]:

```
1 quotes
```

Out[9]:

```
['"The world as we have created it is a process of our thinking. It cannot
be changed without changing our thinking."',
 '"It is our choices, Harry, that show what we truly are, far more than ou
r abilities."',
 '"There are only two ways to live your life. One is as though nothing is
a miracle. The other is as though everything is a miracle."',
 '"The person, be it gentleman or lady, who has not pleasure in a good nov
el, must be intolerably stupid."',
 '"Imperfection is beauty, madness is genius and it's better to be absolut
ely ridiculous than absolutely boring."',
 '"Try not to become a man of success. Rather become a man of value."',
 '"It is better to be hated for what you are than to be loved for what you
are not."',
 '"I have not failed. I've just found 10,000 ways that won't work."',
 '"A woman is like a tea bag; you never know how strong it is until it's i
n hot water."',
 '"A day without sunshine is like, you know, night."']
```

تمرین: سایت را بررسی کنید و از Beautiful Soup برای استخراج ده برچسب برتر از متن درخواست نمایش داده شده در بالا سمت راست صفحه اصلی استفاده کنید (به عنوان مثال Love, Inspirational, Life و غیره). راهنمایی: به خاطر داشته باشید که تحت هر نقل قول همچنین برچسبها وجود دارند، سعی کنید یک کلاس را پیدا کنید که فقط در برچسبهای بالا سمت راست وجود دارد، شاید به span چک کنید.

In [16]:

```
1 # CODE HERE
```

In [10]:

```
1 soup = bs4.BeautifulSoup(res.text, 'lxml')
```

In [11]:

```
1 soup.select('.tag-item')
```

Out[11]:

```
[<span class="tag-item">
  <a class="tag" href="/tag/love/" style="font-size: 28px">love</a>
</span>,
<span class="tag-item">
  <a class="tag" href="/tag/inspirational/" style="font-size: 26px">insp
irational</a>
</span>,
<span class="tag-item">
  <a class="tag" href="/tag/life/" style="font-size: 26px">life</a>
</span>,
<span class="tag-item">
  <a class="tag" href="/tag/humor/" style="font-size: 24px">humor</a>
</span>,
<span class="tag-item">
  <a class="tag" href="/tag/books/" style="font-size: 22px">books</a>
</span>,
<span class="tag-item">
  <a class="tag" href="/tag/reading/" style="font-size: 14px">reading</a>
</span>]
```

In [12]:

```
1 for item in soup.select(".tag-item"):
2     print(item.text)
```

love

inspirational

life

humor

books

reading

friendship

friends

truth

simile

****تمرین:** توجه کنید که بیش از یک صفحه وجود دارد و صفحات بعدی به این صورت به نظر می‌رسند: <http://quotes.toscrape.com/page/2/> (<http://quotes.toscrape.com/page/2/>). با استفاده از آنچه که درباره حلقه for و اتصال رشته می‌دانید، از تمام صفحات عبور کرده و تمام نویسندگان منحصر به فرد در وبسایت را بدست آورید. به خاطر داشته باشید که راه‌های زیادی برای دستیابی به این هدف وجود دارد و همچنین توجه کنید که باید به طریقی بفهمید که حلقه شما در صفحه آخر نقل‌قول‌ها قرار دارد. برای اهداف اشکال

زدایی، به شما اطلاع می‌دهم که تنها 10 صفحه وجود دارد، بنابراین صفحه آخر به این صورت است: <http://quotes.toscrape.com/page/10/> (<http://quotes.toscrape.com/page/10/>). اما سعی کنید یک حلقه‌ی قوی بسازید که به حدی قوی باشد که لازم نباشد تعداد صفحات را به طور پیشین بدانید، شاید از try/except برای این کار استفاده کنید، این به شما بستگی دارد! **

In [22]:

```
1 # CODE HERE
```

روش‌های دیگر زیادی وجود دارند که حتی قوی‌تر و انعطاف‌پذیرتر هستند، اما ایده اصلی همان است، استفاده از یک حلقه while برای گشتن در صفحات مختلف و شرط خروجی بر اساس صفحه نامعتبر است.

In [13]:

```
1 url = 'http://quotes.toscrape.com/page/'
```

راه حل اول

In [14]:

```
1 authors = set()
2
3 for page in range(1,10):
4
5     # Concatenate to get new page URL
6     page_url = url+str(page)
7     # Obtain Request
8     res = requests.get(page_url)
9     # Turn into Soup
10    soup = bs4.BeautifulSoup(res.text, 'lxml')
11    # Add Authors to our set
12    for name in soup.select(".author"):
13        authors.add(name.text)
```

راه حل دیگر

In []:

```
1 page_url = url+str(9999999)
2
3 # Obtain Request
4 res = requests.get(page_url)
5
6 # Turn into Soup
7 soup = bs4.BeautifulSoup(res.text, 'lxml')
```

In []:

```
1 soup
```

In []:

```
1 "No quotes found!" in res.text
```

In []:

```
1 page_still_valid = True
2 authors = set()
3 page = 1
4
5 while page_still_valid:
6
7     # Concatenate to get new page URL
8     page_url = url+str(page)
9
10    # Obtain Request
11    res = requests.get(page_url)
12
13    # Check to see if we're on the last page
14    if "No quotes found!" in res.text:
15        break
16
17    # Turn into Soup
18    soup = bs4.BeautifulSoup(res.text, 'lxml')
19
20    # Add Authors to our set
21    for name in soup.select(".author"):
22        authors.add(name.text)
23
24    # Go to Next Page
25    page += 1
```

In []:

```
1 authors
```

