## 第3章 Python语言入门与Web信息解析

- 3.1 Python语言入门
- 3.2 Python语言进阶
- 3.3 HTTP解析与Python实现
- 3.4 HTML解析与Python实现
- 3.5 Web信息解析

# 3.4 HTML解析与Python实现

- 3.4.1 HTML与CSS要素
- 3.4.2 BeautifulSoup库及对象
- 3.4.3 BeautifulSoup库遍历文档树
- 3.4.4 BeautifulSoup库搜索文档树
- 3.4.5 BeautifulSoup库查找CSS过滤器
- 3.4.6 Python解析网页实例

- Web网页主要由三部分组成:
  - 结构(Structure)
    - 主要包括XHTML和XML,
  - 表现(Presentation)
    - 主要包括CSS
  - 行为(Behavior)。
    - 主要包括对象模型(如W3C DOM)、CMAScript

# 3.4.1 HTML与CS\$

1.文档标记

```
<html>
<head>
  <title>Python爬虫开发与项
  <meta charset="UTF-8">
</head>
<body>
文档设置标记<br>
>这是段落。
这是段落。
这是段落。
<hr>>
<center>居中标记1</center>
<center>居中标记2</center>
<hr>>
```

```
<hr>>
           >
           00:00](music)
           [00:28]你我皆凡人,生在人世间;
           00:35]终日奔波苦,一刻不得闲;
<u1>
Coffee
Milk
type="A">
Coffee
Milk
<d1>
  <dt>计算机</dt>
  <dd>用来计算的仪器 ... ...</dd>
  <dt>显示器</dt>
  <dd>以视觉方式显示信息的装置 ... ...</dd>
  </dl>
  <div >
        <h3>这是标题</h3>
        这是段落。
</div>
```

<

00:00](music)

[00:28]你我皆凡人,生在人世间; [00:35]终日奔波苦,一刻不得闲; [00:43]既然不是仙,难免有杂念;

### 2.图像标记<img>

<img src="路径/文件名.图片格式" width="属性值"
height="属性值" border="属性值" alt="属性值" >

- alt属性有三个作用:
  - 1) 当网页上的图片被加载完成后,鼠标移动到上面去,会显示这个图片指定的属性文字;
  - 2) 如果图像没有下载或者加载失败,会用文字来代替图像显示;
  - 3)搜索引擎可以通过这个属性的文字来抓取图片

- 3.超链接标记
- 爬虫开发中经常需要抽取链接,链接的引用使用的是 <a>标记。
- <a>标记的基本语法:

<a href="链接地址" target="打开方式" name="页面锚 点名称" 《链接文字或者图片</a>

- name属性用来指定页面的锚点名称。
  - \_blank时:在一个新的窗口中打开链接;
  - \_self (默认值): 在当前窗口中打开链接
  - ...

4.表格标记

基本结构包括、<caption>、、、和等

```
学号
 班级
 姓名
 年龄
 籍贯
1500001
 (1)班
 16
 上海
```

- CSS层叠样式表(Cascading Style Sheets)
  - 用来定义如何显示HTML元素,一般和HTML配合使用

```
class:列出一系列以空格分
   隔的CSS类名
      <header class="sohu-head">
         <div class="area sohu-head-box">
100
             <div class="right head-right"</pre>
101
                                           id:将页面唯一标识符
102
                                             附加到某个标签
             </div>
103
         </div>
104
      </header><div class="sohu-ph" id="sohuTopc" style="display:none;">
105
      <div class="ph-link">
106
         <a href="http://news.sohu.com/s2018/guoqing_lindex_shtml" target=" hlank"></a>
107
      </div>
108
                                                      使用id可以快速获取我
      <div id="ph-close" class="ph-close"><a href="javascrip"</pre>
109
                                                      们感兴趣的HTML页面
   </div> <!-- 皮肤 wrapper box -->
                                                           的某些部分
      <div class="theme-skin-wrap" data-spm="top-festival">
111
         <div class="mask"></div>
```

- CSS格式
  - 样式信息被写为冒号分隔的基于键值的语句列表
  - 每个语句本身用分号隔开

color: 'red';

background-color: #ccc;

font-size: 14pt;

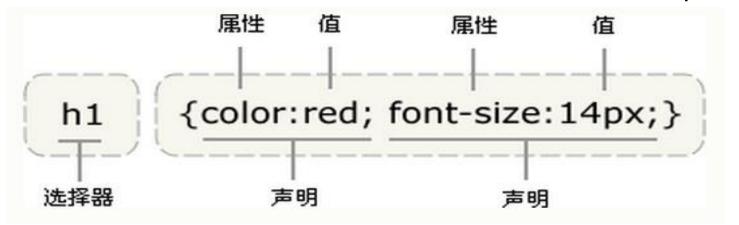
border: 2px solid yellow;

- CSS样式声明有三种方式:
  - 内联样式表:
- CSS代码直接写在现有的HTML标记中,使用style属性改变样式 <body style="background-color:green; margin:0; padding:0;"></body>
  - 嵌入式样式表: 样式应用在哪个标签很清晰
    - CSS样式代码写在<style type="text/css"></style>标记之间
    - 嵌入式CSS样式写在<head></head>之间。
  - 外部样式表:
    - CSS代码写一个单独的外部文件中,以".css"为扩展名
    - 在<head>内(不是在<style>标记内)使用link>标记将CSS 样式文件链接到HTML文件内
    - link rel="StyleSheet" type="text/css" href="style.css"> .

### ■ CSS规则由两个主要的部分构成:

selector {declaration1; declaration2; ... declarationN }

- 选择器:
  - 需要改变样式的HTML元素
  - 三种选择器:标签选择器、类选择器、ID选择器
- 一条或多条声明:
  - 用于定义显示方式,每条声明由一个属性和一个值组成
  - 属性(property)是希望设置的样式属性(style attribute)。
     每个属性有一个值。属性和值由冒号分开
  - 例如:将h1标记中的颜色设置为红色,字体大小为14px



可以任意组合

```
<!DOCTYPE html>
Khtml lang="en" dir="ltr">
  (head)
                                      元素。
    <meta charset="utf-8">
    <title>实例</title>
    <style type='text/css';</pre>
     p, h3 {
       color:red;设置文本颜色
                 为红色
    </style>
  </head>
               标签选择器的使用
  <body>
                               啦啦啦
    <h3>啦啦啦</h3>
                               你好
    你好
    <img src="image/八重樱2.jpg" width:
    〈p〉我的朋友〈/p〉
                               我的朋友
  </body>
 ^{\prime}htm1>
```

样式(color:red;)将应用到 p,h3这两个选择器所引用的

- 标签选择器
  - 选择具有特定标签名称 的所有元素
  - 如、<h3>
  - 代码:

p{font-size:12px;lineheight:1.6em;}

```
<html lang="en">
(head)
   <meta charset="UTF-8">
   <title>实例</title>
   (style)
      .box {
         color: cyan;
                       盒子里有
         font-size: 40px;
                       小动物
   </style>
〈/headXdiv是父元素,span,p,是div的子
〈body〉 继承了父元素的文字颜色与字
                       小猫咪
   <div class="box" >
      〈span〉盒子里有〈/span〉〈
      <span>小动物</span>
      \小猫咪
                        小跳蛙
      \小跳蛙
      </div>
</body>
(/html>
```

类选择器

■ 选择HTML中定义的特定类 的所有元素,即class出现 的位置

■ 语法:

• 类名称{css样式代码;}

```
<!DOCTYPE html>
<html lang="en" dir="ltr">
 (head)
   <meta charset="utf-8">
   <title>实例</title>
   <style type='text/css'>
       #tale{
           color: aqua;
       #tiger{
           color: red;
   </style>
              id选择器的使用
 </head>
 (body)
   (p)
       〈span id="tale">一二三四五
   二三四五 上山打老虎
 </body>
```

|html>

#### id选择器

- 按照id属性值选择匹配的元素
- 任何的标签都可以设置id
- 正确的HTML文档应该确保每 个"id"唯一,并且只给与一个 元素
- 语法:

#intro {font-weight:bold;}

- Python可以使用相同的CSS选择器语法快速查找和检索 HTML页面中的元素
- 例:使用开发者工具(chrome)Elements,选中某元素,点 击右键选择copy→copy selector,可以获得CSS选择器

