

## Day09-正则表示式vs爬虫程序

今日课程学习目标

今日课程大纲

知识点1: 正则表达式功能简介【熟悉】

知识点2: 正则表达式语法-匹配单个字符【熟悉】

知识点3: 正则表达式-匹配多个字符【熟悉】

知识点4: 正则表达式-匹配开头结尾vs其他匹配【熟悉】

知识点5: re模块-match方法的使用【熟悉】

知识点6: re模块-search方法的使用【熟悉】

知识点7: re模块-findall方法的使用【熟悉】

知识点8: re模块-split方法的使用【熟悉】

知识点9: re模块-正则分组操作【常握】

知识点10: re模块-引用正则分组【熟悉】

知识点11: re模块-sub方法使用【熟悉】

知识点12: re模块-正则表示式修饰符【了解】

知识点13: 贪婪模式和非贪婪模式【常握】

知识点14: 浏览器的请求过程-渲染【熟悉】

知识点15: FastAPI-返回图片数据【掌握】

知识点16: FastAPI-提取URL路径数据【掌握】

知识点17: 爬虫的概念及工作过程【了解】

知识点18: requests 模块的功能使用【掌握】

## Day09-正则表示式vs爬虫程序

### 今日课程学习目标

- 1 熟悉正则表达式的基本语法
- 2 熟悉re正则模块中常用方法是使用
- 3 常握使用正则提取指定分组的数据
- 4 理解正则中的贪婪模式和非贪婪模式
- 5 熟悉 FastAPI 提取URL地址数据
- 6 熟悉爬虫的概念和作用

### 今日课程大纲

- 1 # 1. 正则表达式
- 2 正则表达式的作用
- 3 正则表达式的语法
- 4 匹配单个字符
- 5 匹配多个字符
- 6 正则分组
- 7 re正则模块的使用
- 8 贪婪模式vs非贪婪模式
- 9 # 2. 爬虫程序
- 10 FastAPI返回图片数据
- 11 FastAPI提取URL地址数据
- 12 爬虫的概念和作用
- 13 requests简单使用

知识点1：正则表达式功能简介【熟悉】

正则表达式(regular expression)描述了一种字符串匹配的模式，可以用来检查一个串是否含有某种子串、将匹配的子串做替换或者从某个串中取出符合某个条件的子串等。

模式：**一种特定的字符串模式，这个模式是通过一些特殊的符号组成的。**

正则表达式的功能：

- 数据验证（表单验证、如手机、邮箱、IP地址）
- 数据检索（数据检索、数据抓取）
- 数据隐藏（ 135\*\*\*\*6235 王先生）
- 数据过滤（论坛敏感关键词过滤）

正则表达式并不是Python所特有的，在Java、PHP、Go以及JavaScript等语言中都是支持正则表达式的。

知识点2：正则表式语法-匹配单个字符【熟悉】

正则表达式在线练习网址：<https://tool.oschina.net/regex/>

匹配单个字符：

正则语法	描述
.	匹配任意1个字符（除了\n）
[]	匹配[]中列举的字符
\d	匹配数字，即0-9
\D	匹配非数字，即不是数字
\s	匹配空白，即 空格，tab键
\S	匹配非空白
\w	匹配非特殊字符，即a-z、A-Z、0-9、_、汉字
\W	匹配特殊字符，即非字母、非数字、非汉字

练习：

- 在字符串 "abcd123" 中匹配 a： a
- 在字符串 "abcd123" 中匹配任意一个字符： .
- 在字符串 "abcd123" 中匹配 b 或 d： [bd]
- 在字符串 "abcd123" 中匹配数字： \d
- 在字符串 "abcd123" 中匹配非数字内容： \D
- 在字符串 "abcd 123" 中匹配空白字符串： \s
- 在字符串 "abcd 123" 中匹配非空白字符串： \S
- 在字符串 "abcd\_123" 中匹配非特殊字符： \w
- 在字符串 "abcd&%123" 中匹配特殊字符： \W

### 知识点3：正则表达式-匹配多个字符【熟悉】

匹配多个字符：

正则语法	描述
*	匹配前一个字符出现0次或者无限次，即可有可无
+	匹配前一个字符出现1次或者无限次，即至少有1次
?	匹配前一个字符出现1次或者0次，即要么有1次，要么没有
{m}	匹配前一个字符出现m次
{m,n}	匹配前一个字符出现从m到n次

练习：

- 在字符串 "传智播客\_python\_666" 中匹配非数字内容： `\D*` 返回： `传智播客_python_`
- 在字符串 "传智播客\_python\_666" 中匹配数字： `\d+` 返回： `666`
- 在字符串 "传智播客\_python\_666" 中匹配y或py： `[p]?y` 返回： `py`
- 在字符串 "传智播客\_python\_666" 中匹配2个数字： `\d{2}` 返回： `66`
- 在字符串 "传智播客\_python\_666" 中匹配英文字母出现1-3次： `[a-z]{1,3}` 返回： `pyt hon`

### 知识点4：正则表达式-匹配开头结尾vs其他匹配【熟悉】

匹配开头和结尾：

正则语法	描述
^	匹配字符串开头
\$	匹配字符串结尾

练习：

- 在字符串 "abc\_python\_666" 中匹配以a开头： `^a` 返回： `a`
- 在字符串 "abc\_python\_666" 中匹配以数字结尾： `\d$` 返回： `6`

其他匹配：

正则语法	描述
[^指定字符]	匹配除了指定字符以外的所有字符
	匹配左右任意一个表达式

练习：

- 在字符串 "abc\_python\_666" 中匹配除了数字以外的字符： `[^\d]+` 返回： `abc_python_`
- 在字符串 "abc-python-666" 中匹配数字和特殊字符： `\d+|\W+` 返回： `- - 666`

## 知识点5: re模块-match方法的使用【熟悉】

### 函数语法格式:

```
1 match函数: re.match(pattern, string, flags=0)
2 参数:
3     pattern: 匹配的正则表达式
4     string: 要匹配的字符串
5     flags: 标志位(暂时忽略, 后面会介绍)
6
7 功能: 尝试从字符串起始位置匹配一个正则表达式
8     1) 如果不能从起始位置匹配成功, 则返回None;
9     2) 如果能从起始位置匹配成功, 则返回一个匹配的对象
```

### 示例代码:

```
1 my_str = 'abc_123_DFG_456'
2
3 # 匹配字符串bc(注: 从头开始)
4 res = re.match(r'bc', my_str)
5 print(type(res), res) # 注意: 不能从起始位置匹配成功, 则返回None
6
7
8 # 匹配字符串abc(注: 从头开始)
9 res = re.match(r'abc', my_str)
10 print(type(res), res) # 注意: 能从起始位置匹配成功, 则返回一个匹配的对象, Match类实例对象
11
12 # Match对象.group(): 获取正则表达式匹配到的内容
13 print(res.group())
```

## 知识点6: re模块-search方法的使用【熟悉】

### 函数语法格式:

```
1 search函数: re.search(pattern, string, flags=0)
2 参数:
3     pattern: 匹配的正则表达式
4     string: 要匹配的字符串
5     flags: 标志位(暂时忽略, 后面会介绍)
6
7 功能: 根据正则表达式扫描整个字符串, 并返回第一个成功的匹配
8     1) 如果不能匹配成功, 则返回None;
9     2) 如果能匹配成功, 则返回一个匹配对象
```

### 示例代码:

```
1 my_str = 'abc_123_DFG_456'
2
3 # 匹配连续的3位数字
4 # \d: 匹配一位数字
5 res = re.search(r'\d{3}', my_str)
6 print(type(res), res) # 注意: 匹配成功, 则返回一个匹配的对象, Match类实例对象
7
8 # 获取正则表达式匹配到的内容
9 print(res.group())
```

## 知识点7: re模块-findall方法的使用【熟悉】

### 函数语法格式:

```
1  findall函数: re.findall(pattern, string, flags=0)
2  参数:
3      pattern: 匹配的正则表达式
4      string: 要匹配的字符串
5      flags: 标志位(暂时忽略, 后面会介绍)
6
7  功能: 根据正则表达式扫描整个字符串, 并返回所有能成功匹配的子串
8      1) 如果不能匹配成功, 则返回一个空列表;
9      2) 如果能匹配成功, 则返回包含所有匹配子串 of 列表
```

### 示例代码:

```
1  my_str = 'abc_123_DFG_456'
2
3  # 匹配字符串中的所有连续的3位数字
4  # raw string: 原生字符串
5  res = re.findall(r'\d{3}', my_str)
6  print(type(res), res)
```

## 知识点8: re模块-split方法的使用【熟悉】

### 函数语法格式:

```
1  split函数: re.split(pattern, string, maxsplit=0, flags=0)
2  参数:
3      pattern: 匹配的正则表达式
4      string: 要进行分割的字符串
5      maxsplit: 最多分割次数, 默认为0, 表示分割所有
6      flags: 标志位(暂时忽略, 后面会介绍)
7
8  功能: 根据正则表达式匹配的子串对原字符串进行分割, 返回分割后的列表
```

### 示例代码:

```
1  import re
2
3  my_str = '传智播客, Python, 数据分析'
4
5  # 需求: 按照 `, ` 对上面的字符串进行分割
6  res = re.split(r',\s', my_str)
7  print(type(res), res)
8
9  res = re.split(r',\s', my_str, maxsplit=1)
10 print(type(res), res)
11
12 res = my_str.split(', ')
13 print(type(res), res)
14
15 my_str2 = '传智播客,Python;数据分析'
16
17 res = re.split(r'[;,]', my_str2)
18 print(type(res), res)
```

## 知识点9: re模块-正则分组操作【常握】

在使用正则表达式进行匹配操作时, 可以将正则表达式进行分组, 并在匹配之后获取相应分组的匹配数据。

**正则匹配分组语法:**

正则语法	描述
(正则表达式)	将括号中字符作为一个分组
(?P<name>正则表达式)	分组起别名

### 示例1: 正则匹配分组操作

语法: (正则表达式)

```
1 import re
2
3 my_str = '13155667788'
4
5 # 需求: 使用正则提取出手机号的前3位、中间4位以及后 4 位数据
6
7 res = re.match(r'(\d{3})(\d{4})(\d{4})', my_str)
8 print(type(res), res)
9
10 # 获取整个正则表达式匹配的内容
11 print(res.group())
12
13 # 获取正则表达式指定分组匹配到的内容
14 # Match对象.group(组号)
15 print(res.group(1)) # 131
16 print(res.group(2)) # 5566
17 print(res.group(3)) # 7788
```

### 示例2: 给正则分组起别名

语法: (?P<分组别名>正则表达式)

```
1 my_str = '<div><a href="https://www.itcast.cn" target="_blank">传智播客</a><p>Python</p></div>'
2
3 # 需求: 使用正则提取出 my_str 字符串中的 `传智播客` 文本
4
5 res = re.search(r'<a.*?(?P<text>.*)</a>', my_str)
6 print(type(res), res)
7
8 # 获取整个正则表达式匹配到的内容
9 print(res.group())
10
11 # 获取指定分组匹配到的内容
12 print(res.group(1)) # 传智播客
13
14 # 根据分组的名称, 获取指定分组匹配到的内容
15 # Match对象.group(分组名称)
16 print(res.group('text')) # 传智播客
```

## 知识点10: re模块-引用正则分组【熟悉】

在正则表达式中, 放在圆括号 `()` 中表示的是一个组, 当使用 `()` 定义了一个正则表达式组后, 正则引擎会将匹配到的组按照顺序编号, 存入缓存, 这样的话我们想在后面已经匹配过的内容进行引用时, 引用分组的方式如下:

正则语法	描述
<code>\num</code>	引用正则中第 num 个分组匹配到的字符串 例如: <code>\1</code> 表示第一个分组, <code>\2</code> 表示第二个分组.....
<code>(?P=name)</code>	引用正则中别名为 name 分组匹配到的字符串

分组引用示例代码:

```
1  # 需求: 写一个正则表达式, 匹配字符串形如: 'xxx xxx xxx'
2  # 注意: xxx可以是任意多位的数字, 但是这三个xxx必须一致, 比如: '123 123 123', '6666 6666 6666'
3
4  import re
5
6  my_str = '123 123 123'
7
8  # 根据分组序号引用分组
9  res = re.match(r'(\d+)\s\1\s\1', my_str)
10 print(type(res), res)
11 # 获取整个正则表达式匹配到的内容
12 print(res.group())
13
14 print('=' * 20)
15
16 # 根据分组名称引用分组
17 res = re.match(r'(?P<num>\d+)\s(?P=num)\s(?P=num)', my_str)
18 print(type(res), res)
19 # 获取整个正则表达式匹配到的内容
20 print(res.group())
```

## 知识点11: re模块-sub方法使用【熟悉】

函数语法格式:

```
1  sub函数: re.sub(pattern, repl, string, count=0, flags=0)
2
3  参数:
4      pattern: 匹配的正则表达式
5      repl: 替换内容
6      string: 原字符串
7      count: 最多替换次数, 默认为0, 表示替换所有
8      flags: 标志位(暂时忽略, 后面会介绍)
9
10  功能: 根据正则表达式匹配字符串中的所有子串, 然后使用指定内容进行替换
11      1) 函数返回的是替换后的新字符串
```

示例代码1:

```
1  import re
```

```

2
3 my_str = "传智播客-Python-666"
4
5 # 需求: 将字符串中的 - 替换成 _
6 new_str = re.sub(r'-', r'_', my_str)
7 print(type(new_str), new_str)
8
9 new_str = re.sub(r'-', r'_', my_str, count=1)
10 print(type(new_str), new_str)
11
12 new_str = my_str.replace('-', '_')
13 print(type(new_str), new_str)
14
15 my_str = "传智播客,Python;666"
16
17 # 需求: 将字符串中的 ,和; 替换成 _
18 new_str = re.sub(r'[;,]', r'_', my_str)
19 print(type(new_str), new_str)

```

#### 示例代码2:

```

1 import re
2
3 # 需求: 将字符串 `abc.123` 替换为 `123.abc`
4 my_str = 'abc.123'
5
6 new_str = re.sub(r'([a-z]{3})\.(\\d{3})', r'\\2.\\1', my_str)
7 print(new_str)

```

### 知识点12: re模块-正则表示式修饰符【了解】

正则表达式可以包含一些可选标志修饰符来控制匹配的模式。

修饰符被指定为一个可选的标志。多个标志可以通过按位 OR( | ) 它们来指定。如 re.I | re.M 被设置成 I 和 M 标志:

修饰符	描述
re.I	匹配时不区分大小写
re.M	多行匹配, 影响 ^ 和 \$
re.S	使 . 匹配包括换行在内的所有字符

#### 示例代码:

```

1 import re
2
3 my_str = 'aB'
4
5 # re.I: 匹配时不区分大小写
6 res = re.match(r'ab', my_str, flags=re.I)
7 print(bool(res)) # 非None就True, None就是False
8
9
10 print('=' * 20)
11 # 开启多行模式
12 # ^ 可以匹配字符串开头 (字符串的开始位置), 也可以匹配行的开头 (即换行符\n之后的位置)

```



```

13 # $ 可以匹配字符串结尾（字符串的结束位置），也可以匹配行的结尾（即换行符\n之前的位置）
14
15 # 关闭多行模式
16 # ^ 只能匹配字符串开头
17 # $ 只能匹配字符串结尾
18 my_str = 'aabb\nbbcc'
19
20 res = re.findall(r'^[a-z]{4}$', my_str, flags=re.M)
21 print(res)
22 print(bool(res))
23 print('=' * 20)
24
25 my_str = '\nabc'
26 # re.S: 影响 . 符号，设置之后，.符号就能匹配\n了
27 res = re.match(r'.', my_str, flags=re.S)
28 print(bool(res))

```

### 知识点13：贪婪模式和非贪婪模式【常握】

**贪婪模式：**在整个表达式匹配成功的前提下，尽可能多的匹配

**非贪婪模式：**在整个表达式匹配成功的前提下，尽可能少的匹配

正则中的量词包括：{m,n}、?、\*和+，这些量词默认都是贪婪模式的匹配，可以在这些量词后面加?将其变为非贪婪模式。

例如：\d{2,5}

**示例代码：**

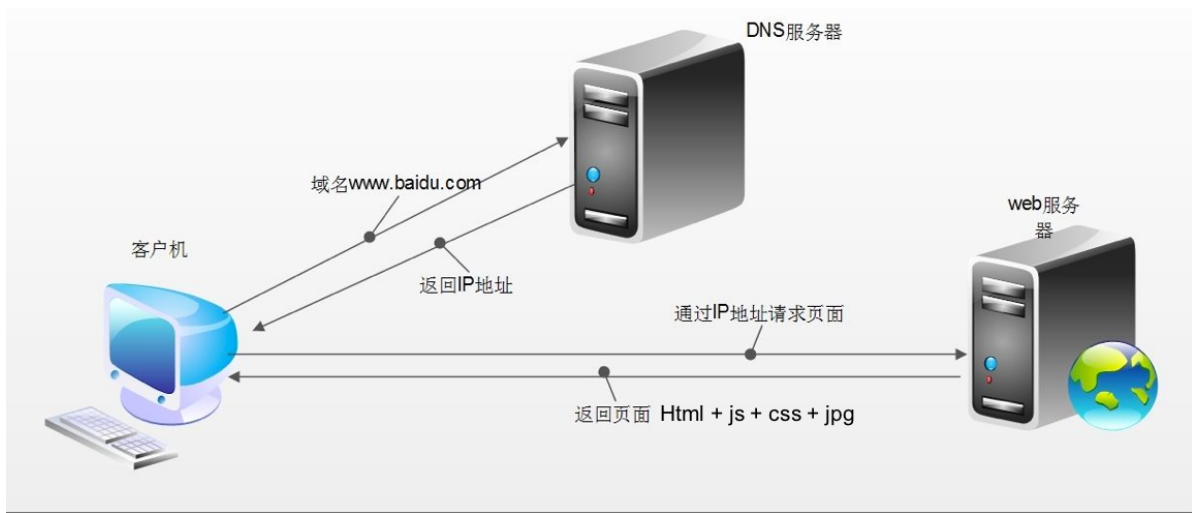
```

1 import re
2
3 my_str = '<div>test1</div><div>test2</div>'
4
5 # 贪婪模式：在整个表达式匹配成功的前提下，尽可能多的匹配
6 re_obj = re.match(r'<div>.*</div>', my_str)
7 print(re_obj)
8 print(re_obj.group()) # 获取整个正则表达式匹配的内容
9
10 # 非贪婪模式：在整个表达式匹配成功的前提下，尽可能少的匹配
11 re_obj = re.match(r'<div>.*?</div>', my_str)
12 print(re_obj)
13 print(re_obj.group()) # 获取整个正则表达式匹配的内容

```

### 知识点14：浏览器的请求过程-渲染【熟悉】

**HTTP单次请求过程：**



### 浏览器请求过程:

1. 浏览器通过域名解析服务器（DNS）获取IP地址
2. 浏览器先向IP发起请求，并获取响应
3. 在返回的响应内容（html）中，可能会带有css、js、图片等url地址，浏览器按照响应内容中的顺序依次发送其他的请求，并获取相应的响应
4. 浏览器每获取一个响应就对展示出的结果进行添加（加载），js、css等内容可能会修改页面的内容，js也可以重新发送请求，获取响应
5. 从获取第一个响应并在浏览器中展示，直到最终获取全部响应，并在展示的结果中添加内容或修改，这个过程叫做浏览器的渲染

### 知识点15：FastAPI-返回图片数据【掌握】

#### 示例代码：

```
1  # 导入 FastAPI 类
2  from fastapi import FastAPI
3  # 导入 uvicorn
4  import uvicorn
5  # 导入 Response 响应类
6  from fastapi import Response
7
8  # 创建 FastAPI 对象
9  app = FastAPI()
10
11
12  # 定义业务处理函数并设置对应的 URL 地址
13  # get: 表示请求方式
14  # /index.html: 表示请求的 URL 地址
15  @app.get('/index.html')
16  def index():
17      with open('./sources/html/index.html', 'r', encoding='utf8') as f:
18          content = f.read()
19
20      # 返回响应对象
21      return Response(content, media_type='html')
22
23
24  @app.get('/gdp.html')
25  def gdp():
26      with open('./sources/html/gdp.html', 'r', encoding='utf8') as f:
27          content = f.read()
28
```

```
29     # 返回响应对象
30     return Response(content, media_type='html')
31
32
33 @app.get('/images/0.jpg')
34 def get_images_0():
35     with open('./sources/images/0.jpg', 'rb') as f:
36         content = f.read()
37
38     # 返回响应对象
39     return Response(content, media_type='jpg')
40
41
42 @app.get('/images/1.jpg')
43 def get_images_1():
44     with open('./sources/images/1.jpg', 'rb') as f:
45         content = f.read()
46
47     # 返回响应对象
48     return Response(content, media_type='jpg')
49
50
51 @app.get('/images/2.jpg')
52 def get_images_2():
53     with open('./sources/images/2.jpg', 'rb') as f:
54         content = f.read()
55
56     # 返回响应对象
57     return Response(content, media_type='jpg')
58
59
60 @app.get('/images/3.jpg')
61 def get_images_3():
62     with open('./sources/images/3.jpg', 'rb') as f:
63         content = f.read()
64
65     # 返回响应对象
66     return Response(content, media_type='jpg')
67
68
69 @app.get('/images/4.jpg')
70 def get_images_4():
71     with open('./sources/images/4.jpg', 'rb') as f:
72         content = f.read()
73
74     # 返回响应对象
75     return Response(content, media_type='jpg')
76
77
78 @app.get('/images/5.jpg')
79 def get_images_5():
80     with open('./sources/images/5.jpg', 'rb') as f:
81         content = f.read()
82
83     # 返回响应对象
84     return Response(content, media_type='jpg')
85
86
```

```

87 @app.get('/images/6.jpg')
88 def get_images_6():
89     with open('./sources/images/6.jpg', 'rb') as f:
90         content = f.read()
91
92     # 返回响应对象
93     return Response(content, media_type='jpg')
94
95
96 if __name__ == '__main__':
97     # 启动 Web 服务器
98     uvicorn.run(app, host='127.0.0.1', port=8080)

```

## 知识点16: FastAPI-提取URL路径数据【掌握】

问题:

1 1. 对于上面的代码, 假如我们有1000张图片要返回, 该怎么操作?

通过观察, 我们可以发现返回图片的代码都是相似的, 只要能从 URL 路径中获得图片的名称, 上面的返回图片处理函数就能融合成一个, 在函数中只要根据图片名称读取对应的图片内容返回就可以了。

1 2. 如何从 URL 地址中提取数据?

FastAPI 可以从 URL 地址中提取数据, 并将提取的数据传递给对应的处理函数, 格式如下:

```

1 from fastapi import FastAPI
2
3 app = FastAPI()
4
5
6 # @app.get('/.../{参数名}')
7 @app.get("/items/{item_id}")
8 def read_item(item_id): # FastAPI 提取了数据之后, 会将提取的数据传递给下面处理函数的对应形参
9     pass

```

动态返回图片数据示例代码:

```

1 # 导入 FastAPI 类
2 from fastapi import FastAPI, Path
3 # 导入 uvicorn
4 import uvicorn
5 # 导入 Response 响应类
6 from fastapi import Response
7
8 # 创建 FastAPI 对象
9 app = FastAPI()
10
11
12 # 定义业务处理函数并设置对应的 URL 地址
13 # get: 表示请求方式
14 # /index.html: 表示请求的 URL 地址
15 @app.get('/index.html')
16 def index():
17     with open('./sources/html/index.html', 'r', encoding='utf8') as f:
18         content = f.read()
19
20     # 返回响应对象

```

```

21     return Response(content, media_type='html')
22
23
24 @app.get('/gdp.html')
25 def gdp():
26     with open('./sources/html/gdp.html', 'r', encoding='utf8') as f:
27         content = f.read()
28
29     # 返回响应对象
30     return Response(content, media_type='html')
31
32
33 # 从 URL 地址中提取图片名称
34 @app.get('/images/{image_name}')
35 def get_images(image_name):
36     print('image_name: ', image_name)
37
38     with open(f'./sources/images/{image_name}', 'rb') as f:
39         content = f.read()
40
41     # 返回响应对象
42     return Response(content, media_type='jpg')
43
44
45 if __name__ == '__main__':
46     # 启动 Web 服务器
47     uvicorn.run(app, host='127.0.0.1', port=8080)

```

## 知识点17: 爬虫的概念及工作过程【了解】

问题:

- 1 1. 什么是爬虫? 爬虫有什么作用?
- 2 2. 爬虫工作过程是什么?

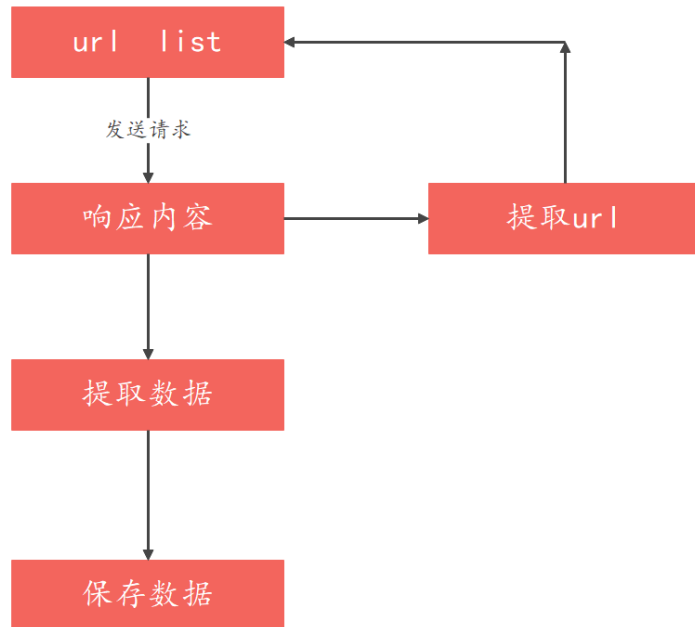
### 什么是爬虫? 爬虫有什么作用?

网络爬虫(又被称为网页蜘蛛, 网络机器人)就是模拟浏览器发送网络请求, 接收请求响应, 一种按照一定的规则, 自动地抓取互联网信息的程序。

- 原则上, 只要是浏览器(客户端)能做的事情, 爬虫都能做。
- 爬虫只能获取到浏览器(客户端)所展示出来的数据。

爬虫的作用: 数据分析中, 进行数据采集的一种方式。

### 爬虫工作过程是什么?



1. 向起始的 url 地址发送请求，并获取响应数据
2. 对响应内容进行提取
3. 如果提取 url，则继续发送请求获取响应
4. 如果提取数据，将数据进行保存

## 知识点18: requests 模块的功能使用【掌握】

问题：

1. requests模块是什么？有什么作用？
2. 如果使用 requests 模块？

### requests模块是什么？有什么作用？

requests 是用 python 语言编写的一个开源的HTTP库，可以通过 requests 库编写 python 代码发送网络请求，其简单易用，是编写爬虫程序时必知必会的一个模块。

中文文档：[https://docs.python-requests.org/zh\\_CN/latest/index.html](https://docs.python-requests.org/zh_CN/latest/index.html)

### 如果使用 requests 模块？

#### 1) 安装

```
1 pip install requests
2 或者
3 pip install requests -i https://pypi.tuna.tsinghua.edu.cn/simple
```

#### 2) 示例：请求百度，并获取响应内容

```
1 """
2 requests模块基本使用
3 学习目标：能够使用 requests 模块请求URL地址并获取响应内容
4 """
5
6 # 导入 requests 包
7 import requests
8
9 # 准备 url 地址
10 url = 'https://www.baidu.com'
```

```
11
12 # 使用 requests 发送 GET 请求
13 # 响应对象
14 response = requests.get(url)
15
16 # 获取响应的内容
17 # response.content: bytes, 服务器返回的原始响应内容
18 # bytes -> str: bytes数据.decode('解码方式')
19 print(response.content.decode())
```