

Day10-爬虫程序vs数据可视化

今日课程学习目标

今日课程大纲

知识点1: 爬虫示例-爬取单张图片数据【掌握】

知识点2: 爬虫示例-爬取多张图片数据【掌握】

知识点3: 爬虫示例-爬取GDP数据【掌握】

知识点4: 爬虫示例-多任务版爬虫【了解】

知识点5: pyecharts 模块简介和安装【了解】

知识点6: pyecharts 绘制饼图-GDP数据可视化【掌握】

知识点7: 程序日志简介及作用【掌握】

知识点8: 程序日志的5个等级说明【掌握】

知识点9: logging模块的基本使用【掌握】

知识点10: 数据埋点介绍【了解】

Day10-爬虫程序vs数据可视化

今日课程学习目标

- 1 掌握爬虫爬取图片数据和GDP数据
- 2 熟悉pyecharts模块绘制饼图
- 3 掌握日志的作用和级别
- 4 熟悉logging模块的基本使用

今日课程大纲

- 1 # 1. 爬虫程序
- 2 爬取单张图片数据
- 3 爬取多张图片数据
- 4 爬取GDP数据
- 5 多任务爬虫
- 6 # 2. 数据可视化
- 7 pyecharts绘制饼图
- 8 # 3. 程序日志记录
- 9 日志的作用和级别
- 10 logging模块的使用
- 11 数据埋点简介

知识点1: 爬虫示例-爬取单张图片数据【掌握】

需求: 使用 requests 编写爬虫程序, 爬取 <http://127.0.0.1:8080/images/1.jpg> 图片数据并保存。

示例代码:

```
1  """
2  爬虫示例-爬虫单张图片数据
3  学习目标: 能够使用 requests 爬取单张图片数据并保存
4  """
5
6  import requests
7
8  # 准备请求的 URL 地址
9  url = 'http://127.0.0.1:8080/images/1.jpg'
```

```

10
11     # 发送请求
12     response = requests.get(url)
13
14     # 获取响应图片内容
15     image_content = response.content
16
17     # 将响应图片内容保存成本地图片文件
18     with open('./spider/1.jpg', 'wb') as f:
19         f.write(image_content)

```

知识点2：爬虫示例-爬取多张图片数据【掌握】

需求：访问 <http://127.0.0.1:8080/index.html> 网址，获取页面上的所有图片保存到本地。

示例代码：

```

1     """
2     爬虫示例-爬取多张图片数据
3     学习目标：能够使用 requests 爬取多张图片数据并保存
4     """
5
6     import requests
7     import re
8
9     # 思路
10    # ① 先请求 http://127.0.0.1:8080/index.html，获取响应内容
11    # ② 从上一步的响应内容中提取所有图片的地址
12    # ③ 遍历每一个图片地址，向每个图片地址发送请求，并将响应的内容保存成图片文件
13
14
15    def get_images():
16        # ① 先请求 http://127.0.0.1:8080/index.html，获取响应内容
17        url = 'http://127.0.0.1:8080/index.html'
18        # 发送请求
19        response = requests.get(url)
20
21        # 获取响应的内容
22        html_str = response.content.decode()
23        print(html_str)
24
25        # ② 从上一步的响应内容中提取所有图片的地址
26        image_url_list = re.findall(r' 网址，提取页面上的国家和GDP数据并保存到本地。

示例代码：

```

1 """
2 爬虫示例-爬取GDP数据
3 学习目标：能够使用 requests 爬取GDP数据并保存
4 """
5
6 # 思路
7 # ① 先请求 http://127.0.0.1:8080/index.html，获取响应内容
8 # ② 使用正则提取页面上的国家和GDP数据
9 # ③ 将提取的 GDP 保存到文件 gdp.txt 中
10
11 import requests
12 import re
13
14
15 def get_gdp_data():
16 # ① 先请求 http://127.0.0.1:8080/gdp.html，获取响应内容
17 url = 'http://127.0.0.1:8080/gdp.html'
18 # 发送请求
19 response = requests.get(url)
20
21 # 获取响应的内容
22 html_str = response.content.decode()
23 print(html_str)
24
25 # ② 使用正则提取页面上的国家和GDP数据
26 gdp_data = re.findall('(.*?).*?¥(.*?)亿元', html_str, flags=re.S)
27 print(gdp_data)
28
29 # ③ 将提取的 GDP 保存到文件 gdp.txt 中
30 with open('./spider/gdp.txt', 'w', encoding='utf8') as f:
31 f.write(str(gdp_data))
32
33 print('保存GDP数据完毕!!!')
34
35
36 if __name__ == '__main__':
37 get_gdp_data()

```

### 知识点4：爬虫示例-多任务版爬虫【了解】

真正的工作环境中，我们爬取的数据可能非常的多，如果还是使用单任务实现，这时候就会让我们爬取数据的时间很长，那么显然使用多任务可以大大提升我们爬取数据的效率。

需求：使用多线程实现分别爬取图片数据和 GDP 数据。

示例代码：

```

1 """
2 爬虫示例-爬虫多任务版
3 学习目标：能够使用多线程的方式执行多任务爬虫
4 """
5
6 import requests
7 import re
8 import threading
9
10
11 def get_images():
12 # ① 先请求 http://127.0.0.1:8080/index.html，获取响应内容
13 url = 'http://127.0.0.1:8080/index.html'
14 # 发送请求
15 response = requests.get(url)
16
17 # 获取响应的内容
18 html_str = response.content.decode()
19 # print(html_str)
20
21 # ② 从上一步的响应内容中提取所有图片的地址
22 image_url_list = re.findall(r'(.*?).*?¥(.*?)亿元', html_str, flags=re.S)
53 # print(gdp_data)
54
55 # ③ 将提取的 GDP 保存到文件 gdp.txt 中
56 with open('./spider/gdp.txt', 'w', encoding='utf8') as f:

```

```

57 f.write(str(gdp_data))
58
59 print('保存GDP数据完毕!!!')
60
61
62 if __name__ == '__main__':
63 # 创建两个线程，分别执行爬取图片和爬取GDP数据的任务
64 image_thread = threading.Thread(target=get_images)
65 gdp_thread = threading.Thread(target=get_gdp_data)
66
67 # 启动线程
68 image_thread.start()
69 gdp_thread.start()

```

## 知识点5: pyecharts 模块简介和安装【了解】

问题:

```
1 1. 什么是 pyecharts? 有什么作用?
```

echarts 是个由百度开源的 javascript 语言实现数据可视化框架，凭借着良好的交互性，精巧的图表设计，得到了众多开发者的认可。pyecharts 是 echarts 的 python 版本的实现，在 python 中，可以使用 pyecharts 进行数据可视化操作。

官网地址: <https://pyecharts.org/#/>

特性:

1. 简洁的API设计，使用如丝滑般流畅，支持链式调用
2. 囊括了30+种常见图表，应有尽有
3. 支持主流Notebook 环境，Jupyter Notebook 和JupyterLab
4. 可轻松集成至Flask, Django等主流Web框架
5. 高度灵活的配置项，可轻松搭配出精美的图表
6. 详细的文档和示例，帮助开发者更快的上手项目
7. 多达400+地图文件以及原生的百度地图，为地理数据可视化提供强有力的支持

安装:

```
1 pip install pyecharts
```

## 知识点6: pyecharts 绘制饼图-GDP数据可视化【掌握】

需求: 加载爬虫抓取的 GDP 数据，使用 pyecharts 绘制饼状图显示GDP前十的国家。

示例代码:

```

1 """
2 pyecharts-GDP数据可视化
3 学习目标: 能够使用 pyecharts 绘制饼图
4 """
5
6 # 思路
7 # ① 从文件中读取 GDP 数据
8 # ② 使用 pyecharts 绘制饼状图显示GDP前十的国家
9
10 # 导入饼图类
11 from pyecharts.charts import Pie
12 # 导入配置选项模块
13 import pyecharts.options as opts

```

```

14
15
16 def data_view_pie():
17 # ① 从文件中读取 GDP 数据
18 gdp_data = []
19
20 with open('./spider/gdp.txt', 'r', encoding='utf8') as f:
21 content = f.read()
22 gdp_data = eval(content)
23
24 # ② 使用 pyecharts 绘制饼状图显示GDP前十的国家
25 # 获取 GDP 前 10 的国家数据
26 gdp_top10 = gdp_data[:10]
27
28 # 创建饼图
29 pie = Pie(init_opts=opts.InitOpts(width="1400px", height="800px"))
30 # 给饼图添加数据
31 pie.add(
32 "GDP",
33 gdp_top10,
34 label_opts=opts.LabelOpts(formatter='{b}:{d}%')
35)
36 # 给饼图设置标题
37 pie.set_global_opts(title_opts=opts.TitleOpts(title="2020年世界GDP排名",
38 subtitle="美元"))
39 # 保存结果，默认保存到 render.html 文件
40 pie.render()
41
42 if __name__ == '__main__':
43 data_view_pie()

```

## 知识点7：程序日志简介及作用【掌握】

问题：

- 1 1. 什么是程序日志？
- 2 2. 日志有什么作用？
- 3 3. python中如何记录日志？

### 什么是程序日志？

程序的日志是记录程序在运行过程中，一些关键性的信息，比如：xx用户做了xx操作，服务器出现了xx错误等。

### 日志有什么作用？

- 1) 可以很方便的了解程序的运行情况
- 2) 方便开发人员检查bug
- 3) 可以分析用户的操作行为、喜好等信息

### python中如何记录日志？

python中，可以使用 logging 模块记录日志。

## 知识点8：程序日志的5个等级说明【掌握】

同其他信息一样，日志信息也是分重要程度的。

日志按重要程度从低到高，分成如下 5 个等级：

- 1) DEBUG：程序调试bug时使用
- 2) INFO：程序正常运行时使用
- 3) WARNING：程序未按预期运行时使用，但并不是错误，如：用户登录密码错误
- 4) ERROR：程序出错误时使用，如：IO操作失败
- 5) CRITICAL：特别严重的问题，导致程序不能再继续运行时使用，如：磁盘空间为空，一般很少使用

从低到高的：DEBUG < INFO < WARNING < ERROR < CRITICAL

## 知识点9：logging模块的基本使用【掌握】

logging 模块记录日志时，可以设置日志的等级，默认日志的等级是 WARNING。

记录日志信息时，<= 设置日志等级的信息，会直接被忽略(不做记录)。

示例代码：

```
1 """
2 logging模块记录日志
3 学习目标：能够使用 logging 模块记录日志
4 """
5
6 import logging
7
8 # 设置日志的等级、设置日志的记录格式
9 # level=logging.DEBUG：将日志等级设置成DEBUG
10 logging.basicConfig(level=logging.DEBUG,
11 format='%(asctime)s - %(filename)s[line:%(lineno)d] - %
12 (levelname)s: %(message)s',
13 filename='log.txt',
14 filemode='a')
15
16 # logging记录日志的默认级别是：WARNING
17 logging.debug('这是一个debug级别的日志')
18 logging.info('这是一个info级别的日志')
19 logging.warning('这是一个warning级别的日志')
20 logging.error('这是一个error级别的日志')
21 logging.critical('这是一个critical级别的日志')
```

## 知识点10：数据埋点介绍【了解】

问题：

1. 什么是数据埋点？数据埋点有什么作用？
2. 如何做数据埋点？

### 什么是数据埋点？数据埋点有什么作用？

数据埋点是数据采集的一种方式，是针对特定业务场景进行数据采集和上报的技术方案。简单来说就是根据具体的业务分析需求，将关心的用户行为信息记录都日志中，采集数据以便进行后续的分析操作。

## 如何做数据埋点？

- 1) 代码埋点：在产品源代码中添加记录日志的代码（需要开发人员配合）
- 2) 可视化埋点：通过第三方产品，配置记录关心的用户行为数据（业务人员自己即可配置）
- 3) 无埋点(全埋点)：通过第三方产品，配置记录用户的所有行为数据（业务人员自己即可配置）