

数据持久化

- CSV

```
1 import csv
2 with open('xxx.csv','w',encoding='utf-8',newline='') as f:
3     writer = csv.writer(f)
4     writer.writerow([])
```

- MongoDB

```
1 import pymongo
2
3 # __init__(self):
4
5     self.conn = pymongo.MongoClient('IP',27017)
6     self.db = self.conn['cardb']
7     self.myset = self.db['car_set']
8
9 # save_html(self,r_list):
10
11     self.myset.insert_one(dict)
```

MySQL

- pymysql

```
1 import pymysql
2
3 db = pymysql.connect('localhost','root','123456','maoyandb',charset='utf8')
4 cursor = db.cursor()
5
6 ins = 'insert into filmtab values(%s,%s,%s)'
7 cursor.execute(ins,['霸王别姬','张国荣','1993'])
8
9 db.commit()
10 cursor.close()
11 db.close()
```

```

1  # __init__(self):
2      self.db = pymysql.connect('IP',... ..)
3      self.cursor = self.db.cursor()
4
5  # save_html(self,r_list):
6      self.cursor.execute('sql',[data1])
7      self.db.commit()
8
9  # run(self):
10     self.cursor.close()
11     self.db.close()

```

- 练习 - 将电影信息存入MySQL数据库

```

1  【1】提前建库建表
2  mysql -h127.0.0.1 -uroot -p123456
3  create database maoyandb charset utf8;
4  use maoyandb;
5  create table maoyantab(
6  name varchar(100),
7  star varchar(300),
8  time varchar(100)
9  ) charset=utf8;
10
11 【2】使用execute()方法将数据存入数据库思路
12     2.1) 在 __init__() 中连接数据库并创建游标对象
13     2.2) 在 save_html() 中将所抓取的数据处理成列表，使用execute()方法写入
14     2.3) 在 run() 中等数据抓取完成后关闭游标及断开数据库连接

```

- 汽车之家二手车信息抓取

```

1  【1】URL地址
2      进入汽车之家官网，点击 二手车
3      即: https://www.che168.com/beijing/a0_0msdgscncgpi11to1cspexx0/
4
5  【2】抓取目标
6      每辆汽车的
7      2.1) 汽车名称
8      2.2) 行驶里程
9      2.3) 城市
10     2.4) 个人还是商家
11     2.5) 价格
12
13 【3】抓取前5页

```

- 参考答案

```

1  import requests
2  import re
3  import time

```

```

4  import random
5
6  class CarSpider:
7      def __init__(self):
8          self.url =
          'https://www.chel68.com/beijing/a0_0msdgsncgpilltolcsp{}exx0/?
          pvareaid=102179#currngpostion'
9          self.headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; WOW64)
          AppleWebKit/537.36 (KHTML, like Gecko) Chrome/81.0.4044.138 Safari/537.36'}
10
11     def get_html(self, url):
12         html = requests.get(url=url,
13         headers=self.headers).content.decode('gb2312', 'ignore')
14         self.parse_html(html)
15
16     def parse_html(self, html):
17         pattern = re.compile('<li class="cards-li list-photo-li".*?<div
18         class="cards-bottom">.*?<h4 class="card-name">(.*)</h4>.*?<p class="cards-unit">
19         (.*)</p>.*?<span class="pirce"><em>(.*)</em>', re.S)
20         car_list = pattern.findall(html)
21         self.save_html(car_list)
22
23     def save_html(self, car_list):
24         for car in car_list:
25             print(car)
26
27     def run(self):
28         for i in range(1,6):
29             page_url = self.url.format(i)
30             self.get_html(page_url)
31             time.sleep(random.randint(1,2))
32
33 if __name__ == '__main__':
34     spider = CarSpider()
35     spider.run()

```

MongoDB

- MongoDB特点

- 1 【1】非关系型数据库,数据以键值对方式存储,端口27017
- 2 【2】MongoDB基于磁盘存储
- 3 【3】MongoDB数据类型单一,值为JSON文档,而Redis基于内存,
- 4 3.1> MySQL数据类型:数值类型、字符类型、日期时间类型、枚举类型
- 5 3.2> Redis数据类型:字符串、列表、哈希、集合、有序集合
- 6 3.3> MongoDB数据类型:值为JSON文档
- 7 【4】MongoDB: 库 -> 集合 -> 文档
- 8 MySQL : 库 -> 表 -> 表记录

- MongoDB常用命令

```
1  Linux进入: mongo
2  >show dbs                - 查看所有库
3  >use 库名                 - 切换库
4  >show collections         - 查看当前库中所有集合
5  >db.集合名.find().pretty() - 查看集合中文档
6  >db.集合名.count()        - 统计文档条数
7  >db.集合名.drop()         - 删除集合
8  >db.dropDatabase()        - 删除当前库
9  # MongoDB - Command - 库->集合->文档
10
11  mongo
12
13  > show dbs
14  > use db_name
15  > show collections
16  > db.集合名.find().pretty()
17  > db.集合名.count()
18  > db.集合名.drop()
19  > db.dropDatabase()
```

- pymongo模块使用

```
1  import pymongo
2
3  # 1.连接对象
4  conn = pymongo.MongoClient(host = 'localhost',port = 27017)
5  # 2.库对象
6  db = conn['maoyandb']
7  # 3.集合对象
8  myset = db['maoyanset']
9  # 4.插入数据库
10 myset.insert_one({'name':'赵敏'})
```

- 练习 - 将电影信息存入MongoDB数据库

```
1  """
2  猫眼电影top100抓取 (电影名称、主演、上映时间)
3  存入mongodb数据库中
4  """
5  import requests
6  import re
7  import time
8  import random
9  import pymongo
10
11  class MaoyanSpider:
12      def __init__(self):
13          self.url = 'https://maoyan.com/board/4?offset={}'
```

```

14         self.headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; WOW64)
AppleWebKit/537.36 (KHTML, like Gecko) Chrome/81.0.4044.113 Safari/537.36'}
15         # 三个对象: 连接对象、库对象、集合对象
16         self.conn = pymongo.MongoClient('127.0.0.1', 27017)
17         self.db = self.conn['maoyandb']
18         self.myset = self.db['maoyanset2']
19
20     def get_html(self, url):
21         html = requests.get(url=url, headers=self.headers).text
22         # 直接调用解析函数
23         self.parse_html(html)
24
25     def parse_html(self, html):
26         """解析提取数据"""
27         regex = '<div class="movie-item-info">.*?title="(.*?)".*?<p class="star">
(.*?)</p>.*?<p class="releasetime">(.*?)</p>'
28         pattern = re.compile(regex, re.S)
29         r_list = pattern.findall(html)
30         # r_list: [('活着', '牛犇', '2000-01-01'), (), (), ..., ())
31         self.save_html(r_list)
32
33     def save_html(self, r_list):
34         """数据处理函数"""
35         for r in r_list:
36             item = {}
37             item['name'] = r[0].strip()
38             item['star'] = r[1].strip()
39             item['time'] = r[2].strip()
40             print(item)
41             # 存入到mongodb数据库
42             self.myset.insert_one(item)
43
44     def run(self):
45         """程序入口函数"""
46         for offset in range(0, 91, 10):
47             url = self.url.format(offset)
48             self.get_html(url=url)
49             # 控制数据抓取频率: uniform() 生成指定范围内的浮点数
50             time.sleep(random.uniform(0, 1))
51
52 if __name__ == '__main__':
53     spider = MaoyanSpider()
54     spider.run()

```

CSV

- csv描述

```
1  【1】作用
2      将爬取的数据存放到本地的csv文件中
3
```

```
4  【2】使用流程
5      2.1> 打开csv文件
6      2.2> 初始化写入对象
7      2.3> 写入数据(参数为列表)
8
```

```
9  【3】示例代码
10     import csv
11     with open('sky.csv','w') as f:
12         writer = csv.writer(f)
13         writer.writerow([])
```

- 示例

```
1  【1】题目描述
2      创建 test.csv 文件，在文件中写入数据
3
4  【2】数据写入 - writerow([])方法
5      import csv
6      with open('test.csv','w') as f:  # with open('test.csv','w',newline='') as f:-
      ---->windows里面的写法，因为再wiondows中每条数据会有一个空行
7          writer = csv.writer(f)
8          writer.writerow(['超哥哥','25'])
9
```

- 练习 - 使用 writerow() 方法将猫眼电影数据存入本地 maoyan.csv 文件

```
1  【1】在 __init__() 中打开csv文件，因为csv文件只需要打开和关闭1次即可
2  【2】在 save_html() 中将所抓取的数据处理成列表，使用writerow()方法写入
3  【3】在run() 中等数据抓取完成后关闭文件
```

- 代码实现

```
1  """
2  猫眼电影top100抓取 (电影名称、主演、上映时间)
3  存入csv文件,使用writerow()方法
4  """
5  import requests
6  import re
7  import time
8  import random
9  import csv
10
11  class MaoyanSpider:
12      def __init__(self):
13          self.url = 'https://maoyan.com/board/4?offset={}'
```

```

14         self.headers = {'User-Agent': 'Mozilla/5.0 (compatible; MSIE 9.0; Windows
NT 6.1; Win64; x64; Trident/5.0; .NET CLR 2.0.50727; SLCC2; .NET CLR 3.5.30729;
.NET CLR 3.0.30729; Media Center PC 6.0; InfoPath.3; .NET4.0C; Tablet PC 2.0;
.NET4.0E)'}
15         # 打开文件,初始化写入对象
16         self.f = open('maoyan.csv', 'w', newline='', encoding='utf-8')
17         self.writer = csv.writer(self.f)
18
19     def get_html(self, url):
20         html = requests.get(url=url, headers=self.headers).text
21         # 直接调用解析函数
22         self.parse_html(html)
23
24     def parse_html(self, html):
25         """解析提取数据"""
26         regex = '<div class="movie-item-info">.*?title="(.*?)".*?<p class="star">
(.*?)</p>.*?<p class="releasetime">(.*?)</p>'
27         pattern = re.compile(regex, re.S)
28         r_list = pattern.findall(html)
29         # r_list: [('活着', '牛犇', '2000-01-01'), (), (), ..., ())
30         self.save_html(r_list)
31
32     def save_html(self, r_list):
33         """数据处理函数"""
34         for r in r_list:
35             li = [ r[0].strip(), r[1].strip(), r[2].strip() ]
36             self.writer.writerow(li)
37             print(li)
38
39     def run(self):
40         """程序入口函数"""
41         for offset in range(0, 91, 10):
42             url = self.url.format(offset)
43             self.get_html(url=url)
44             # 控制数据抓取频率:uniform()生成指定范围内的浮点数
45             time.sleep(random.uniform(1,2))
46
47         # 所有数据抓取并写入完成后关闭文件
48         self.f.close()
49
50 if __name__ == '__main__':
51     spider = MaoyanSpider()
52     spider.run()

```