

爬虫期中作业：

姓名：李雪萌

班级:16 计算机科学与技术三班

学号：16219111337

项目内容：

Django 项目+爬虫代码+数据库+实验报告

涉及知识：

Python 语言的基本语法以及第三方库的使用；

静态网页的爬取；

网页的解析；

动态网页的爬取；

模拟无头浏览器爬取动态网页；

爬虫与反爬虫

涉及其他：

正则表达式

Django 的 mvc 模型

Selenium 的安装调用

豆瓣爬取：

```
1 import requests
2 from lxml import etree
3 import MySQLdb
4
5
6 con= MySQLdb.connect(host='localhost',user='root',password='root',db='test',charset="utf8")
7 cursor=con.cursor()
8 cursor.execute("create table test(No varchar(32) ,content char(255))")
9
10 #静态网页抓取：request爬虫实践：top250电影数据
11 def get_page(start_num):
12     url='https://movie.douban.com/top250?start=%s&filter=' %start_num
13     print(url)
14
15     response=requests.get(url)
16     tree=etree.HTML(response.text)
17     title=tree.xpath('//span[@class="title"][1]/text()')
18     return title
19
20 def get_all_page(start,end):
21     result=[]
22     for i in range(start,end-start):
23         title_list=get_page(i*25)
24         result+=title_list
25         print(result)
26     return result
27
28 if __name__=="__main__":
29     result=get_all_page(0,10)
30
31     for i in range(len(result)):
32         cursor.execute("INSERT INTO test(No,content) values(%d,'%s')" %(i+1,result[i]))
33     cursor.close()
34     con.commit()
35     con.close()
```

天气爬取：

```

from bs4 import BeautifulSoup
from bs4 import UnicodeDammit
import urllib.request
import MySQLdb

class weatherDB:
    def openDB(self):
        self.con= MySQLdb.connect(host='localhost',user='root',password='root',db='test',charset="utf8")
        self.cursor=self.con.cursor()
        try:
            self.cursor.execute("create table weathers1(wCity varchar(16),wDate varchar(16),wWeather varchar(64),wTemp varchar(32),constraint pk_we
        except:
            self.cursor.execute("delete from weathers")

    def closeDB(self):
        self.con.commit()
        self.con.close()

    def insert(self,city,date,weather,temp):
        try:
            self.cursor.execute("insert into weathers1(wCity,wDate,wWeather,wTemp) values(?,?,?,?)",(city,date,weather,temp))
        except Exception as err:
            print(err)

    def show(self):
        self.cursor.execute("select * from weathers1")
        rows=self.cursor.fetchall()
        print("%-16s%-16s%-32s%-16s"%(city,"date","weather","temp"))
        for row in rows:
            print("%-16s%-16s%-32s%-16s"%(row[0],row[1],row[2],row[3]))

class weatherForecast:
    def __init__(self):
        self.headers={"User-Agent":"Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/73.0.3683.103 Safari/53
        self.cityCode={"北京":"101010100","上海":"101020100","广州":"101280101","深圳":"101280601"}

    def forecastCity(self,city):
        if city not in self.cityCode.keys():
            print(city+"code cannot be found")
            return

```

```

url="http://www.weather.com.cn/weather/"+self.cityCode[city]+".shtml"
try:
    req=urllib.request.Request(url,headers=self.headers)
    data=urllib.request.urlopen(req)
    data=data.read()
    dammit=UnicodeDammit(data,["utf-8","gbk"])
    data=dammit unicode_markup
    soup=BeautifulSoup(data,"lxml")
    lis=soup.select("ul[class='t clearfix'] li")
    for li in lis:
        try:
            date=li.select('h1')[0].text
            weather=li.select('p[class="wea"]')[0].text
            temp=li.select('p[class="tem"] span')[0].text+"/"+li.select('p[class="tem"] i')[0].text
            print(city,date,weather,temp)
            self.db.insert(city,date,weather,temp)
        except Exception as err:
            print(err)
    except Exception as err:
        print(err)

def process(self,cities):
    self.db=weatherDB()
    self.db.openDB()
    for city in cities:
        self.forecastCity(city)

    #self.db.show()
    self.db.closeDB()

ws=weatherForecast()
ws.process(["北京","上海","深圳","广州"])
print("completed")

```

京东爬取：

```

1 from selenium import webdriver
2 from selenium.webdriver.chrome.options import Options
3 import urllib.request
4 import threading
5 import MySQLdb
6 import os
7 import datetime
8
9
10 class MySpider:
11     header={"User-Agent: Mozilla/5 (Windows NT 10 ; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/73.0.3683.103 Safari/537.36"}
12     imagePath="download"
13
14     def startUp(self,url,key):
15         try:
16             self.con=MySQLdb.connect(host='localhost',user='root',password='root',db='test',charset="utf8")
17             self.cursor=self.con.cursor()
18             try:
19                 self.cursor.execute("drop table phones")
20             except:
21                 pass
22             try:
23                 sql="create table phones(mNo varchar(32) primary key,mMark varchar(256),mPrice varchar(32),mNote varchar(1024),mFile varchar(256))"
24                 self.cursor.execute(sql)
25             except:
26                 pass
27         except Exception as err:
28             print(err)
29
30         chrome_option=Options()
31         chrome_option.add_argument('--headless')
32         chrome_option.add_argument('--disable-gpu')
33
34         self.driver=webdriver.Chrome(chrome_options=chrome_option)
35         self.threads=[]
36         self.No=0
37         self.imgNo=0
38
39         try:
40             if not os.path.exists(MySpider.imagePath):
41                 os.mkdir(MySpider.imagePath)
42             images=os.listdir(MySpider.imagePath)
43             for img in images:

```

```

44                 for img in images:
45                     s=os.path.join(MySpider.imagePath,img)
46                     os.remove(s)
47
48         except Exception as err:
49             print(err)
50
51         try:
52             if not os.path.exists(MySpider.imagePath):
53                 os.mkdir(MySpider.imagePath)
54             images=os.listdir(MySpider.imagePath)
55             for img in images:
56                 s=os.path.join(MySpider.imagePath,img)
57                 os.remove(s)
58         except Exception as err:
59             print(err)
60
61         self.driver.get(url)
62         keyInput=self.driver.find_element_by_id("key")
63         keyInput.send_keys(key)
64         keyInput.send_keys(Keys.ENTER)
65
66     def closeUp(self):
67         try:
68             self.con.commit()
69             self.con.close()
70             self.driver.close()
71         except Exception as err:
72             print(err)
73
74     def insertDB(self,mNo,mMark,mPrice,mNote,mFile):
75         try:
76             sql="insert into phones(mNo,mMark,mPrice,mNote,mFile)values(?,?,?,?,?)"
77             self.cursor.execute(sql,(mNo,mMark,mPrice,mNote,mFile))
78         except Exception as err:
79             print(err)
80
81     def showDB(self):
82         try:
83             con= MySQLdb.connect(host='localhost',user='root',password='root',db='test',charset="utf8")
84             cursor=con.cursor()
85             print("%-8s %-16s %-8s %-16s %s"%("No","Mark","Price","Image","Note"))
86             cursor.execute("select mNo,mMark,mPrice,mNote from phones order by mNo")
87             rows=cursor.fetchall()
88             for row in rows:
89                 print("%-8s %-16s %-8s %-16s %s"%(row[0],row[1],row[2],row[3],row[4]))
90             con.close()

```

```

        except Exception as err:
            print(err)
    def download(self,src1,src2,mFile):
        data=None
        if src1:
            try:
                req=urllib.request.Request(src1,headers=MySpider.headers)
                resp=urllib.request.urlopen(req,timeout=400)
                data=resp.read()
            except:
                pass
        if not data and src2:
            try:
                req=urllib.request.Request(src2,headers=MySpider.headers)
                resp=urllib.request.urlopen(req,timeout=400)
                data=resp.read()
            except:
                pass
        if data:
            fobj=open(MySpider.imagePath+"\\\\"+mFile,"wb")
            fobj.write(data)
            fobj.close()
            print("download",mFile)

    def processSpider(self):
        try:
            time.sleep(10)
            print(self.driver.current_url)
            lis=self.driver.find_element_by_xpath("//div[@id='J_goodsList']/li[@class='gl-item']")
            for li in lis:
                try:
                    src1=li.find_element_by_xpath(".\\div[@class='p-img']/a/img").get_attribute("src")
                except:
                    src1=""
                try:
                    src2=li.find_element_by_xpath(".\\div[@class='p-img']/a/img").get_attribute("data-lazy-img")
                except:
                    price="0"
                try:
                    note=li.find_element_by_xpath(".//div[@class='p-price']/i").text
                    mark=note.split(" ")[0]
                    mark=mark.replace("爱心东东\\n","")
                    mark=mark.replace(",","")
                    note=note.replace("爱心东东\\n","")
                    note=note.replace(",","")

```

```

                    note=note.replace(",","")
                except:
                    note=""
                    mark=""
            self.No=self.No+1
            no=str(self.No)
            while len(no)<6:
                no="0"+no
            print(no,mark,price)
            if src1:
                src1=urllib.request.urljoin(self.driver.current_url,src1)
                p=src1.rfind(".")
                mFile=no+src1[p:]
            elif src2:
                src2=urllib.request.urljoin(self.driver.current_url,src2)
                p=src2.rfind(".")
                mFile=no+src2[p:]
            if src1 or src2:
                T=threading.Thread(target=self.download,args=(src1,src2,mFile))
                T.setDaemon(False)
                T.start()
                self.threads.append(T)
            else:
                mFile=""
                self.insertDB(no,mark,price,note,mFile)
            try:
                self.driver.find_element_by_xpath("//span[@class='p-num']/a[@class='pn-next disabled']")
            except:
                nextPage=self.driver.find_element_by_xpath("//span[@class='p-num']/a[@class='pn-next']")
                nextPage.click()
                self.processSpider()
        except Exception as err:
            print(err)

    def executeSpider(self,url,key):
        starttime=datetime.datetime.now()
        print("Spider starting.....")
        self.startUp(url,key)
        self.processSpider()
        self.closeUp()
        for t in self.threads:
            t.join()
        print("Spider completed.....")
        endtime=datetime.datetime.now()
        elapsed=(endtime-starttime).seconds
        print("Total",elapsed,"second elapsed")

```

```

        t.join()
    print("Spider completed.....")
    endtime=datetime.datetime.now()
    elapsed=(endtime-starttime).seconds
    print("Total",elapsed,"second elapsed")

url='http://www.jd.com'
spider=MySpider()
while True:
    print("1.爬取")
    print("2.显示")
    print("3.退出")
    s=input("请输入选择 (1, 2, 3) :")
    if s=="1":
        spider.executeSpider(url,"手机")
    elif s=="2":
        spider.showDB()
    elif s=="3":
        break

```

淘宝爬取：

```

import requests
import re
import MySQLdb

def getHTMLText(url):
    try:
        r=requests.get(url,timeout=30)
        r.raise_for_status()
        r.encoding=r.apparent_encoding
        return r.text
    except:
        return ""

def parsePage(ilt,html):
    try:
        plt=re.findall(r'view_price"\:\:("[\d\.]*)"',html)
        tlt=re.findall(r'raw_title"\:\:("[\d\.\s]*"',html)
        for i in range(len(plt)):
            price=eval(plt[i].split(':')[1])
            title=eval(tlt[i].split(':')[1])
            ilt.append([price,title])
    except:
        print("")

def printGoodList(ilt):
    con= MySQLdb.connect(host='localhost',user='root',password='root',db='test',charset="utf8")
    cursor=con.cursor()
    cursor.execute("select * from shubao")
    tplt="{:4}\t{:8}\t{:16}"
    print(tpl.format("序号","价格","商品名称"))
    count=0
    for g in ilt:
        count=count+1
        print(tpl.format(count,g[0],g[1]))
        cursor.execute("insert into shubao(id,Price,Title) values('"+count+"','"+g[0]+"','"+g[1]+"')")
    cursor.close()
    con.commit()
    con.close()

def main():
    goods='书包'
    depth=3

```

```

1 depth=3
2 start_url='https://s.taobao.com/search?q='+goods
3 infoList=[]
4 for i in range(depth):
5     try:
6         url=start_url+'&s='+str(44*i)
7         html=getHTMLText(url)
8         parsePage(infoList,html)
9     except:
10        continue
11 printGoodList(infoList)
12
13 main()

```

爬取结果展示：


```
e:\爬虫\16219111337_lixuemeng\pythonfile>python tq.py
```

```
list index out of range
```

北京 26日(明天) 晴转多云 21°C/10°C

北京 27日(后天) 多云 18°C~7°C

北京 28日(周日) 多云 21°C/9°C

北京 29日(周一) 多云转小雨 24°C/13°C

北京 30日(周二) 多云 26°C/14°C

北京 1日(周三) 多云 24°C/14°C

```
list index out of range
```

上海 26日(明天) 多云转晴 18°C/12°C

上海 27日(后天) 多云 18°C/14°C

上海 28日(周日) 小雨转多云 23°C/17°C

上海 29日(周一) 中雨转阴 24°C/17°C

上海 30日(周二) 阴转小雨 22°C/17°C

上海 1日(周三) 小雨转多云 23°C/17°C

```
list index out of range
```

广州 26日(明天) 中雨转中到大雨 28°C/24°C

广州 27日(后天) 中到大雨转雷阵雨 28°C/24°C

广州 28日(周日) 雷阵雨 29°C/25°C

广州 29日(周一) 雷阵雨转中雨 30°C/25°C

广州 30日(周二) 中雨转大到暴雨 30°C/22°C

广州 1日(周三) 大到暴雨转多云 26°C/20°C

```
list index out of range
```

深圳 26日(明天) 中雨转大雨 29°/24°C

深圳 27日(后天) 大雨转雷阵雨 27°C/23°C

深圳 28日(周日) 雷阵雨 28°C/24°C

深圳 29日(周一) 雷阵雨 30°C/25°C

深圳 30日(周二) 雷阵雨转暴雨 30°C/23°C

深圳 1日(周三) 暴雨转阵雨 27°C/23°C

完成！

27	69.00	牛牛发卖双页日本2019新款韩版时尚男女生书包电脑包双肩书包
28	209.00	双页双页书包2019新款韩版牛牛发卖书包时尚百搭双肩书包2019
49	69.00	双页双页书包2019新款韩版牛牛发卖书包时尚百搭双肩书包2019
50	69.00	双页双页韩版时尚男女生书包电脑包双肩书包2019新款韩版牛牛
51	69.00	地球书包2019新款韩版牛牛发卖书包时尚百搭双肩书包2019
52	128.00	双肩双肩书包日本2019新款ins风双肩书包书包书包ins风alleno
53	149.00	韩版男式双肩双肩书包日本2019新款韩版时尚双肩书包时尚韩版书包
54	169.00	韩版男式双肩双肩书包日本2019新款韩版时尚双肩书包时尚韩版书包
55	69.00	韩版男式双肩双肩书包日本2019新款韩版时尚双肩书包时尚韩版书包
56	169.00	七色狼双肩书包 新款韩版男式双肩书包女初中生大容量书包学生书
57	139.00	韩版双肩书包日本2019新款韩版时尚双肩书包时尚韩版书包
58	129.00	韩版双肩双肩书包日本2019新款韩版时尚双肩书包时尚韩版书包
59	139.00	ukel双肩双肩男式男生双肩1-6年级韩版时尚双肩书包1-12年级韩版
60	228.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
61	69.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
62	69.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
63	149.00	七色狼双肩书包 新款韩版男式双肩书包女初中生大容量书包学生书
64	158.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
65	85.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
66	116.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
67	149.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
68	298.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
69	79.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
70	80.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
71	139.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
72	358.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
73	149.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
74	158.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
75	169.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
76	249.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
77	124.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
78	69.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
79	69.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
80	49.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
81	169.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
82	29.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
83	168.00	FAL双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
84	129.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
85	138.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
86	169.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
87	25.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
88	38.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
89	69.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
90	99.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
91	274.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
92	58.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
93	169.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
94	169.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
95	168.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
96	169.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
97	169.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
98	169.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
99	169.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩
100	169.00	双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩双肩

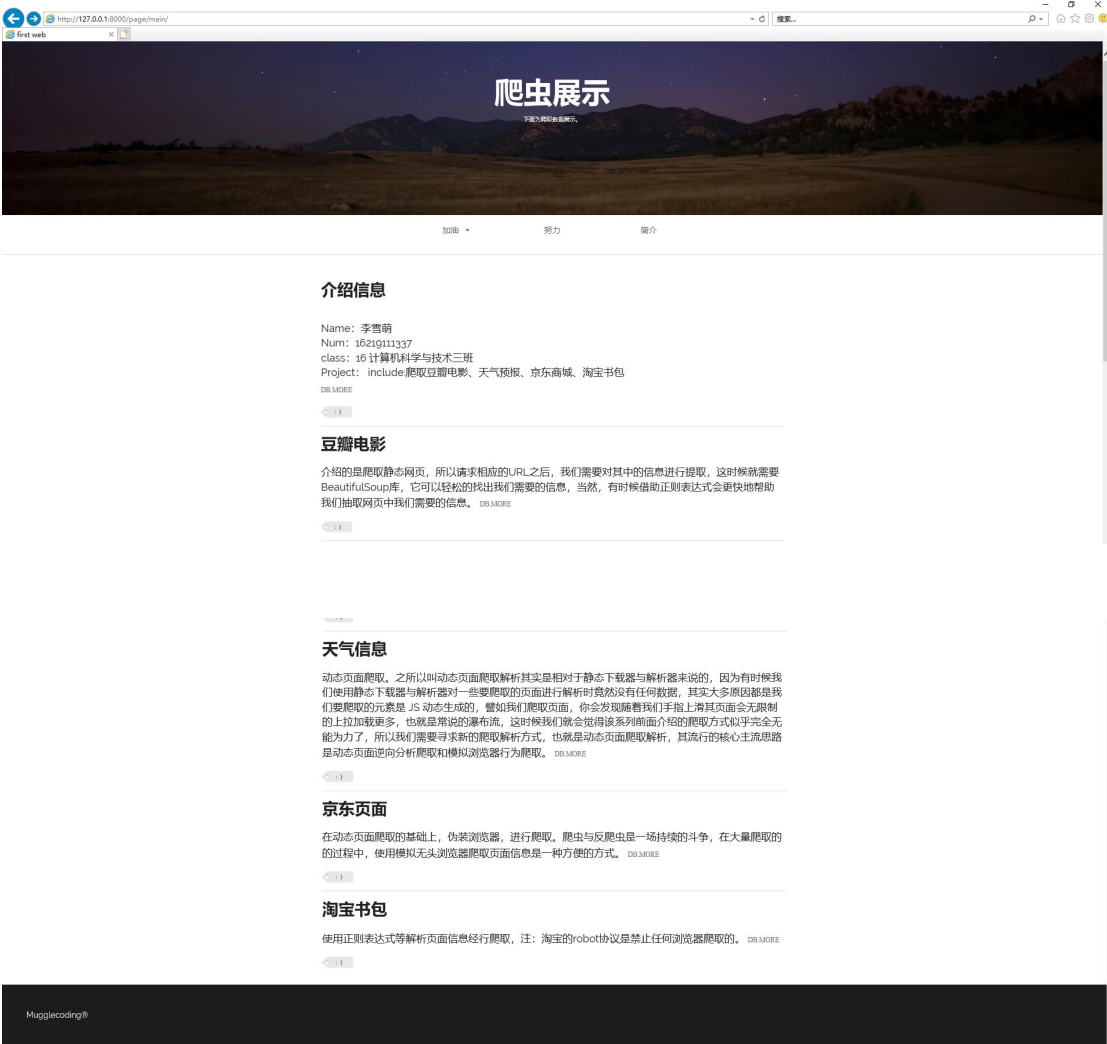
网站思想：

1.运通 mvc 框架（简单用 mvt）

2.Template 模板页用于显示 view 层的信息

3.Modle 层是 django 特有的数据库形式，以类的方式引用

4.页面使用 html+css（semantic ui）+简单 js+django 特有编辑语言



代码：

web：

```

1 <!DOCTYPE html>
2 {% load staticfiles %}
3 <html>
4 <head>
5 <meta charset="utf-8">
6 <title>first web</title>
7 <link rel="stylesheet" href="{% static 'css/semantic.css' %}" media="screen" title="no title" charset="utf-8">
8 <link href="https://fonts.googleapis.com/css?family=Oswald|Raleway" rel="stylesheet">
9
10
11 <style type="text/css">
12     h1 {
13         font-family: 'Oswald', sans-serif!important;
14         font-size: 40px;
15     }
16
17     body {
18         font-family: 'Raleway', sans-serif;
19     }
20     p {
21         font-family: 'Raleway', sans-serif;
22         font-size: 18px;
23     }
24     .ui.vertical.segment.masthead {
25         height: 300px;
26         background-image: url("{% static 'images/star_banner.jpg' %}");
27         background-size: cover;
28         background-position: 100% 80%;
29     }
30
31     .ui.container.segment {
32         width: 800px;
33     }
34
35     .ui.center.aligned.header.blogslogon {
36         margin-top: 40px;
37     }
38
39     .ui.center.aligned.header.blogslogon p {
40         margin-top: 10px;
41         color: white;
42         font-size: 10px;
43     }
44     .ui.container.nav {
45         width: 500px;
46     }
47 </style>

```

```

{% block css %}{% endblock%}

</head>
<body>
    <div class="ui inverted vertical segment masthead">
        <h1 class="ui center aligned header blogslogon" style="font-size:50px;font-family: 'Raleway', sans-serif!important;">
            爬虫展示
            <p class="ui sub header">
                下面为爬取数据展示。
            </p>
        </h1>
    </div>
    <div class="ui container nav">
        <div class="ui borderless text three item menu">
            <div class="ui simple dropdown item">
                加油
                <i class="dropdown icon"></i>
                <div class="menu">
                    <a class="item" href="">简介</a>
                    <a class="item" href="">简介</a>
                </div>
            </div>
            <a class="item">
                努力
            </a>
            <a class="item">
                简介
            </a>
        </div>
    </div>
    <div class="ui divider"></div>

{% block header %}{% endblock%}

<div class="ui vertical segment">
    <div class="ui container vertical segment">
        <a href="#">
            <h1 class="ui header">
                介绍信息
            </h1>
        </a>
    </div>

```

```

</a>

<i class="icon grey small unhide"></i>
<p>
    <form >
        <div class="six wide column">
            <p><label >Name: 李雪萌</label></p>
            </div>
            <div >
                <p><label >Num: 16219111337</label></p>
            </div>
            <div >
                <p><label >class: 16 计算机科学与技术三班</label></p>
            </div>
            <div >
                <p>
                    <label>
                        Project:
                        | include:爬取豆瓣电影、天气预报、京东商城、淘宝书包
                    </label>
                </p>
            </div>
        </form>
        <a href="#">
            <i class="angle tiny double grey right icon">DB.MORE</i>
        </a>
    </p>

    <div class="ui mini tag label">
        | : )
    </div>
</div>
<div class="ui container vertical segment">
    <a href="#">
        <h1 class="ui header">
            | 豆瓣电影
        </h1>
    </a>

    <i class="icon grey small unhide"></i>
    <p>
        介绍的是爬取静态网页，所以请求相应的URL之后，我们需要对其中的信息进行提取，这时候就需要BeautifulSoup库，它可以轻松的找出我们需要的信息
        <a href="http://127.0.0.1:8000/page/db">
            <i class="angle tiny double grey right icon">DB.MORE</i>
        </a>
    </p>

```

```

    </p>

    <div class="ui mini tag label">
        | : )
    </div>
</div>

<div class="ui container vertical segment">
    <a href="#">
        <h1 class="ui header">
            | 天气信息
        </h1>
    </a>

    <i class="icon grey small unhide"></i>
    <p>
        动态页面爬取。之所以叫动态页面爬取解析其实是相对于静态下载器与解析器来说的，因为有时候我们使用静态下载器与解析器对一些要爬取的页面进行爬
        <a href="http://127.0.0.1:8000/page/tq">
            <i class="angle tiny double grey right icon">DB.MORE</i>
        </a>
    </p>

    <div class="ui mini tag label">
        | : )
    </div>
</div>
<div class="ui container vertical segment">
    <a href="#">
        <h1 class="ui header">
            | 京东页面
        </h1>
    </a>

    <i class="icon grey small unhide"></i>
    <p>
        在动态页面爬取的基础上，伪装浏览器，进行爬取。爬虫与反爬虫是一场持续的斗争，在大量爬取的过程中，使用模拟无头浏览器爬取页面信息是一种方
        <a href="http://127.0.0.1:8000/page/jd">
            <i class="angle tiny double grey right icon">DB.MORE</i>
        </a>
    </p>

    <div class="ui mini tag label">
        | : )
    </div>
</div>
<div class="ui container vertical segment">

```

```

        : )
    </div>
    <div class="ui container vertical segment">
        <a href="#">
            <h1 class="ui header">
                | 淘宝书包
            </h1>
        </a>

        <i class="icon grey small unhide"></i>
        <p>
            使用正则表达式等解析页面信息行爬取，注：淘宝的robot协议是禁止任何浏览器爬取的。
            <a href="http://127.0.0.1:8080/page/tb">
                | <i class="angle tiny double grey right icon">DB.MORE</i>
            </a>
        </p>

        <div class="ui mini tag label">
            | : )
        </div>
    </div>

</div>

{% block body %}{% endblock%}

<div class="ui inverted vertical very padded segment">
    | Mugglecoding®
</div>

{% block foot %}{% endblock%}

</body>
</html>

```

web_db:

```
{% load staticfiles %}
<html>
  <head>
    <meta charset="utf-8">
    <title>first web</title>
    <link rel="stylesheet" href="{% static 'css/semantic.css' %}" media="screen" title="no title" charset="utf-8">
    <link href="https://fonts.googleapis.com/css?family=Oswald|Raleway" rel="stylesheet">

  <style type="text/css">
    h1 {
      font-family: 'Oswald', sans-serif!important;
      font-size: 40px;
    }

    body {
      font-family: 'Raleway', sans-serif;
    }
    p {
      font-family: 'Raleway', sans-serif;
      font-size: 18px;
    }
    .ui.vertical.segment.masthead {
      height: 300px;
      background-image: url("{% static 'images/star_banner.jpg' %}");
      background-size: cover;
      background-position: 100% 80%;
    }

    .ui.container.segment {
      width: 800px;
    }

    .ui.center.aligned.header.blogslogon {
      margin-top: 40px;
    }

    .ui.center.aligned.header.blogslogon p {
      margin-top: 10px;
      color: □white;
      font-size: 10px;
    }

    .ui.container.nav {
      width: 500px;
    }
  </style>
</head>
  <body>
    <div class="ui container">
      <div class="ui vertical segment masthead">
        <img alt="Star Banner" data-bbox="115 115 885 300"/>
      </div>
      <div class="ui container segment">
        <div class="ui center aligned header blogslogon">
          <p>
            <span class="ui icon"></span>
            <span>Blogs Logon</span>
          </p>
          <p>
            <span class="ui icon"></span>
            <span>Blogs Logon</span>
          </p>
        </div>
        <div class="ui container nav">
          <div class="ui button">
            <span class="ui icon"></span>
            <span>Blogs Logon</span>
          </div>
          <div class="ui button">
            <span class="ui icon"></span>
            <span>Blogs Logon</span>
          </div>
        </div>
      </div>
    </div>
  </body>
</html>
```



```

    }

</style>
{% block css %}{% endblock%}

</head>
<body>
    <div class="ui inverted vertical segment masthead">

        <h1 class="ui center aligned header blogslogon" style="font-size:50px;font-family: 'Raleway', sans-serif!important;">
            爬虫展示
        <p class="ui sub header">
            下面为爬取数据展示。
        </p>

    </h1>
</div>
<div class="ui container nav">
    <div class="ui borderless text three item menu ">
        <div class="ui simple dropdown item">
            加油
            <i class="dropdown icon"></i>
            <div class="menu">
                <a class="item" href="">简介</a>
                <a class="item" href="">简介</a>
            </div>
        </div>
        <a class="item">
            努力
        </a>
        <a class="item">
            简介
        </a>
    </div>
</div>

<div class="ui divider"></div>

{% block header %}{% endblock%}

<div class="ui container nav" >

    <br>

    <br>

```

```

    <br>

    <h1>豆瓣电影</h1>

    <br>

    <br>

    <h3><a href='http://127.0.0.1:8000/page/main'>回到页面</a></h3>

    <br>

    <br>

    {% for content in content_list %}

    <h3>
        编号: {{content.id}} <br>
        名字: {{content.content}}<br>
    </h3>

    <br>

    <br><br><br>

    {% endfor %}

    {% if data %}

    <ul id="pages" class="pagination pagination-sm pagination-xs">

        {% if data.first %}

        <li><a href="?page=1">1</a></li>

        {% endif %}

        {% if data.left %}

        {% if data.left_has_more %}

```

```

        {% if data.left_has_more %}

        <li><span>...</span></li>

        {% endif %}

        {% for i in data.left %}

        <li><a href="?page={{i}}">{{i}}</a></li>

        {% endfor %}

        {% endif %}

        <li class="active"><a href="?page={{data.page}}">{{data.page}}</a></li>

        {% if data.right %}

        {% for i in data.right %}

        <li><a href="?page={{i}}">{{i}}</a></li>

        {% endfor %}

        {% if data.right_has_more %}

        <li><span>...</span></li>

        {% endif %}

        {% endif %}

        {% if data.last %}

        <li><a href="?page={{data.total_pages}}">{{data.total_pages}}</a></li>

        {% endif %}

    </ul>

    {% endif %}

</div>

{% block body %}{% endblock%}

<div class="ui inverted vertical very padded segment">

</div>

{% block body %}{% endblock%}

<div class="ui inverted vertical very padded segment">
|   Mugglecoding®
</div>

{% block foot %}{% endblock%}

</body>
</html>

```

web_tq:

```

!DOCTYPE html>
{% load staticfiles %}
<html>
  <head>
    <meta charset="utf-8">
    <title>first web</title>
    <link rel="stylesheet" href="{% static 'css/semantic.css' %}" media="screen" title="no title" charset="utf-8">
    <link href="https://fonts.googleapis.com/css?family=Oswald|Raleway" rel="stylesheet">

    <style type="text/css">
      h1 {
        font-family: 'Oswald', sans-serif!important;
        font-size: 40px;
      }

      body {
        font-family: 'Raleway', sans-serif;
      }

      p {
        font-family: 'Raleway', sans-serif;
        font-size: 18px;
      }

      .ui.vertical.segment.masthead {
        height: 300px;
        background-image: url("{% static 'images/star_banner.jpg' %}");
        background-size: cover;
        background-position: 100% 80%;
      }

      .ui.container.segment {
        width: 800px;
      }

      .ui.center.aligned.header.blogslogon {
        margin-top: 40px;
      }

      .ui.center.aligned.header.blogslogon p {
        margin-top: 10px;
        color: white;
        font-size: 10px;
      }

      .ui.container.nav {
        width: 500px;
      }
    </style>
  </head>
  <body>
    <div class="ui inverted vertical segment masthead">
      <h1 class="ui center aligned header blogslogon" style="font-size: 50px; font-family: 'Raleway', sans-serif!important;">
        爬虫展示
        <p class="ui sub header">
          下面为爬取数据展示。
        </p>
      </h1>
    </div>
    <div class="ui container nav">
      <div class="ui borderless text three item menu">
        <div class="ui simple dropdown item">
          加油
          <i class="dropdown icon"></i>
          <div class="menu">
            <a class="item" href="">简介</a>
            <a class="item" href="">简介</a>
          </div>
        </div>
        <a class="item">
          努力
        </a>
        <a class="item">
          简介
        </a>
      </div>
    </div>
    <div class="ui divider"></div>
    {% block header %}{% endblock%}

    <div class="ui container nav">
      <br>
      <br>
    </div>
  </body>
</html>

```

```

    </div>
  </body>
</html>

```

[illegible]

```
{% if data.left_has_more %}

<li><span>...</span></li>

{% endif %}

{% for i in data.left %}

<li><a href="?page={{i}}">{{i}}</a></li>

{% endfor %}

{% endif %}

<li class="active"><a href="?page={{data.page}}">{{data.page}}</a></li>

{% if data.right %}

{% for i in data.right %}

<li><a href="?page={{i}}">{{i}}</a></li>

{% endfor %}

{% if data.right_has_more %}

<li><span>...</span></li>

{% endif %}

{% endif %}

{% if data.last %}

<li><a href="?page={{data.total_pages}}">{{data.total_pages}}</a></li>

{% endif %}

</ul>

{% endif %}

</div>

{% block body %}{% endblock %}

<div class="ui inverted vertical very padded segment">
```



```
</div>

{% block body %}{% endblock%}

<div class="ui inverted vertical very padded segment">
|   Mugglecoding®
</div>

{% block foot %}{% endblock%}

</body>
</html>
```

Web_jd:

```
{% load staticfiles %}

<html>

<head>

<meta charset="utf-8">
<title>first web</title>
<link rel="stylesheet" href="{% static 'css/semantic.css'%}" media="screen" title="no title" charset="utf-8">
<link href="https://fonts.googleapis.com/css?family=Oswald|Raleway" rel="stylesheet">

<style type="text/css">
h1 {
    font-family: 'Oswald', sans-serif!important;
    font-size: 40px;
}

body {
    font-family: 'Raleway', sans-serif;
}
p {
    font-family: 'Raleway', sans-serif;
    font-size: 18px;
}

.ui.vertical.segment.masthead {
    height: 300px;
    background-image: url("{% static 'images/star_banner.jpg'%}");
    background-size: cover;
    background-position: 100% 80%;
}

.ui.container.segment {
    width: 800px;
}

.ui.center.aligned.header.blogslogon {
    margin-top: 40px;
}

.ui.center.aligned.header.blogslogon p {
    margin-top: 10px;
    color: white;
    font-size: 10px;
}

.ui.container.nav {
    width: 500px;
}
```

```

    }

</style>
{% block css %}{% endblock%}

</head>
<body>
    <div class="ui inverted vertical segment masthead">

        <h1 class="ui center aligned header blogslogon" style="font-size:50px;font-family: 'Raleway', sans-serif!important;">
            爬虫展示
        <p class="ui sub header">
            下面为爬取数据展示。
        </p>

    </h1>
</div>
<div class="ui container nav">
    <div class="ui borderless text three item menu ">
        <div class="ui simple dropdown item">
            加油
            <i class="dropdown icon"></i>
            <div class="menu">
                <a class="item" href="">简介</a>
                <a class="item" href="">简介</a>
            </div>
        </div>
        <a class="item">
            努力
        </a>
        <a class="item">
            简介
        </a>
    </div>
</div>

<div class="ui divider"></div>

{% block header %}{% endblock%}

<div class="ui container nav" >

    <br>

    <br>

```

```

<h1>京东手机</h1>

<br>

<br>

<h3><a href='http://127.0.0.1:8000/page/main'>回到页面</a></h3>

<br>

<br>

{% for content in content_list %}

    <h4>

        编号: {{content.mNo}}<br>

        图片: <br>

        <br>

        品牌: {{content.mMark}}<br>

        价格: {{content.mPrice}}<br>

        简介: {{content.mNote}}<br>

    </h4>

    <br><br><br>

{% endfor %}

{% if data %}
<ul id="pages" class="pagination pagination-sm pagination-xs">

    {% if data.first %}

        <li><a href="?page=1">1</a></li>

    {% endif %}

    {% if data.left %}

        {% if data.left has more %}

```

```

    {% if data.left_has_more %}

    <li><span>...</span></li>

    {% endif %}

    {% for i in data.left %}

    <li><a href="?page={{i}}">{{i}}</a></li>

    {% endfor %}

    {% endif %}

    <li class="active"><a href="?page={{data.page}}">{{data.page}}</a></li>

    {% if data.right %}

    {% for i in data.right %}

    <li><a href="?page={{i}}">{{i}}</a></li>

    {% endfor %}

    {% if data.right_has_more %}

    <li><span>...</span></li>

    {% endif %}

    {% endif %}

    {% if data.last %}

    <li><a href="?page={{data.total_pages}}">{{data.total_pages}}</a></li>

    {% endif %}

</ul>

{% endif %}

</div>

{% block body %}{% endblock%}

<div class="ui inverted vertical very padded segment">

</div>

{% block body %}{% endblock%}

<div class="ui inverted vertical very padded segment">
|  Mugglecoding®
</div>

{% block foot %}{% endblock%}

</body>
</html>

```

Web_tb:

```

!DOCTYPE html>
{% load staticfiles %}
<html>
  <head>
    <meta charset="utf-8">
    <title>first web</title>
    <link rel="stylesheet" href="{% static 'css/semantic.css'%}" media="screen" title="no title" charset="utf-8">
    <link href="https://fonts.googleapis.com/css?family=Oswald|Raleway" rel="stylesheet">

    <style type="text/css">
      h1 {
        font-family: 'Oswald', sans-serif!important;
        font-size: 40px;
      }

      body {
        font-family: 'Raleway', sans-serif;
      }

      p {
        font-family: 'Raleway', sans-serif;
        font-size: 18px;
      }

      .ui.vertical.segment.masthead {
        height: 300px;
        background-image: url("{% static 'images/star_banner.jpg'%}");
        background-size: cover;
        background-position: 100% 80%;
      }

      .ui.container.segment {
        width: 800px;
      }

      .ui.center.aligned.header.blogslogon {
        margin-top: 40px;
      }

      .ui.center.aligned.header.blogslogon p {
        margin-top: 10px;
        color: white;
        font-size: 10px;
      }

      .ui.container.nav {
        width: 500px;
      }
    </style>
  </head>
  <body>
    <div class="ui inverted vertical segment masthead">
      <h1 class="ui center aligned header blogslogon" style="font-size: 50px; font-family: 'Raleway', sans-serif!important;">
        爬虫展示
        <p class="ui sub header">
          下面为爬取数据展示。
        </p>
      </h1>
    </div>
    <div class="ui container nav">
      <div class="ui borderless text three item menu">
        <div class="ui simple dropdown item">
          加油
          <i class="dropdown icon"></i>
          <div class="menu">
            <a class="item" href="">简介</a>
            <a class="item" href="">简介</a>
          </div>
        </div>
        <a class="item">
          努力
        </a>
        <a class="item">
          简介
        </a>
      </div>
    </div>
    <div class="ui divider"></div>
  </body>
</html>

```

```

    <div class="ui container nav">
      <div class="ui borderless text three item menu">
        <div class="ui simple dropdown item">
          加油
          <i class="dropdown icon"></i>
          <div class="menu">
            <a class="item" href="">简介</a>
            <a class="item" href="">简介</a>
          </div>
        </div>
        <a class="item">
          努力
        </a>
        <a class="item">
          简介
        </a>
      </div>
    </div>
    <div class="ui divider"></div>
  </body>
</html>

```



```

<div class="ui container nav" >
  <br>
  <br>
  <h1>淘宝书包</h1>
  <br>
  <br>
  <h3><a href="http://127.0.0.1:8000/page/main">回到页面</a></h3>
  <br>
  <br>
  {% for shubao in content_list %}
    <h3>
      编号: {{shubao.id}}<br>
      价格: {{shubao.Price}}<br>
      简介: {{shubao.Title}}</h3> |
    <br><br>
  {% endfor %}
{% if data %}

<ul id="pages" class="pagination pagination-sm pagination-xs">

  {% if data.first %}
    <li><a href="?page=1">1</a></li>

  {% endif %}

  {% if data.left %}
    {% if data.left_has_more %}

    <li><span>...</span></li>

    {% endif %}

    {% for i in data.left %}

    <li><a href="?page={{i}}">{{i}}</a></li>

    {% endfor %}

    {% endif %}

    <li class="active"><a href="?page={{data.page}}">{{data.page}}</a></li>

```

```

    {% if data.left_has_more %}

    <li><span>...</span></li>

    {% endif %}

    {% for i in data.left %}

    <li><a href="?page={{i}}">{{i}}</a></li>

    {% endfor %}

    {% endif %}

    <li class="active"><a href="?page={{data.page}}">{{data.page}}</a></li>

    {% if data.right %}

    {% for i in data.right %}

    <li><a href="?page={{i}}">{{i}}</a></li>

    {% endfor %}

    {% if data.right_has_more %}

    <li><span>...</span></li>

    {% endif %}

    {% endif %}

    {% if data.last %}

    <li><a href="?page={{data.total_pages}}">{{data.total_pages}}</a></li>

    {% endif %}

  </ul>

  {% endif %}
</div>

{% block body %}{% endblock %}

<div class="ui inverted vertical very padded segment">

```

```

</div>

{% block body %}{% endblock%}

<div class="ui inverted vertical very padded segment">
  | Mugglecoding®
</div>

{% block foot %}{% endblock%}

</body>
</html>

```

数据库：

NO.	Content
1	肖申克的救赎
2	霸王别姬
3	这个杀手不太冷
4	阿甘正传
5	美丽人生
6	泰坦尼克号
7	千与千寻
8	辛德勒的名单
9	盗梦空间
10	忠犬八公的故事
11	机器人总动员
12	三傻大闹宝莱坞
13	海上钢琴师
14	放牛班的春天
15	楚门的世界
16	大话西游之大圣娶亲
17	星际穿越
18	龙猫
19	教父
20	熔炉
21	无间道
22	疯狂动物城
23	当幸福来敲门
24	怦然心动
25	触不可及
26	蝙蝠侠：黑暗骑士
27	乱世佳人
28	活着
29	少年派的奇幻漂流

ormation_schema
sql
rformance_schema
;
t
表
auth_group
auth_group_permissions
auth_permission
auth_user
auth_user_groups
auth_user_user_permissions
django_admin_log
django_content_type
django_migrations
django_session
firstapp_test
showapp_movie
showapp_phones
showapp_shubao
showapp_test1
showapp_weathers
视图
函数
事件
查询
报表
备份
lost_3306

wCity	wDate	wWeather	wTemp
上海	1日 (周三)	小雨转多云	23°C/17°C
上海	26日 (明天)	多云转晴	18°C/12°C
上海	27日 (后天)	多云	18°C/14°C
上海	28日 (周日)	小雨转多云	23°C/17°C
上海	29日 (周一)	中雨转阴	24°C/17°C
上海	30日 (周二)	阴转小雨	22°C/17°C
北京	1日 (周三)	多云	24°C/14°C
北京	26日 (明天)	晴转多云	21°C/10°C
北京	27日 (后天)	多云	18°C/7°C
北京	28日 (周日)	多云	21°C/9°C
北京	29日 (周一)	多云转小雨	24°C/13°C
北京	30日 (周二)	多云	26°C/14°C
广州	1日 (周三)	大到暴雨转多云	26°C/20°C
广州	26日 (明天)	中雨转中到大雨	28°C/24°C
广州	27日 (后天)	中到大雨转雷阵雨	28°C/24°C
广州	28日 (周日)	雷阵雨	29°C/25°C
广州	29日 (周一)	雷阵雨转中雨	30°C/25°C
广州	30日 (周二)	中雨转大到暴雨	30°C/22°C
深圳	1日 (周三)	暴雨转阵雨	27°C/23°C
深圳	26日 (明天)	中雨转大雨	29°C/24°C
深圳	27日 (后天)	大雨转雷阵雨	27°C/23°C
深圳	28日 (周日)	雷阵雨	28°C/24°C
深圳	29日 (周一)	雷阵雨	30°C/25°C
深圳	30日 (周二)	雷阵雨转暴雨	30°C/23°C

id	Price	Title
1	45.80	小学生书包男生1-3-4-6年级6-12周岁儿童
2	39.90	迪卡侬双肩包运动背包男女健身包书包儿童学生户外旅行包KIPSTA
3	119.00	kk树书包小学生女孩6-12周岁儿童1-3-6年级女童双肩背包护脊减负
4	499.00	Fjallraven/北极狐双肩包kanken classic书包女户外旅行背包23510
5	129.00	小米双肩包简约休闲多功能书包男女笔记本电脑包时尚潮流旅行背包
6	258.00	电视剧款JanSport旗舰店官网杰斯伯双肩包时尚女书包背包男大容量
7	348.00	爆款anello官方旗舰店日本ins潮风双肩女背包男离家出走书包
8	199.00	小米 米兔儿童书包 6-12岁男女小学生潮双肩背包幼儿园大容量背包
9	79.00	双肩包男士背包大容量旅行包电脑休闲女时尚潮流高中初中学生书包
10	109.00	七匹狼商务双肩包男书包中学生女电脑包旅行包休闲男士背包大容量
11	148.00	佑一良品男士背包双肩包男韩版大学生书包男时尚潮流大容量旅行包
▶ 12	69.00	巴布豆旗舰店书包1-3年级护脊减负儿童书包男4-6小学生书包轻便
13	299.00	BOPAI博牌电脑背包男户外旅行休闲双肩包商务书包出差多功能男包
14	49.00	小学生书包6-12周岁 女儿童双肩包 3-5年级女童背包 1-3年级女孩
15	45.80	儿童书包小学生男童1-3年级6-12周岁4-6年级男孩双肩背包轻便减负
16	59.80	商务背包男士双肩包韩版潮流旅行包休闲女学生书包简约时尚电脑包
17	168.00	双肩包男书包男士时尚潮流青年休闲简约潮牌旅行背包大学生电脑包
18	69.00	迪士尼书包小学生男女1-3-4-6年级米奇减负背包儿童书包8-10-12岁
19	119.00	巴朗商务双肩包休闲时尚潮流大学生书包15.6寸电脑包男士背包男潮
20	99.00	米熙休闲运动背包双肩包女书包中学生男韩版时尚大容量旅游旅行包
21	195.02	国家地理背包女运动户外时尚双肩包男牛津布旅行防水学生情侣书包
22	79.00	双肩包女士2019新款小韩版百搭时尚书包2019旅行防盗牛津帆布背包
23	89.00	迪士尼书包小学生男童1-3-4-5年级6-12周岁女孩儿童减负双肩背包
24	69.00	双肩包女2019新款潮牌韩版时尚百搭女士休闲pu软皮小背包旅行书包
25	138.00	瑞士军刀双肩包男 背包休闲商务旅行大容量瑞士书包电脑男士户外
26	139.00	花花公子男士双肩包时尚潮流休闲电脑旅行书包潮牌大学生帆布背包
27	89.00	HK双肩包男简约个性书包韩版时尚潮流休闲电脑包户外旅行轻便背包
28	49.00	背包男双肩包休闲大容量旅行包时尚潮流韩版高中生初中学生书包男
29	69.00	休闲双肩包男士韩版简约电脑旅行背包女时尚潮流初中高中学生书包

id	mNo	mMark	mPrice	mNote	mFile
1	000001	OPPO	3599.00	OPPO Reno 全面屏拍照手机	000001.jpg
2	000002	Apple	5899.00	Apple iPhone XR (A2108)	000002.jpg
3	000003	【KPL官方比赛用机】vivo	3298.00	【KPL官方比赛用机】vivo i	000003.jpg
4	000004	荣耀8X	1299.00	荣耀8X 千元屏霸 91%屏占比	000004.jpg
5	000005	荣耀10青春版	1299.00	荣耀10青春版 幻彩渐变 24C	000005.jpg
6	000006	vivo	799.00	vivo U1 水滴全面屏 AI智慧	000006.jpg
7	000007	荣耀V20	2799.00	荣耀V20 胡歌同款 麒麟980	000007.jpg
8	000008	vivo	3598.00	vivo X27 8GB+256GB大内存	000008.jpg
9	000009	OPPO	2999.00	OPPO Reno手机 新品 全面	000009.jpg
10	000010	小米	1199.00	小米 红米Redmi Note7 幻彩	000010.jpg
11	000011	荣耀畅玩8C两天一充	899.00	荣耀畅玩8C两天一充 莱茵护	000011.jpg
12	000012	小米	799.00	小米 红米6 4GB+64GB 铂晶	000012.jpg
13	000013	黑鲨游戏手机2	3499.00	黑鲨游戏手机2 8GB+128GB	000013.jpg
14	000014	Apple	6199.00	Apple iPhone X (A1865) 6	000014.jpg
15	000015	小米	799.00	小米 红米Redmi 7 幻彩渐变	000015.jpg
16	000016	小米8SE	1599.00	小米8SE 全面屏智能游戏拍	000016.jpg
17	000017	三星	6999.00	三星 Galaxy S10+ 8GB+12	000017.jpg
18	000018	Apple	9699.00	Apple iPhone XS Max (A2	000018.jpg
19	000019	Apple	3799.00	Apple iPhone 8 (A1863) 6	000019.jpg
20	000020	vivo	1598.00	vivo Z3 6GB+64GB 极光蓝	000020.jpg
21	000021	华为	3999.00	华为 HUAWEI Mate 20 麒麟	000021.jpg
22	000022	vivo	2298.00	vivo S1 6GB+128GB 宠爱精	000022.jpg
23	000023	Apple	4699.00	Apple iPhone 8 Plus (A18	000023.jpg
24	000024	小米8青春版	1499.00	小米8青春版 镜面渐变AI双摄	000024.jpg
25	000025	三星	1448.00	三星 Galaxy A6s 6GB+64G	000025.jpg
26	000026	荣耀10	2199.00	荣耀10 GT游戏加速 AIS手机	000026.jpg
27	000027	诺基亚	1099.00	诺基亚 NOKIA X6 6GB+64	000027.jpg
28	000028	华为	1999.00	华为 HUAWEI nova 4e 32C	000028.jpg

补充:

多进程:

```
#!/usr/bin/env python
# -*- coding=utf-8 -*-

from multiprocessing import Process, Queue

import time
from lxml import etree
import requests

class DouBanSpider(Process):
    def __init__(self, url, q):
        # 重写父类的__init__方法
        super(DouBanSpider, self).__init__()
        self.url = url
        self.q = q
        self.headers = {
            'Host': 'movie.douban.com',
            'Referer': 'https://movie.douban.com/top250?start=225&filter=',
            'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/59.0.3071.104 Safari/537.36',
        }

    def run(self):
        self.parse_page()

    def send_request(self, url):
        """
        用来发送请求的方法
        :return: 返回网页源码
        """
        # 请求出错时，重复请求3次，
        i = 0
        while i <= 3:
            try:
                print (u"[INFO]请求url:" + url)
                return requests.get(url=url, headers=self.headers).content
            except Exception as e:
                print (u'[INFO] %s%s' % (e, url))
                i += 1

    def parse_page(self):
        """
        解析网站源码，并采用xpath提取 电影名称和评分放到队列中
        :return:
        """
        response = self.send_request(self.url)
```

```
        response = self.send_request(self.url)
        html = etree.HTML(response)
        # 获取到一页的电影数据
        node_list = html.xpath("//div[@class='info']")
        for move in node_list:
            # 电影名称
            title = move.xpath('.//a/span/text()')[0]
            # 评分
            score = move.xpath('.//div[@class="bd"]//span[@class="rating_num"]/text()')[0]

            # 将每一部电影的名称跟评分加入到队列
            self.q.put(score + "\t" + title)

def main():
    # 创建一个队列用来保存进程获取到的数据
    q = Queue()
    base_url = 'https://movie.douban.com/top250?start='
    # 构造所有url
    url_list = [base_url + str(num) for num in range(0, 225 + 1, 25)]

    # 保存进程
    Process_list = []
    # 创建并启动进程
    for url in url_list:
        p = DouBanSpider(url, q)
        p.start()
        Process_list.append(p)

    # 让主进程等待子进程执行完成
    for i in Process_list:
        i.join()

    while not q.empty():
        print (q.get())

if __name__ == "__main__":
    start = time.time()
    main()
    print ('[info]耗时: %s'%(time.time()-start))
```

多线程:

```
#!/usr/bin/env python
# -*- coding=utf-8 -*-

from threading import Thread
from queue import Queue
import time
from lxml import etree
import requests

class DoubanSpider(Thread):
    def __init__(self, url, q):
        super(DoubanSpider, self).__init__()
        self.url = url
        self.q = q
        self.headers = {
            'Cookie': 'll="118282"; bid=ctyiEarSLfw; ps=y; __yadk_uid=05r85yZ9d4bEeLKhv4w36950F0PoedzC; dbc12="155150959:0Eu4dds1G1o"; as="https://',
            'Host': 'movie.douban.com',
            'Referer': 'https://movie.douban.com/top250?start=225&filter=',
            'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/59.0.3071.104 Safari/537.36',
        }

    def run(self):
        self.parse_page()

    def send_request(self, url):
        """
        用来发送请求的方法
        :return: 返回网页源码
        """
        # 请求出错时，重复请求3次，
        i = 0
        while i <= 3:
            try:
                print (u"[INFO]请求url:" + url)
                html = requests.get(url=url, headers=self.headers).content
            except Exception as e:
                print (u'[INFO] %s' % (e, url))
                i += 1
            else:
                return html

    def parse_page(self):
        """
        解析网站源码，并采用xpath提取 电影名称和评分放到队列中
        :return:
        """
        response = self.send_request(self.url)
```

```
        response = self.send_request(self.url)
        html = etree.HTML(response)
        # 获取到一页的电影数据
        node_list = html.xpath("//div[@class='info']")
        for move in node_list:
            # 电影名称
            title = move.xpath('.//a/span/text()')[0]
            # 评分
            score = move.xpath('.//div[@class="bd"]//span[@class="rating_num"]/text()')[0]

            # 将每一部电影的名称跟评分加入到队列
            self.q.put(score + "\t" + title)

def main():
    # 创建一个队列用来保存进程获取到的数据
    q = Queue()
    base_url = 'https://movie.douban.com/top250?start='
    # 构造所有url
    url_list = [base_url + str(num) for num in range(0, 225 + 1, 25)]

    # 保存线程
    Thread_list = []
    # 创建并启动线程
    for url in url_list:
        p = DoubanSpider(url, q)
        p.start()
        Thread_list.append(p)

    # 让主线程等待子线程执行完成
    for i in Thread_list:
        i.join()

    while not q.empty():
        print (q.get())

if __name__ == "__main__":
    start = time.time()
    main()
    print ('[info]耗时: %s' % (time.time() - start))
```

总结：

- 1.正确爬取爬虫代码，并将其存储在数据库中。
- 2.使用 `django`，设计首页，包括样式，链接跳转，子页面。
- 3.子页面中`{{classname: idname}}`显示数据库调取的信息。
- 4.其中，我使用了 `semantic ui` 框架来作为我页面的美化。
- 5.但是在分页面中我是用`{% extends 'xxx.html'%}`来继承主页面的样式时由于版本的原因，不能经行覆盖。
- 6.代码上传，编写报告，完成。

