# 温州大学瓯江学院

## WENZHOU UNIVERSITY OUJIANG COLLEGE

# 《爬虫综合实验》

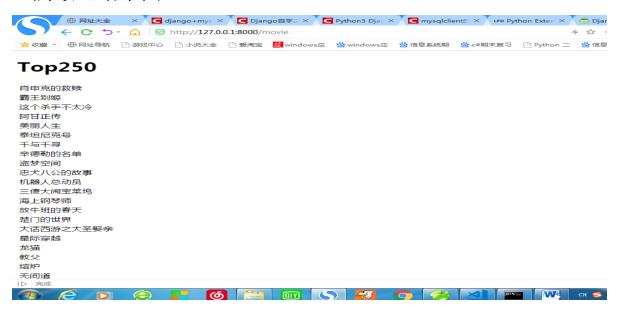题　　目：　　　爬虫实验合集　　　

分　　院：　　理工分院　　　　　　

班　　级：　　16 计算机科学与技术本 3

姓　　名：　　王家跃　　　　　　　

学　　号：　　16219111301　　　　

完成日期：　　2019 年 6 月 18 日　

温州大学瓯江学院教务部

二〇一二年十一月制

# 一、期中作业

## 豆瓣最终结果图：



## 爬取京东手机 Django 截图：

（只爬取 2 页的数据 ）



Django　（爬取豆瓣 TOP250 和京东手机的价格）

1. 安装 django　pip install Django
2. 创建第一个项目：django-admin　startproject　JDdjango
3. 启动服务器：python manage.py runserver 0.0.0.0:8000
4. 新建 view.py
5. 数据库配置：

```
DATABASES = {
    'default': {
        'ENGINE': 'django.db.backends.mysql',
        'NAME': 'mypachou', #数据库名字
        'USER': 'root1',
        'PASSWORD': 'wangjiayue',
        'HOST':'localhost',
        'PORT':'3306',
    }
}
```

6. 创建一个 TestModel 的 APP: Django-admin.py startapp TestModel

7. TestModel/model.py 代码：建立一个 movie 表和一个 test 表（京东手机）



8. 在 setting.py 中修改代码：

```
INSTALLED_APPS = [
    'django.contrib.admin',
    'django.contrib.auth',
    'django.contrib.contenttypes',
    'django.contrib.sessions',
    'django.contrib.messages',
    'django.contrib.staticfiles',
    'TestModel', #添加此项
]
```
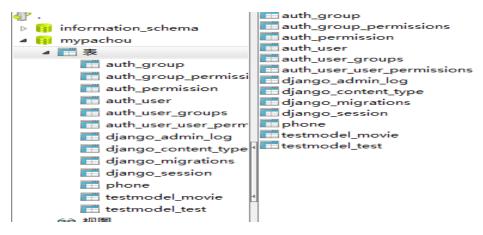
9. Cmd 运行：

Python manage.py migrate

Python manage.py makemigrations TestModel

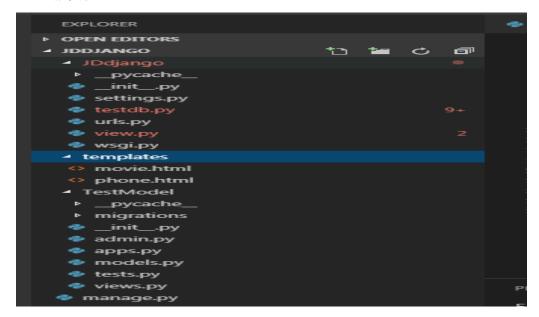Python manage.py migrate TestModel

三句运行成功后在 test 数据库下会生成 testmodel_test ,testmodel_movie 表
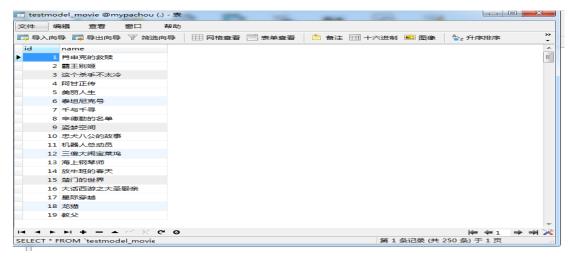
10. 在添加 testdb.py 代码来获取爬虫写入数据库的数据：

代码见源码

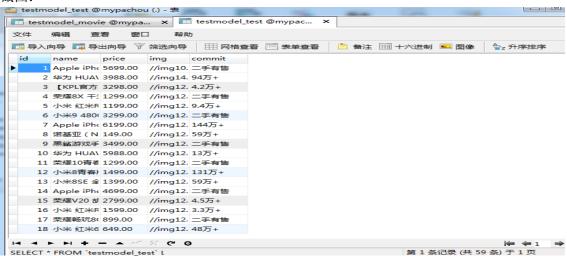11. 目录：（导入 models 中的 Test 和 Movie 表，一直报错，不存在 objects 成员，没解决但是能实现）



在 python 分别运行爬取豆瓣和爬取京东手机的代码，将爬取的数据插入表中：

代码见源码

结果截图：

（京东手机只爬取了两页内容，代码以实现全部爬取：）

截图：



新建 templates 文件夹并在文件夹下新建 phone.html movie.html:

```html
<!DOCTYPE html>
<html lang="en">
<head>
    <meta charset="UTF-8">
    <title>Title</title>
</head>
<body>
    <h1>Top250</h1>
    <table>
    {% for line in movie_list %}
    <tr>
        <td>{{line.name}}</td>

    </tr>
    {% endfor %}
    </table>
</body>
</html>
```

```html
<!DOCTYPE html>
<html>
<head>
    <meta charset="UTF-8">
    <title>Title</title>
</head>
<body>
    <h1>京东手机</h1>
    <table>
    {% for line in phone_list %}
    <tr>
        <td>{{line.name}}</td>
        <td>{{line.price}}</td>

    </tr>
    {% endfor %}
    </table>
</body>
</html>
```

修改 View.py：

```
from django.shortcuts import render
import pymysql
from TestModel import models


def hello(request):
    phone_list = models.Test.objects.all()
    return render(request, 'phone.html', {"phone_list":phone_list})


def movie(request):
    movie_list = models.Movie.objects.all()
    return render(request, 'movie.html', {"movie_list":movie_list})
```

以及绑定到 urls.py:

```
from JDdjango import testdb
from django.conf.urls import url
from JDdjango import view

urlpatterns = [
    url(r'^phone$', view.hello),
    url(r'^movie$', view.movie),
    url(r'^testdb$', testdb.testdb),
    url(r'^testdb$', testdb.testmovie),
]
```

# 二、期末作业

## 12306 自动登陆代码（基本实现全自动化）：

```
import re
import time
import base64
import requests
import sys
from selenium import webdriver
from selenium.webdriver.support.ui import WebDriverWait
import urllib.request
from selenium.webdriver.support import expected_conditions as EC
from selenium.webdriver.common.by import By
from selenium.webdriver.common.keys import Keys
from selenium.webdriver.common.action_chains import ActionChains

class Login(object):
    def __init__(self):
        self.login_url = "https://kyfw.12306.cn/otn/resources/login.html"
        self.geturl = "http://littlebigluo.qicp.net:47720/"
        self.totalFlush = 0
        self.startTime = time.time()
```

```python
        # self.driver = '' #驱动 chrome 浏览器进行操作
        driver = webdriver.Chrome()
        self.driver = driver #驱动 chrome 浏览器进行操作


    def login_input(self):
        self.driver.get(self.login_url)
        time.sleep(0.2)
        account = self.driver.find_element_by_class_name("login-hd-account")
        account.click()
        userName = self.driver.find_element_by_id("J-userName")
        userName.send_keys("15606879517")    # 12306 账号
        password = self.driver.find_element_by_id("J-password")
        password.send_keys("64890110wjy")    # 12306 密码
        time.sleep(5)


    def getImage(self):
        img_element                                            =WebDriverWait(self.driver,
100).until(EC.presence_of_element_located((By.ID, "J-loginImg")))
        try:
            img_element
        except Exception as e:
            print("出错了，请稍后再试")
        base64_str=img_element.get_attribute("src").split(",")[-1]
        imgdata=base64.b64decode(base64_str)
        with open('mylogin.jpg','wb') as file:
            file.write(imgdata)
        self.img_element=img_element
        time.sleep(2)


    def getResult(self):
        driver = webdriver.Chrome()
        driver.get(self.geturl)
        time.sleep(5)
        files ="C:\\Users\\Administrator\\Desktop\\爬虫期中作业\\mylogin.jpg"
        onput = driver.find_element_by_name("pic_xxfile")
        onput.send_keys(files)
        myclick = driver.find_element_by_xpath("/html/body/form/input[2]").click()
        result = driver.find_element_by_xpath("/html/body/p[1]/font/font")
        result=result.text.split(" ")
        print(result)
        self.result=result
```

```python
            time.sleep(2)
            driver.quit()



    def Click(self):
        js=""
        for i in self.result:
            if i=='1':
                js=js+'<div randcode="42,50" class="lgcode-active" style="top: 66px; left:
29px;"></div>'
            elif i=='2':
                js=js+'<div randcode="109,43" class="lgcode-active" style="top: 59px; left:
96px;"></div>'
            elif i=='3':
                js=js+'<div randcode="186,36" class="lgcode-active" style="top: 52px; left:
178px;"></div>'
            elif i=='4':
                js=js+'<div randcode="256,50" class="lgcode-active" style="top: 62px; left:
241px;"></div>'
            elif i=='5':
                js=js+'<div randcode="31,124" class="lgcode-active" style="top: 140px; left:
18px"></div>'
            elif i=='6':
                js=js+'<div randcode="119,110" class="lgcode-active" style="top: 130px; left:
103px;"></div>'
            elif i=='7':
                js=js+'<div randcode="182,110" class="lgcode-active" style="top: 124px; left:
169px;"></div>'
            else :
                js=js+'<div randcode="258,110" class="lgcode-active" style="top: 128px; left:
245px;"></div>'
        resultjs='document.getElementById("J-passCodeCoin").innerHTML='+'"'+js+'"'
        print(resultjs)
        self.driver.execute_script(resultjs)
        time.sleep(5)



    def submit(self):
        self.driver.find_element_by_id("J-login").click()
```

```
if __name__ == '__main__':
    spider = Login()
    spider.login_input()
    spider.getImage()
    spider.getResult()
    spider.Click()
    spider.submit()
```

## 成功截图：