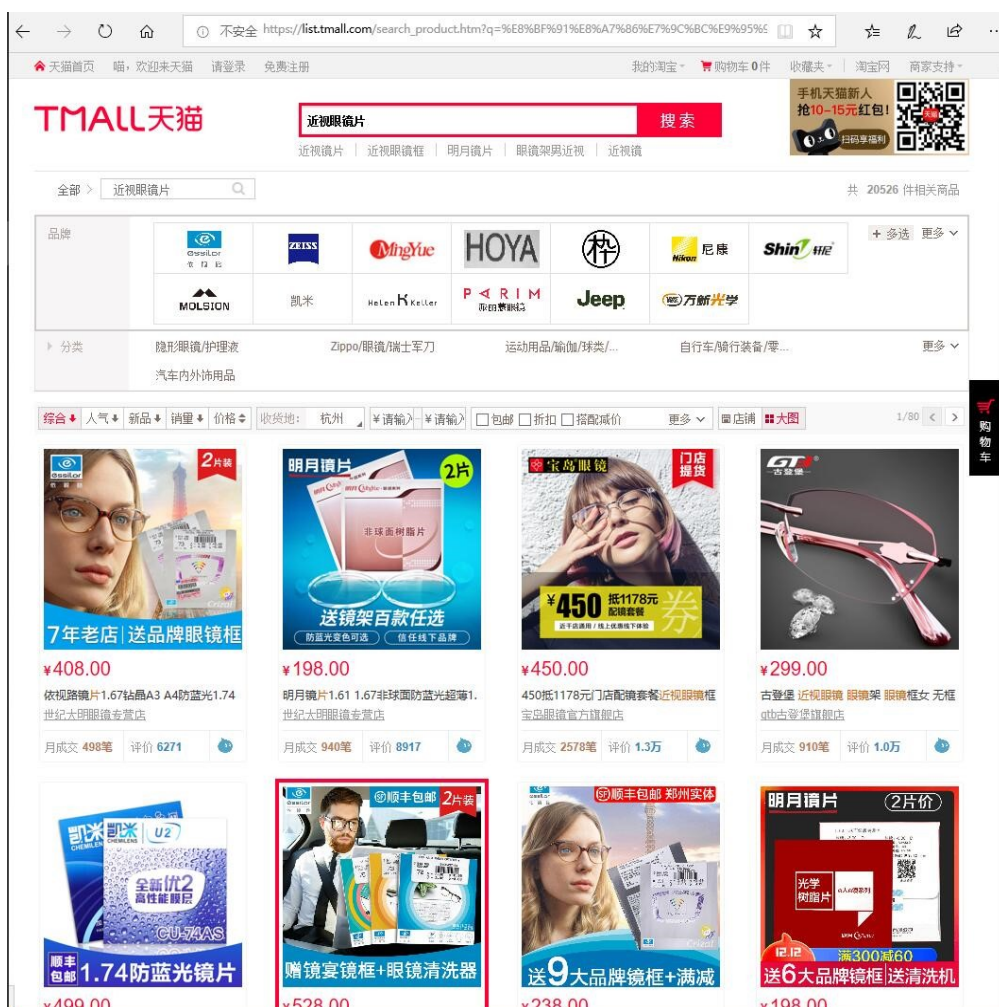


使用selenium爬取天猫近视眼镜片评论和销售情况并分析

内容简介：先进行数据采集，获取商品数据后再根据商品url拼接评论内容url，通过评论url获得评论数据，最后进行数据分析

1. 数据收集

采用selenium



得到的商品url

[https://list.tmall.com/search_product.htm?
q=%BD%FC%CA%D3%D1%DB%BE%B5%C6%AC&type=p&spm=a220m.1000858.a2227oh.d100&from=.list.pc_1_searchbutton](https://list.tmall.com/search_product.htm?q=%BD%FC%CA%D3%D1%DB%BE%B5%C6%AC&type=p&spm=a220m.1000858.a2227oh.d100&from=.list.pc_1_searchbutton)
([https://list.tmall.com/search_product.htm?
q=%BD%FC%CA%D3%D1%DB%BE%B5%C6%AC&type=p&spm=a220m.1000858.a2227oh.d100&from=.list.pc_1_searchbutton](https://list.tmall.com/search_product.htm?q=%BD%FC%CA%D3%D1%DB%BE%B5%C6%AC&type=p&spm=a220m.1000858.a2227oh.d100&from=.list.pc_1_searchbutton))

换页条



点击第三页后的商品url

[https://list.tmall.com/search_product.htm?
spm=a220m.1000858.0.0.3888589de9xg9o&s=120&q=%BD%FC%CA%D3%D1%DB%BE%B5%C6%AC&sort=s&style=g&from=mallfp..pc_1_searchbutton
&type=pc#J_Filter](https://list.tmall.com/search_product.htm?spm=a220m.1000858.0.0.3888589de9xg9o&s=120&q=%BD%FC%CA%D3%D1%DB%BE%B5%C6%AC&sort=s&style=g&from=mallfp..pc_1_searchbutton&type=pc#J_Filter) ([https://list.tmall.com/search_product.htm?
spm=a220m.1000858.0.0.3888589de9xg9o&s=120&q=%BD%FC%CA%D3%D1%DB%BE%B5%C6%AC&sort=s&style=g&from=mallfp..pc_1_searchbutton
&type=pc#J_Filter](https://list.tmall.com/search_product.htm?spm=a220m.1000858.0.0.3888589de9xg9o&s=120&q=%BD%FC%CA%D3%D1%DB%BE%B5%C6%AC&sort=s&style=g&from=mallfp..pc_1_searchbutton&type=pc#J_Filter))

分析：相比第一页多了一个s参数，s=120，接着跳转到80页，得到的url中参数s=4740

所以可以通过s参数换页获取到不同商品信息， $s = (\text{页数} - 1) \times 60$

```
for page in range(1, 78):
    s = (page - 1)*60
    search_url = "https://list.tmall.com/search_product.htm" + "?"
    cat=50041298&s=" + str(s)\
    +
    "&q=%BD%FC%CA%D3%D1%DB%BE%B5%C6%AC&sort=s&style=g&from=mallfp..pc_1_s
    earchbutton&active=2" \

    "&industryCatId=50041298&spm=a220m.1000858.0.0.2f4e58981G3iTJ&type=pc
    #J_Filter "
    cat_spider = CatSpider(url=search_url)
    cat_spider.get_goods_info(page=page)
    print(search_url)
    sleep(6 + randint(1, 4))
```

为了防止触发天猫的反爬虫，加了一个随机等待时间

在获取到商品数据后，根据商品地址拼接商品评论url

商品url如下

<http://detail.tmall.com/item.htm?>

[id=565334203629&skuld=4063956831708&areald=430100&user_id=3355869104&cat_id=50041298&is_b=1&rn=1f4a2c261ac23e5987794fcc03e7a9c7](http://detail.tmall.com/item.htm?id=565334203629&skuld=4063956831708&areald=430100&user_id=3355869104&cat_id=50041298&is_b=1&rn=1f4a2c261ac23e5987794fcc03e7a9c7) (http://detail.tmall.com/item.htm?id=565334203629&skuld=4063956831708&areald=430100&user_id=3355869104&cat_id=50041298&is_b=1&rn=1f4a2c261ac23e5987794fcc03e7a9c7)

[id=565334203629&skuld=4063956831708&areald=430100&user_id=3355869104&cat_id=50041298&is_b=1&rn=1f4a2c261ac23e5987794fcc03e7a9c7](http://detail.tmall.com/item.htm?id=565334203629&skuld=4063956831708&areald=430100&user_id=3355869104&cat_id=50041298&is_b=1&rn=1f4a2c261ac23e5987794fcc03e7a9c7))

分析

id: 商品id

user_id: 卖家id

当点击商品评论的下一页时，商品的url并没有变化，说明评论数据的更新是通过Ajax获取

利用开发者工具查看Ajax请求的地址

名称	协议	方法	结果	标头	正文	参数	Cookie	计时
taobao.live.video.pause?gmkey=CLK&gokey=app...	HTTP/2	GET	200	请求 URL: https://rate.tmall.com/list_detail_rate.htm?itemid=39310588779&spuld=0&sellerid=775323974&order=3¤tPage=1&append=1&content=1&tagid=&posi=&picture=0&groupid=&ua=098%23E1h				
tmallrate.6.2.7?logtype=2&&_tm_cache=15454724...	HTTP/2	GET	200	请求方法: GET				
list_detail_rate.htm?itemid=39310588779&spuld=0...	HTTP/2	GET	200	状态代码: 200 /				
punish?5secdata=5e0de1365474455070961b803...	HTTP/2	GET	200	请求标头				
TB10tqLMMPMelJy1XbXcwwVXa-694-685.png	HTTPS	GET	200	Accept: */*				
flexible.js	HTTPS	GET	200	Accept-Encoding: gzip, deflate, br				

评论请求的url

标头	正文	参数	Cookie	计时
请求 URL: https://rate.tmall.com/list_detail_rate.htm?itemid=39310588779&spuld=0&sellerid=775323974&order=3¤tPage=1&append=1&content=1&tagid=&posi=&picture=0&groupid=&ua=098%23E1hTv9LvQvUvCkvvvvgjPR25pgDn2FWQJ02Pm...				

根据网上的资料分析：

评论数据url [https://rate.tmall.com/list_detail_rate.htm?](https://rate.tmall.com/list_detail_rate.htm?itemId=39310588779&sellerId=775323974¤tPage=1&callback=jsonp1210)

[itemId=39310588779&sellerId=775323974¤tPage=1&callback=jsonp1210](https://rate.tmall.com/list_detail_rate.htm?itemId=39310588779&sellerId=775323974¤tPage=1&callback=jsonp1210) ([https://rate.tmall.com/list_detail_rate.htm?](https://rate.tmall.com/list_detail_rate.htm?itemId=39310588779&sellerId=775323974¤tPage=1&callback=jsonp1210)

[itemId=39310588779&sellerId=775323974¤tPage=1&callback=jsonp1210](https://rate.tmall.com/list_detail_rate.htm?itemId=39310588779&sellerId=775323974¤tPage=1&callback=jsonp1210))

其中

itemId: 商品id 商品url 的 id

sellerId: 卖家id 商品url中的user_id

currentPage: 页码

callback:

作为回调函数的一部分，返回的是json格式数据，“callback=jsonp”这部分是固定不变的，后面的数字利用random函数生成一个随机数拼接上去

利用parse解析出商品url中的参数

```
url_param = good_info.split('?')
res = parse.parse_qs(url_param[1])
url_str = 'https://rate.tmall.com/list_detail_rate.htm?itemId=' +
''.join(res['id']) + '&sellerId=' \
        + ''.join(res['user_id']) + '&currentPage=1&callback=jsonp'
+ str(randint(1000, 20000))
urls_str.append(url_str)
```

访问100多个页面后，会遇到访问验证，此时服务器返回的json是固定的，包含极验验证的地址，获取地址后，跳转到该地址，获得滑块对象，通过开发者工具可以很容易看到滑动条长度为300px，接着杜撰一个滑动轨迹，（这里我是参考了《python网络爬虫实战》里的模拟验证，但其实没什么用，最后还是自己根据我验证通过的经验，模拟了一个轨迹），开始验证直至验证通过(最开始成功过几次，后来再也没成功过)。

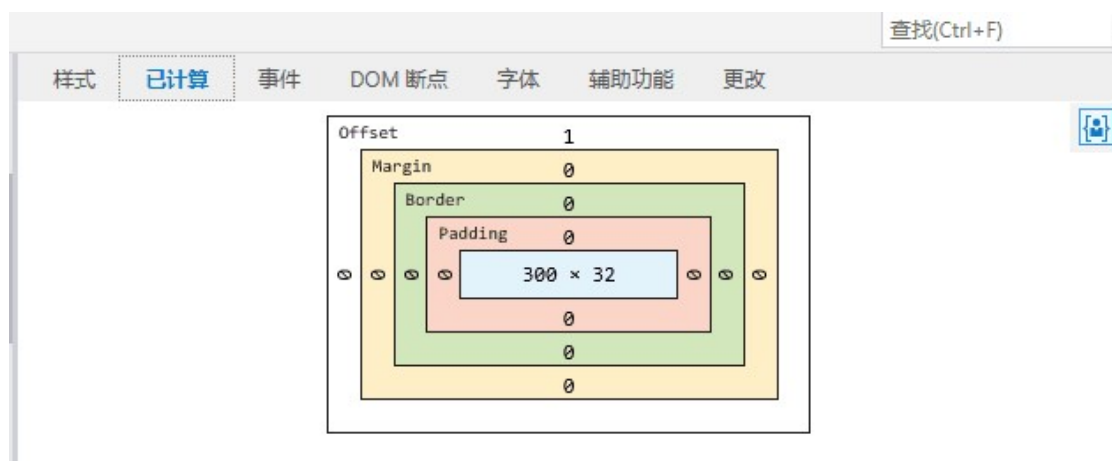


休息会呗，坐下来喝口水，
我们马上回来。

为保证您的正常访问请进行验证

>>

请按住滑块，拖动到最右边



遇到的问题：

多次访问后，（大概访问100多个页面）就会触发阿里的反爬虫，第一次需要登录后，再滑动验证，天猫和淘宝采用的都是第三代极验验证，号称采用机器学习识别轨迹，实测真的很坑，在前几次验证还可以通过，验证多次后感觉难度越来越大，基本不可能通过，就算是人工验证也不鸟你。只有等过几个小时，才会解封。

2. 数据处理

爬取到的商品数据

X14																	
	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q
1	店名	付款人数	标价	评论数	商品地址												
2	至美上品鞋 558双	499	60	6298	\\detail.tmall.com/item.htm?d=39310588779&skuld=3642495193234&areald=430100&user_id=775323974&cat_id=50041298&is_b=1&m=14a2c261												
3	明秀高鞋 635双	398	60	7323	\\detail.tmall.com/item.htm?d=25687452603&skuld=79384653841&areald=430100&user_id=165255874&cat_id=50041298&is_b=1&m=14a2c261												
4	至美上品鞋 74双	368	60	922	\\detail.tmall.com/item.htm?d=534620921543&skuld=3469740570161&areald=430100&user_id=775323974&cat_id=50041298&is_b=1&m=14a2c261												
5	明秀高鞋 429双	198	60	723	\\detail.tmall.com/item.htm?d=5622734225698&skuld=372949090080&areald=430100&user_id=165255874&cat_id=50041298&is_b=1&m=14a2c261												
6	世纪大明鞋 505双	498	60	6307	\\detail.tmall.com/item.htm?d=5227273232098&skuld=394967379593&areald=430100&user_id=796527718&cat_id=50041298&is_b=1&m=14a2c261												
7	宝岛眼镜 356双	290	1513		\\detail.tmall.com/item.htm?d=4268772334&skuld=46116860611115123&areald=430100&user_id=746173362&cat_id=50041298&is_b=1&m=14a2c261												
8	世纪大明鞋 418双	772	2710		\\detail.tmall.com/item.htm?d=41204976563&skuld=392841187457&areald=430100&user_id=796527718&cat_id=50041298&is_b=1&m=14a2c261												
9	世纪精典鞋 683双	203	3123		\\detail.tmall.com/item.htm?d=54391462881&skuld=3504056681218&areald=430100&user_id=3021903329&cat_id=50041298&is_b=1&m=14a2c261												
10	博士眼镜 406双	498	602		\\detail.tmall.com/item.htm?d=524117067957&skuld=4005772738051&areald=430100&user_id=1148656175&cat_id=50041298&is_b=1&m=14a2c261												
11	世纪精典鞋 539双	204	1529		\\detail.tmall.com/item.htm?d=54395086089&skuld=401736408750&areald=430100&user_id=3021903329&cat_id=50041298&is_b=1&m=14a2c261												
12	宝岛眼镜 62双	598	338		\\detail.tmall.com/item.htm?d=56288502606&skuld=381071917086&areald=430100&user_id=746173362&cat_id=50041298&is_b=1&m=14a2c261												
13	大明眼镜 227双	272	1648		\\detail.tmall.com/item.htm?d=42743672391&skuld=3917246887493&areald=430100&user_id=2300311521&cat_id=50041298&is_b=1&m=14a2c261												
14	可得光 137双	616	343		\\detail.tmall.com/item.htm?d=565348476757&skuld=35826367971182&areald=430100&user_id=3170525383&cat_id=50041298&is_b=1&m=14a2c261												
15	卫光眼镜 292双	202	402		\\detail.tmall.com/item.htm?d=56379594283&skuld=4034511447749&areald=430100&user_id=3355869104&cat_id=50041298&is_b=1&m=14a2c261												
16	卫光眼镜 68双	499	211		\\detail.tmall.com/item.htm?d=5653344036298&skuld=4063956831708&areald=430100&user_id=3355869104&cat_id=50041298&is_b=1&m=14a2c261												
17	大明眼镜 265双	198	1656		\\detail.tmall.com/item.htm?d=528876473966&skuld=3976753211941&areald=430100&user_id=2300311521&cat_id=50041298&is_b=1&m=14a2c261												
18	康乐眼镜 342双	198	3535		\\detail.tmall.com/item.htm?d=528992051223&skuld=332501260205&areald=430100&user_id=268627496&cat_id=50041298&is_b=1&m=14a2c261												
19	pulais眼镜 1416	220	25000		\\detail.tmall.com/item.htm?d=35277146667&skuld=318776783198&areald=430100&user_id=594992150&cat_id=50041298&is_b=1&m=14a2c261												
20	大明眼镜 68双	138	678		\\detail.tmall.com/item.htm?d=44262761793&skuld=375034861261&areald=430100&user_id=82881337&cat_id=50041298&is_b=1&m=14a2c261												
21	博士眼镜 40双	268	51		\\detail.tmall.com/item.htm?d=56873114558&skuld=38488984143&areald=430100&user_id=1148656175&cat_id=50041298&is_b=1&m=14a2c261												
22	卫光眼镜 802双	190	8739		\\detail.tmall.com/item.htm?d=526692282851&skuld=324453324539&areald=430100&user_id=594992150&cat_id=50041298&is_b=1&m=14a2c261												
23	宝岛眼镜 82双	398	4985		\\detail.tmall.com/item.htm?d=56067220259&skuld=3509087908260&areald=430100&user_id=55558604&cat_id=50041298&is_b=1&m=14a2c261												
24	木九十方 144双	599	101		\\detail.tmall.com/item.htm?d=56938116638&skuld=365813671720&areald=430100&user_id=278295698&cat_id=50041298&is_b=1&m=14a2c261												
25	木九十方 191双	598	132		\\detail.tmall.com/item.htm?d=55419311656&skuld=35046528825&areald=430100&user_id=278295698&cat_id=50041298&is_b=1&m=14a2c261												
26	雷蒙方 73双	699	100		\\detail.tmall.com/item.htm?d=56278641156&skuld=3763589312703&areald=430100&user_id=666872058&cat_id=50041298&is_b=1&m=14a2c261												
27	雷蒙方 93双	300	1188		\\detail.tmall.com/item.htm?d=38930405282&skuld=3521266637074&areald=430100&user_id=2069747240&cat_id=50041298&is_b=1&m=14a2c261												
28	loh眼镜 113双	800	509		\\detail.tmall.com/item.htm?d=559821894623&skuld=350716596019&areald=430100&user_id=2370381708&cat_id=50041298&is_b=1&m=14a2c261												
29	loh眼镜 151双	300	3419		\\detail.tmall.com/item.htm?d=4387640183&skuld=368260513639&areald=430100&user_id=2370381708&cat_id=50041298&is_b=1&m=14a2c261												
30	雷蒙方 63双	1105	40		\\detail.tmall.com/item.htm?d=5696436496298&skuld=391878006866&areald=430100&user_id=666872058&cat_id=50041298&is_b=1&m=14a2c261												
31	美杜古旗 153双	598	39		\\detail.tmall.com/item.htm?d=574750433271&skuld=389484738465&areald=430100&user_id=1710682596&cat_id=50041298&is_b=1&m=14a2c261												
32	百视达眼镜 7双	398	148		\\detail.tmall.com/item.htm?d=4289591860&skuld=91092612619&areald=430100&user_id=227115016&cat_id=50041298&is_b=1&m=14a2c261												
33	JINS眼镜 128双	500	636		\\detail.tmall.com/item.htm?d=54116815718&skuld=46116855959541622&areald=430100&user_id=1658173778&cat_id=50041298&is_b=1&m=14a2c261												
34	视光眼镜 69双	368	3449		\\detail.tmall.com/item.htm?d=53389253122&skuld=3767817839110&areald=430100&user_id=552201090&cat_id=50041298&is_b=1&m=14a2c261												
35	视光眼镜 67双	178	598		\\detail.tmall.com/item.htm?d=56267445924&skuld=371696592719&areald=430100&user_id=55558604&cat_id=50041298&is_b=1&m=14a2c261												
36	亿超眼镜 55双	138	371		\\detail.tmall.com/item.htm?d=57393255959&skuld=39157830312&areald=430100&user_id=190884141&cat_id=50041298&is_b=1&m=14a2c261												
37	JINS眼镜 21双	300	1929		\\detail.tmall.com/item.htm?d=54116788738&skuld=461168559604275642&areald=430100&user_id=1658173778&cat_id=50041298&is_b=1&m=14a2c261												
38	哆啦眼镜 33双	168	1408		\\detail.tmall.com/item.htm?d=541129513123&skuld=324767200675&areald=430100&user_id=300425476&cat_id=50041298&is_b=1&m=14a2c261												
39	木九十眼镜 54双	150	1094		\\detail.tmall.com/item.htm?d=54930520340&skuld=35030533625&areald=430100&user_id=304003141&cat_id=50041298&is_b=1&m=14a2c261												
40	comani眼镜 53双	275	126		\\detail.tmall.com/item.htm?d=43321702258&skuld=373795166934&areald=430100&user_id=6644840&cat_id=50041298&is_b=1&m=14a2c261												
41	视光眼镜 4双	1080	23		\\detail.tmall.com/item.htm?d=534137860169&skuld=318562392319&areald=430100&user_id=221501090&cat_id=50041298&is_b=1&m=14a2c261												
42	innix眼镜 142双	98	2358		\\detail.tmall.com/item.htm?d=21194036289&skuld=3126546650829&areald=430100&user_id=506921075&cat_id=50041298&is_b=1&m=14a2c261												
43	ago眼镜 15双	299	22		\\detail.tmall.com/item.htm?d=5749610488&skuld=39398568184&areald=430100&user_id=300600339&cat_id=50041298&is_b=1&m=14a2c261												
44	巨匠眼镜 5双	399	302		\\detail.tmall.com/item.htm?d=3610741582&skuld=37182970155&areald=430100&user_id=166669227&cat_id=50041298&is_b=1&m=14a2c261												
45	innix眼镜 77双	198	1398		\\detail.tmall.com/item.htm?d=401427719&skuld=312644285966&areald=430100&user_id=506921075&cat_id=50041298&is_b=1&m=14a2c261												
46	曼丝眼镜 66双	500	91		\\detail.tmall.com/item.htm?d=5766272136&skuld=390835201373&areald=430100&user_id=311002743&cat_id=50041298&is_b=1&m=14a2c261												
47	汉源眼镜 74双	150	1328		\\detail.tmall.com/item.htm?d=5491326784&skuld=46116865706066250&areald=430100&user_id=223849392&cat_id=50041298&is_b=1&m=14a2c261												
48	派丽豪方 53双	199	101		\\detail.tmall.com/item.htm?d=537102919084&skuld=320660100325&areald=430100&user_id=928190478&cat_id=50041298&is_b=1&m=14a2c261												
49	大明眼镜 18双	188	146		\\detail.tmall.com/item.htm?d=550374807907&skuld=3354110881205&areald=430100&user_id=828813377&cat_id=50041298&is_b=1&m=14a2c261												
>>> 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 ...																	

爬取到的商品评论数据

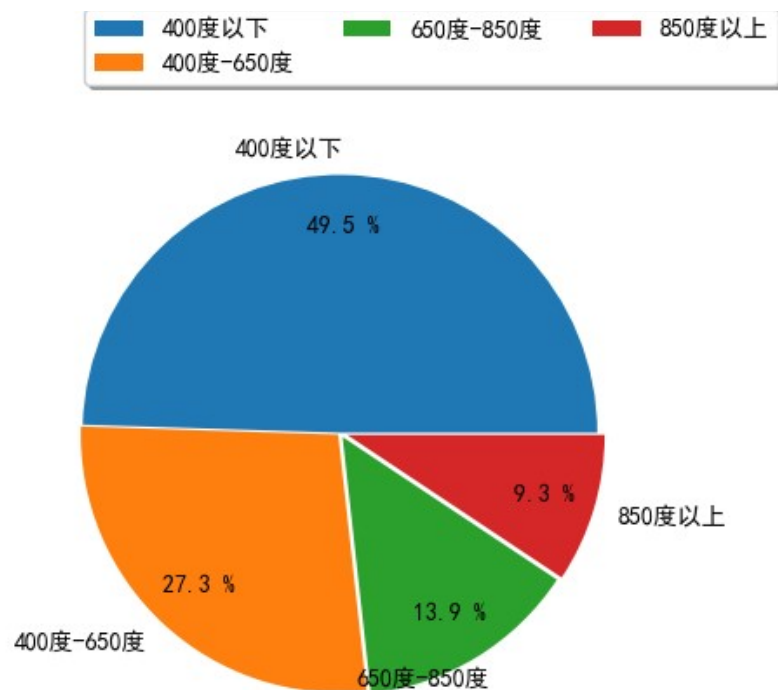

```

try:
    dict_data = json.loads(json_data, strict=False)
    for json_item in dict_data['rateDetail']['rateList']:
        if json_item['appendComment'] is not None:
            append_comment_count = append_comment_count + 1
            f.write(json_item['appendComment']['content'] + "\n")
            comment_time = json_item['rateDate'].split()[0] # 评论日期
            dushu = json_item['auctionSku'][json_item['auctionSku'].index
('镜片适合度数') + 7:]
            dushu = re.findall(r"\d+\.\d*", dushu) # 提取度数
            dushus.append(dushu)
            comment_time_list.append(comment_time)
            count = count + 1
except Exception:
    continue

```

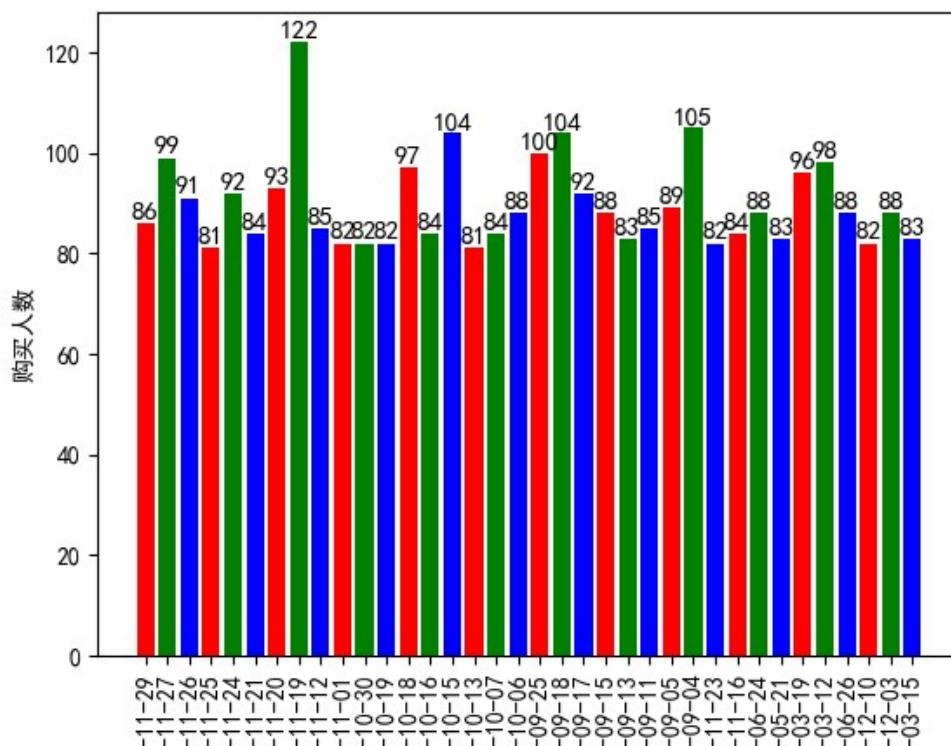
dushus 保存的是所有眼镜片销售的度数规格

comment_time_list 保存的是所有评论的日期



1. 评论日期分析

分析网购的规律



分析可以发现销量最高一天为11月19号，这个比较符合实际规律，双11销量比较高（考虑8天的快递时间）。

1. 追加评论情感分析(snowNLP)

从追评的角度分析网购眼镜片的可靠性

Features

- * 中文分词 ([Character-Based Generative Model](http://aclweb.org/anthology/Y/Y09/Y09-2047.pdf))
- * 词性标注 ([TnT](http://aclweb.org/anthology/A/A00/A00-1031.pdf) 3-gram 隐马)
- * 情感分析 (现在训练数据主要是买卖东西时的评价，所以对其他的一些可能效果不是很好，待解决)
- * 文本分类 (Naive Bayes)
- * 转换成拼音 (Trie树实现的最大匹配)
- * 繁体转简体 (Trie树实现的最大匹配)
- * 提取文本关键词 ([TextRank](http://acl.ldc.upenn.edu/acl2004/emnlp/pdf/Mihalcea.pdf)算法)
- * 提取文本摘要 ([TextRank](http://acl.ldc.upenn.edu/acl2004/emnlp/pdf/Mihalcea.pdf)算法)
- * tf, idf
- * Tokenization (分割成句子)
- * 文本相似 ([BM25](http://en.wikipedia.org/wiki/Okapi_BM25))
- * 支持python3 (感谢[erning](https://github.com/erning))

如果有追加评论的话，那么aooendComment字段不为None
在json字符串解析时顺便将追评保存到AppendComment.txt

```
if json_item['appendComment'] is not None:
    append_comment_count = append_comment_count + 1
    f.write(json_item['appendComment']['content'] + "\n")
```

AppendComment.txt



代码：

```

from snownlp import SnowNLP
import matplotlib.pyplot as plt

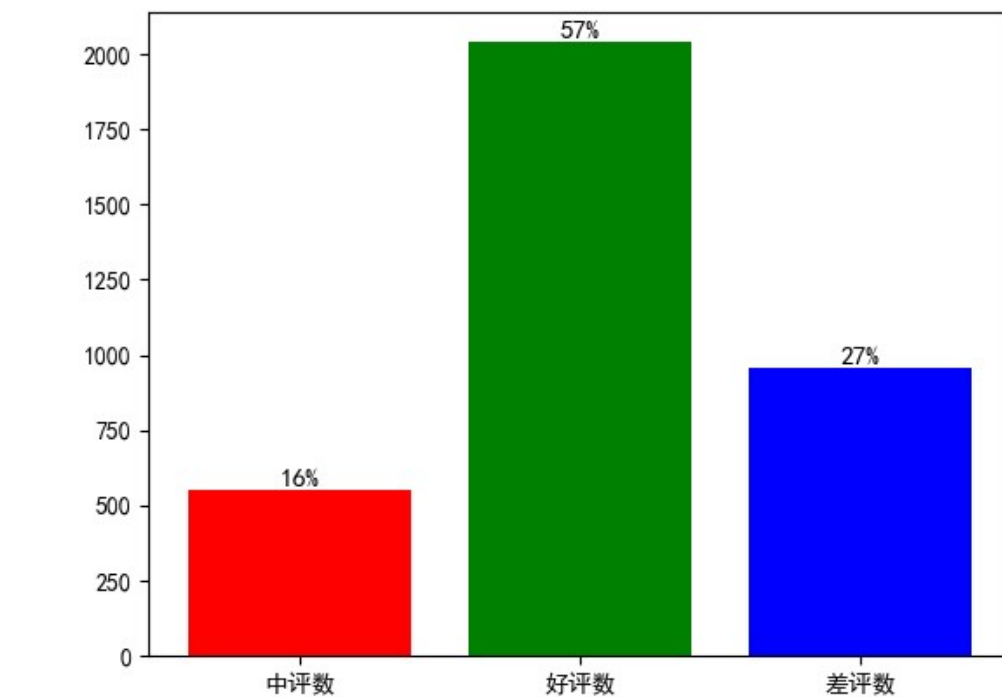
file = "appendComment.txt" # 评论数据文件
bad = "bad.txt" # 差评保存文件
good = "good.txt"
comment_dict = {}
z = 0 # 总数
with open(file, "r", encoding="gbk", errors='ignore') as text:
    bad = open(bad, "w", encoding="utf-8")
    good = open(good, 'w', encoding='utf-8')
    for comment in text:
        z += 1
        s = SnowNLP(comment) # 文本分析
        s = s.sentiments # 情感系数
        if s >= 0.66:
            good.write(comment)
            comment_dict['好评数'] = comment_dict.get('好评数', 0) + 1
        elif 0.66 > s > 0.33:
            comment_dict['中评数'] = comment_dict.get('中评数', 0) + 1
        else:
            bad.write(comment) # 写入差评数
            comment_dict['差评数'] = comment_dict.get('差评数', 0) + 1
    bad.close()
    good.close()
for a, b in comment_dict.items():
    pctb = b/z*100
    # ha 文字指定在柱体中间, va指定文字位置 fontsize指定文字体大小
    plt.text(a, b + 0.05, '%.0f%%' % pctb, ha='center', va='bottom',
             fontsize=11)
# 设置X轴Y轴数据, 两者都可以是list或者tuple
x_axis = tuple(comment_dict.keys())
y_axis = tuple(comment_dict.values())
plt.bar(x_axis, y_axis, color='rgb') # 如果不指定color, 所有的柱体都会
是一个颜色

plt.savefig("cat_spider_img\\append_comment_analyze.png") # 保存为图
片
plt.show()
print(comment_dict)

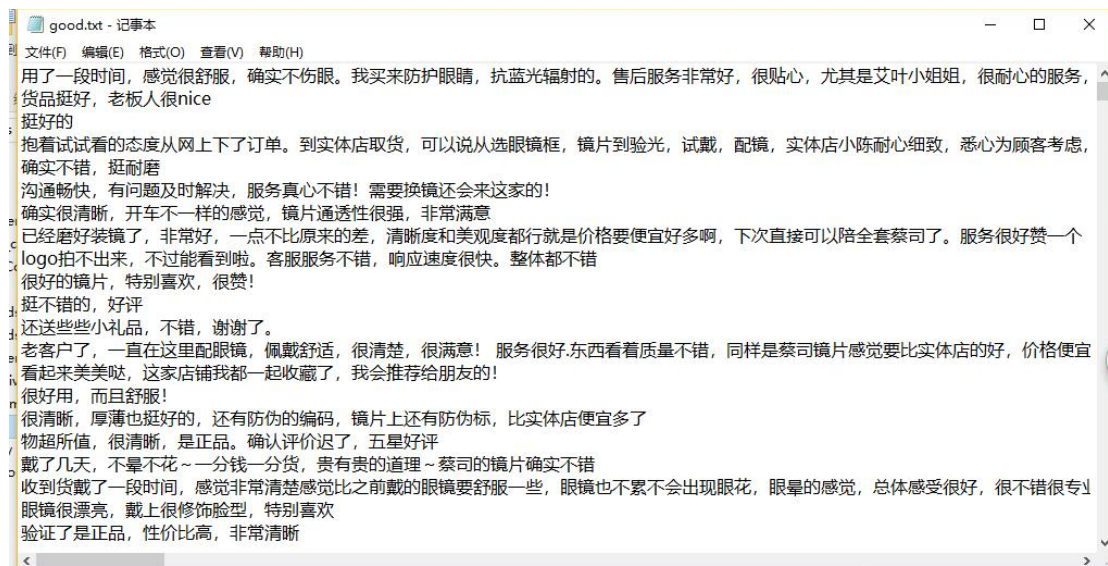
```

最后得到评论总数 37309 追评数 3638

情感分析结果



好评结果



差评结果

