# Application Research of Crawler and Data Analysis Based on Python

Wu Hejing

East University of Heilongjiang

Heilongjiang, 150086

E-mail: 499917928@qq.com,

Liu Fang

East University of Heilongjiang

Heilongjiang, 150086）

Zhao Long

East University of Heilongjiang

Heilongjiang, 150086

Shao Yabin

East University of Heilongjiang

Heilongjiang, 150086

Cui Ran

East University of Heilongjiang

Heilongjiang, 150086

*Abstract*—**Combined with the actual situation, this paper explores how to develop a crawler method based on the specific framework for the complete interface of steam manufacturers and stores, which should be able to automatically and efficiently crawl the data of specific targets, analyze the dynamic pages, and complete the data cleaning, downloading, saving and other operations, explore the methods of general data analysis, and Analyze the downloaded data, extract useful information from it, analyze and summarize the specific crawler method and data analysis method through practical application.**

*Keywords-Python; Scrapy; Selenium; BeautifulSoup*

## I.    INTRODUCTION

The 21st century is a book written by information. With the rapid development of information technology, today's society has become a huge information polymer, and there are various kinds of data in this huge polymer. Data is a kind of embodiment of information. In this era of information explosion, how to efficiently find the data we want from all kinds of miscellaneous data and extract them from the network in batches has become a key problem. However, sometimes the unprocessed data itself may be confusing for people. How to process the huge and complex data obtained through what kind of technical means, and finally become an intuitive number, or trend, and become the information that people can obtain intuitively is also a very important topic to be studied in this data age.

## II.    STATISTICAL INVESTIGATION ON THE PREFERENCE SALES VOLUME

In this project, the American Steam online game platform mall is selected as the research object of the crawler. By setting a specific game company as a search keyword in steam's online mall, the data of all works of the company in steam mall are crawled, and the useful information is extracted by analyzing the basic data of each manufacturer's preference for game production type, series sales volume, and praise In addition, the game manufacturers are comprehensively scored and evaluated.

## III. RELEVANT TECHNOLOGY AND FRAMEWORK

This project will use the scrapy framework based on Python language to crawl steam website. Python as a language has the advantages of lightweight, simplicity, wide range of application and so on. At present, various crawler frameworks and application libraries based on Python have been very mature, among which the crawler framework is very popular in the application of general web crawlers. Its first version was released in 2008, and now it is quite mature as a crawler framework. The basicprinciple of the scrapy framework is shown in Figure 1.
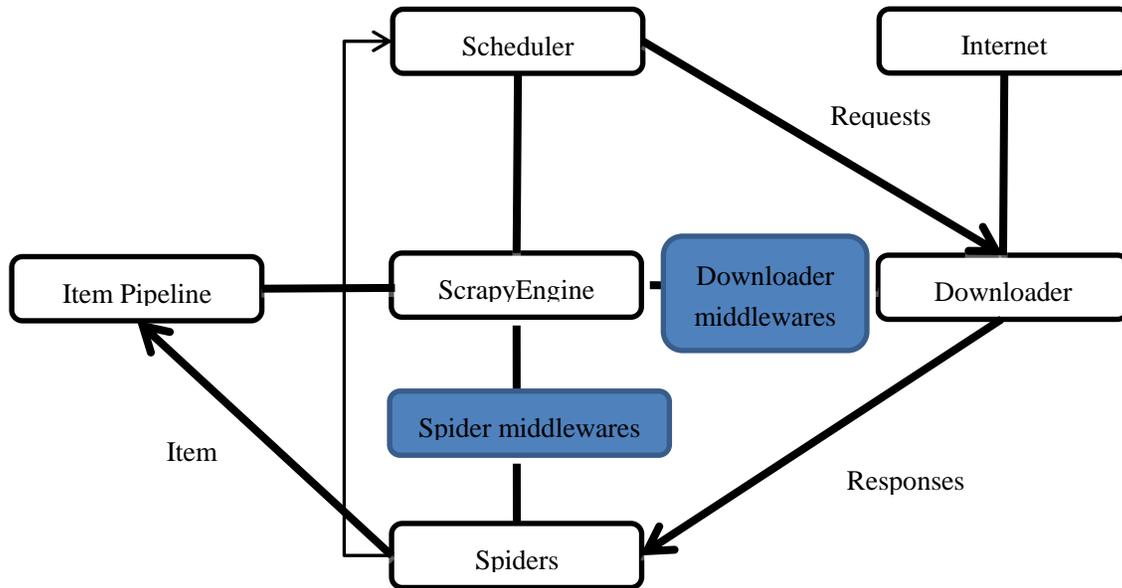


Figure 1.    Basic principles of Scrapy frame

## IV. DESIGN OF CRAWLER

### A. General design idea

The process of crawler itself is actually to simulate the user's operation on the browser with a program. First of all, the starting point and range of crawling need to be specified. As the target of crawling is for manufacturers and their works, the interface of manufacturers is taken as the starting point. For example, the page of paradox, a manufacturer, first analyzes the entire manufacturer's page, and finds that the page links and information of all games or game related DLC downloads of the manufacturer are stored in the recommendation div framework of each sub recommendation of recommendations rows, as shown in Figure 2

### B. Design and implementation of reptile functions

The crawler architecture is composed of items, spiders, piplings and middleware. Among them, items are mainly used to define the items to be crawled, spiders are responsible for defining the whole process of crawling, what means to crawl, pipes are responsible for the basic operations such as data cleaning and saving, middleware can be responsible for the bridge service of scratch and other plug-ins or architectures.
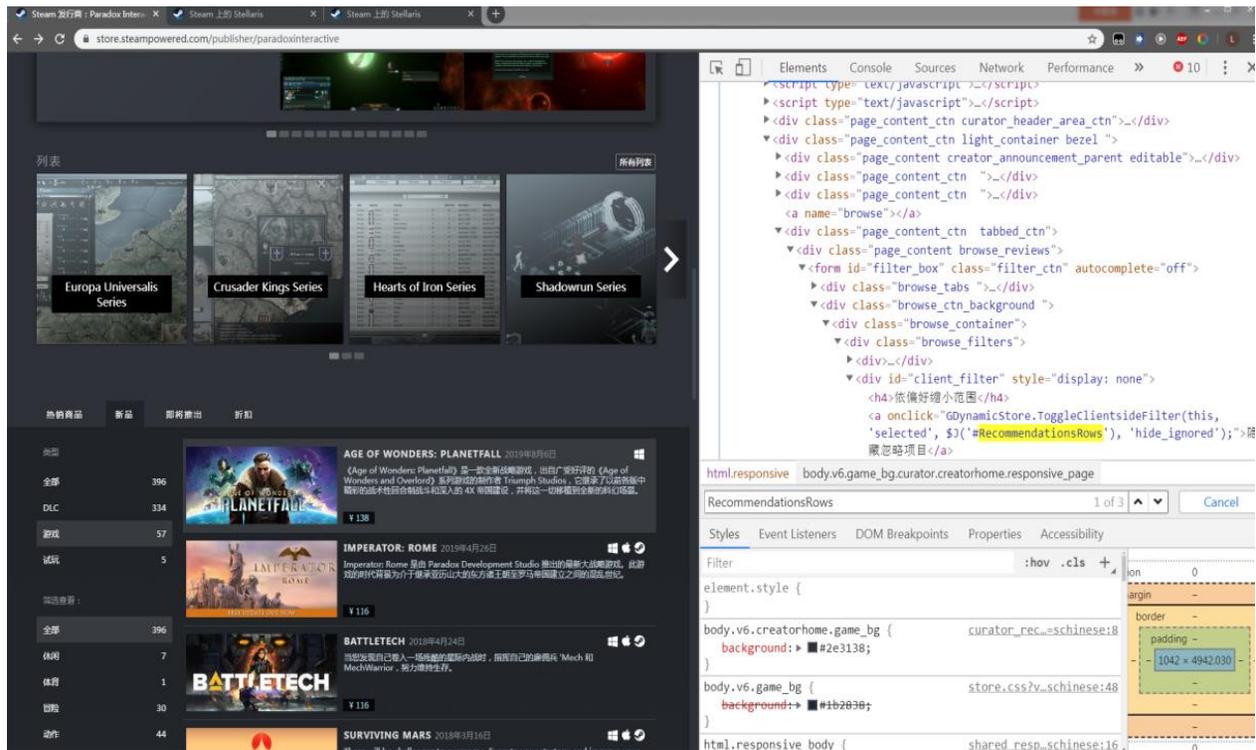
Figure 2.    Investigation of HTML page structure of steam manufacturers by using viewers

First, the items to be crawled are defined in the items file. Finally, these items may be submitted to the analysis part for data analysis. The specific design and implementation code is:

```
import scrapy

class SteamDevItem(scrapy.Item):

    # define the fields for your item here like:

    # name = scrapy.Field()

    qry_nam = scrapy.Field()

    if_dev = scrapy.Field()

    pub_sum = scrapy.Field()

    pub_gam_sum = scrapy.Field()

    pub_dlc_sum = scrapy.Field()

    dev_nam = scrapy.Field()

    pub_nam = scrapy.Field()

    gam_title = scrapy.Field()

    res_date = scrapy.Field()

    gam_type = scrapy.Field()

    gam_tag = scrapy.Field()

    if_muti = scrapy.Field()

    gam_score = scrapy.Field()

    gam_score_sum = scrapy.Field()

    gam_score_ratio = scrapy.Field()

pass
```

## C. Spider design

The design of spider is the key point of this project. Whether the initial dynamic page connection or the last static page information crawling mode will be defined

in this file. In this project, spider will be named steam, and some key implementation codes will be pasted here, with running results and some notes attached. First, introduce start_ the design method of dynamic page crawling of selenium in requests method:

```
chrome_opt = webdriver.ChromeOptions()

    prefs = {


"profile.managed_default_content_settings.images": 2,

        'permissions.default.stylesheet': 2

    }

    chrome_opt.add_experimental_option("prefs",
prefs)

    browser                                    =
webdriver.Chrome(options=chrome_opt)
```

```
    browser.get("https://store.steampowered.com/"
+ Qry_sta + "/" + Qry_Target)

    bs     =     BeautifulSoup(browser.page_source,
'html.parser')        #Beautiful Soup
```

The specific store connections of each product exist in the a anchor label of each entry, and these connections are read to the defined links using the loop_ In the list list, crawling of the list is completed, but sometimes the text and picture in the entry may contain a tag, and they all point to the same page. If direct application may cause repeated crawling, a loop is used here, and if not in statement is used to de duplicate the list.

After using the print statement to verify the function of the module, the verification results are shown in Figure 3.
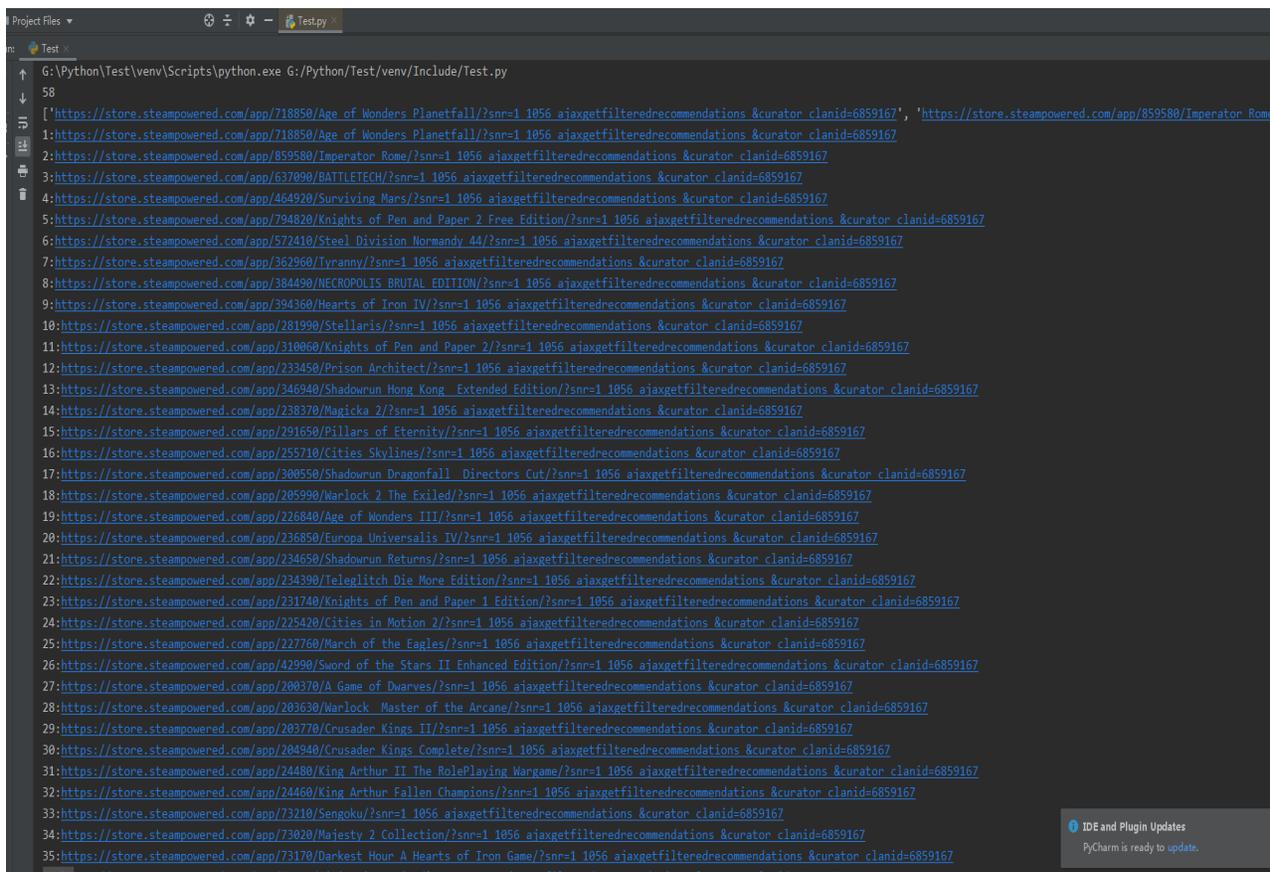


Figure 3.    List of URLs obtained by selenium and beautiful soup

*D. Start directional climbing*

After designing and debugging the spider, run the CMD command window of the system, open the root directory of the crawler file, and input the crawler stream-o SteamDev.csv , crawl the target website.

Input - O SteamDev.csv The purpose is to let the crawler save the last crawled data in the form of CSV table. The saved data appears in the project root. See Figure 4 for the climbing process.



Figure 4.    Executing the start request method selenium pop-up browser to crawl the dynamic page

## V.  DATA ANALYSIS

Next, we will perform basic visual operations on the crawled data in the form of operation tables. In the crawler project, we crawled for the Paradox Interactive publisher. The crawled data is presented in the form of CSV tables, as shown in Figure 5.

Through the use of spreadsheets and further collation of the crawled data, the following data are obtained: the publisher has published 396 works in steam platform, of which the majority of DLC has published 334 DLC, most of the games published are single player games, and each game published in its mall has an average of 6800 reviews, of which the proportion of favorable reviews is about 76.4 8%, see the chart below for detailed visual analysis.

| dev_nam | gam_score | gam_score | gam_score | gam_tag | gam_title | gam_type | if_dev | if_muti | pub_dlc_s | pub_gam_s | pub_nam | pub_sum | qry_nam | res_date | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | 发行商 | | 334 | 57 | | 396 | Paradox Interactive - Official | | |
| Triumph S | 1845 | | 7 | Very Posi | Strategy, | Age of Wc | Strategy | | Single-player | | | Paradox Interactive | | 6 Aug, 2019 | |
| Paradox L | 11187 | | 6 | Mostly Pc | Strategy, | Imperator | Simulation, Strateg | Single-player | | | Paradox Interactive | | 25 Apr, 2019 | |
| Harebrain | 10594 | | 7 | Very Posi | Mechs, Str | BATTLETEC | Action, Adventure, S | Single-player | | | Paradox Interactive | | 24 Apr, 2018 | |
| Haemimont | 7060 | | 9 | Very Posi | Colony Si | Surviving | Simulation, Strateg | Single-player | | | Paradox Interactive | | 15 Mar, 2018 | |
| Kyy Games | 116 | | 6 | Mixed | RPG, Indie | Knights c | Indie, RPG, Simulati | Single-player | | | Paradox Interactive | | 21 Feb, 2018 | |
| Eugen Sys | 3661 | | 7 | Very Posi | World War | Steel Div | Action, Simulation, | Single-player | | | Paradox Interactive | | 23 May, 2017 | |
| Obsidian | 5949 | | 9 | Very Posi | RPG, Story | Tyranny | Adventure, RPG | Single-player | | | Paradox Interactive | | 10 Nov, 2016 | |
| Harebrain | 3384 | | 6 | Mostly Pc | Souls-lik | NECROPOLI | Action, Adventure, I | Single-player | | | Paradox Interactive | | 16-Jul | |
| Paradox L | 47255 | | 9 | Very Posi | Space, Str | Stellaris | Simulation, Strateg | Single-player | | | Paradox Interactive | | 9 May, 2016 | |
| Kyy Games | 1234 | | 9 | Very Posi | RPG, Adver | Knights c | Adventure, Indie, RF | Single-player | | | Paradox Interactive | | 20 Oct, 2015 | |
| Double El | 34388 | | 9 | Very Posi | Simulatic | Prison Ar | Indie, Simulation, S | Single-player | | | Paradox Interactive | | 6 Oct, 2015 | |
| Harebrain | 2106 | | 9 | Very Posi | RPG, Cyber | Shadowrur | Adventure, Indie, RF | Single-player | | | Paradox Interactive | | 20 Aug, 2015 | |
| Pieces Ir | 5097 | | 7 | Very Posi | Magic, Co- | Magicka 2 | Action, Adventure | Single-player | | | Paradox Interactive | | 26 May, 2015 | |
| Obsidian | 9125 | | 9 | Very Posi | RPG, Fanta | Pillars c | RPG | Single-player | | | Paradox Interactive | | 26 Mar, 2015 | |
| Colossal | 74207 | | 9 | Very Posi | City Buil | Cities: S | Simulation, Strateg | Single-player | | | Paradox Interactive | | 10 Mar, 2015 | |
| Harebrain | 3520 | | 9 | Very Posi | RPG, Cyber | Shadowrur | Adventure, Indie, RF | Single-player | | | Paradox Interactive | | 18 Sep, 2014 | |
| Ino-Co Pl | 697 | | 6 | Very Posi | Strategy, | Warlock 2 | Strategy | Single-player | | | Paradox Interactive | | 10 Apr, 2014 | |
| Triumph S | 5144 | | 9 | Very Posi | Strategy, | Age of Wc | RPG, Strategy | Single-player | | | Paradox Interactive | | 31 Mar, 2014 | |
| Paradox L | 47223 | | 9 | Very Posi | Grand Str | Europa Ur | RPG, Cyber | Simulation, Strateg | Single-player | | | Paradox Interactive | | 13 Aug, 2013 | |
| Harebrain | 7411 | | 9 | Very Posi | RPG, Cyber | Shadowrur | Adventure, Indie, RF | Single-player | | | Paradox Interactive | | 25 Jul, 2013 | |
| Test3 Pro | 819 | | 9 | Very Posi | Action Rc | Teleglitc | Action, Indie | Single-player | | | Paradox Interactive | | 24 Jul, 2013 | |
| Behold St | 1541 | | 9 | Very Posi | RPG, Turn- | Knights c | Indie, RPG | Single-player | | | Paradox Interactive | | 18 Jun, 2013 | |
| Colossal | 1016 | | 6 | Mixed | Simulatic | Cities ir | Simulation, Strateg | Single-player | | | Paradox Interactive | | 2 Apr, 2013 | |
| Paradox L | 158 | | 6 | Mixed | Strategy, | March of | Simulation, Strateg | Single-player | | | Paradox Interactive | | 18 Feb, 2013 | |
| Zeal Game | 305 | | 6 | Mixed | Strategy, | A Game of | Casual, Strategy | Single-player | | | Paradox Interactive | | 23 Oct, 2012 | |
| Ino-Co Pl | 738 | | 9 | Very Posi | Strategy, | Warlock - | Strategy | Single-player | | | Paradox Interactive | | 8 May, 2012 | |
| Paradox L | 49512 | | 9 | Very Posi | Grand Str | Crusader | Free to Play, RPG, S | Single-player | | | Paradox Interactive | | 14 Feb, 2012 | |
| Paradox L | 111 | | 6 | Mixed | Strategy, | Crusader | Strategy | Single-player | | | Paradox Interactive | | 14 Feb, 2012 | |
| NeocoreGa | 233 | | 6 | Mixed | Strategy, | King Arth | RPG, Strategy | Single-player | | | Paradox Interactive | | 27 Jan, 2012 | |
| NeocoreGa | 27 | | 6 | Mixed | Strategy, | King Arth | RPG, Strategy | Single-player | | | Paradox Interactive | | 16 Sep, 2011 | |
| Paradox L | 268 | | 6 | Mixed | Strategy, | Sengoku | RPG, Simulation, Str | Single-player | | | Paradox Interactive | | 15 Sep, 2011 | |
| 1C:InoCo | 701 | | 7 | Very Posi | Strategy, | Majesty 2 | Strategy | Single-player | | | Paradox Interactive | | 19 Apr, 2011 | |

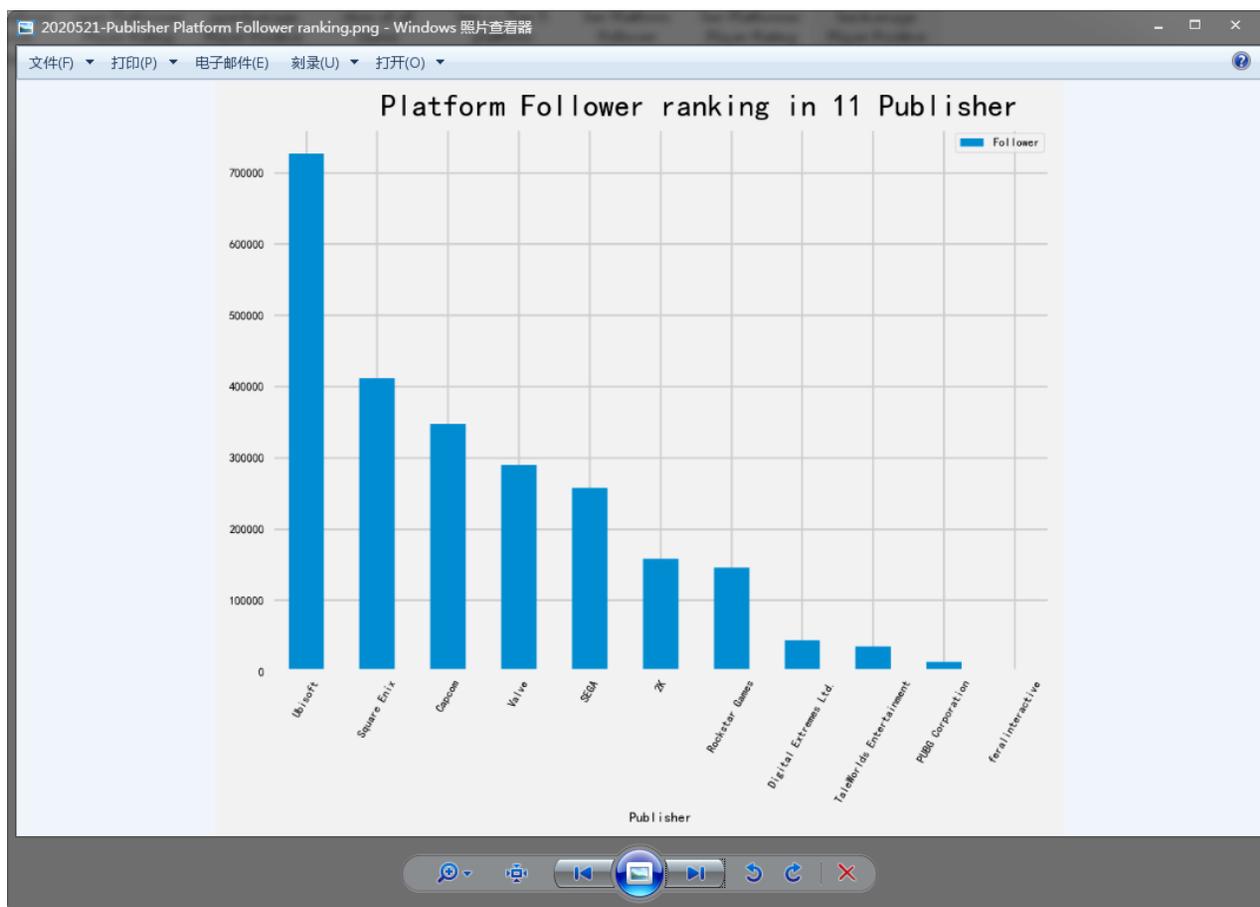Figure 5.    Crawled data list



Figure 6.     Output the publisher platform follower ranking chart

## VI. CONCLUSION

Through demonstration and part of practice, this paper explores the process of data crawling and basic data analysis of dynamic pages by combining the general Python's story framework with selenium + beautiful soup through crawling the steam online game mall website.

The crawler has good scalability. For example, if you want to compare the crawling data of multiple game manufacturers, you can write a query manufacturer list to get the product URL list from the dynamic web page of the manufacturer list first. In terms of anti-crawler, selenium itself has a very good anti crawler ability. If you want to further anti crawler, you can also expand multiple cookies, and even establish a proxy IP pool.

## ACKNOWLEDGMENT

## REFERENCE

[1] Yuhao Fan. Design and implementation of distributed crawler system based on scrapy[J].IOP Conference Series: Earth and Environmental Science,2018,108(4):2-8.

[2] Jing Wang, Yuchun Guo. Scrapy-based crawling and user-behavior characteristics analysis on taobao[P]. Cyber-Enabled Distributed Computing and Knowledge Discovery (CyberC), 2012 International Conference on, 20120:1-5.

[3] Ryan Mitchell. Python web crawler authority Guide (Second Edition) [M]. Beijing: People's post and Telecommunications Press, 2019:57-70.

[4] Wei Chengcheng. Data information crawler technology based on Python [J]. Electronic world, 2018 (11): 208-209.

[5] Mark.Lutz . Python learning manual (Fifth Edition, Volume I) [M]. Beijing: Mechanical Industry Press, 2019:1-2.

[6] Fan Chuanhui. Python reptile development and project practice [M]. Beijing: Mechanical Industry Press, 2017 (3): 69-72.

[7] Song Yongsheng, Huang Rongmei, Wang Jun. research on Python based data analysis and visualization platform [J]. Modern information technology: 2019 (21): 1-4.

[8] Liu Yuke, Wang Ping. Statistics and graph output of student achievement data based on Python + pandas + Matplotlib [J]. Fujian computer. 2017 (11): 2-6.

[9] Liu Yuke, Wang Ping. Statistics and graph output of student achievement data based on Python + pandas + Matplotlib [J]. Fujian computer. 2017 (11): 2-6.

[10] Long Hu, Yang Hui. Data analysis and visualization in the context of big data [J]. Journal of Kaili University. 2016 (03): 1-3.