



Property Tycoon

Team members:

马志宇、刘润东、潘乔、吴铮、汪啸麟

CONTENTS

A person is silhouetted against a night sky filled with stars and the Milky Way. The person is standing on a dark, rocky hill or ridge. The sky is a deep purple and blue, with the Milky Way appearing as a bright, hazy band of light stretching across the upper half of the image. The overall mood is contemplative and vast.

1.Introduction

2.Data collection

3.Difficulties

4.Data analysis

5.Advantages and disadvantages

6.Summary

A person stands in silhouette on a dark, rocky hill under a vast night sky. The Milky Way galaxy is visible as a bright, hazy band of light stretching diagonally across the frame. The sky is filled with numerous stars, and the overall color palette is dominated by deep blues, purples, and blacks.

1

PART

Introduction

INTRODUCTION



1

REQUESTS

2

PANDAS

3

NUMPY

4

BEAUTIFULSOUP4

2 PART

Data collection



1.preparation

```
import request
import time
from bs4 import BeautifulSoup
import pandas as pd
import numpy as np
```


2.url

A Uniform Resource Locator (URL), is a reference to a web resource that specifies its location on a computer network and a mechanism for retrieving it. URLs occur most commonly to reference web pages (http), but are also used for file transfer (ftp), email (mailto), database access (JDBC), and many other applications

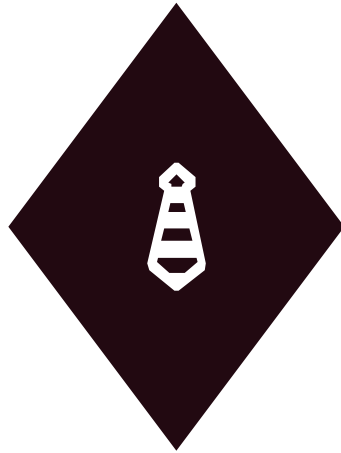
`url='http://bj.lianjia.com/ershoufang/pg'` Here we use 'url' as a name of our variable.

3.Headers

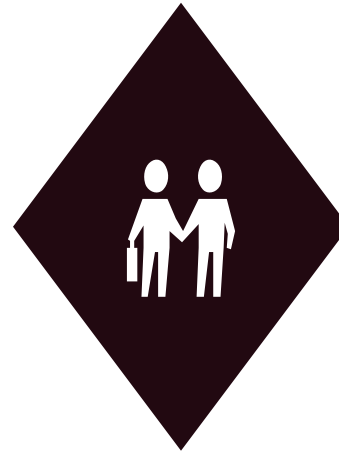
```
headers={'User-Agent':'Mozilla/5.0 (compatible; MSIE  
10.0; Windows NT 6.1; WOW64; Trident/6.0; SLCC2;.NET  
CLR 2.0.50727; .NET CLR 3.5.30729; .NET CLR 3.0.30729;  
InfoPath.3; .NET4.0C; .NET4.0E)',  
        'Accept':'image/webp,image/*,*/*;q=0.8',  
        'Referer':'http://bj.lianjia.com/ershoufang/pg9/',  
        'Accept-Encoding':'gzip, deflate',  
        'Connection':'keep-alive' }
```




3.1 Some web sites often judged by the UserAgent to the operating system



3.2 Accept



3.3 Referer



3.4 Accept-Encoding

4 .Use the for loop to generate 1-100 of the numbers, convert the format to the previous URL fixed part, and spell the URL you want to get. Here we set 0.5 pages per two seconds apart. Then the pages are saved in html.

```
for i in range(1,100):
```

```
    if i == 1:
```

```
        i=str(i)
```

```
html=requests.get(url=url+i+'/',headers=headers).content
```

```
    else:
```

```
        i=str(i)
```

```
html2=requests.get(url=url+i+'/',headers=headers).content
```

```
    html=html+html2
```

```
    time.sleep(0.5) Time.sleep
```

Defer the specified time to run the thread. Unit is' per second'

5. Parse pages and extract information



原生墅 低楼层通透三居 税费低 业主诚售 随时可看 房主自荐

原生墅 | 3室1厅 | 144.67平米 | 东 西 北 | 简装 | 无电梯

中楼层(共6层)2007年建板楼 - 枣园

205人关注 / 共79次带看 / 29天以前发布

距离4号线高米店南站917米

链家优选

房本满五年

随时看房

730万

单价50460元/平米

```
price=lj.find_all('div',attrs={'class':'priceInfo'})
tp=[]
for a in price:
    totalPrice=a.span.string
    tp.append(totalPrice)
```

#提取房源信息

```
houseInfo=lj.find_all('div',attrs={'class':'houseInfo'})
hi=[]
for b in houseInfo:
    house=b.get_text()
    hi.append(house)
```

#提取房源关注度

```
followInfo=lj.find_all('div',attrs={'class':'followInfo'})
fi=[]
for c in followInfo:
    follow=c.get_text()
    fi.append(follow)
```

6.Create data tables and clean data

Import the pandas library, collect the listings, the total price and the attention before, and then generate the data table. Easy to analyze later.

```
house=pd.DataFrame({'totalprice':tp,'houseinfo':hi,'followinfo':fi})

houseinfo_split=pd.DataFrame((x.split('|') for x in house.houseinfo),\
                               index=house.index,\
                               columns=['小区','户型','mianji','朝向','装修','电梯'])

house=pd.merge(house,houseinfo_split,right_index=True,left_index=True)

followinfo_split=pd.DataFrame((x.split('/') for x in house.followinfo),\
                               index=house.index,\
                               columns=['guanzhu','热度','日期'])

house=pd.merge(house,followinfo_split,right_index=True,left_index=True)

print(house)
```

Export data to CSV files and perform further analysis in excel

```
house.to_csv('final.csv')
```


A person stands in silhouette on a dark, rocky hill under a vast night sky. The Milky Way galaxy is visible as a bright, hazy band of light stretching diagonally across the frame. The sky is filled with numerous stars, and the overall color palette is dark with hints of purple and blue from the galaxy's light.

3

PART

Difficulties



1. learning of different libraries.

requests, beautifulsoup4, pandas and numpy

convert the form of data into the attribute of function

2. The locked IP addresses of our computers

lianjia.com have already locked our IP address

have changed several laptops to run the code

Finally, we got a csv file with over 2,000 pieces of data

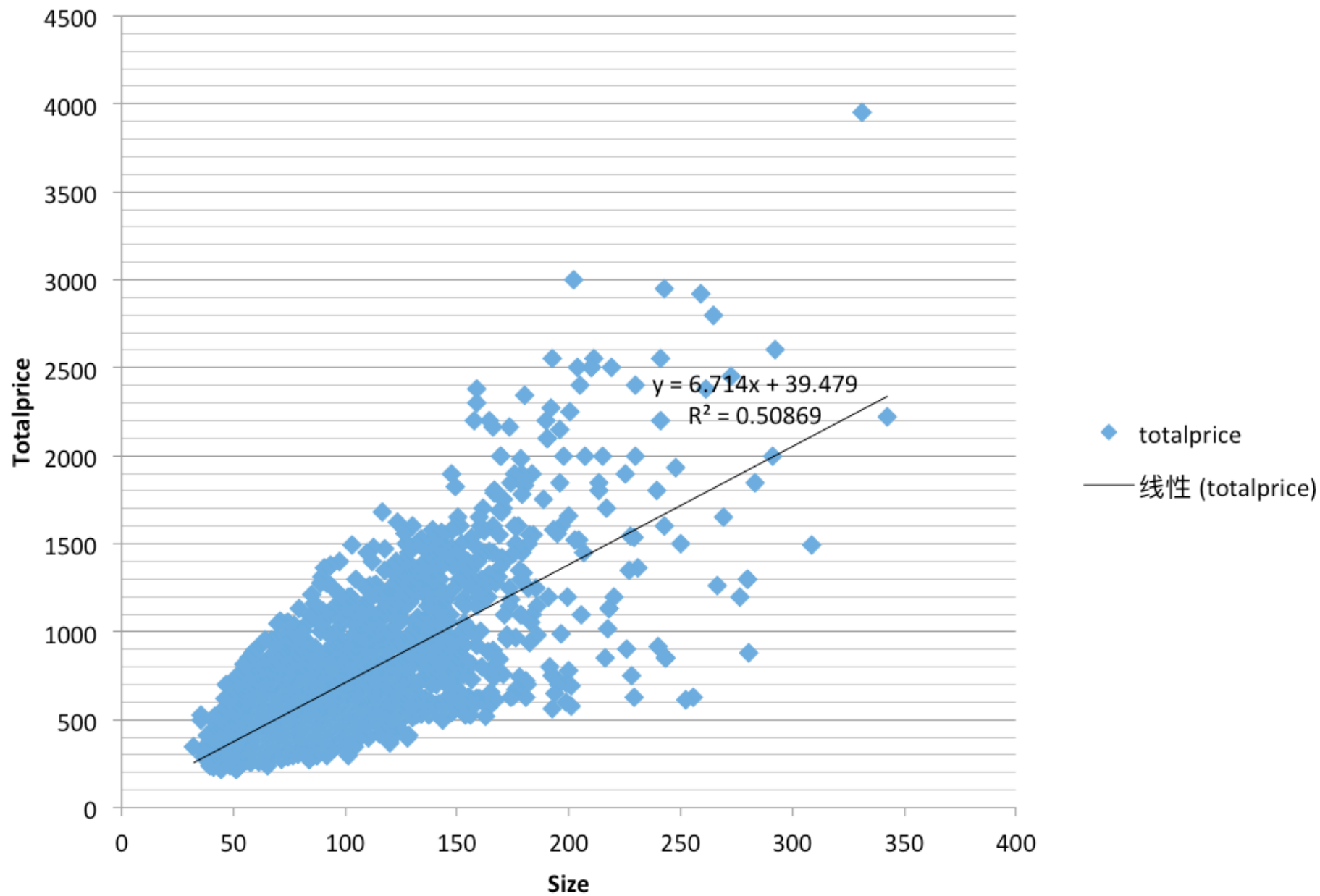
A person stands on a dark, silhouetted hill under a vast night sky filled with stars and the Milky Way. A thin white line extends from the bottom of the number '4' towards the person.

4

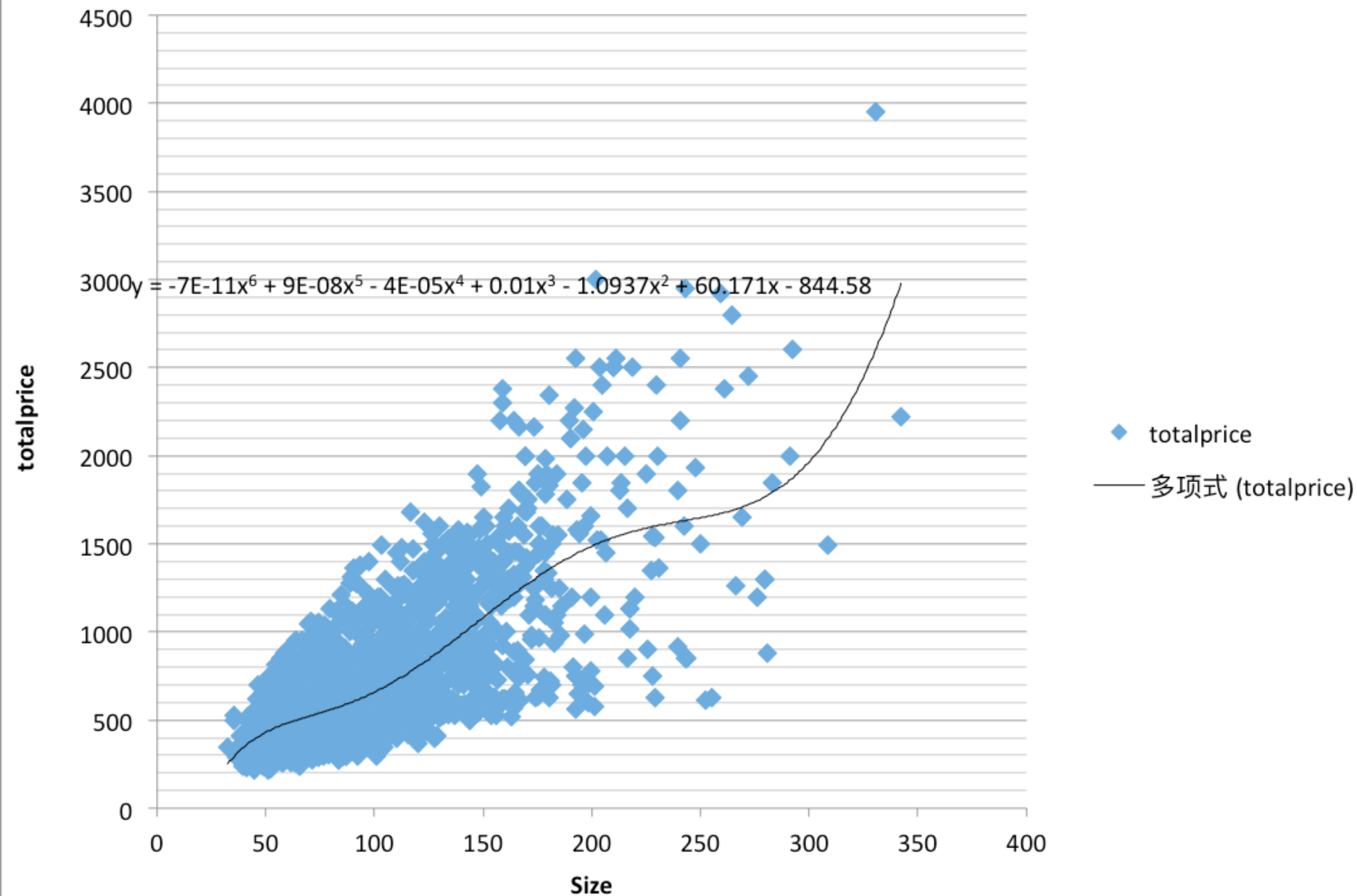
PART

Data analysis

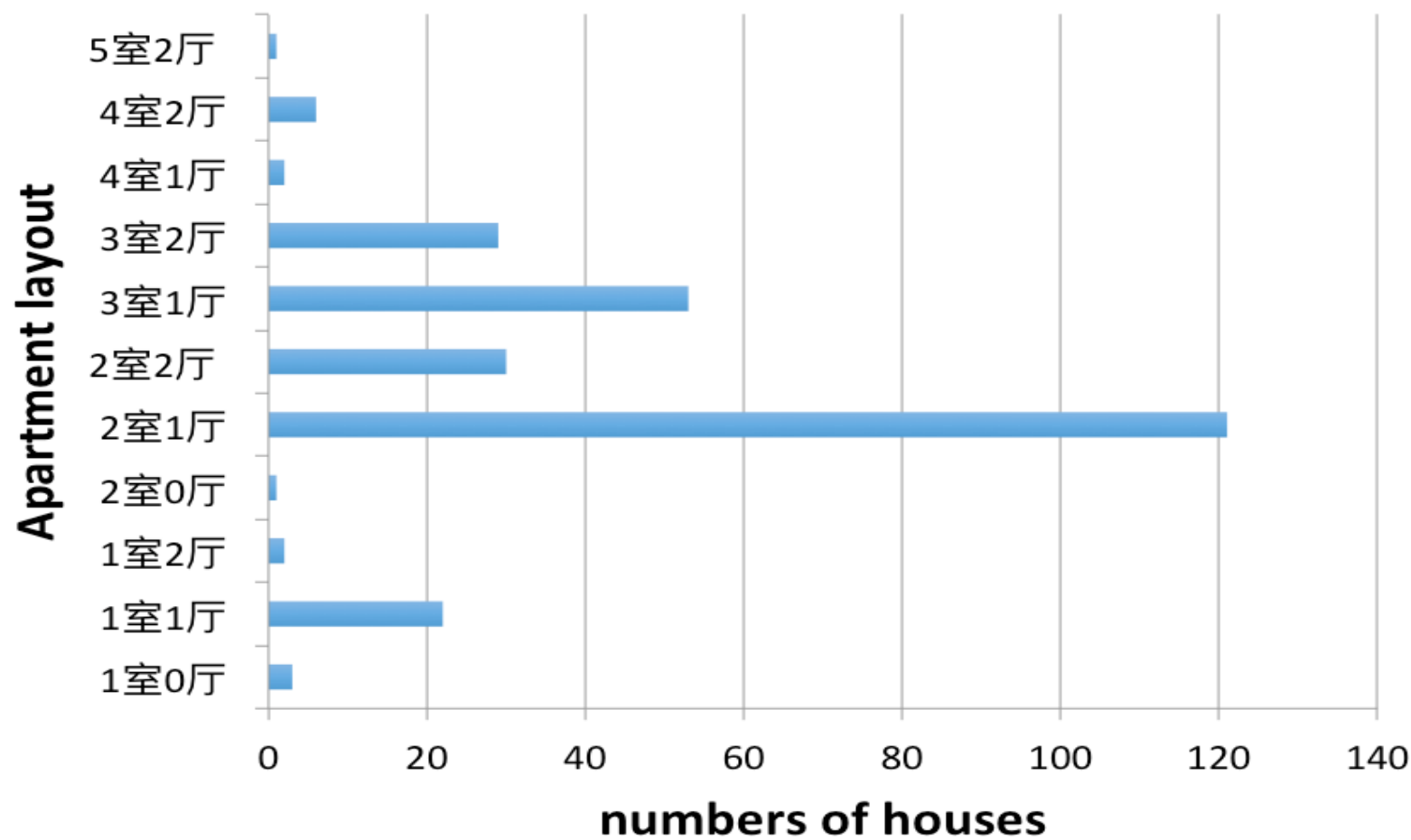
Relationship between totalprice & size



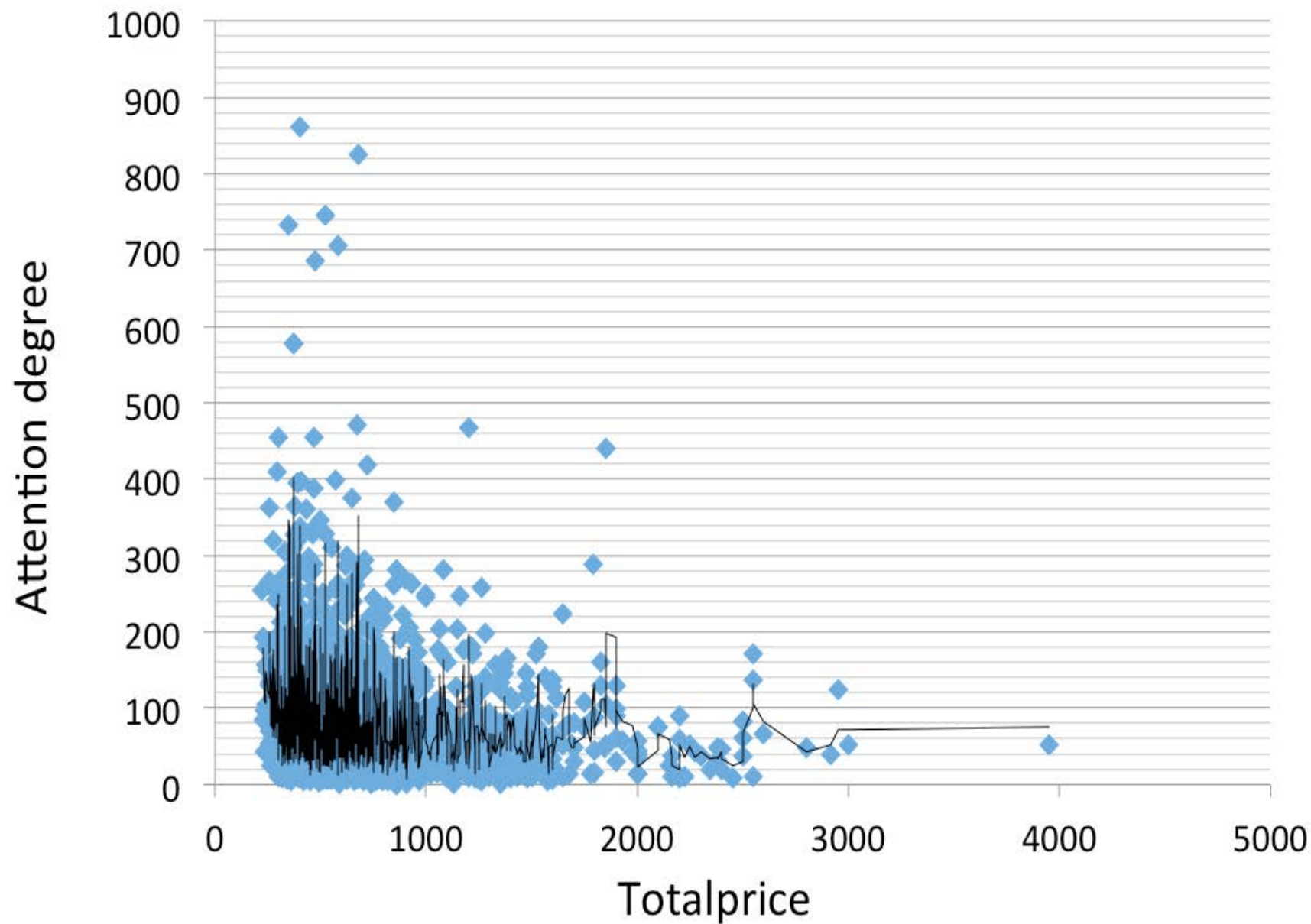
Relationship between totalprice & size



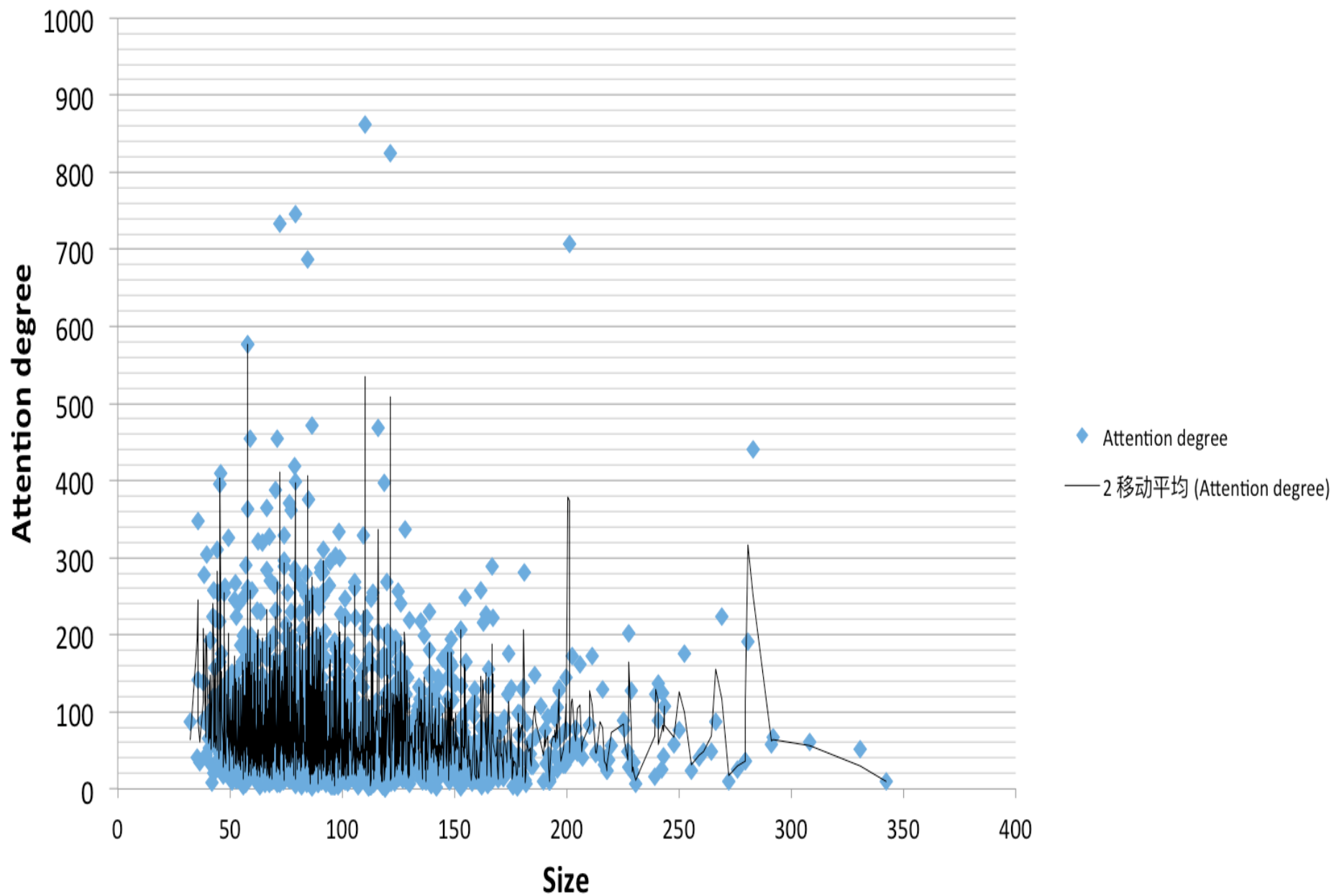
Relationship between apartment layout & numbers of houses



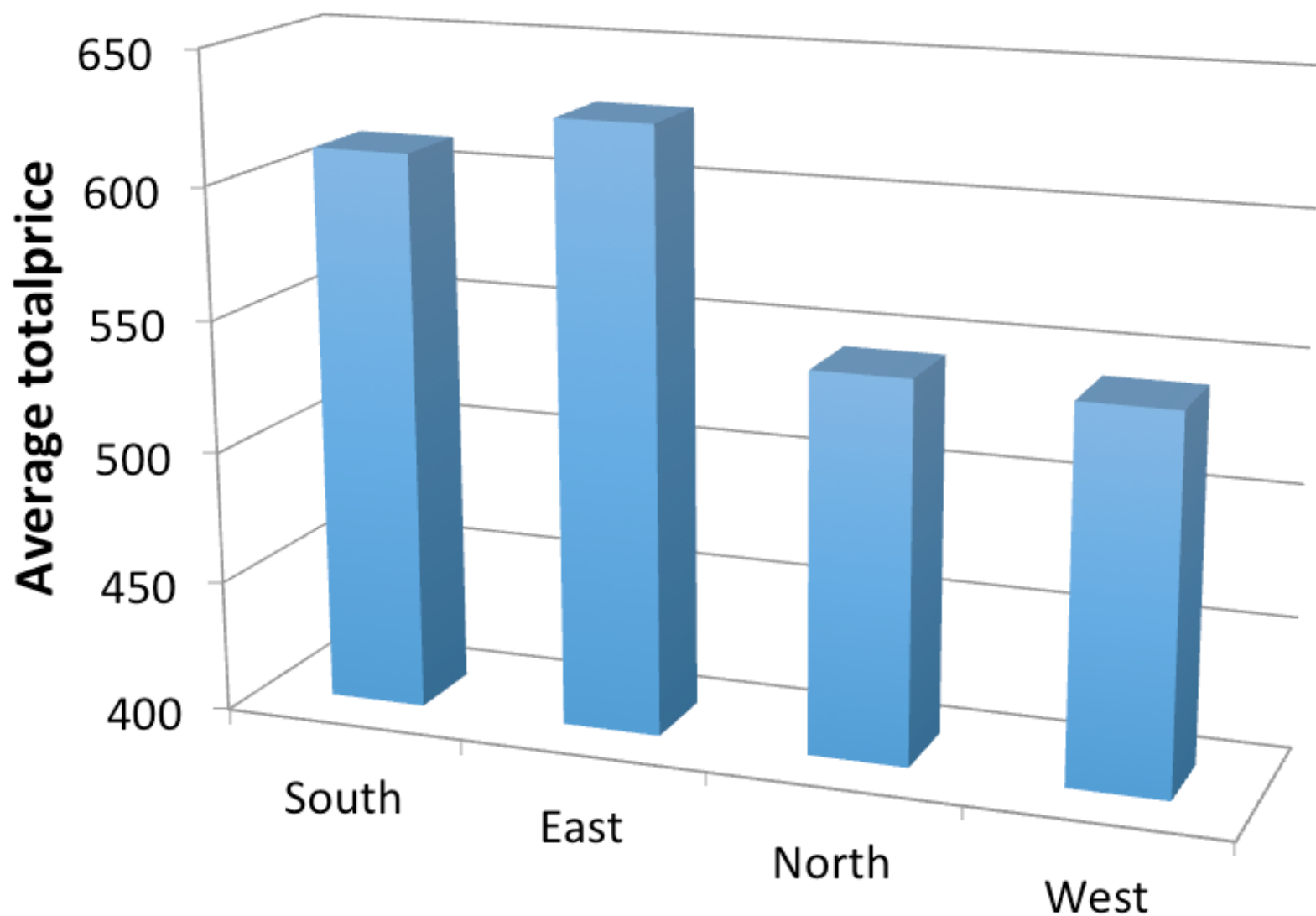
Relationship between attention degree & totalprice



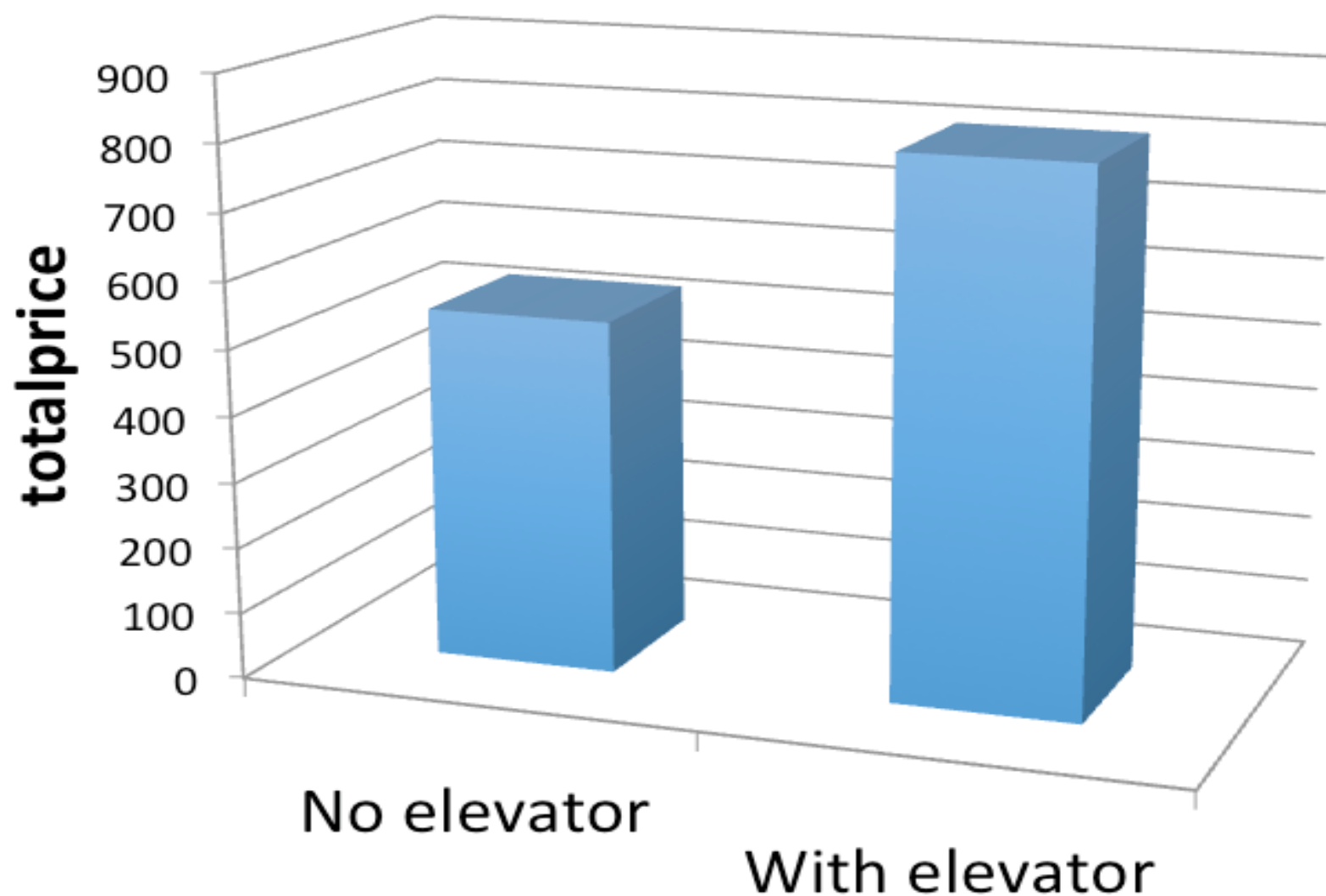
Relationship between attention degree & size



Relationship between totalprice & the direction of the house



The relationship between totalprice & whether there is elevator





5

PART

Advantages
and
disadvantages

Advantages of our project

Since we can collect all the data on the Lianjia website, up to 2372 different apartments, with this code in a few minutes. It saves quite a lot of energy on collecting data. Otherwise, people have to collect thousands of data by hand, which is quite a heavy work.

Disadvantages of our project

The web scrapping program can only collect the data on the main list. It cannot get into every page of each apartment. As a result, we can only collect a few data of a specific apartment, which is not enough for further study.

Summary:

To begin with, we import all kinds of libraries, such as `ch`, `requests`, `time`, `pandas`, `numpy`, `matplotlib.pyplot` and `beautifulsoup`.

Then, we conduct the website-choosing and data collection. Firstly, we observe the target page structure, especially the format of URL. Through the for cycle, we can switch the URL variable part of the page code, and set the header information in a HTTP request, grasp list page, and save in the HTML. Secondly, through the `pandas` library, we create, clean, and generate the data tables. Finally, we export it to the CSV file with `pandas` for later processing.

The last process is the visualization of the collected data. We use `Matplotlib` to draw the relevant chart in order to study the relationships between different type of values.

A full-page background image featuring a person standing on a dark, rocky ridge, silhouetted against a vast night sky. The Milky Way galaxy is visible, stretching from the bottom left towards the top center, with a vibrant pinkish-purple hue. The sky is filled with numerous stars, and the overall color palette transitions from deep blue and black at the top to warm orange and yellow near the horizon.

THANK

YOU