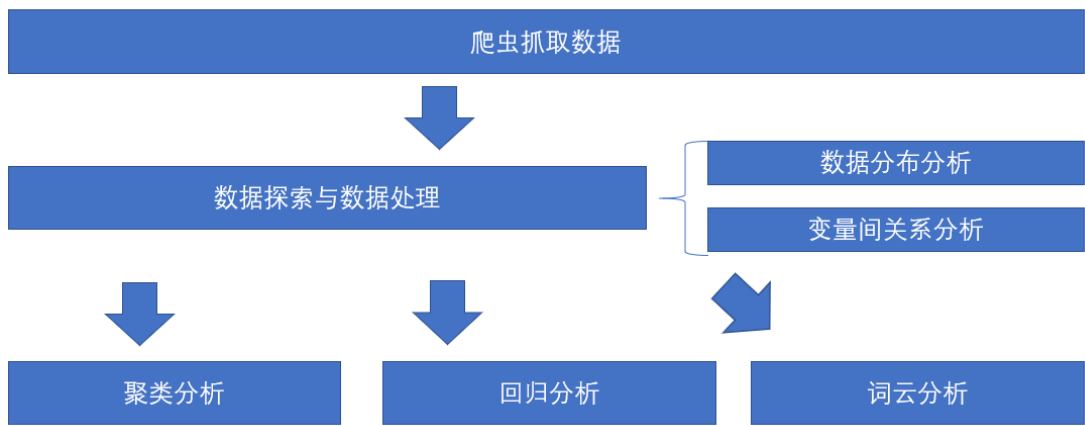


# 上海链家网租房数据分析报告

## 1 分析简介

上海链家网租房数据分析框架如下图所示



本文首先使用爬虫程序抓取房源数据，之后对数据进行了探索与处理，最后对进行了聚类分析、回归分析以及词云分析。

## 2 数据获取

本文使用 python scrapy 框架对上海链家网租房(<https://sh.lianjia.com/zufang>)数据进行抓取。

网页页面如下图所示，

链家网上海站 > 上海租房

按区域 ^ 按地铁线 v

不限 静安 徐汇 黄浦 长宁 普陀 浦东 宝山 闸北 虹口 杨浦 闵行 金山 嘉定 崇明 奉贤 松江 青浦

方式 不限 整租 合租

租金 ☐ ≤1000元 ☐ 1000-1500元 ☐ 1500-2000元 ☐ 2000-3000元 ☐ 3000-5000元 ☐ 5000-8000元 ☐ ≥8000元  -  元 确定

户型 ☐ 一居 ☐ 两居 ☐ 三居 ☐ 四居+

朝向 ☐ 东 ☐ 西 ☐ 南 ☐ 北 ☐ 南北

更多 v

已为您找到 20168 套上海租房 清空条件

综合排序 最新上架 价格 面积

下载链家APP

整租·保平小区 2室1厅 南

闸北-彭浦 / 54㎡ / 南 / 2室1厅1卫

链家

23天前发布

近地铁

3200 元/月

整租·南丹小区 1室1厅 东南/南

徐汇-徐家汇 / 37㎡ / 东南 南 / 1室1厅1卫

贝壳优选

1个月前发布

近地铁 随时看房

4850 元/月

扫描二维码 随时查看新房源 了解更多 >

租房网站上，已知的租房房源大约在 20000 套，而该站点一次只显示 3000 条，单一页面下，无法一次性抓取，因此需要采取别的措施。本文通过价格区间筛选，构造了初始链接，对 20000 余套房源进行了抓取。

整租·马当小区 1室1厅 东南

黄浦-新天地 / 43㎡ / 东南 / 1室1厅1卫

链家

1个月前发布

限女生 近地铁 精装

7300 元/月

整租·上海阳城 2室2厅 南

闵行-梅陇 / 93㎡ / 南 / 2室2厅1卫

链家

5天前发布

近地铁 新上

7300 元/月

一次显示3000套

1 2 3 ... 100 下一页

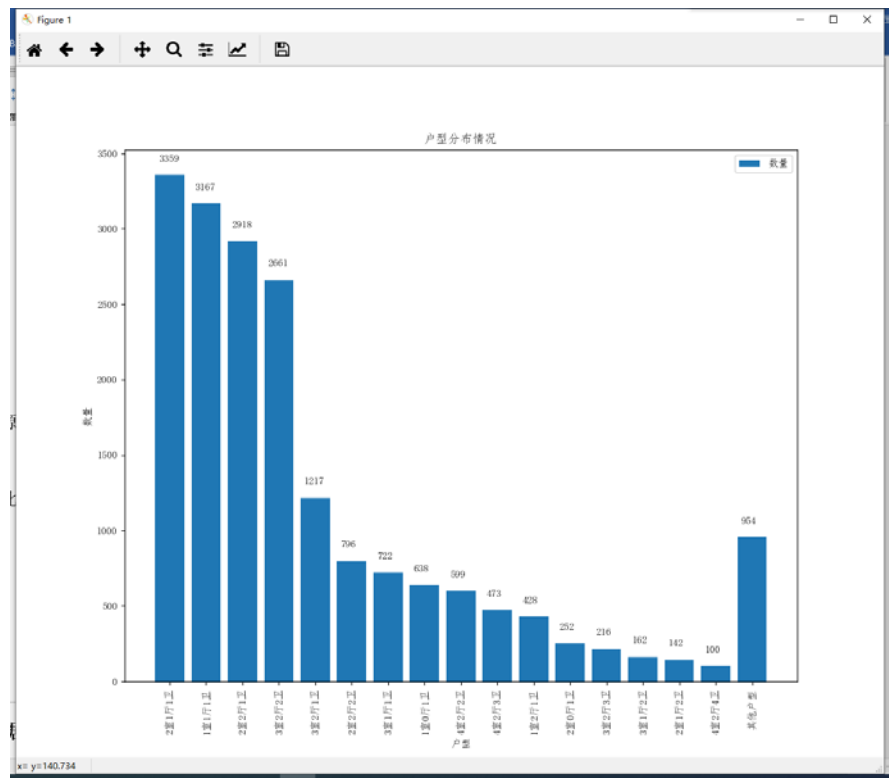
根据网站上的信息，本文抓取了每套房源的 id(房源编号)、area(区域)、adress(地址)、community(小区)、shelf\_time(上架时间)、rent(租金)、house\_type(户型)、square(面积)、towards(朝向)、floor(楼层)、describe(介绍)、subway(地铁信息)、img\_url(图片链接)、url(租房房源链接)、download\_time(采集时间)。

最终爬取 18805 条房源数据，如下图所示，大致接近所有租房房源数据，可以进行整体分析。爬虫代码见 lianjiaSpider 文件夹。



3.1.2 户型分布

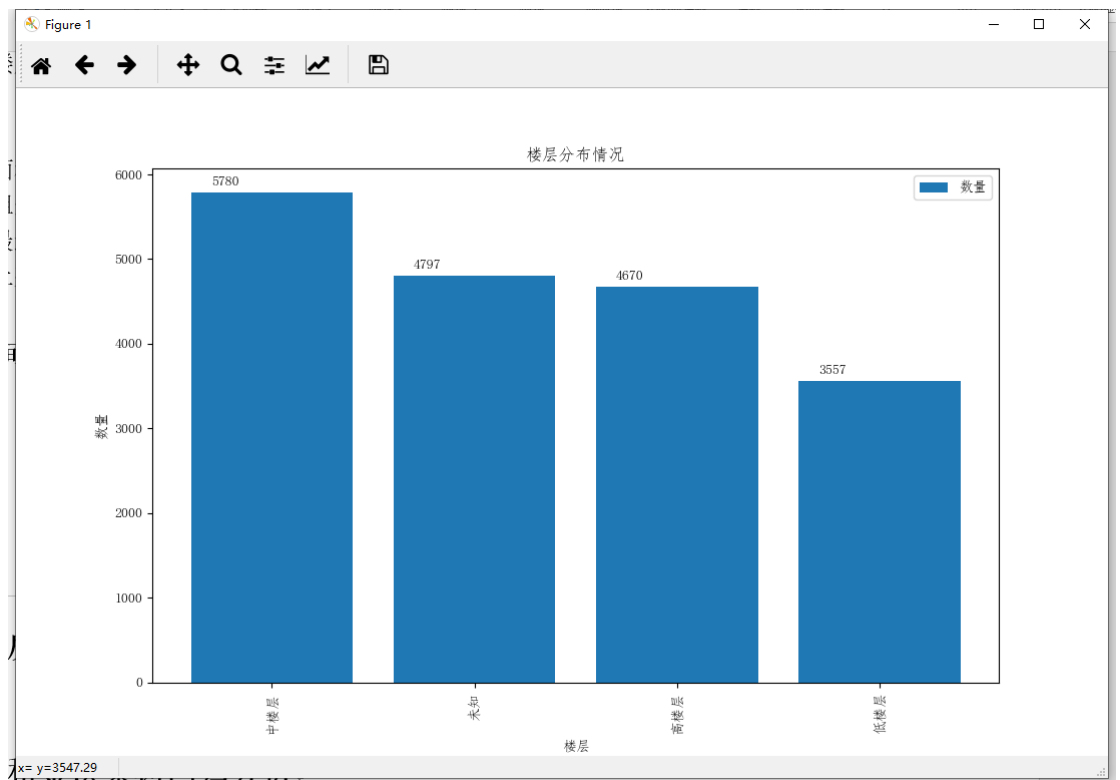
统计发现户型的种类多达 122 种,而大部分户型市场上的数量很少甚至只有 1 个,因此考虑将户型数量小于 100 的归为其他户型,之后绘制分布图如下。



从图中可以看到，市场上数量较多的户型为 2 室 1 厅 1 卫、1 室 1 厅 1 卫、2 室 2 厅 1 卫、3 室 2 厅 2 卫，可见这些户型较受市场欢迎。

3.1.3 楼层分布

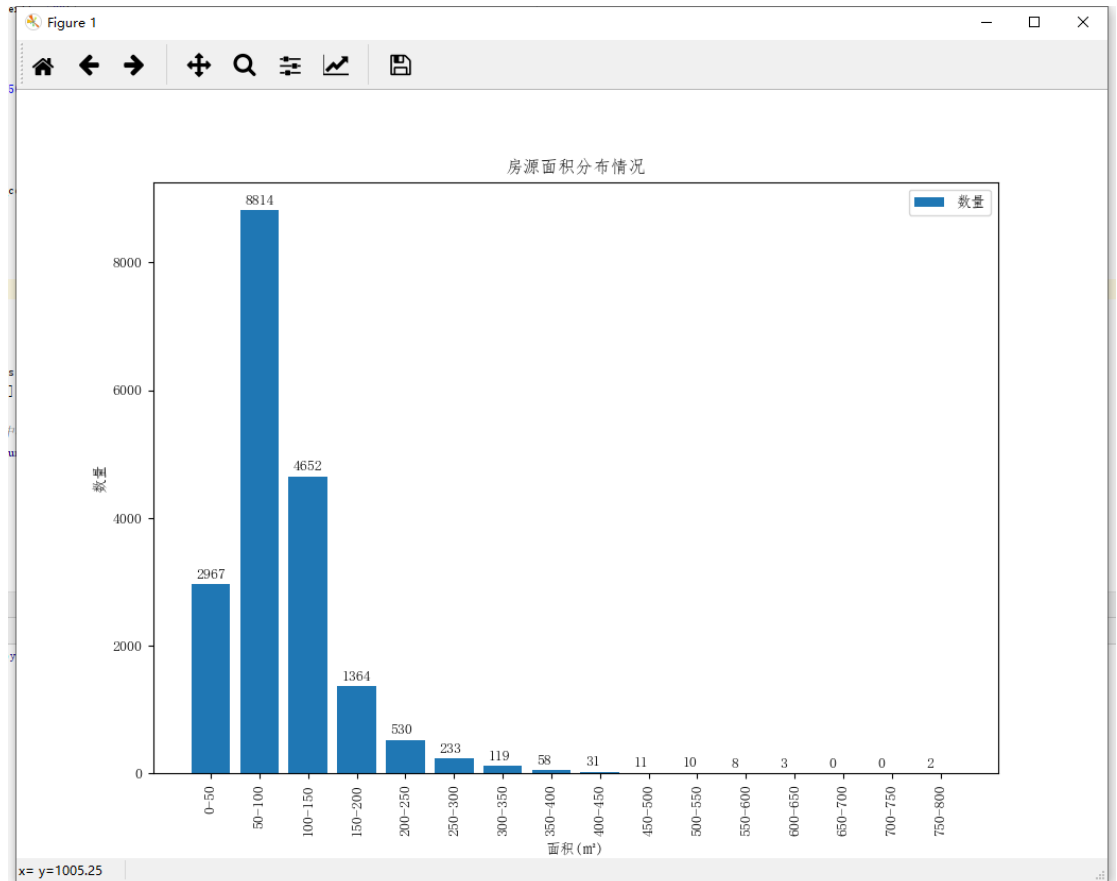
统计发现市场上的楼层描述为楼层类型/总楼层，因而楼层的水平多达 171 种，为了便于分析，统一将楼层划分为低楼层、中楼层、高楼层三类，剩余的归为未知楼层，分布如下图所示。



从图中可以看出，市场上各个楼层均有一定的数量分布，能够满足多样化需求。

### 3.1.4 面积分布

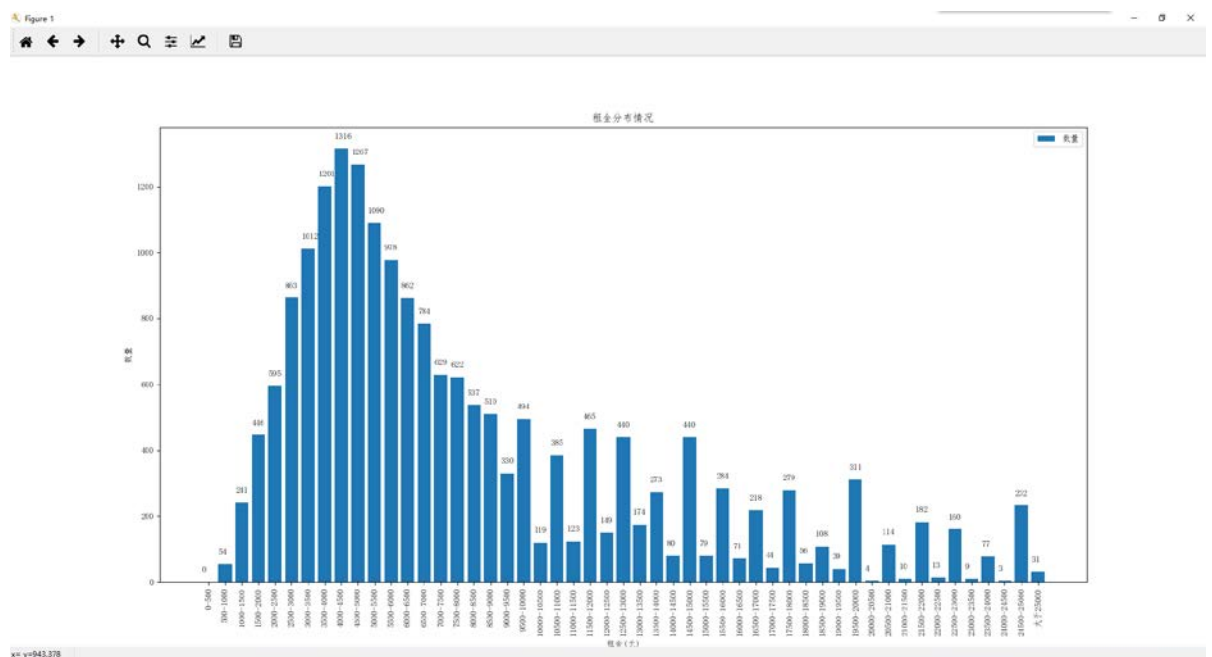
面积数据属于连续数据，需要进行离散化处理才能绘制分布图，本文将面积数据按照一定的区间进行划分，之后绘制分布图，如下所示。



可以看到，租房房源面积集中分布于 0-150 m<sup>2</sup>，而超过 300 m<sup>2</sup>的房源数量较少，房源面积分布整体呈右偏分布。

### 3.1.5 租金分布

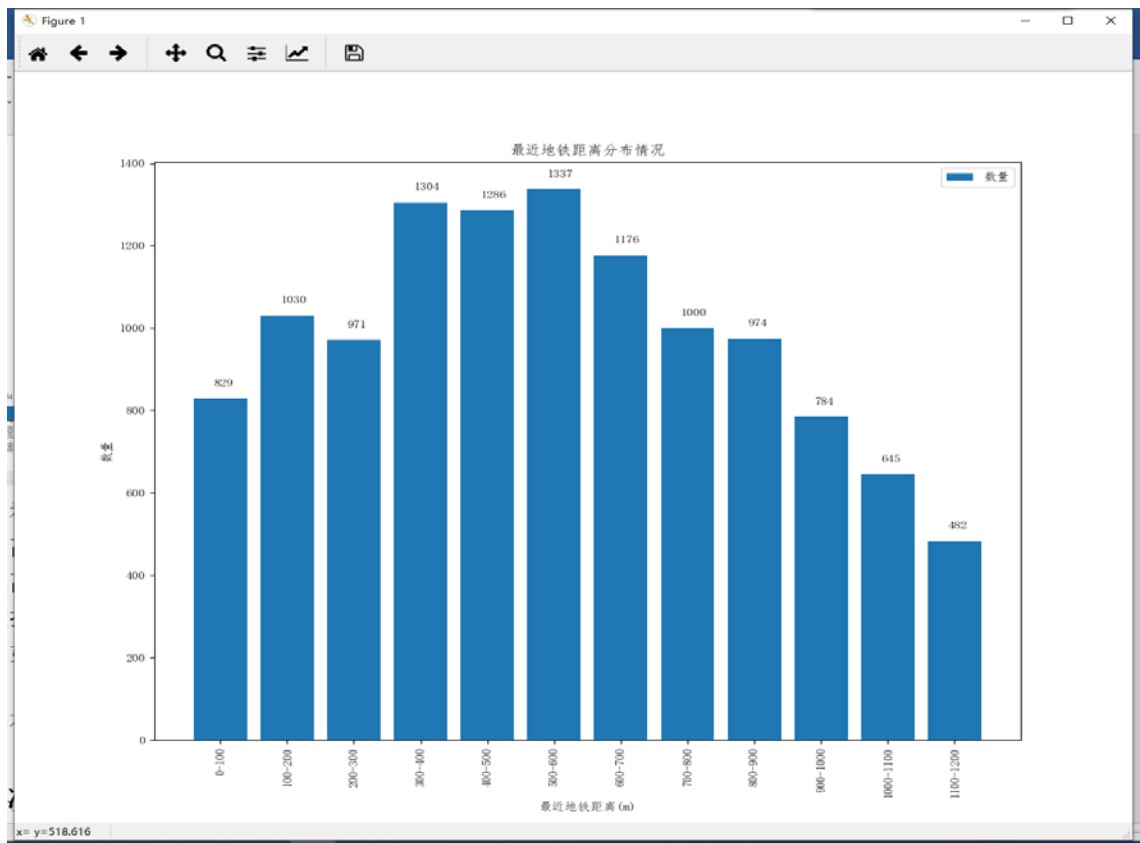
租金数据也为连续数据，因此也需要离散化处理，分布图如下。



从图中可以看到，租金价格主要集中于 2500-6500 元/月。0-500 元/月这一区间没有相应租房房源，而在高价位仍有房源分布，可见上海租房市场的繁荣。

3.1.6 最近地铁距离分布

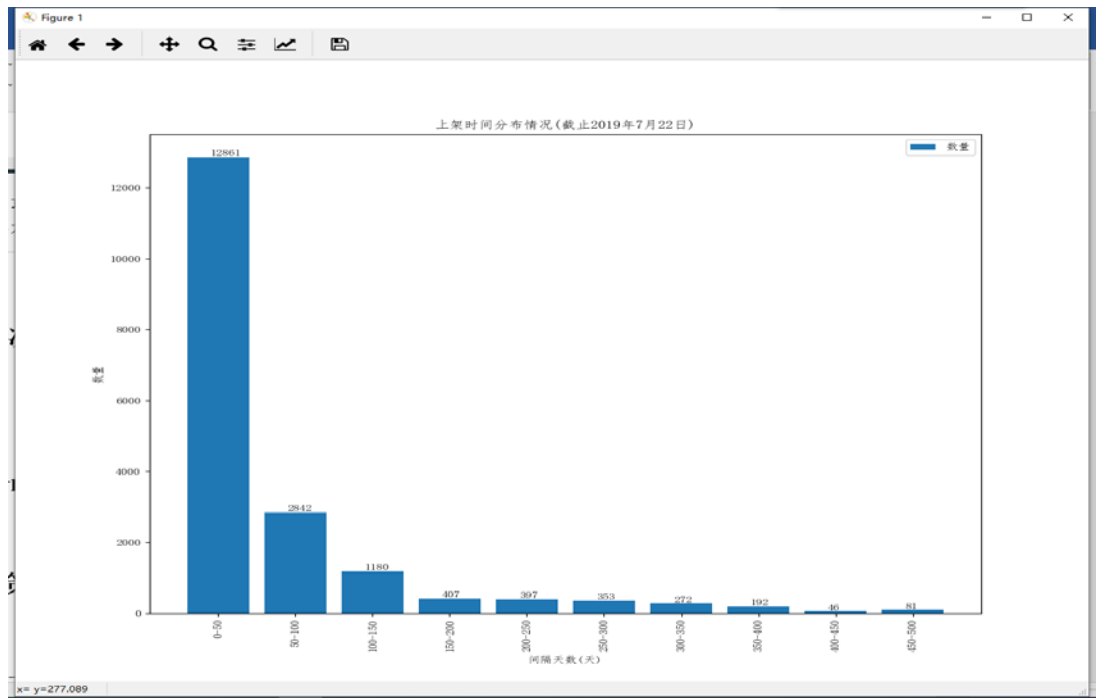
部分房源数据没有地铁相应数据，因此本文只对有地铁数据的房源进行了最近地铁距离分布分析。本文首先对抓取得地铁数据进行了最近距离提取，随后绘制了分布图如下。



可以看到，最近地铁距离在各个区间均有分布，无明显特征。

3.1.7 上架时长分布

本文根据采集时间以及房源发布时间进行了上架时长计算，之后绘制了分布图如下。

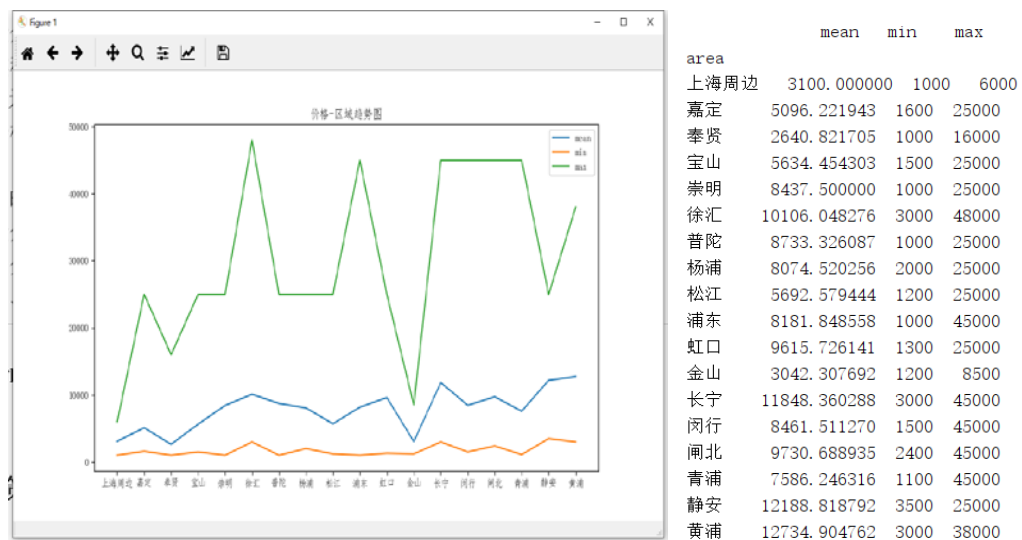


可以看到，上架时长呈分布呈递减趋势，时间越短数量越多。

## 3.2 数据间相关情况分析

### 3.2.1 区域与租金分析

分析各区域租金时发现，某些地区的部分房源租金非常之高，例如长宁区有套房源租金高达 42 万元/月([链接](#))属于极端数据，为了便于分析，本文只保留了租金小于 50000 元/月的数据进行分析，下图是各区域的租金情况趋势图以及相对应表格。

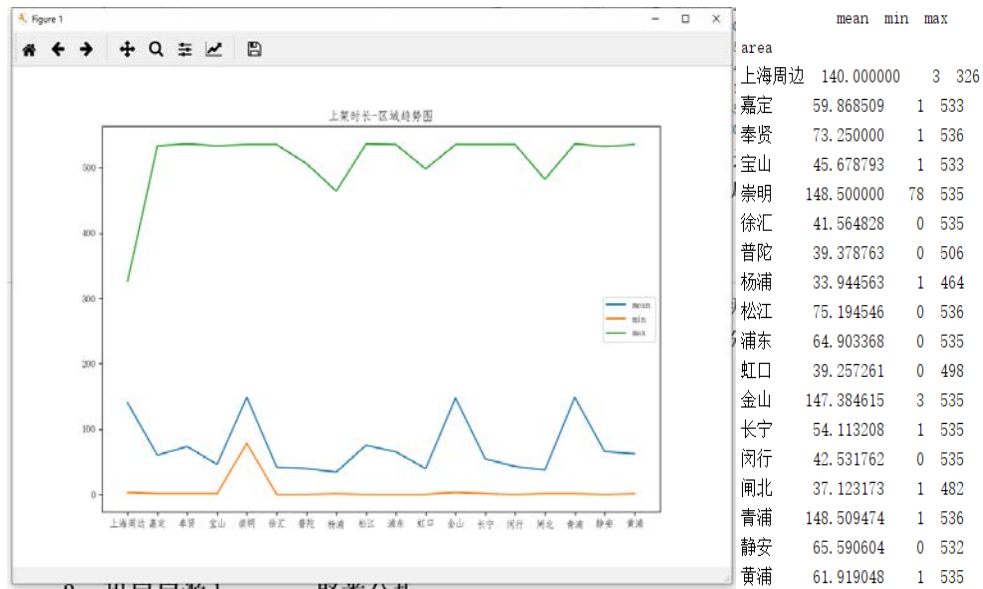


可以看到徐汇、长宁、静安、黄浦的租房价格相对较高，平均价格超过了 10000 元/月，而房源数量最多的浦东地区的租房均价则位于中游。从上表也可以看到，上海地区的租房最低为 1000 元/月。



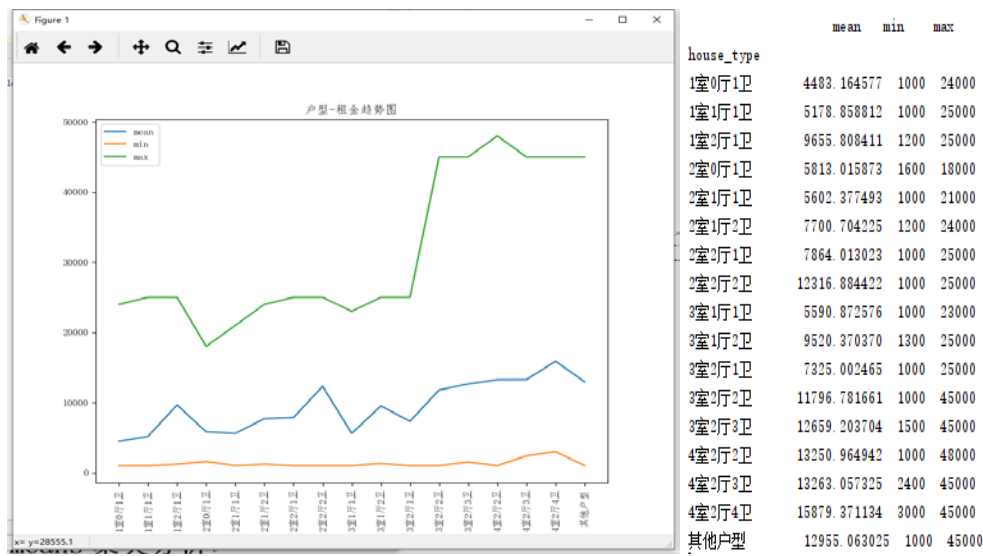
3.2.2 区域与上架时长分析

分析上架时长时发现，有套房源(链接)的发布时间竟然为 1970 年，初步断定为录入错误数据，因此本文将其进行了剔除。随后对各区域的上架时长进行了分析，如下所示。



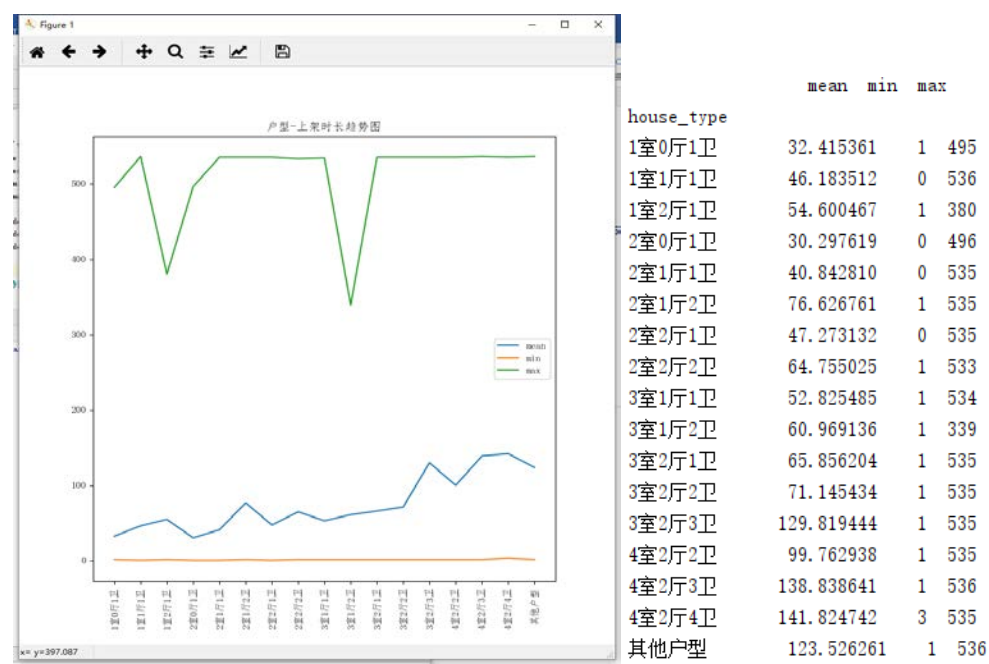
从表中可以看到，崇明这一区域的房源上架时长最久，可能是过于偏僻。其余地区的房源上架时长差异相对不明显。

3.2.3 户型与租金分析



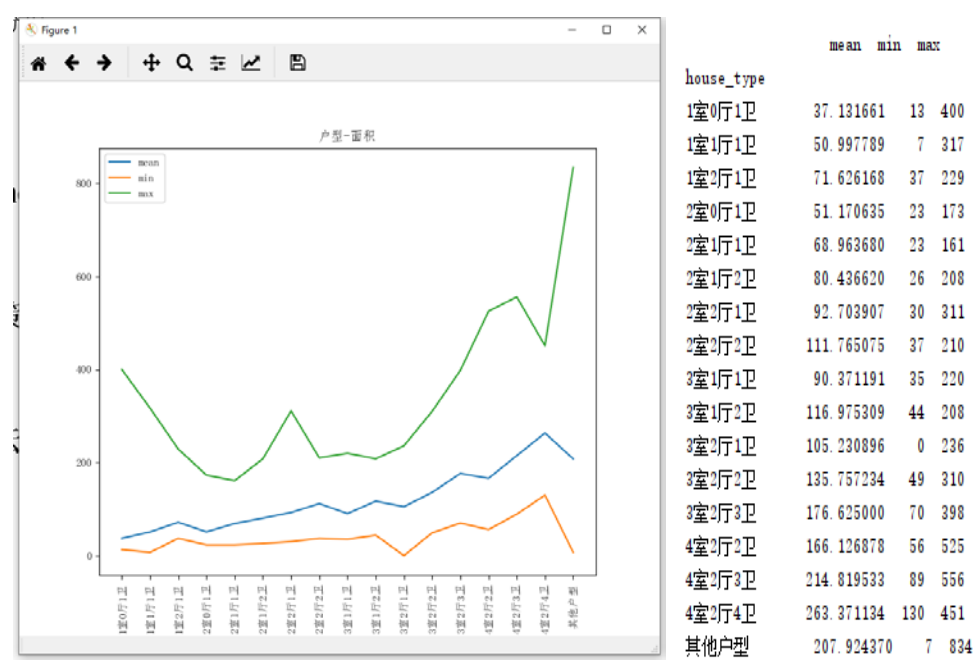
从中可以看出，随着房间数量的增加，租金有上升趋势。而常见户型（2室1厅1卫、1室1厅1卫）平均租金也超过了 5000 元/月。

3.2.4 户型与上架时长分析



从中可以看到，大户型的平均上架时长相对较长，而像常见户型（2室1厅1卫、1室1厅1卫）以及其他小户型则平均上架时长较短。

3.2.5 户型与面积分析



显然，随着房间数量的增加，面积也增加。然而对于户型与面积存在对应关系，例如2室1厅1卫的平均面积为68.96 m²。

3.2.6 租金、面积、上架时长、最近地铁距离相关分析

对于这四种连续变量，本文首先计算了它们的相关系数矩阵，如下表所示。

	rent	square	time_delta	subway
rent	1.000000	0.649142	0.178392	-0.088873
square	0.649142	1.000000	0.279384	0.034736
time_delta	0.178392	0.279384	1.000000	-0.006909
subway	-0.088873	0.034736	-0.006909	1.000000

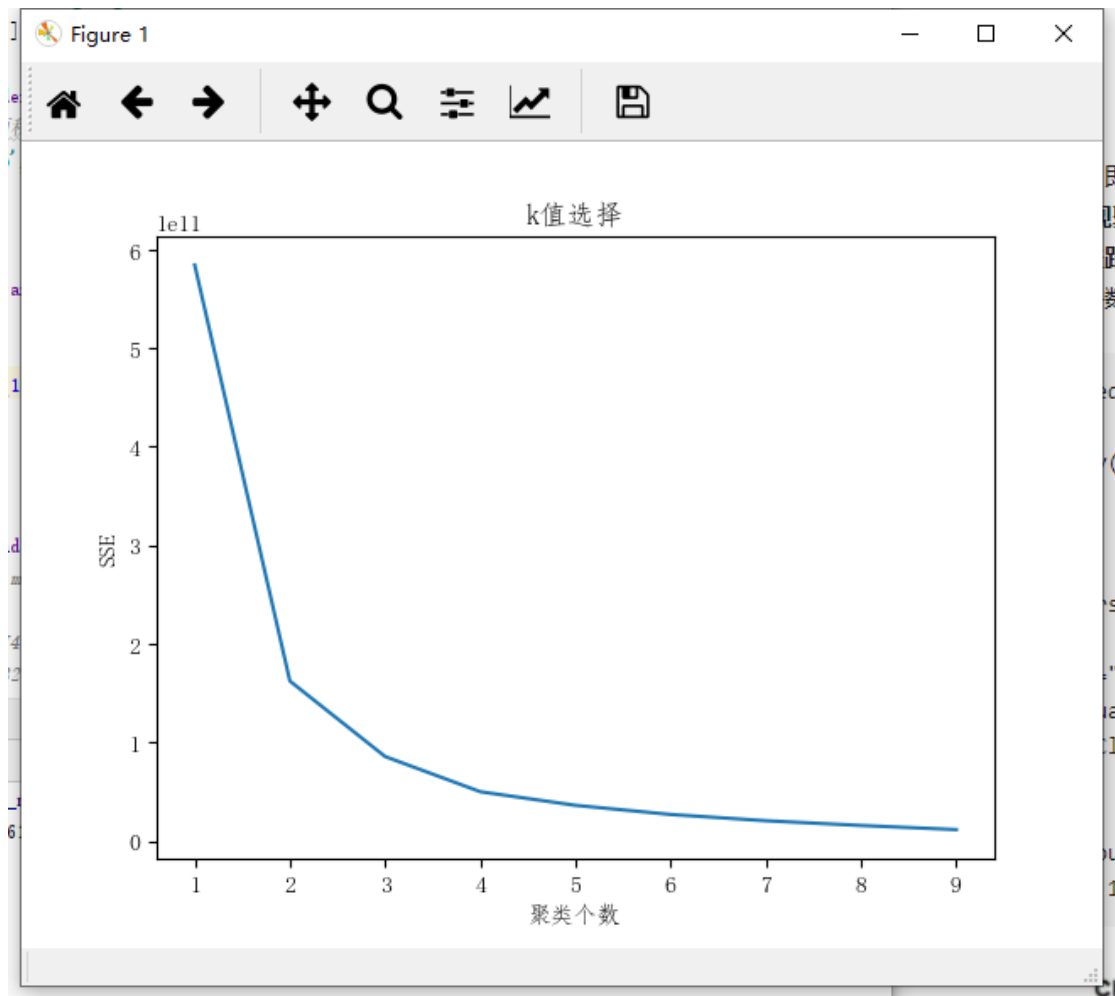
从表中可以看到，square（面积）与 rent（租金）呈正相关，square（面积）也与 time\_delta（上架时长）呈弱正相关，而最近地铁距离与其余三个变量几乎不相关。因此为了便于后续分析，将地铁信息这一变量进行二分类化，代表有无地铁信息。

4 租房房源 k-means 聚类分析

本文选取了数值型变量租金(rent)、面积(square)、上架时间间隔(time\_delta)进行了 k-means 聚类分析。

4.1k 值选择

在聚类之前需要确定**聚为几类**，即 k 值。根据聚类原则：组内差距要小，组间差距要大，本文绘制了不同类簇下的组内离差平方和图。



从图中可以看到，聚类个数达 4 个时，SSE 取值接近平缓，因此 k 值选择为 4。

## 4.2 聚类结果统计

聚类后各个类别的数量如下表。

0	1	2	3
8528	3183	1664	5421

聚类中心的结果如下表所示。

类别	租金	面积	上架时长
0	3918.60096154	72.89774859	54.47678236
1	13626.01916431	132.40182218	69.71190701
2	21248.48617788	179.64963942	97.20552885
3	7756.37114923	96.8544549	53.36081904

结合前文 3.2.5 户型与面积的关系推测，这四个中心可以分别代表 1 室 2 厅 1 卫、3 室 2 厅 2 卫、3 室 2 厅 3 卫、2 室 2 厅 1 卫这 4 种户型。

区域与聚类结果数量统计

cls_result	0	1	2	3	
area					
上海周边		8	0	0	1
嘉定	922	38	21	213	
奉贤	497	5	0	14	
宝山	720	48	9	350	
崇明	5	0	2	1	
徐汇	420	388	178	464	
普陀	423	274	77	422	
杨浦	192	73	30	174	
松江	1308	106	54	439	
浦东	2035	822	404	1282	
虹口	161	103	59	159	
金山	23	0	0	3	
长宁	188	348	218	359	
闵行	833	297	186	636	
闸北	143	108	53	175	
青浦	511	98	87	254	
静安	69	256	147	273	
黄浦	70	219	139	202	

有无地铁（0 表示无地铁信息）与聚类结果数量统计

cls_result	0	1	2	3
subway				
0	4554	562	349	1519
1	3974	2621	1315	3902

### 4.3 聚类结果解释

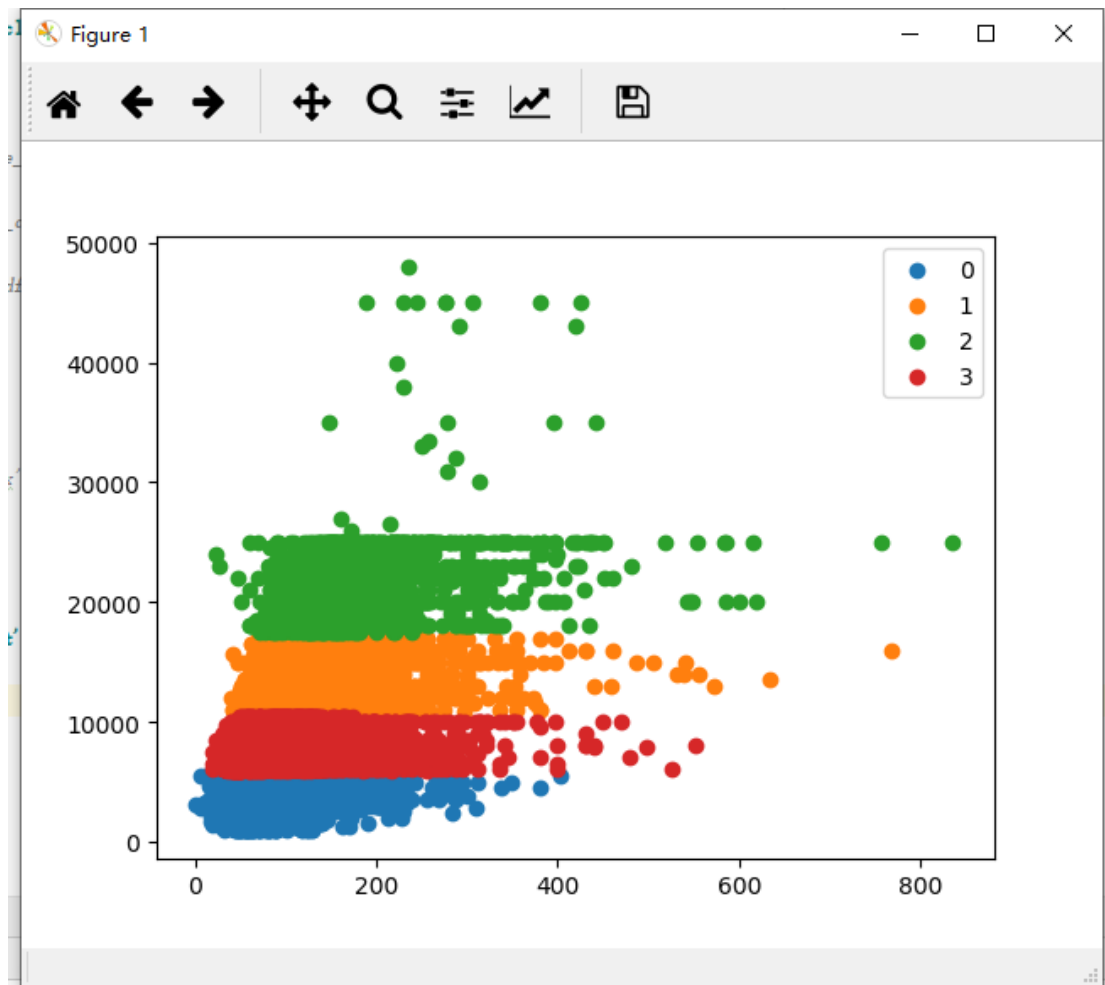
整体来看，可以将聚类结果合成如下：

**大户型**（3室2厅2卫、3室2厅3卫），属于第1、2类，平均面积超过130 m<sup>2</sup>，平均租金也较高，且上架时长较长，分布于浦东、长宁等地。

**经济型**（1室2厅1卫），属于第0类，平均租金较低，且周边大多无地铁，分布于浦东、松江等地。

**中间型**（2室2厅1卫），属于第3类，平均租金介于大户型与经济型之间，房源较多，平均上架时间相对较短，分布于浦东、闵行等地。

最后绘制聚类结果同租金与面积散点图如下



## 5 租房租金决策树回归分析

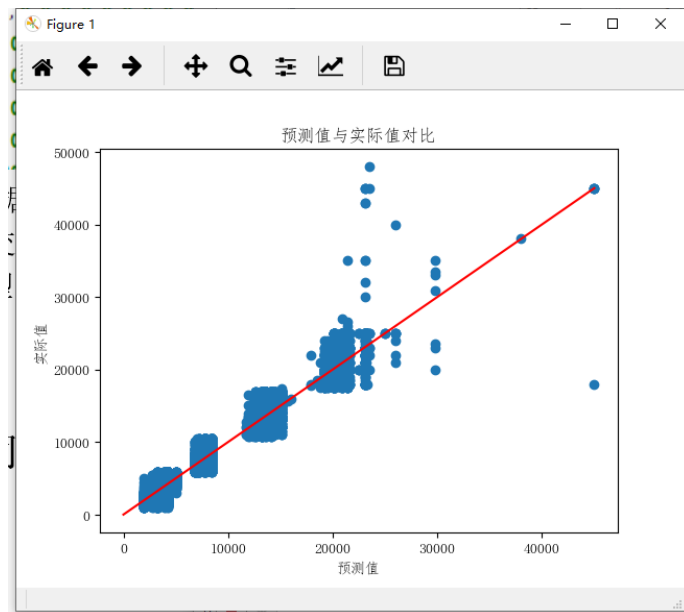
对于租金价格预测，本文选择了（面积、地铁、区域、户型、楼层、聚类水平）这些变量作为特征，构建了决策树回归模型。

### 5.1 预处理编码

首先需要对非数值变量（区域、楼层、户型）进行编码，编码字典如下表所示。

{'静安': 0, '黄浦': 1, '闸北': 2, '嘉定': 3, '宝山': 4, '上海周边': 5, '普陀': 6, '徐汇': 7, '长宁': 8, '松江': 9, '崇明': 10, '奉贤': 11, '杨浦': 12, '浦东': 13, '闵行': 14, '青浦': 15, '金山': 16, '虹口': 17}
{'中楼层': 0, '未知': 1, '高楼层': 2, '低楼层': 3}
{'3室2厅1卫': 0, '3室1厅1卫': 1, '2室2厅1卫': 2, '3室1厅2卫': 3, '2室1厅2卫': 4, '其他户型': 5, '4室2厅3卫': 6, '2室2厅2卫': 7, '3室2厅3卫': 8, '1室1厅1卫': 9, '4室2厅4卫': 10, '3室2厅2卫': 11, '1室0厅1卫': 12, '2室0厅1卫': 13, '1室2厅1卫': 14, '4室2厅2卫': 15, '2室1厅1卫': 16}





可以看到,模型的输出拟合效果也较好,模型可以初步用于上海租房房源市场租金预测。

## 6 租房介绍词云分析

链家网上的租房介绍（房源描述）如下图所示。

### 房源描述



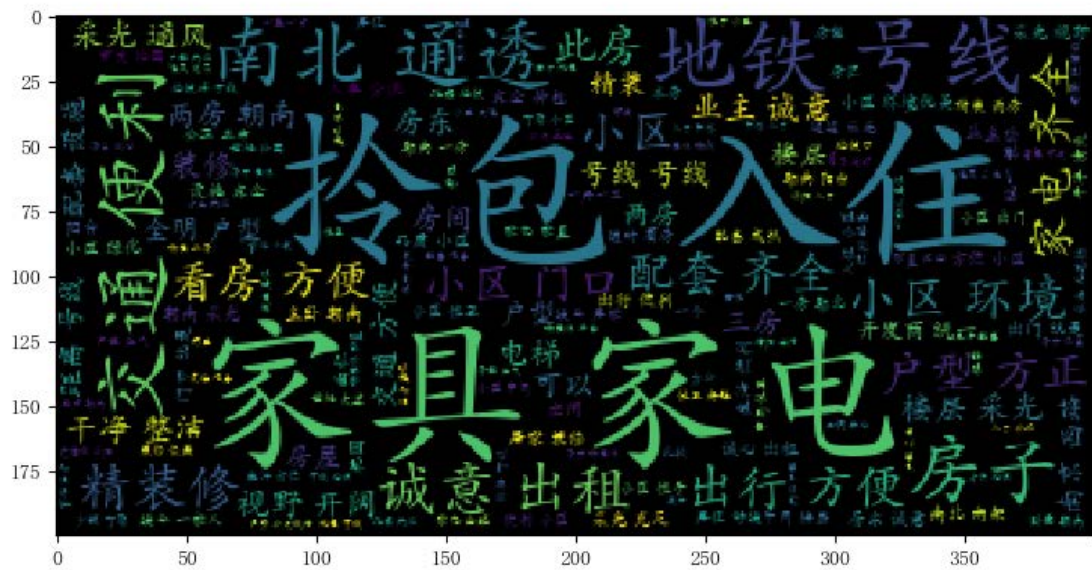
张亚萍 在线咨询  
链家 经纪人

4000134576转2481

【房源亮点】 本房位于紫晶南园,紫晶南园小区位于剑川路上,105平三房两厅三开间,两房朝南一房朝北  
【交通出行】 出小区门向左向右都有公交站台,729,958,956,闵吴线等等公交都在此公交站台停靠  
【出租原因】 房子目前空置状态,半年短租,房东其他地方有房子住,房子可接受一家人

本文对上述文本信息进行了词云分析。首先对房源亮点、装修描述、出租原因、交通出行、户型介绍这些无效频繁词进行了剔除,并且对文本中的标点、空格等符号也进行了剔除,之后利用 jieba 对文本进行了分词,最后生成词云,如下图所示。





从图中可以看到，家具、家电、拎包入住成为词云的亮点高频词，可以推测租房房源市场可能对快速入住有较大需求，因而在租房介绍上重点描述了这方面的信息。

## 7 总结

本文基于上述分析，初步得出上海租房市场如下结论：

1. 房源集中分布于浦东、闵行、松江等地；
2. 房源户型、楼层种类很多，能够满足市场多样化需求；
3. 经济型房源大多无周边地铁信息；
4. 上架时长久的房源普遍为大户型房源；
5. 租金主要与面积有关，与区域、有无地铁、户型也有一定关系，而与楼层几乎没有关系；
6. 租房市场对快速入住有较高的需求。