

大数据在房产领域的实践



蔡白银

QCon

2016.10.20~22

上海·宝华万豪酒店

全球软件开发大会 2016

[上海站]



购票热线: 010-64738142

会务咨询: qcon@cn.infoq.com

赞助咨询: sponsor@cn.infoq.com

议题提交: speakers@cn.infoq.com

在线咨询(QQ): 1173834688

团·购·享·受·更·多·优·惠

7折 优惠(截至06月21日)
现在报名, 立省2040元/张

自我介绍

- 蔡白银

毕业于北京大学，在大数据数据挖掘领域有多年的经验, 目前就任链家网大数据架构师，负责链家网大数据体系的建设，运用大数据挖掘大数据价值助力房产领域的O2O，提升房屋买卖体验，使买卖房屋不再难

自我介绍

- 蔡白银

毕业于北京大学，在大数据数据挖掘领域有多年的经验，目前就任链家网大数据架构师，负责链家网大数据体系的建设，运用大数据挖掘大数据价值助力房产领域的O2O，提升房屋买卖体验，使买卖房屋不再难

我在链家网做大数据



你去链家当中介了？！

就一中介，还大数据（翻白眼）？！！



友谊的小船说翻就翻



提纲

- 蜀道难难于上青天

- 行困难而正确之事

- 往事可鉴未来可追



蜀道难难于上青天

- 客少、物少——数据来源少
- 买卖行为少周期长——行为数据稀少
- 线下行为重容易分流——线上线下难打通
- 业务复杂性——分析挖掘无坦途



提纲

- 蜀道难难于上青天
- 行困难而正确之事
- 往事可鉴未来可追



提升服务品质的环节

- 房源真实无虚假
- 合适的房屋给合适的人
- 房屋买卖不再难
- 缩短周期见效率
- 减少资源浪费

 新浪财经 生活 > 正文

演员方子哥买房被骗393万

2015年09月24日 02:40 北京晨报

2012年6月，方子哥的妻子颜女士为儿子看中了位于朝阳区千鹤家园的一套房屋，后通过 ~~某房产中介~~，约定以420万元的价格购买，并支付居间服务费92400元以及向担保公司支付17600元。此后，颜女士陆续向卖家李佳霖个人支付了393万元。她没有想到，李佳霖隐瞒已将房屋委托他人出售并办理最高额抵押贷款的事实，使用伪造的房屋所有权证骗取房款。后李佳霖因诈骗罪、合同诈骗罪被判处无期徒刑。

效果概述

1000万/天

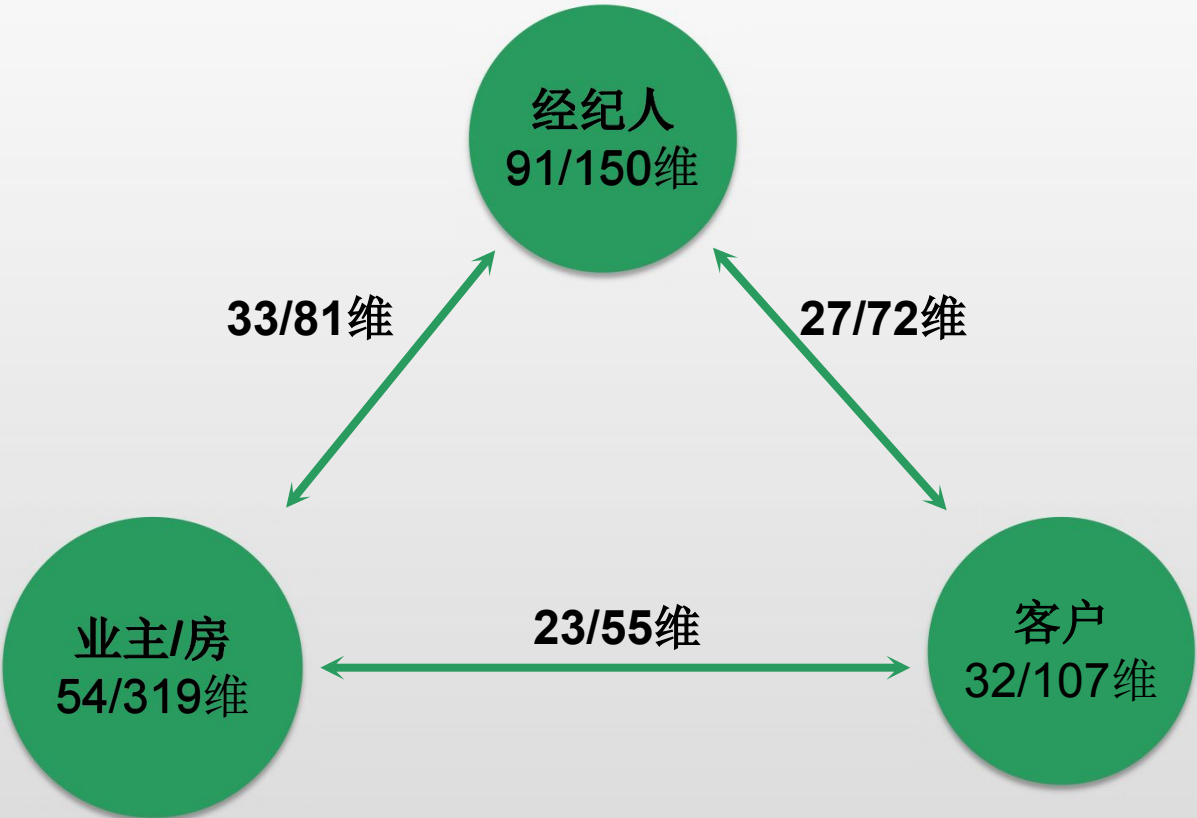
数百万/天

数T级别/天

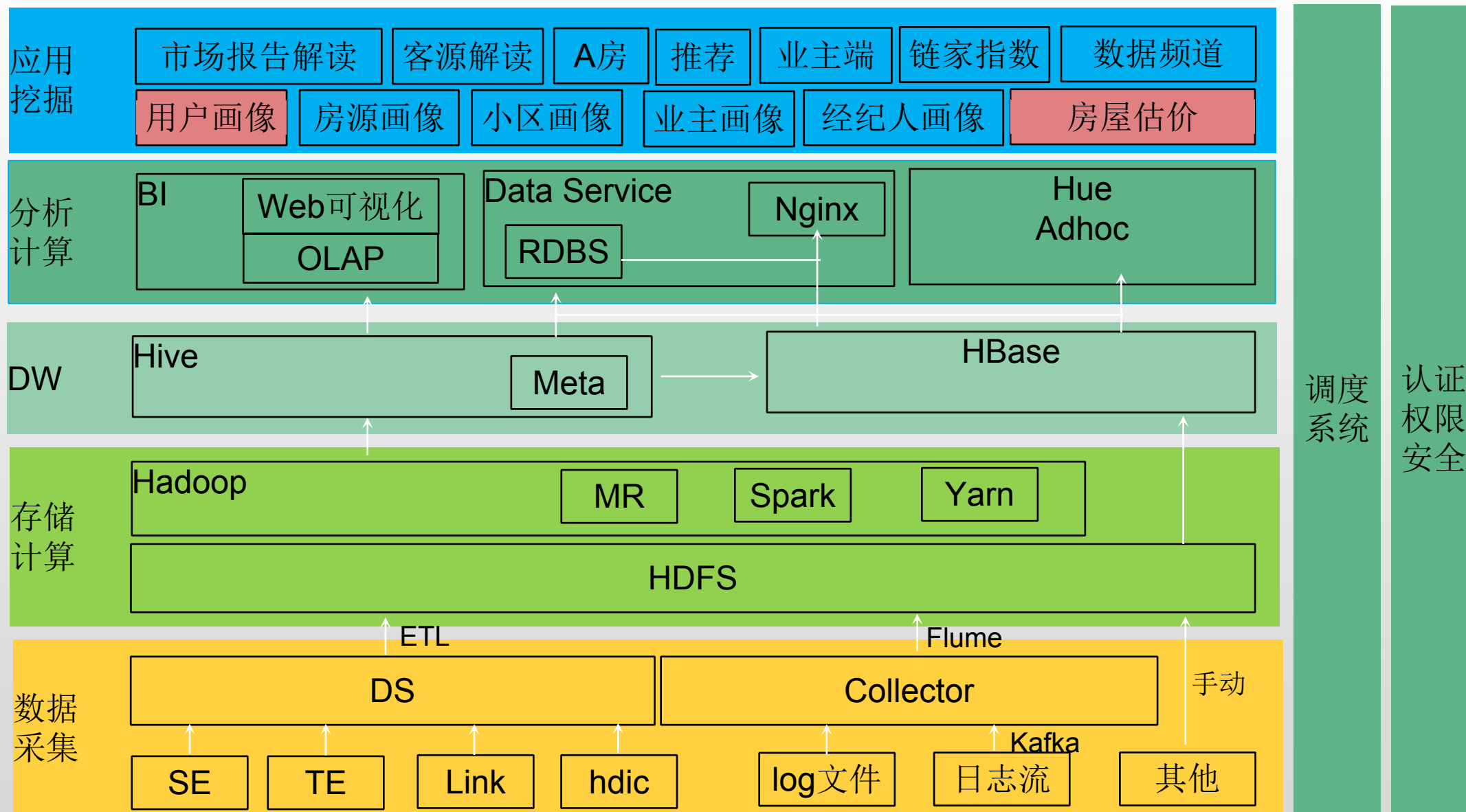
6000万

2300万

效果概述



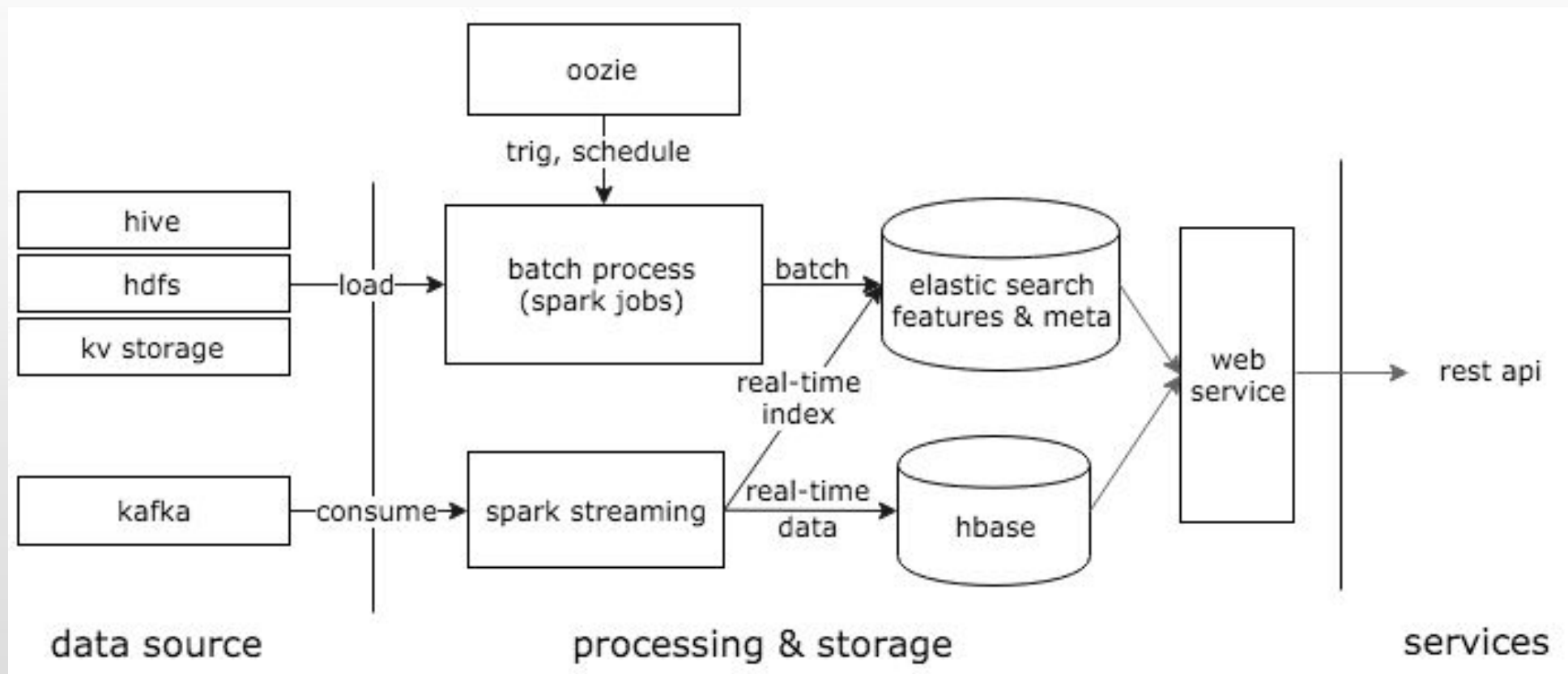
技术架构



用户画像

- *elasticsearch*, *hbase*, *spark*等成熟的开源数据存储、处理系统上
- *elasticsearch*存储、索引融合层全量数据，线上用户行为数据全量索引以及热数据
- *hbase*存储线上用户行为数据
- *spark*完成批量和流式数据处理，包括线下全量/增量数据导入，线上日志流处理并传送至*elasticsearch*集群。

用户画像

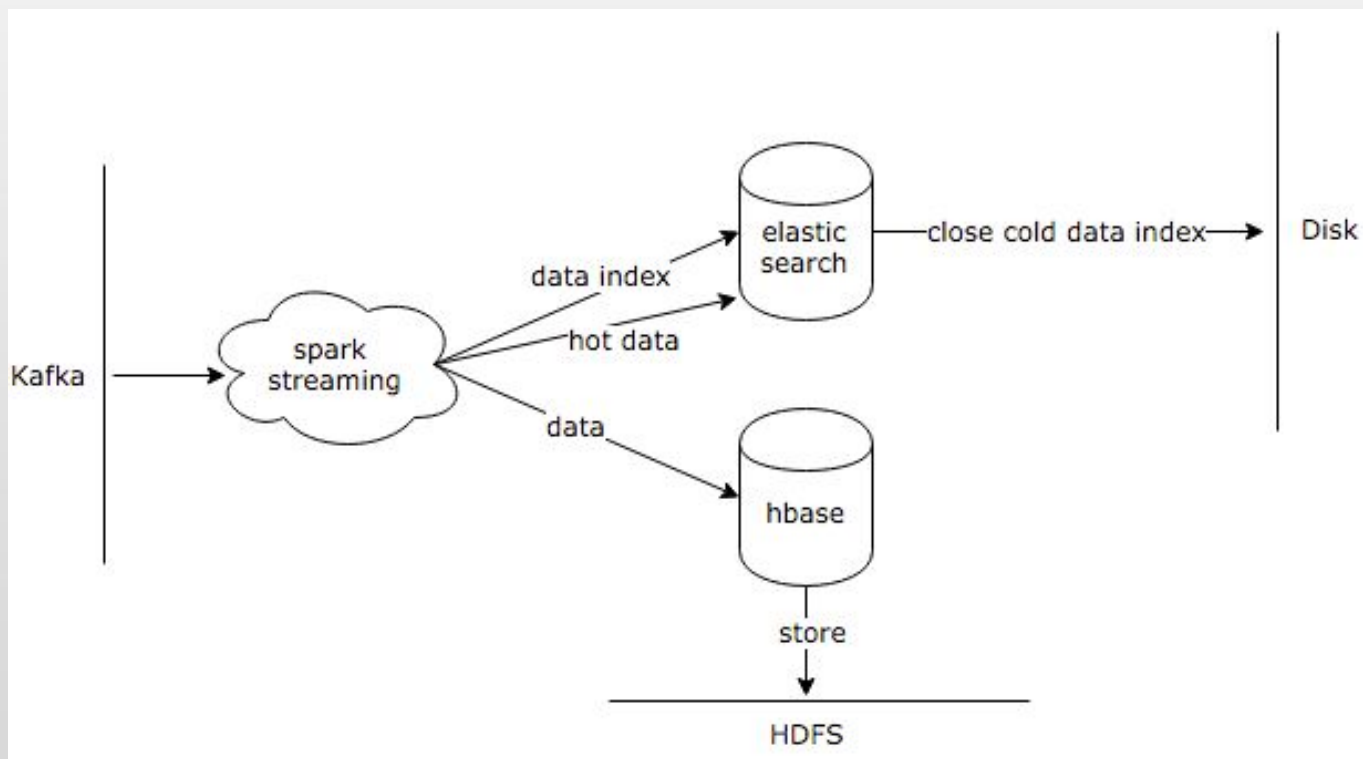


用户画像

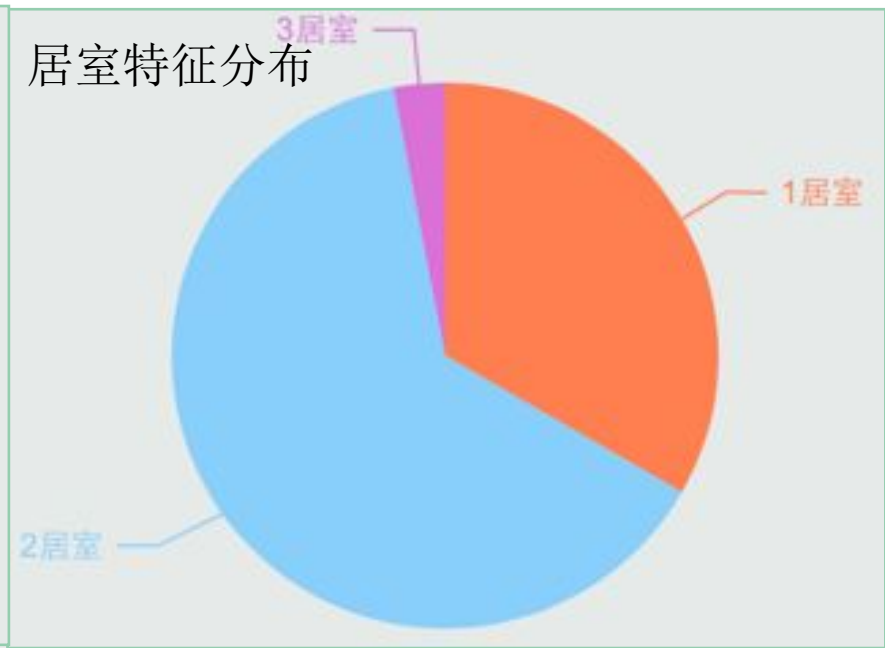
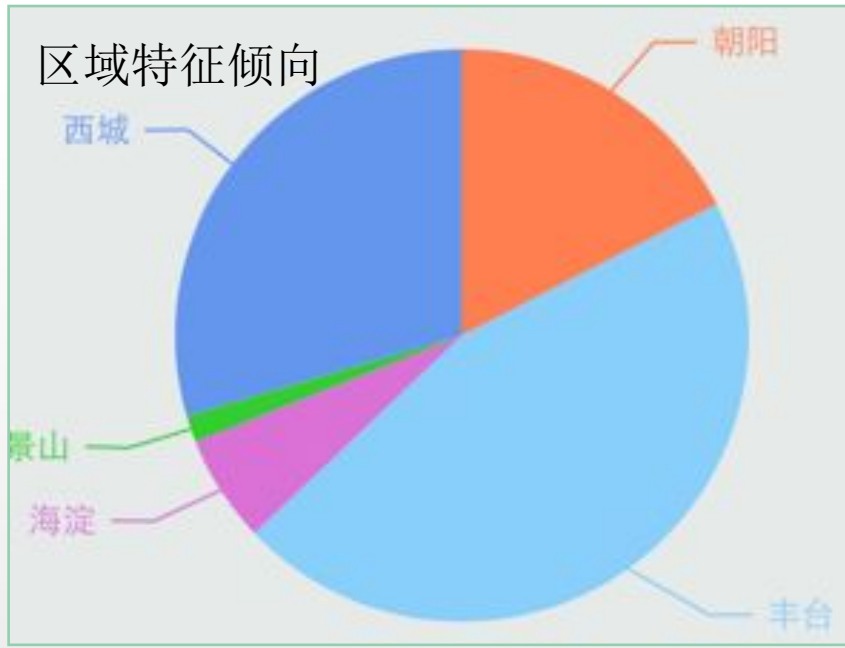
KV查询: 如通过手机号查询客源的一切数据

数据筛选: 如筛选西山商圈, 近三个月新增房源的小区名、挂牌价和房屋状态, 要求房屋必须是精装修或大于3居室

OLAP查询: 如查询海淀区2015年不同月份客源带看次数的分布



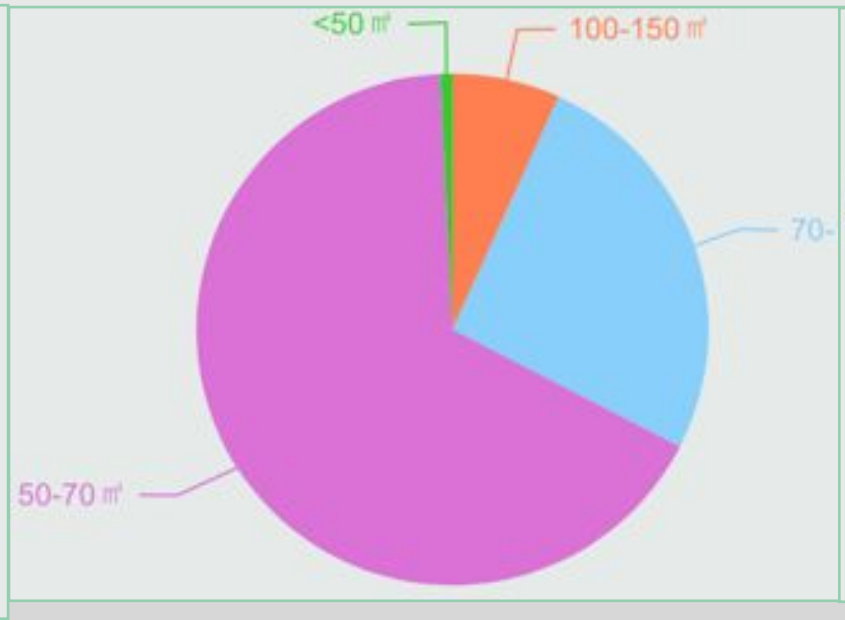
用户画像



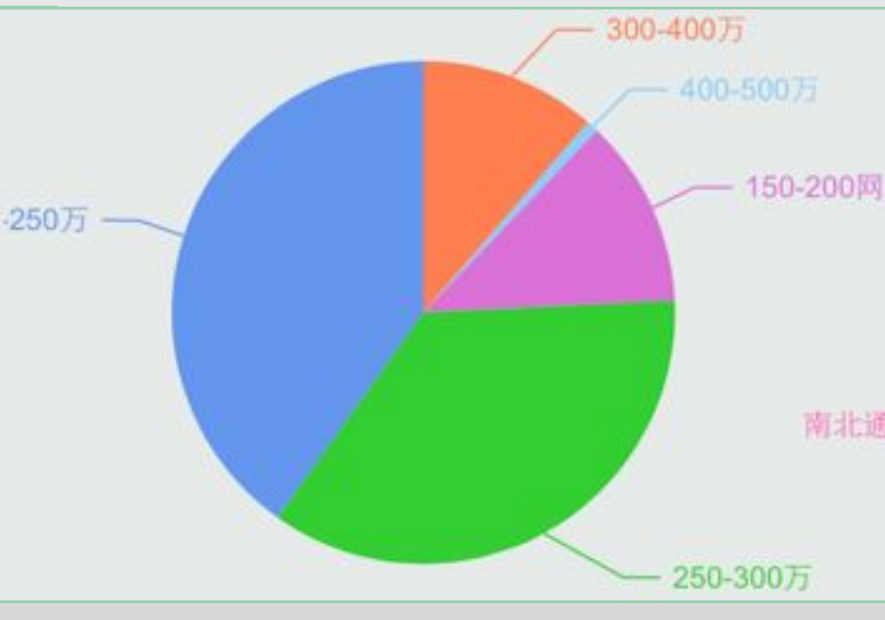
用户特征倾向



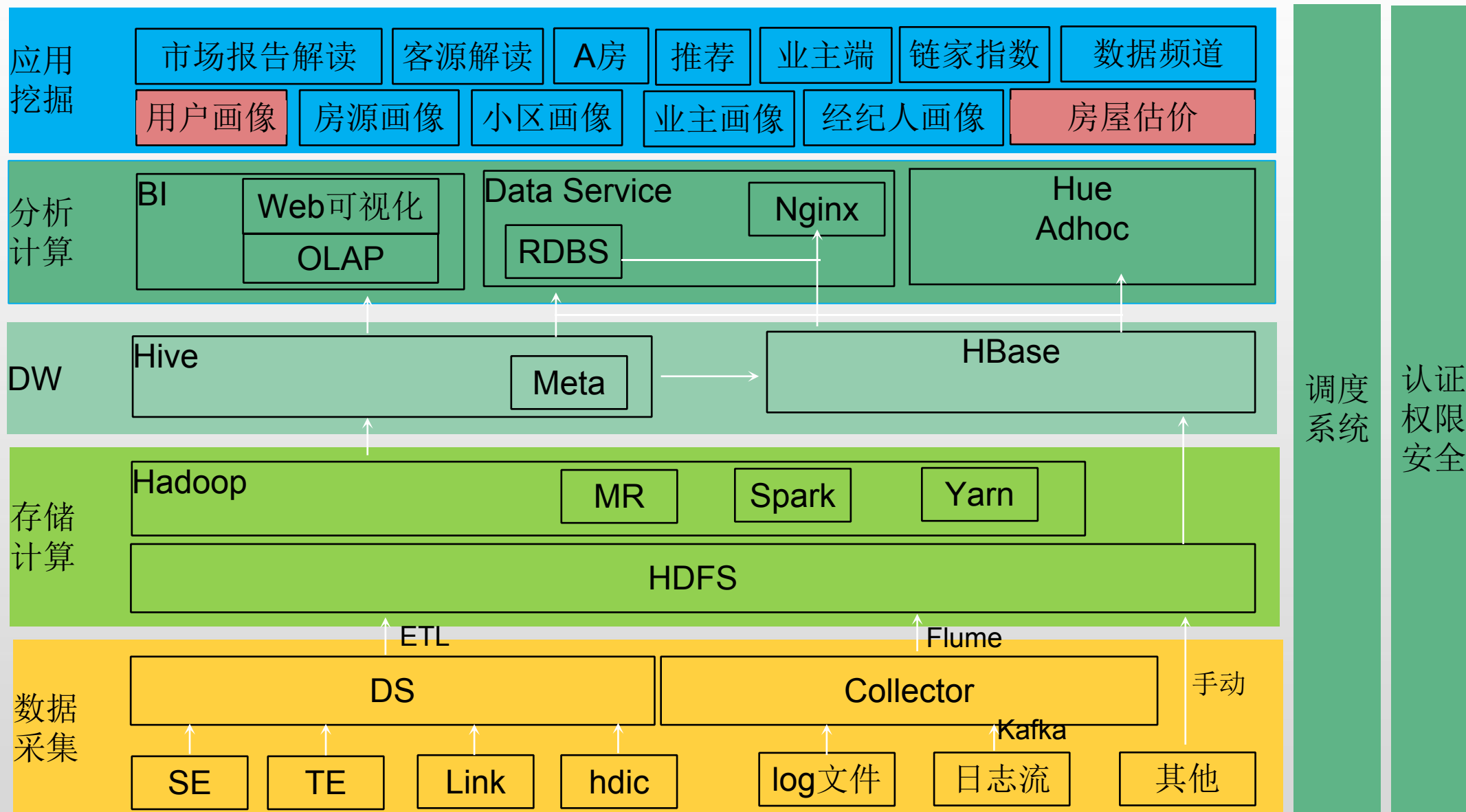
面积特征倾向



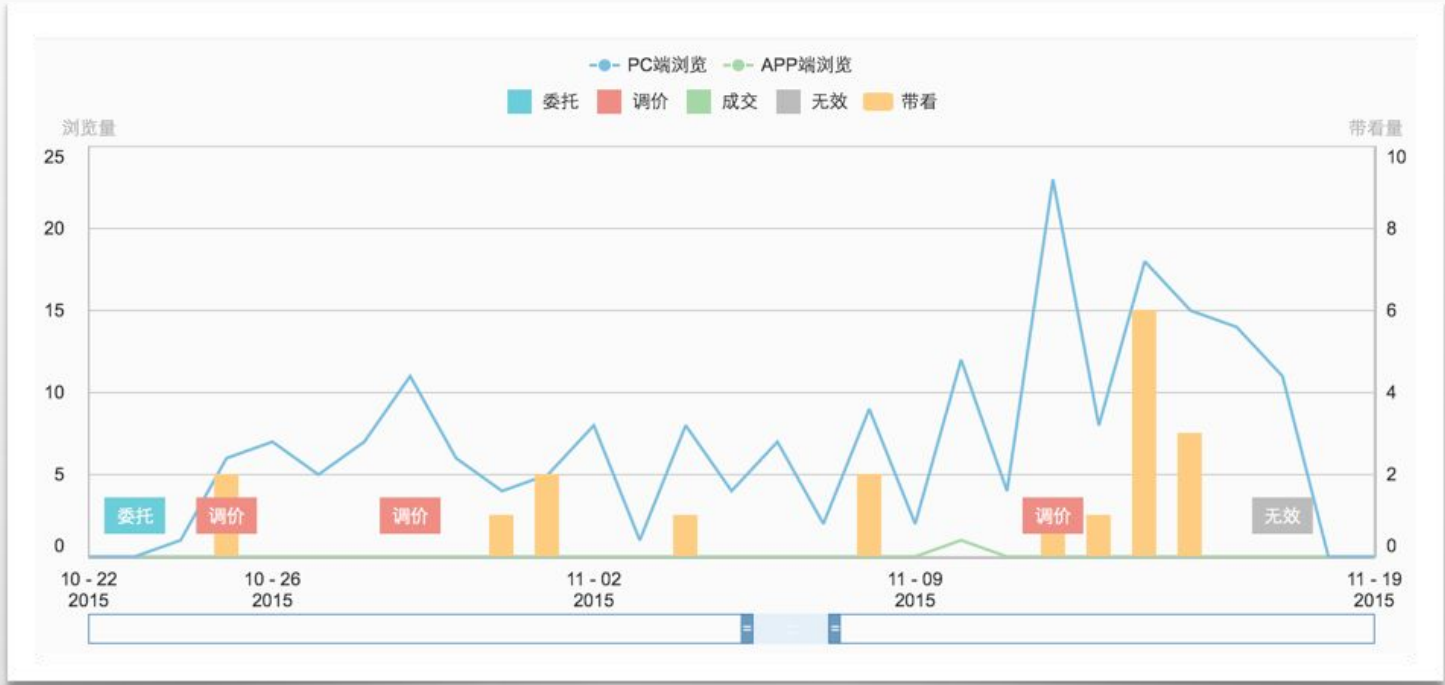
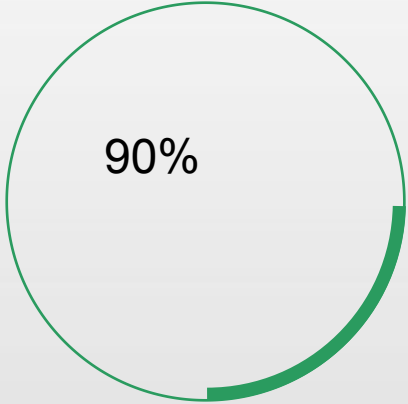
价格特征分布



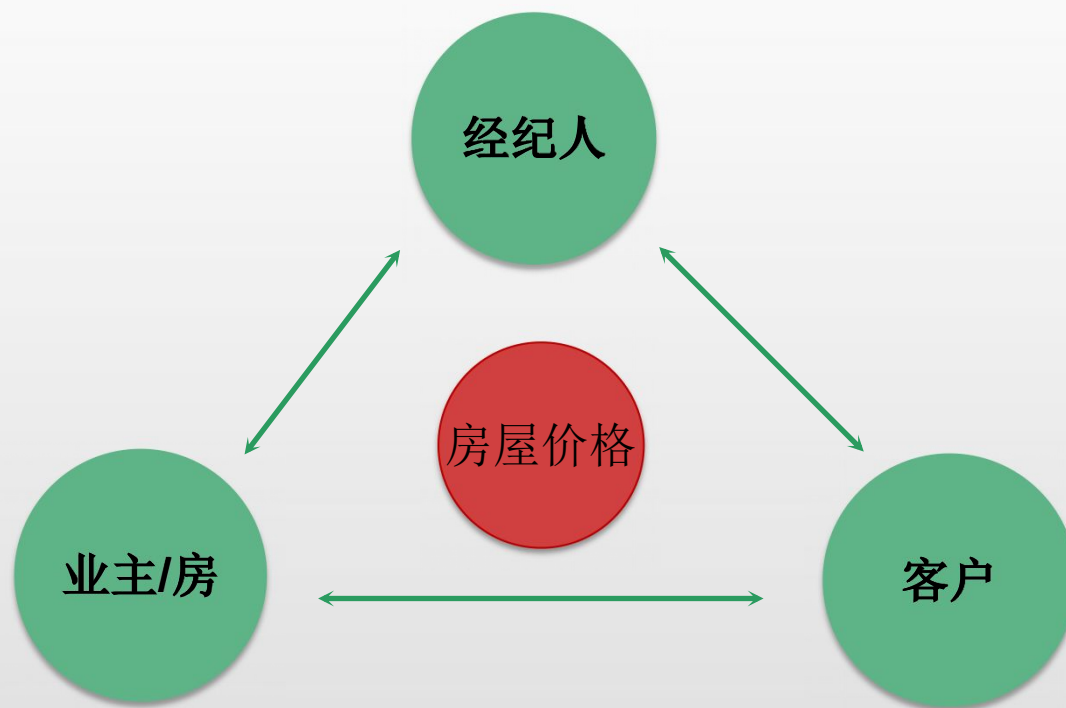
行困难而正确之事



房屋估价



房屋估价



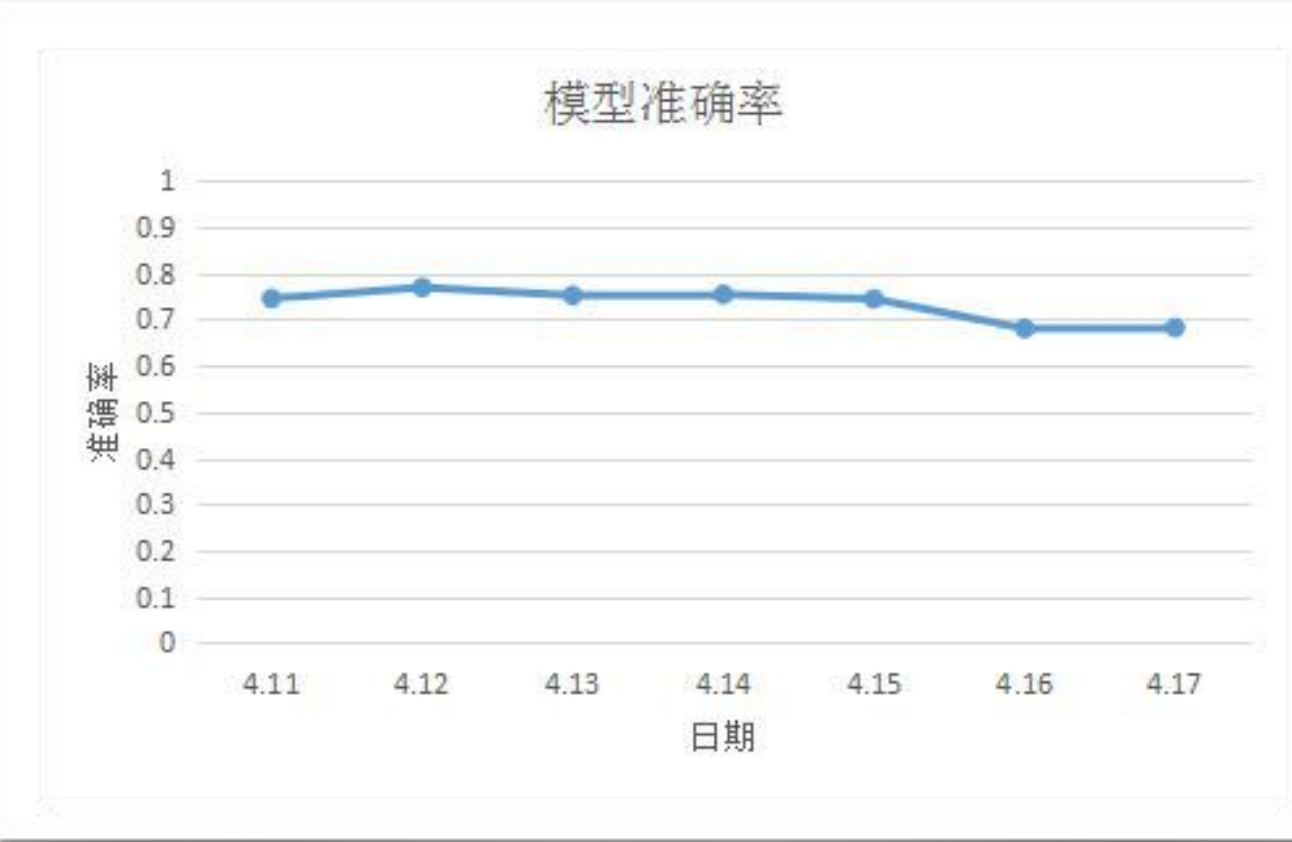
房屋估价



房屋估价

$$diff = \frac{|real_money - predict_money|}{real_money}$$

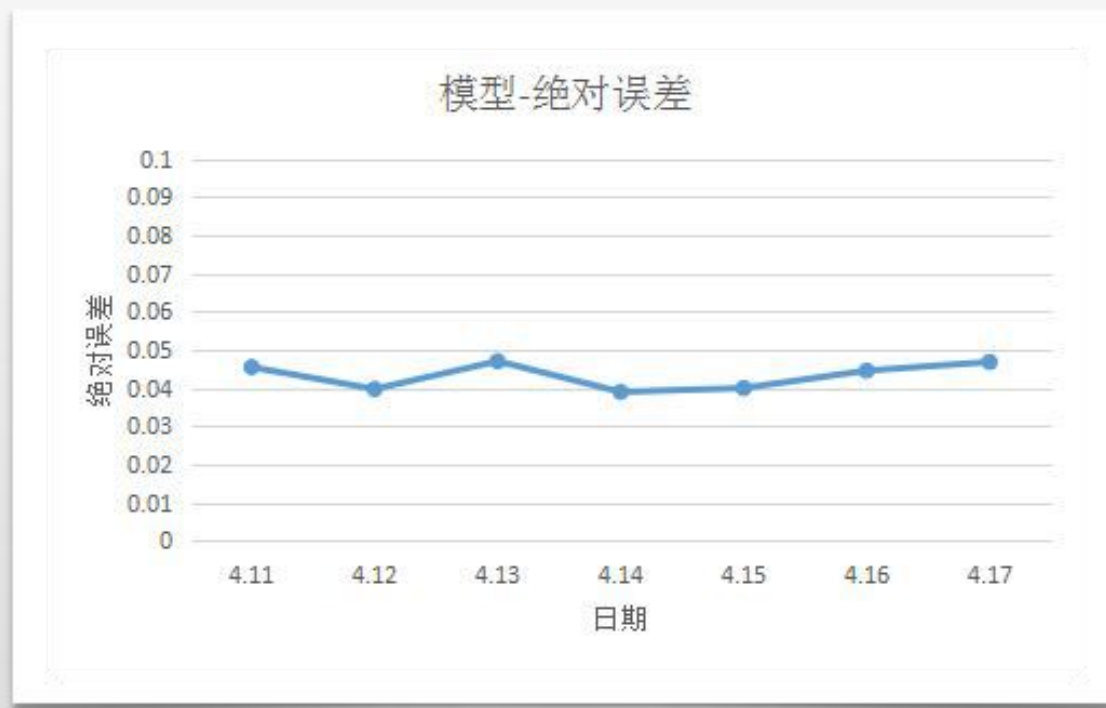
估价准确: $diff \leq 5\%$ $right_precent = \frac{right_amount}{total_amount}$



Data Coverage and Zestimate Accuracy Table
Choose a location type below to change data:
[Top Metro Areas](#)
[States/Countries*](#)
[National](#)

| | Zestimate Accuracy | Homes on Zillow | Homes With Zestimates | Within 5% of Sale Price | Within 10% of Sale Price | Within 20% of Sale Price | Median Error |
|---------------------------|--------------------|-----------------|-----------------------|-------------------------|--------------------------|--------------------------|--------------|
| Denver, CO | ★★★★★ | 957.7K | 887.3K | 40.5% | 66.9% | 88.8% | 6.5% |
| Detroit, MI | ★★★★ | 1.8M | 1.7M | 36.3% | 60.5% | 81.4% | 7.5% |
| Houston, TX | ★ | 2.2M | 1.9M | -- | -- | -- | -- |
| Kansas City, MO | ★ | 756.0K | 692.9K | -- | -- | -- | -- |
| Miami-Fort Lauderdale, FL | ★★ | 2.5M | 2.4M | 31.1% | 54.3% | 79.4% | 8.9% |
| Minneapolis-St Paul, MN | ★★★★ | 1.2M | 1.1M | 35.8% | 60.5% | 84.6% | 7.6% |
| New York, NY | ★★★★ | 5.3M | 4.9M | 32.4% | 55.3% | 77.0% | 8.6% |
| Orlando, FL | ★★★★ | 875.5K | 803.0K | 37.1% | 61.9% | 83.0% | 7.2% |
| Philadelphia, PA | ★★★★ | 2.1M | 2.0M | 35.5% | 58.1% | 79.1% | 7.8% |
| Phoenix, AZ | ★★★★★ | 1.7M | 1.5M | 43.0% | 68.0% | 87.8% | 6.2% |

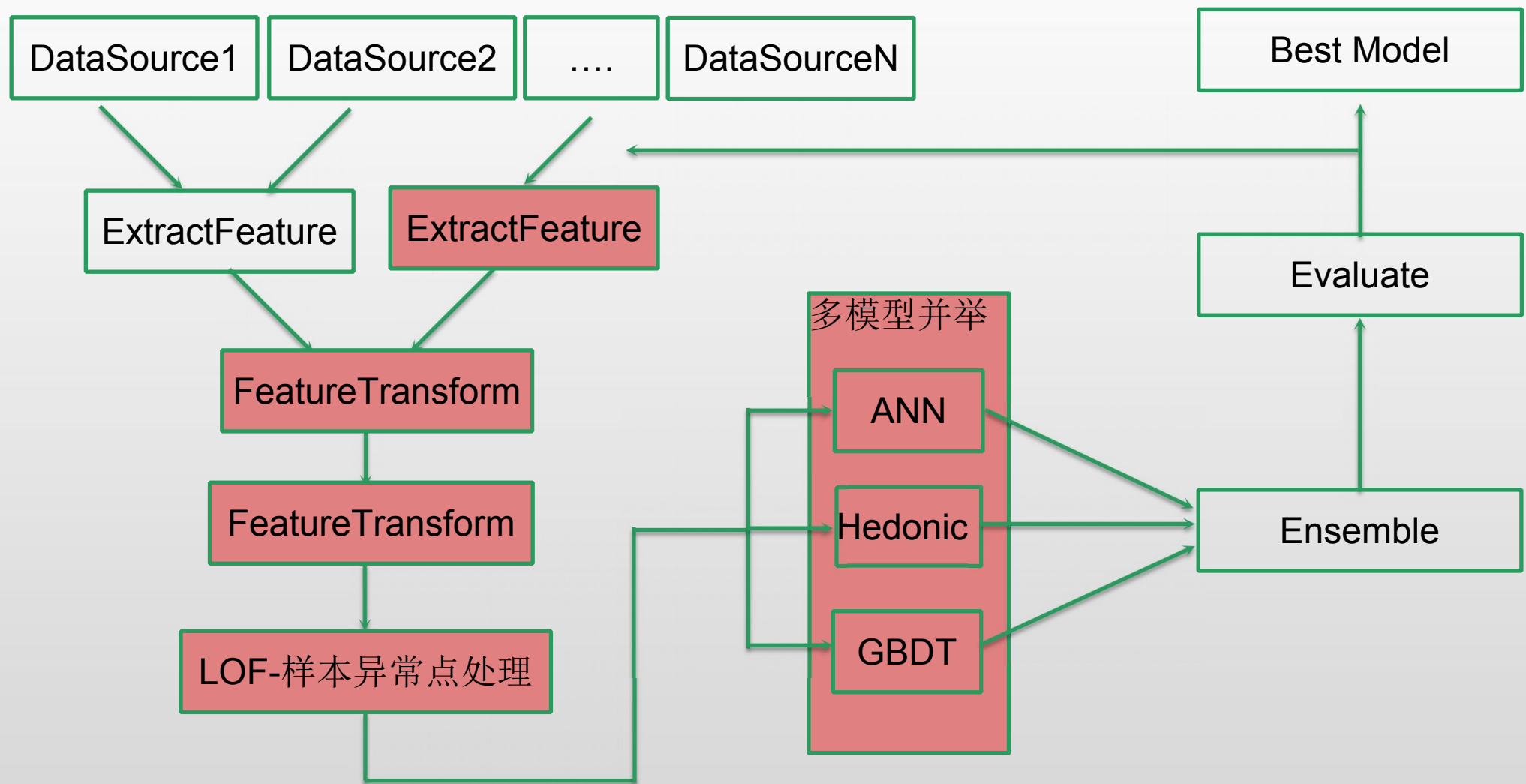
房屋估价



房屋估价



房屋估价



房屋估价

- LOF(Local Outlier Factor)算法是一种机遇密度的异常检测算法，
- 通过计算每个实例相对于其邻居的孤立情况来判断这个实例是否为离群点
- 为每一个每个实例计算一个异常分数，这个分数称为实例的局部离群因子（LOF）
- 较高的LOF值指示这个实例可能是异常的，较低的LOF值指示这个实例可能是正常的

$$LOF_k(p) = \frac{\sum_{o \in N_k(p)} \frac{lrd_k(o)}{lrd_k(p)}}{|N_k(p)|}$$

$$lrd_k(p) = \left[\frac{\sum_{o \in N_k(p)} reach-dist_k(p, o)}{|N_k(p)|} \right]^{-1}$$

提纲

- 蜀道难难于上青天
- 行困难而正确之事

往事可鉴未来可追

往事可鉴未来可追

链家金融

智能家电

链家装修

家政服务

亿万
房产
O2O
服务
平台
打造
住的
入口

3D看房

VR看房

社区服务

海外置业

The background image is a dimly lit interior of a modern home. In the foreground, a long wooden dining table is set with a bowl of fruit. To the right, a white kitchen island is visible. In the background, large glass doors and windows offer a view of a garden and a neighboring house. The overall atmosphere is calm and minimalist.

Thanks