

# 期末模型展示

欧阳语博 2022201462

中国人民大学

指标	Hackthon 分数	训练集	测试	删除异常值后测试集的总数 N
Xgboost	84.067	67990	444328	未删除异常值，测试集使用 10% 的总样本

2025 年 6 月 5 日

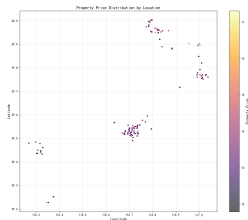
## 一、总体

舍弃缺失比例（0.5 以上）过高的数据；推断缺失值

## 二、环线位置



(a) 城市 0 价格热力图



(b) 城市 1 价格热力图

## 三、构建新的特征

比率：楼栋密度、卫室比；时间特征：交易年份、交易季度、挂牌时间

## 一、树的特性

树的每一次分裂是为了将不同的样本分开

## 二、分箱

处理离散数据

小区房屋、小区楼栋

## 三、处理高基数数据

设置一个阈值，防止过拟合：`threshold=10`

开发商、物业公司

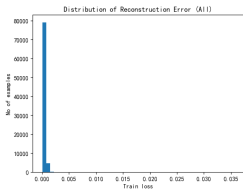
构造一个小型的神经网络，拟合特征之间的关系。以自身为训练目标

## 一、重构误差

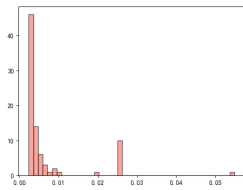
## 二、回顾特征工程

是否有特征被错误处理

## 三、检测异常值



(a) 重建误差分布图



(b) 重建误差极端值图

固定参数，调整树的数量

## 一、XGboost

learning rate、num boost round

## 二、随机森林

max features、max depth、n estimators

## 三、集成

加权、Stacking