

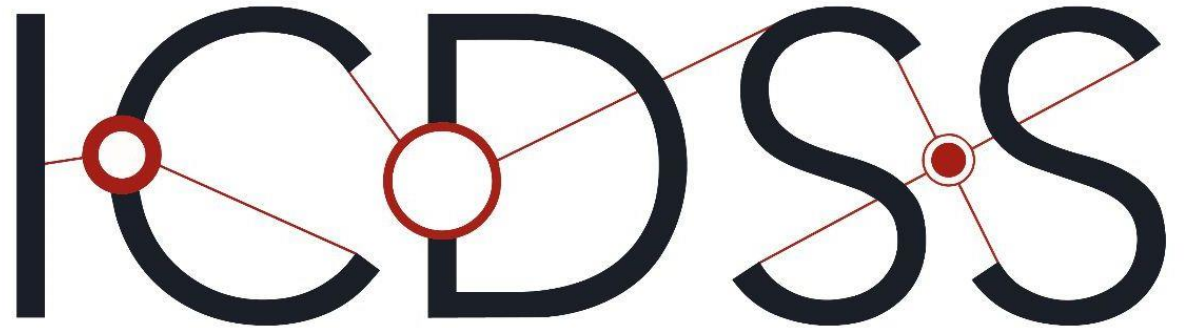


Imperial College
London

Boston Housing Challenge

Team Name: The Outliers

gh: <https://github.com/SamtheSaint/AIHACK21>



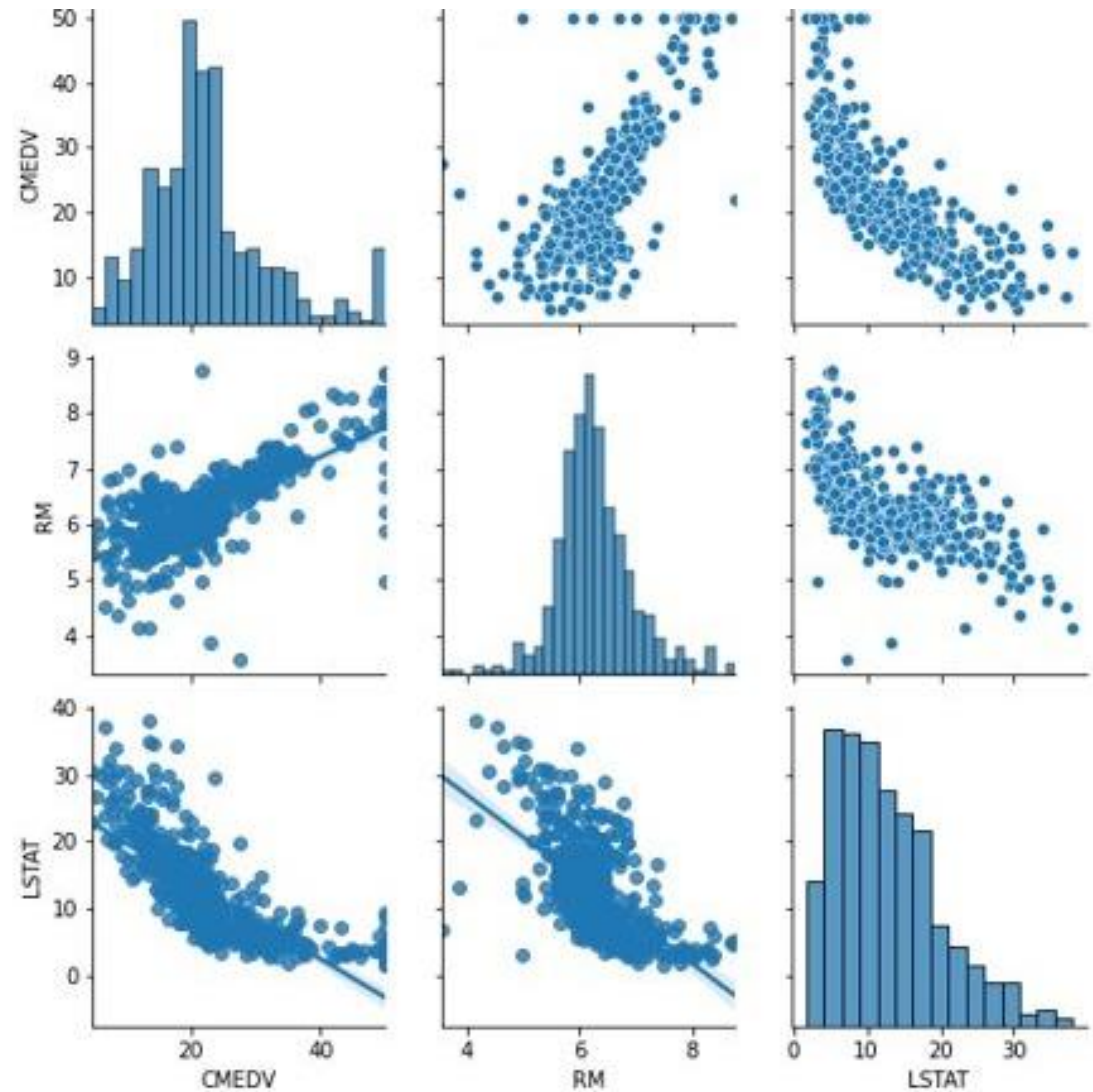
Initial Exploration of Dataset

- Computed correlation coefficients
- With regards to price the number of rooms and percentage of lower status population in the area are the most strongly correlated

TRACT	1	-0.22	-0.23	0.43	0.43	-0.55	0.37	-0.58	0.041	-0.57	0.31	-0.49	0.5	-0.83	-0.79	-0.53	0.37	-0.52
LON	-0.22	1	0.14	-0.32	-0.32	0.065	-0.22	0.063	-0.18	0.16	-0.26	0.2	-0.011	0.034	0.051	0.31	-0.018	0.2
LAT	-0.23	0.14	1	0.0097	0.0068	-0.084	-0.13	-0.041	-0.045	-0.069	-0.069	0.079	-0.083	-0.21	-0.17	-0.0045	0.11	0.046
MEDV	0.43	-0.32	0.0097	1	1	-0.39	0.36	-0.48	0.18	-0.43	0.7	-0.38	0.25	-0.38	-0.47	-0.51	0.33	-0.74
CMEDV	0.43	-0.32	0.0068	1	1	-0.39	0.36	-0.48	0.18	-0.43	0.7	-0.38	0.25	-0.38	-0.47	-0.51	0.33	-0.74
CRIM	-0.55	0.065	-0.084	-0.39	-0.39	1	-0.2	0.41	-0.056	0.42	-0.22	0.35	-0.38	0.63	0.58	0.29	-0.39	0.46
ZN	0.37	-0.22	-0.13	0.36	0.36	-0.2	1	-0.53	-0.043	-0.52	0.31	-0.57	0.66	-0.31	-0.31	-0.39	0.18	-0.41
INDUS	-0.58	0.063	-0.041	-0.48	-0.48	0.41	-0.53	1	0.063	0.76	-0.39	0.64	-0.71	0.6	0.72	0.38	-0.36	0.6
CHAS	0.041	-0.18	-0.045	0.18	0.18	-0.056	-0.043	0.063	1	0.091	0.091	0.087	-0.099	-0.0074	-0.036	-0.12	0.049	-0.054
NOX	-0.57	0.16	-0.069	-0.43	-0.43	0.42	-0.52	0.76	0.091	1	-0.3	0.73	-0.77	0.61	0.67	0.19	-0.38	0.59
RM	0.31	-0.26	-0.069	0.7	0.7	-0.22	0.31	-0.39	0.091	-0.3	1	-0.24	0.21	-0.21	-0.29	-0.36	0.13	-0.61
AGE	-0.49	0.2	0.079	-0.38	-0.38	0.35	-0.57	0.64	0.087	0.73	-0.24	1	-0.75	0.46	0.51	0.26	-0.27	0.6
DIS	0.5	-0.011	-0.083	0.25	0.25	-0.38	0.66	-0.71	-0.099	-0.77	0.21	-0.75	1	-0.49	-0.53	-0.23	0.29	-0.5
RAD	-0.83	0.034	-0.21	-0.38	-0.38	0.63	-0.31	0.6	-0.0074	0.61	-0.21	0.46	-0.49	1	0.91	0.46	-0.44	0.49
TAX	-0.79	0.051	-0.17	-0.47	-0.47	0.58	-0.31	0.72	-0.036	0.67	-0.29	0.51	-0.53	0.91	1	0.46	-0.44	0.54
PTRATIO	-0.53	0.31	-0.0045	-0.51	-0.51	0.29	-0.39	0.38	-0.12	0.19	-0.36	0.26	-0.23	0.46	0.46	1	-0.18	0.37
B	0.37	-0.018	0.11	0.33	0.33	-0.39	0.18	-0.36	0.049	-0.38	0.13	-0.27	0.29	-0.44	-0.44	-0.18	1	-0.37
LSTAT	-0.52	0.2	0.046	-0.74	-0.74	0.46	-0.41	0.6	-0.054	0.59	-0.61	0.6	-0.5	0.49	0.54	0.37	-0.37	1
	TRACT	LON	LAT	MEDV	CMEDV	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT

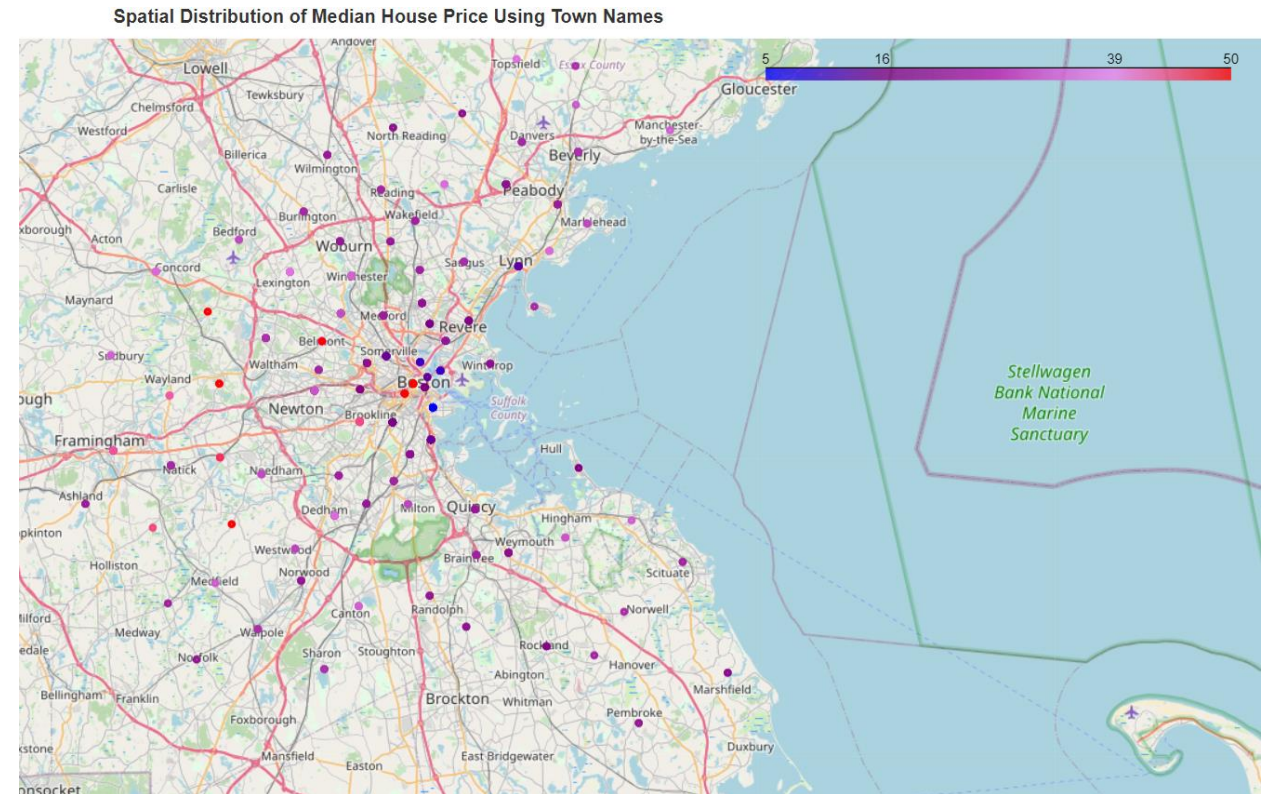
Initial Exploration of Dataset

- More evidence for the strong correlation between price and the number of rooms and the percentage of lower status population
- Notes on distribution of price and related factors

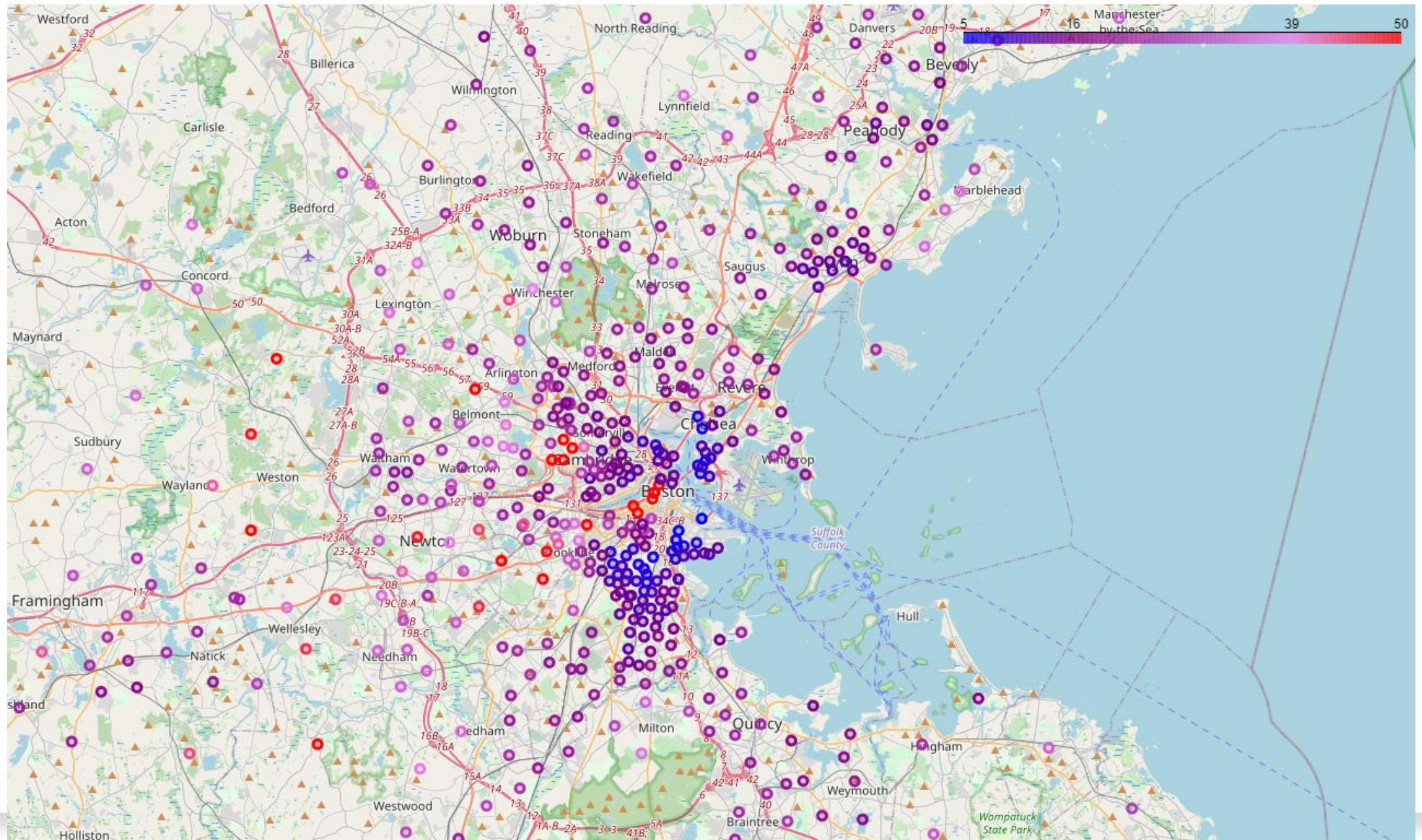


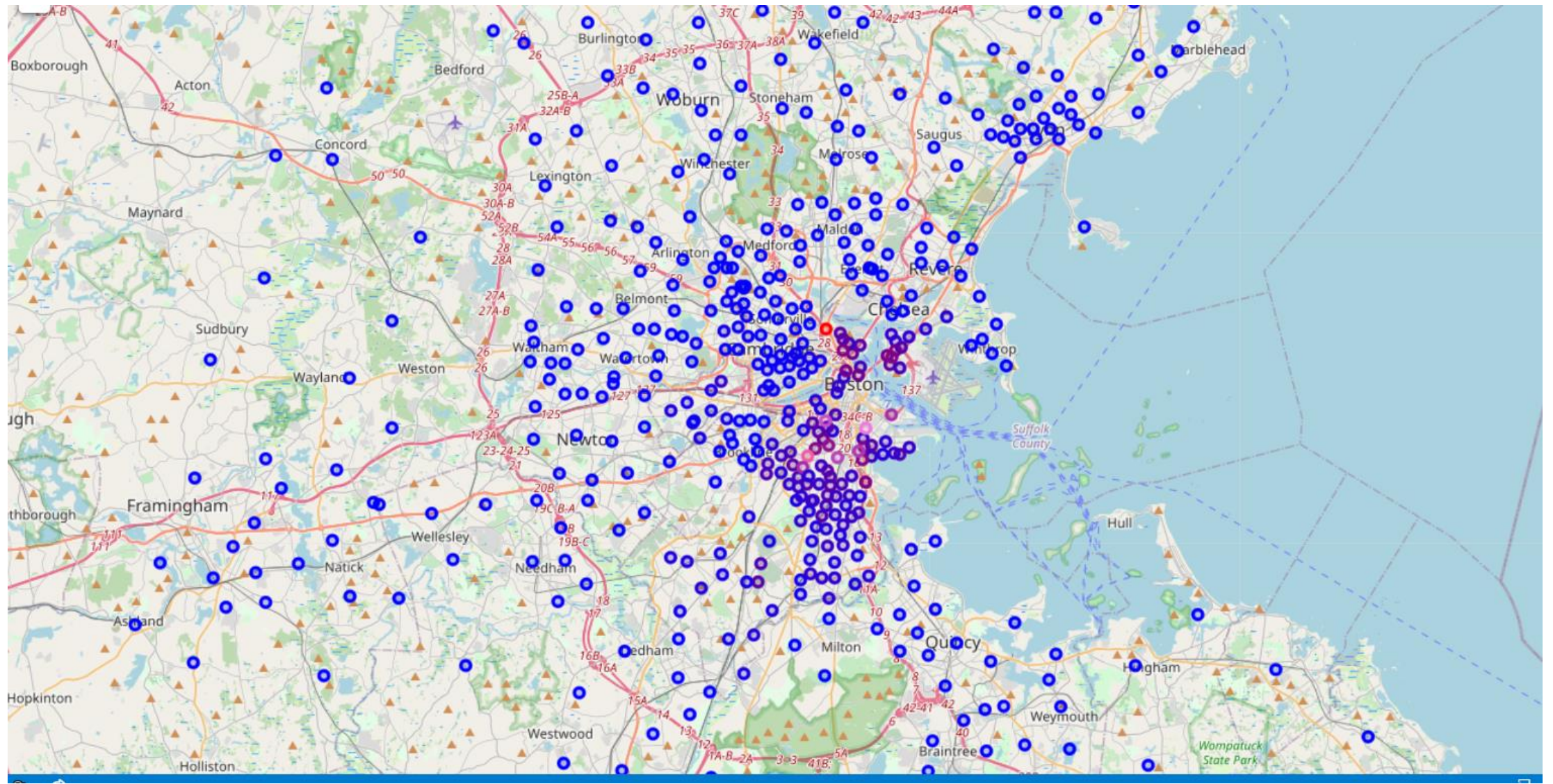
Spatial Analysis

- Initially provided longitude and latitude data was invalid, town names were used to plot distribution of metrics
- With Longitude/Latitude data provided the following day, the spatial distribution map was updated as shown in the next slide
- The effect of 'House Price Spatial Spillover' could be seen, and this was explored through clustering



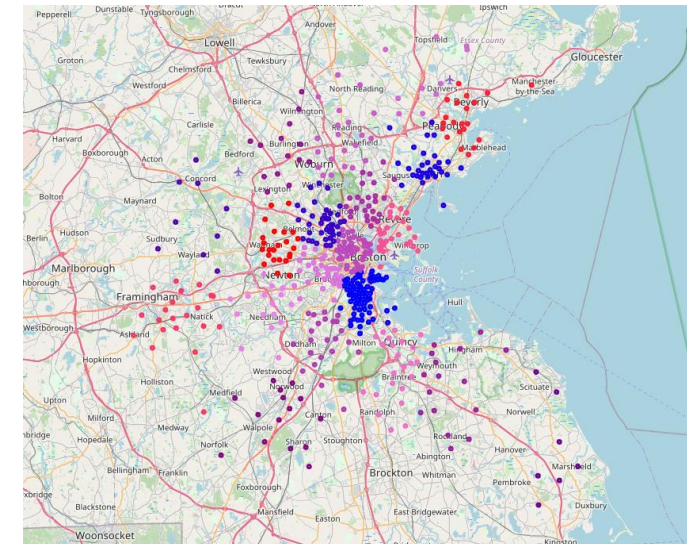
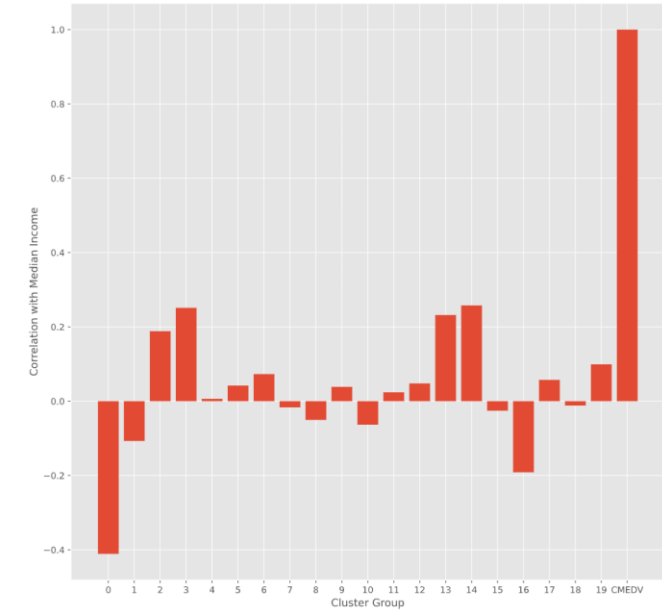
Spatial Distribution of Median House Price Using Lat/Lon Data





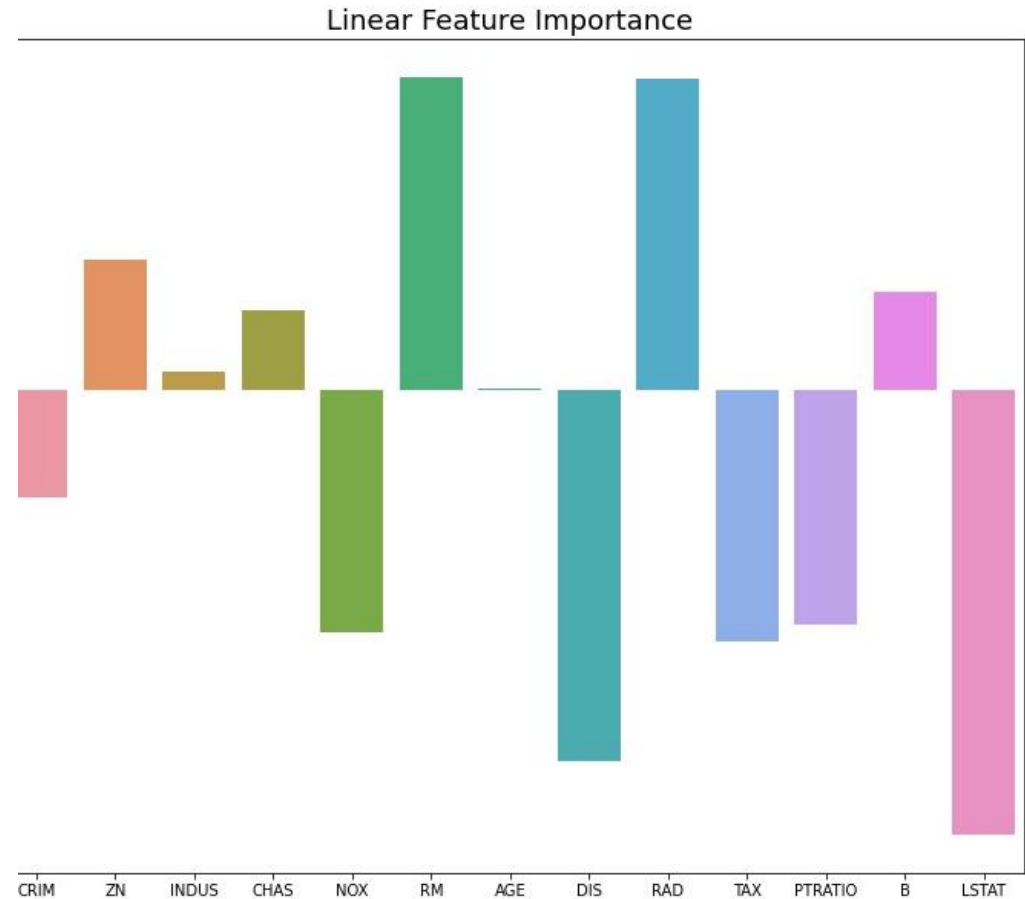
Clustering

- Using K-Means clustering to organize the spatial data.
- r^2 value of 41 %



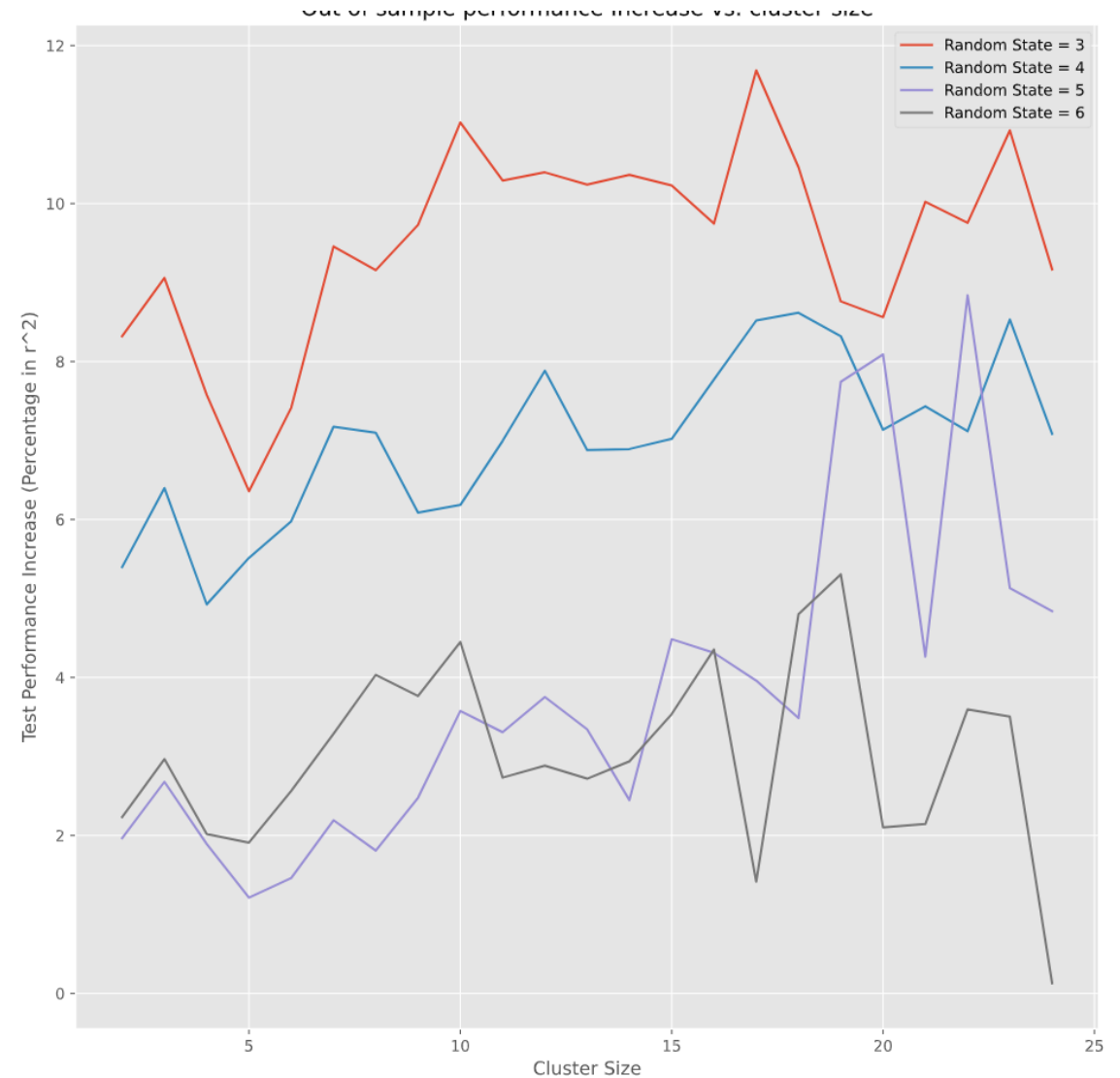
Regression on Price – Linear Regression

- Uses correlations between the features and the price as weights in the model.
- Results:
 - MSE: 22.95671894828107
 - R2: 0.7178446686547839



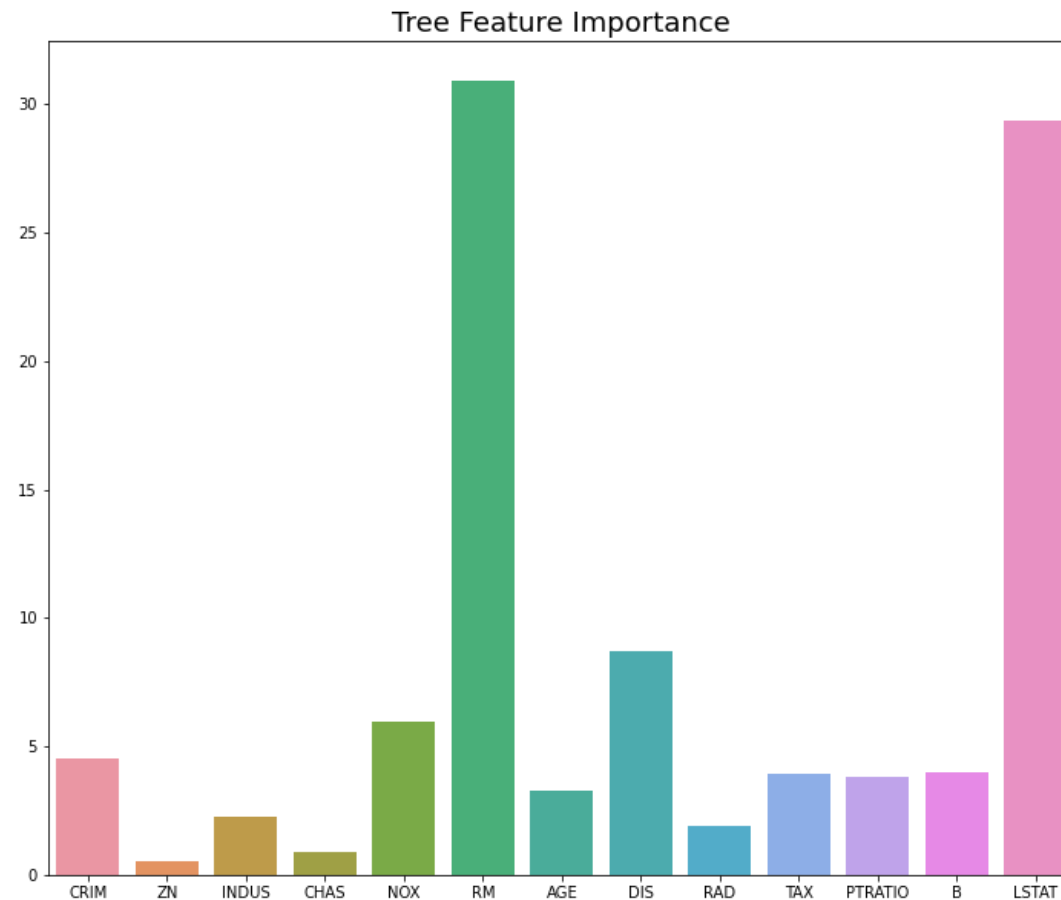
Using Cluster Group In Linear Regression Model

- With a Linear Regression model, from a base r^2 value of 68% using an 80-20% test- train split, by adding the cluster groups, performances of up to 11% could be achieved on the test set. The performance increase was measured for different cluster sizes and random states of the test train split to assess robustness



Regression on Price – Ensemble Methods

- Improved on linear regression using different Gradient Boosting trees: AdaBoost, XGBoost and CatBoost
- CatBoost had the best performance:
 - MSE: 9.112841668593683
 - R2: 0.8900240391522004
- Optimal hyper parameters are exposed in the notebook

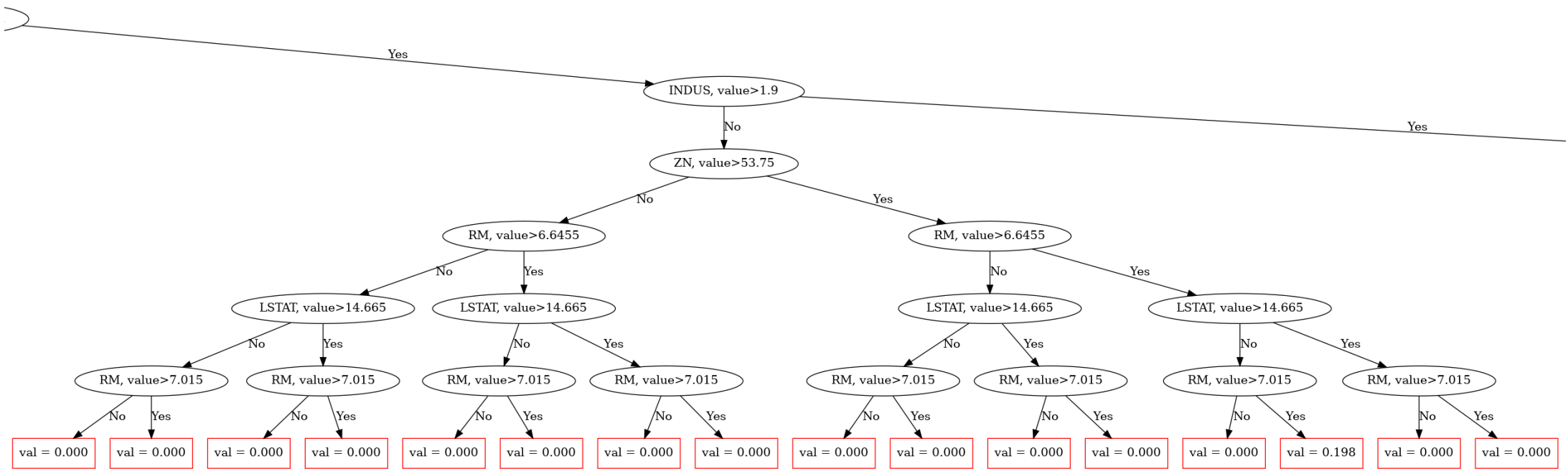


Hyper parameter search - Optuna

- We used optuna to search for optimal hyper parameters to use in the regression model coming up.
- Led to a x% increase in model performance



OPTUNA



Regression on Price – Ensemble Methods

- Extract from one of the trees used in the CatBoost Regressor
- Splits are made using most important features

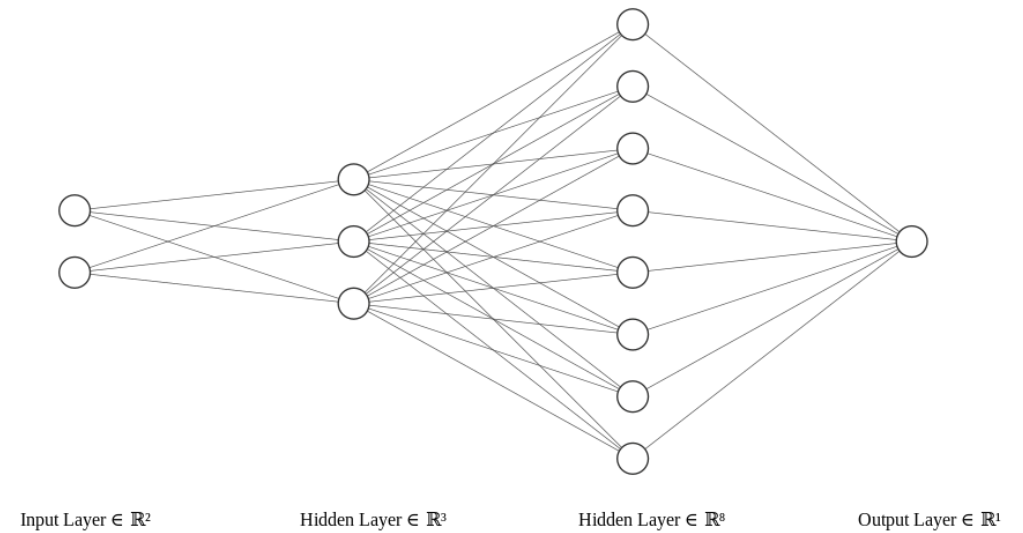
Adding the engineered cluster feature to the CatBoost model resulted in a performance boost

Regression on Price – Clustering Feature

CatBoost Model	R2	MSE
Without engineered spatial feature	0.8900240391522004	9.112841668593683
With engineered spatial feature	0.8934898067515067	8.725613355102025
% Improvement	0.389	4.25

Regression on Price – Neural Networks

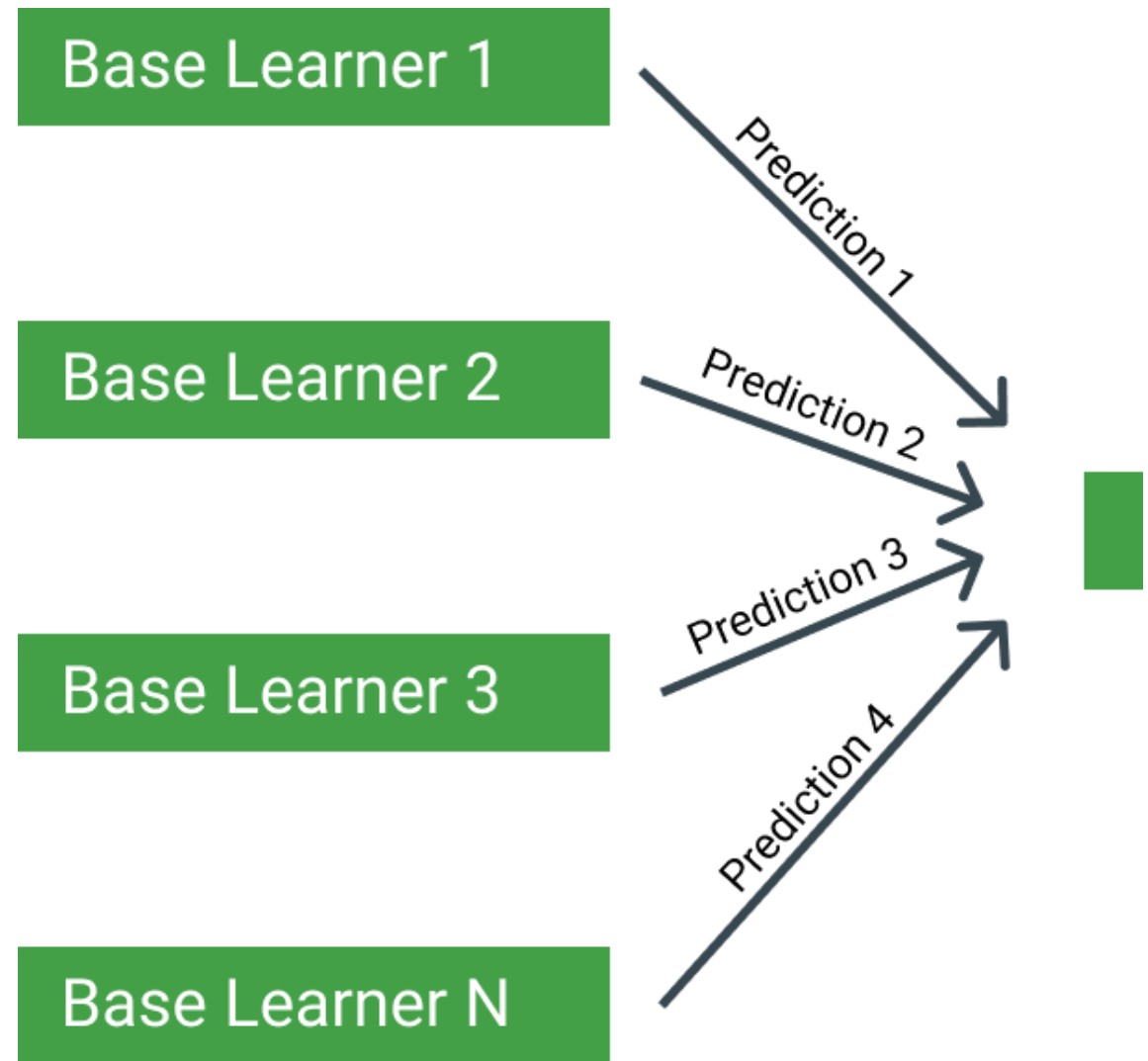
- Used simple artificial neural network structure
- Results:
 - MSE: 25.921140647451026
 - R2: 0.677052182685883
- Not enough data for deep learning methods
- Optimal hyper parameters are exposed in the notebook



Scaled down version of neural network architecture

Regression on Price – Stacking Regression

- Used results from XGBoost and CatBoost as base learners then a simpler Decision Tree model as a meta learner.
- Results:
 - MSE: 15.515438358099788
 - R2: 0.814410486439004
- Not too far off performance but very computationally expensive



Regression on Price - Summary

Model	R2 score	MSE score	Hyperparameter search time (s)	Fit time (s)
Linear Regressor	0.71784	22.957	n/a	0.00457
AdaBoost	0.83157	13.760	82	0.10824
XGBoost	0.88116	9.8720	82	0.27811
CatBoost	0.89002	9.1128	n/a	1.43539
Feedforward Neural Network (MLP)	0.67705	25.921	472	0.65697
Stacking Regressor	0.81441	15.515	n/a	20.229



Policy Based Implications

- Based on the correlation between NOx pollution, crime rates, % of lower status individuals, and house prices, there are very worrying trends shown in the inner city, particularly towards the south side of the city.
 - The government should work towards developing these areas – as gentrification takes place, these areas of higher crime will only move elsewhere and wont be solved.
-



Interpretation and Conclusions from results

- Percentage of lower status population is strongly negatively correlated with median house price
 - Existence of geographic spill-over in house prices shown through spatial regression
 - Strong predictive performance of CatBoost model when combined with spatial data!
-