

# Homework 3

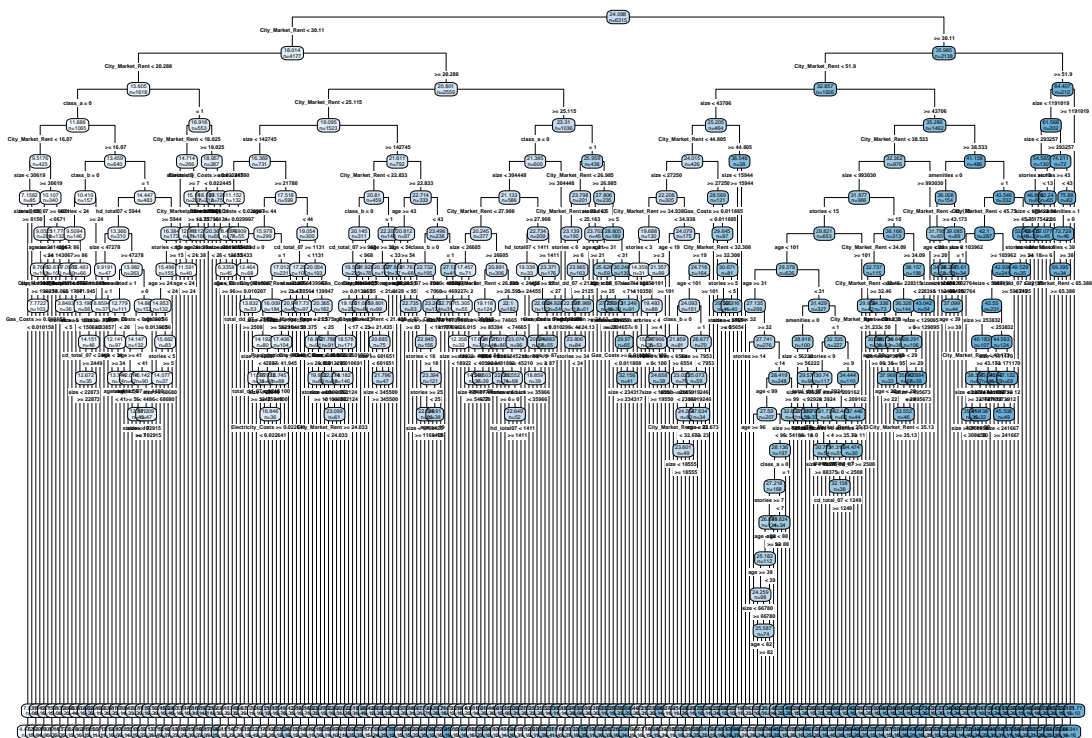
Patrick Massey

4/4/2022

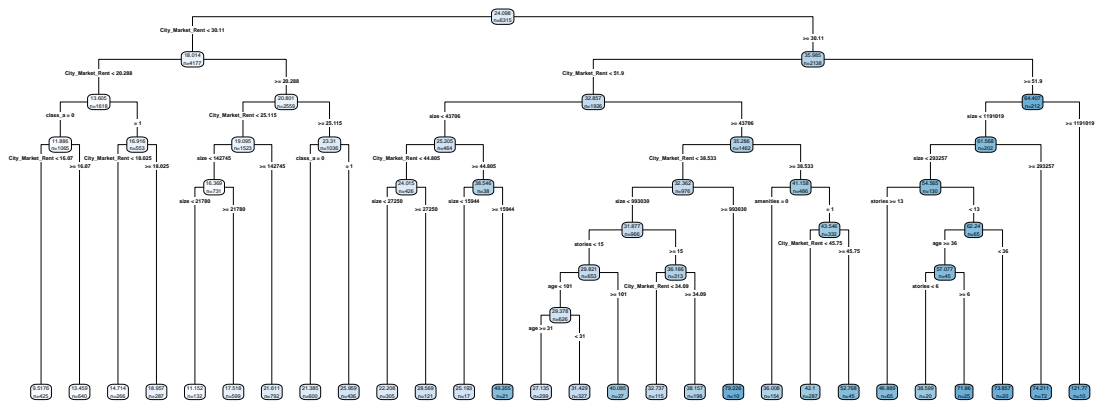
Before developing any models we first begin by performing some feature engineering. The first feature we engineer is the outcome variable of interest revenue which represents the revenue per square foot per calendar year. In order to create this feature we first scale down `leasing_rate` to a percentage by dividing by 100, and then multiplying that by the rent. We also create a new feature called `utility_cost` which is the sum of gas and electricity costs for rents that are quoted on a net contract basis. The purpose of this new feature is to capture the costs associated with a rental offered on a net contract basis. We then create a training set and a testing set with a split of 80/20. This gives us 6315 observations in our training set and 1579 observations in our testing set.

To begin developing our model we start with a linear model using all features of the data set excluding, `CS_PropertyID`, `cluster`, `leasing_rate`, `Rent`, `LEED`, and `Energystar`. We remove `CS_PropertyID` as it is just a unique building ID, and for similar reasons we remove `cluster`. We remove `leasing_rate` and `Rent` since these variables directly calculate our outcome variable. Lastly we remove `LEED` and `Energystar` because we are only concerned if a building is green certified or not, and not what kind of green certification a building may have. We capture this with the `green_rating` feature.

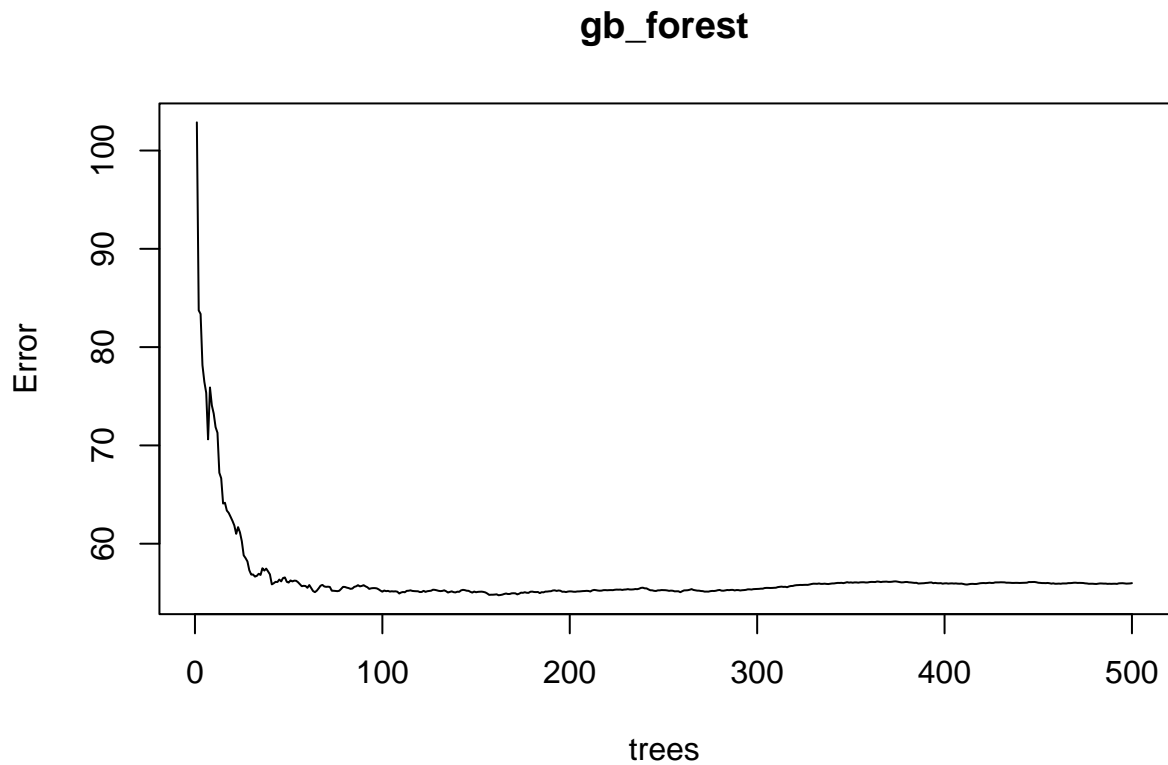
After getting a baseline model we then moved onto predicting using a tree model. The initial tree model generated, shown below, was extremely complex and not readable. This indicated that there might be some overfitting happening.



We then pruned our tree using the 1se method which generated the much simpler decision tree shown below. This tree sacrifices a marginal amount of performance for a much simpler tree.



The visualization of the tree really highlighted the interactions that were not included in our baseline linear model. Naturally after seeing the performance of the tree as compared to the linear model we wanted to see if it could be improved upon using a random forest.



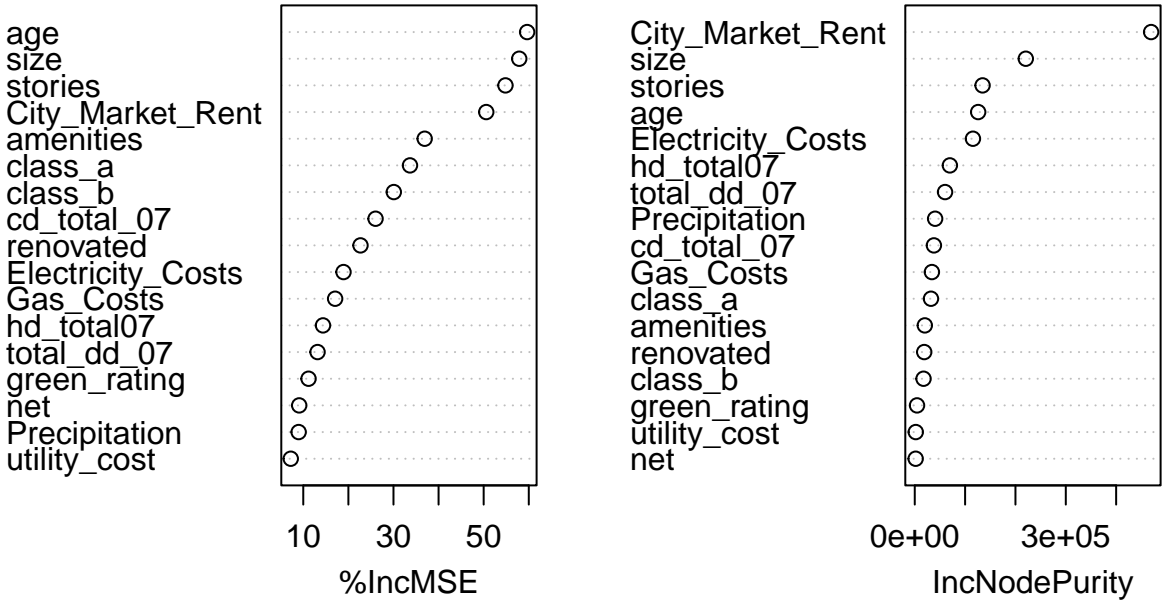
We see that our error really starts to bottom out around 100 trees. The performance of our models is shown below.

Model	RMSE
Linear	10.546099
CART	9.250752
Pruned Tree	9.904348
Random Forest	7.146455

The random forest provides a significant reduction in RMSE as compared to our baseline linear model.

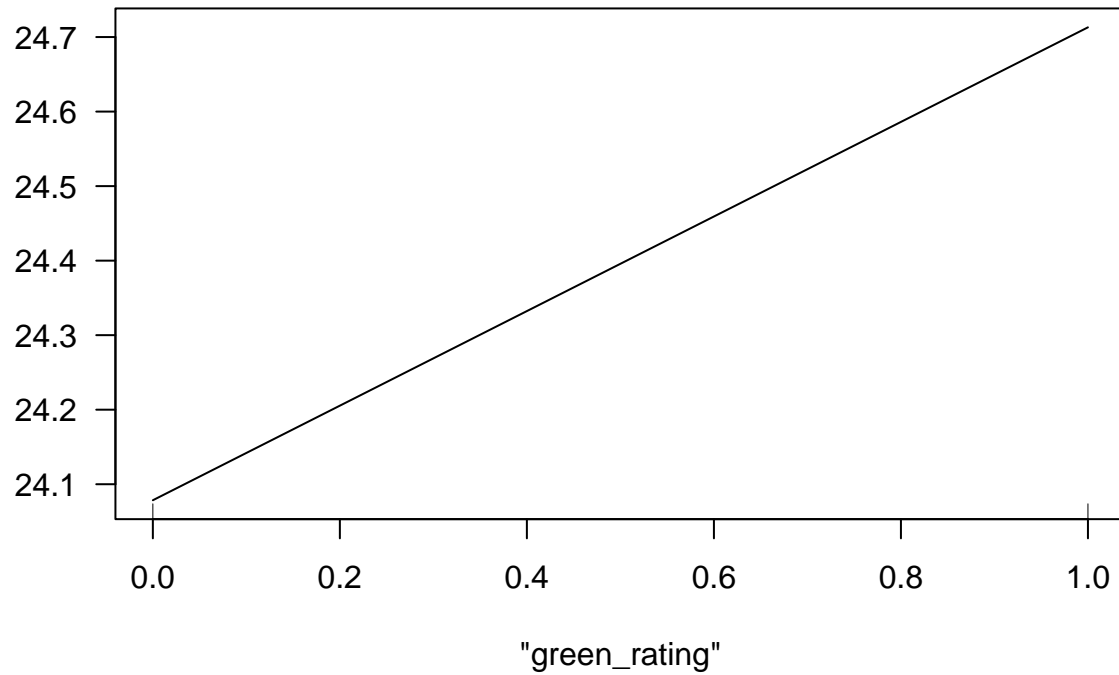
Now that we have developed a model for predicting the revenue generated from an building we will look at the importance of the variables we have used in our model.

# Variable Importance Plot



We can see that from a prediction point of view the green rating of a building does provide a large (>10%) increase in RMSE performance. Now lets look at the dollar increase in revenue from a building that has a green rating by creating a partial dependence plot shown below.

### Partial Dependence on "green\_rating"



From the plot we can see that there is a small marginal improvement in the expected revenue for a green building versus a non-green building. In fact going from a non-green building to a green building will give a revenue increase of 0.6345277 which in percentage terms leads to a 2.635% increase in revenue.