

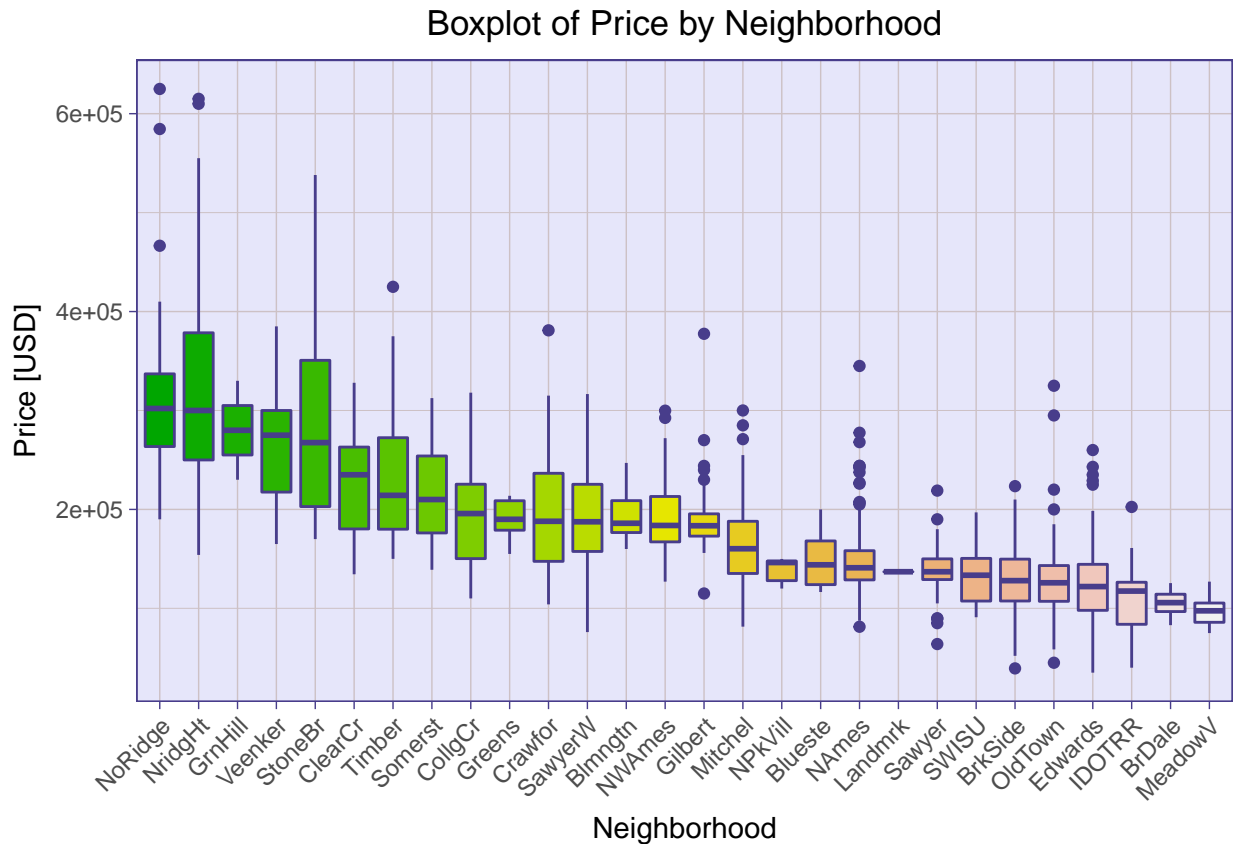
Final Project Writeup

Thomas Fleming, Eden Huang, Blaire Li, Marc Ryser

1. Exploratory data analysis

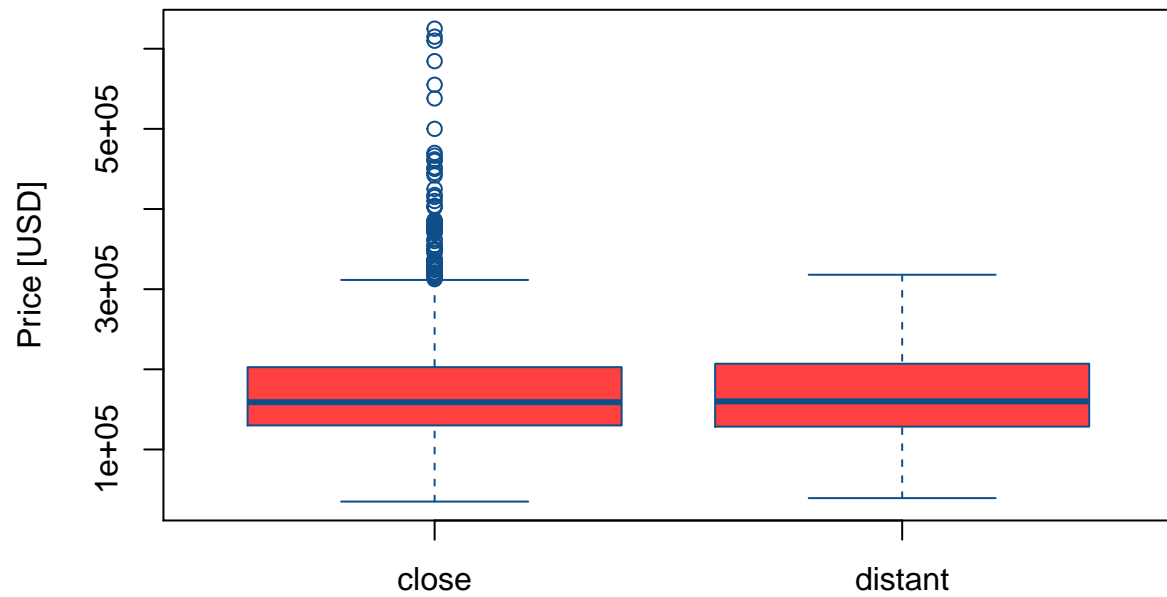
We first performed basic checks on the training data to identify predictors with (truly) missing entries and potential abnormalities. The variable lot frontage was identified to have 282 or 18.8% missing entries. We removed this predictor, and verified at a later stage that adding it in would not lead to better model predictions.

The mantra in real estate appears to be “location, location, location.” This begs for a simple visualization of house price distribution by neighborhood.



An interesting observation is that there is a wider dispersion among the more affluent neighborhoods based on relative inter-quartile ranges, whereas the neighborhoods with cheaper housing tend to be more concentrated around their medians. In other words, homoscedasticity is violated in this data set as variation increases with sale price. This is an important observation, as it indicates that we would be wise to transform the response variable (which we eventually did).

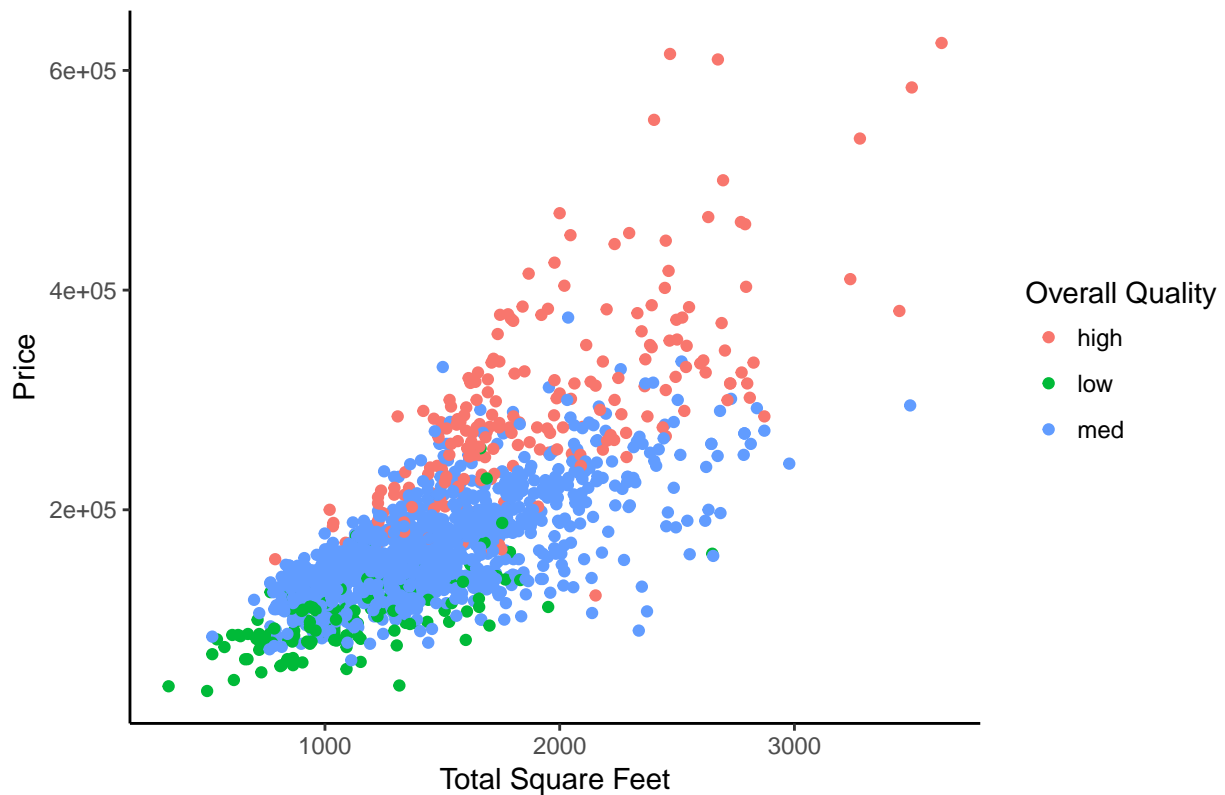
With respect to location, we note that Ames is a college town built around Iowa State. With almost 16,000 employees, the university is by far the biggest employer in town. Using Google maps, we recorded the respective distances between the university to the different neighborhoods. Beyond a univariable analysis, we could not find an association between price and distance to university. This is illustrated in the following boxplot where we compare prices between close and distant houses with respect to the university (defined as within 1.5 miles to Parks Library at ISU, measured using Google Maps). However, the proximity to ISU doesn't seem to affect property prices as expected.



Location with respect to university

Our third plot shows the relationship between sales price and total square footage, stratified by overall quality of the house. We see that there are different slopes in the relationship between sales price and total square footage, indicating the need for interaction terms between these variables (indeed, the more complex model contains such interactions).

Sale Price versus Total Square Feet by Overall Quality



2. Development and assessment of an initial model from Part I

Initial Model and Model Selection

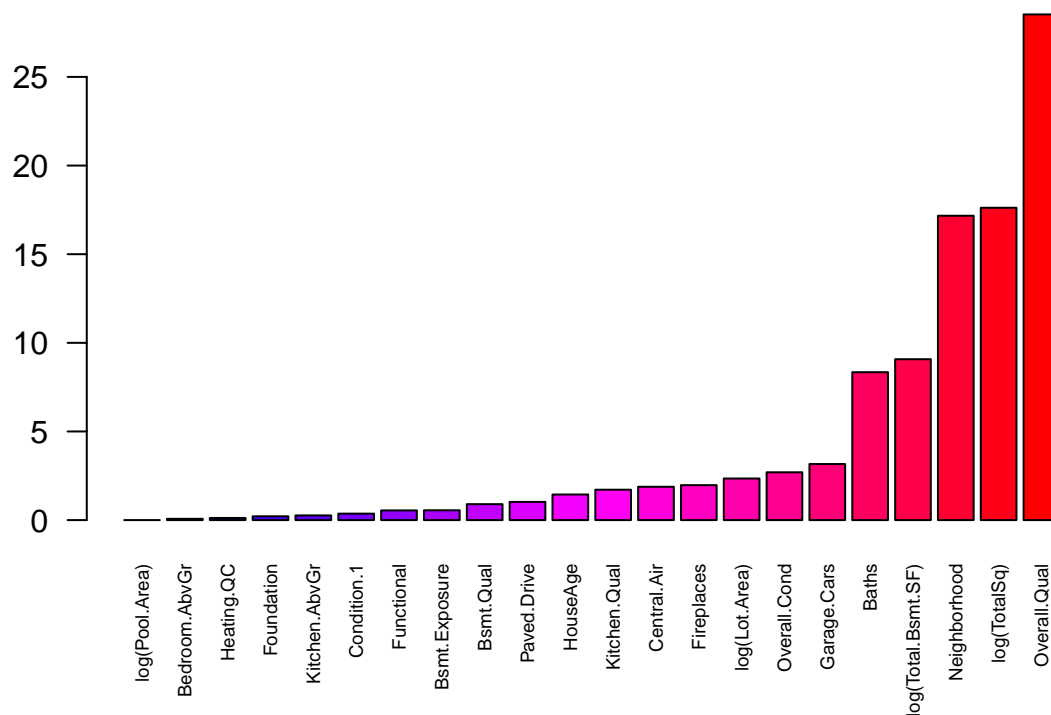
Prior to selecting the variables for our simple model, we cleaned the data (big time). We used a variety of techniques to do this, including converting variables to factors to account for non-linearities (e.g. MS.Subclass, Alley, Bsmt.Qual, Bsmt.Cond, etc.), aggregating like variables, such as combining porch square footage for different types of porches and creating a variable for total number of baths, accounting for the NA's by adding levels named "None" where appropriate, creating a variable for house age- calculated by subtracting the max of year built and year remodeled from the year sold- and filtering out NA's for a few particular variables with few NA's.

Following our data cleaning, we performed exploratory analyses. We first fit a full linear model, including all available variables, and examined its residual plots for indication as to the changes we should make. We detected some non-linearity in the data, seeing a trend in the residual and standardized residual plots, and decided to make a log transformation of our response variable, which improved the model fit. Taking a note from Appendix A in Gelman's book, we also logged all continuous explanatory variables, as this helps provide a multiplicative effect to the model when transformed to the original scale. The BoxCox procedure also indicated that a log-transformation for the response is reasonable.

Having improved our model fit through simple variable transformations, we then went about variable selection. We found that running a step function evaluated using the Bayesian Information Criterion brought us down close to 20 predictors, giving us 22. We then used a boosted tree model (depth 1, no interactions) to evaluate the relative importance of these 22 variables, and removed log(Pool.Area), Bedrooms.AbvGr, and Heating.QC, as these variables had the lowest relative importance. This brought us to our finalized simple model, with 19 predictors.

Interpreting the coefficients of our simple model, we see an intercept of ~ 7.59 , indicating that this would be the predicted housing price given a model where we have the base categorical level for all factors and 0 for numerical variables. However, this is not particularly meaningful given that the majority of our predictions were in the 11 to 13 range. For our variables that are factors, we would include the coefficient in a particular regression given that the observation fell into that particular category for the variable. Most of these coefficients fell between 0 and 1. We also have integer and numerical variables. Holding all else constant, a one-unit increase in Overall.Qual translates to a .05 increase in $\log(\text{price})$; a one-unit increase in Overall.Cond translates to a .35 increase in $\log(\text{price})$; the $\log(\text{price})$ decreases by .00095 with each additional year of house age; a one-unit increase in $\log(\text{Total.Bsmt.SF})$ translates to an increase of .119 in $\log(\text{price})$; each additional bath translates to a .047 increase in $\log(\text{price})$; each additional kitchen above ground translates to a decrease in $\log(\text{price})$ of .10; each unit increase in kitchen quality translates to a .038 increase in $\log(\text{price})$; each additional fireplace translates to a .03 increase in $\log(\text{price})$; each additional car a garage can hold translates to a .04 increase in $\log(\text{price})$; and each additional unit in $\log(\text{TotalSq})$ translates to a .3439479 increase in $\log(\text{price})$.

Relative influence



```
##
##               var      rel.inf
## Overall.Qual      Overall.Qual 28.52418935
## log(TotalSq)      log(TotalSq) 17.61747566
## Neighborhood      Neighborhood 17.16875214
## log(Total.Bsmt.SF) log(Total.Bsmt.SF) 9.07543469
## Baths              Baths      8.34268575
## Garage.Cars        Garage.Cars 3.16191190
## Overall.Cond       Overall.Cond 2.69666562
## log(Lot.Area)      log(Lot.Area) 2.34680989
## Fireplaces         Fireplaces 1.97228410
## Central.Air        Central.Air 1.88031329
## Kitchen.Qual       Kitchen.Qual 1.71313499
## HouseAge           HouseAge 1.44359467
## Paved.Drive        Paved.Drive 1.02695486
## Bsmt.Qual          Bsmt.Qual 0.89967679
## Bsmt.Exposure      Bsmt.Exposure 0.55283394
## Functional         Functional 0.54544829
## Condition.1        Condition.1 0.36368635
## Kitchen.AbvGr      Kitchen.AbvGr 0.26267515
## Foundation         Foundation 0.21371747
## Heating.QC         Heating.QC 0.11987758
## Bedroom.AbvGr      Bedroom.AbvGr 0.07187751
## log(Pool.Area)     log(Pool.Area) 0.00000000
##
## Call:
## lm(formula = log(price) ~ log(Lot.Area) + Neighborhood + Condition.1 +
##     Overall.Qual + Overall.Cond + HouseAge + Foundation + Bsmt.Qual +
##     Bsmt.Exposure + log(Total.Bsmt.SF) + Central.Air + Baths +
```

```

##      Kitchen.AbvGr + Kitchen.Qual + Functional + Fireplaces +
##      Garage.Cars + Paved.Drive + log(TotalSq), data = data_train)
##
## Residuals:
##      Min        1Q      Median        3Q        Max
## -0.66880 -0.05342 -0.00040  0.05532  0.34125
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    7.5928042  0.1034783  73.376 < 2e-16 ***
## log(Lot.Area)    0.1005736  0.0075956  13.241 < 2e-16 ***
## NeighborhoodBlueste  0.0038079  0.0422185   0.090 0.928145
## NeighborhoodBrDale -0.0773350  0.0348923  -2.216 0.026822 *
## NeighborhoodBrkSide -0.0297326  0.0309012  -0.962 0.336120
## NeighborhoodClearCr -0.0032135  0.0340622  -0.094 0.924851
## NeighborhoodCollgCr -0.0225456  0.0276495  -0.815 0.414975
## NeighborhoodCrawfor  0.0440361  0.0313361   1.405 0.160155
## NeighborhoodEdwards -0.1090553  0.0295961  -3.685 0.000237 ***
## NeighborhoodGilbert -0.0503690  0.0287855  -1.750 0.080367 .
## NeighborhoodGreens   0.0254955  0.0444436   0.574 0.566287
## NeighborhoodGrnHill  0.4479441  0.0714376   6.270 4.76e-10 ***
## NeighborhoodIDOTRR -0.1533929  0.0326944  -4.692 2.97e-06 ***
## NeighborhoodLandmrk -0.0763339  0.0957014  -0.798 0.425220
## NeighborhoodMeadowV -0.1420774  0.0384490  -3.695 0.000228 ***
## NeighborhoodMitchel -0.0461630  0.0296321  -1.558 0.119485
## NeighborhoodNames   -0.0670485  0.0286918  -2.337 0.019584 *
## NeighborhoodNoRidge  0.0541880  0.0304830   1.778 0.075674 .
## NeighborhoodNPkVill -0.0388906  0.0385797  -1.008 0.313597
## NeighborhoodNridgHt  0.0455037  0.0292245   1.557 0.119682
## NeighborhoodNWames  -0.0669261  0.0298203  -2.244 0.024965 *
## NeighborhoodOldTown -0.1235192  0.0296508  -4.166 3.29e-05 ***
## NeighborhoodSawyer  -0.0386395  0.0299358  -1.291 0.197000
## NeighborhoodSawyerW -0.0603434  0.0290947  -2.074 0.038255 *
## NeighborhoodSomerst  0.0500211  0.0278563   1.796 0.072757 .
## NeighborhoodStoneBr  0.0703212  0.0323876   2.171 0.030077 *
## NeighborhoodSWISU   -0.0749894  0.0341474  -2.196 0.028249 *
## NeighborhoodTimber  -0.0418053  0.0319528  -1.308 0.190967
## NeighborhoodVeenker -0.0129656  0.0378365  -0.343 0.731894
## Condition.1Artery   -0.0768610  0.0152515  -5.040 5.26e-07 ***
## Condition.1Feedr    -0.0749669  0.0114239  -6.562 7.39e-11 ***
## Condition.1Park      0.0063233  0.0183326   0.345 0.730204
## Condition.1Rail     -0.0441457  0.0149789  -2.947 0.003259 **
## Overall.Qual         0.0509388  0.0034351  14.829 < 2e-16 ***
## Overall.Cond         0.0346458  0.0027915  12.411 < 2e-16 ***
## HouseAge            -0.0009529  0.0001880  -5.069 4.52e-07 ***
## FoundationCBlock     0.0576301  0.0105526   5.461 5.57e-08 ***
## FoundationPConc      0.0749570  0.0115916   6.466 1.37e-10 ***
## FoundationSlab       0.0647669  0.0293347   2.208 0.027413 *
## FoundationStone      0.0002816  0.0393425   0.007 0.994291
## FoundationWood       0.0525257  0.0559323   0.939 0.347841
## Bsmt.QualEx          -0.7533650  0.1145897  -6.574 6.83e-11 ***
## Bsmt.QualFa          -0.8320710  0.1133422  -7.341 3.54e-13 ***
## Bsmt.QualGd          -0.8380252  0.1131896  -7.404 2.25e-13 ***
## Bsmt.QualPo          -0.8410027  0.1454982  -5.780 9.15e-09 ***

```

```

## Bsmt.QualTA      -0.8503317  0.1130739  -7.520  9.63e-14 ***
## Bsmt.ExposureAv  0.1560803  0.0926024   1.685  0.092113 .
## Bsmt.ExposureGd  0.1933630  0.0929084   2.081  0.037591 *
## Bsmt.ExposureMn  0.1117113  0.0928128   1.204  0.228936
## Bsmt.ExposureNo  0.1272957  0.0925440   1.376  0.169187
## log(Total.Bsmt.SF) 0.1196200  0.0092381  12.949  < 2e-16 ***
## Central.AirY     0.0634765  0.0119725   5.302  1.33e-07 ***
## Baths            0.0473048  0.0045589  10.376  < 2e-16 ***
## Kitchen.AbvGr    -0.1041428  0.0140674  -7.403  2.26e-13 ***
## Kitchen.Qual     0.0380916  0.0055894   6.815  1.39e-11 ***
## FunctionalMaj2    -0.1668154  0.0522781  -3.191  0.001449 **
## FunctionalMin1     0.0144202  0.0348213   0.414  0.678848
## FunctionalMin2     0.0342439  0.0344569   0.994  0.320480
## FunctionalMod     -0.0001033  0.0374934  -0.003  0.997801
## FunctionalTyp     0.0846081  0.0315133   2.685  0.007341 **
## Fireplaces        0.0313138  0.0047078   6.651  4.12e-11 ***
## Garage.Cars       0.0395438  0.0047321   8.356  < 2e-16 ***
## Paved.DriveP      0.0070829  0.0182491   0.388  0.697982
## Paved.DriveY      0.0575194  0.0115044   5.000  6.45e-07 ***
## log(TotalSq)      0.3439479  0.0130974  26.261  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.09198 on 1428 degrees of freedom
## Multiple R-squared:  0.9415, Adjusted R-squared:  0.9389
## F-statistic: 359 on 64 and 1428 DF, p-value: < 2.2e-16

```

Residual Plots

The residual plots for our simple model display favorable results. The residual vs. fitted plot shows little to no pattern and while the scale-location plot has a slight pattern, it is not severe. While there are a few outliers displayed on the Normal Q-Q plot, most observations fall within 4 standard deviations, with the majority falling very close to or on the one-to-one line. Finally, after having the function remove observations with leverage 1 (of which there were only 3), we observe a favorable leverage plot, with no observations exceeding a Cook's distance of 0.5.

RMSE

The RMSE for the simple model evaluated on the test was 15,477, which corresponds to approximately 8% of the mean house price. This implies that the RMSE is fairly small.

Model Testing

Beyond the RMSE, we were very pleased to see a bias relatively closer to zero than our fellow teams, at -165.05. We found a maximum deviation of 66,474.27, a mean deviation of 11,458.19, and coverage of 96.2%. The coverage is very favorable and, given that the mean deviation is less than 10% of the mean housing price for the training set, we are please with these results.

In our hand-calculated model check, we evaluated our model on the first observation of the training set, as well as the first observation of the training set. For the training set evaluation, we get 11.73495 translating to \$124,860.20, whereas the true price listing for this observation was \$137,000. In our model check on the test

set, we get a predicted value of 12.22475 or \$203,770.57, whereas the true value was \$192,100. These seem to be reasonable estimates.

	bias	max.dev	mean.dev	RMSE	Coverage
Training	799.9486	152265.48	11858.23	17070.81	0.9564635
Test	-165.0479	66474.27	11458.19	15477.42	0.9620000

3. Development of the final model

We tried a range of different approaches for a more complex model. Among others, we evaluated tree models, bagging, boosting and random forests, as well as Lasso and Ridge. Based on RMSE however, none of these options was able to outperform a linear model with interaction terms. To develop the final linear model, we proceeded as follows:

- We grew a deep tree
- We pruned the tree using cross-validation
- We read out the interactions from the final tree
- We inserted these interactions into the full linear model
- We ran a stepwise BIC

The reason for not using a tree-based model for the final version was that the RMSE of these approaches remained substantially higher compared to the final OLS model.

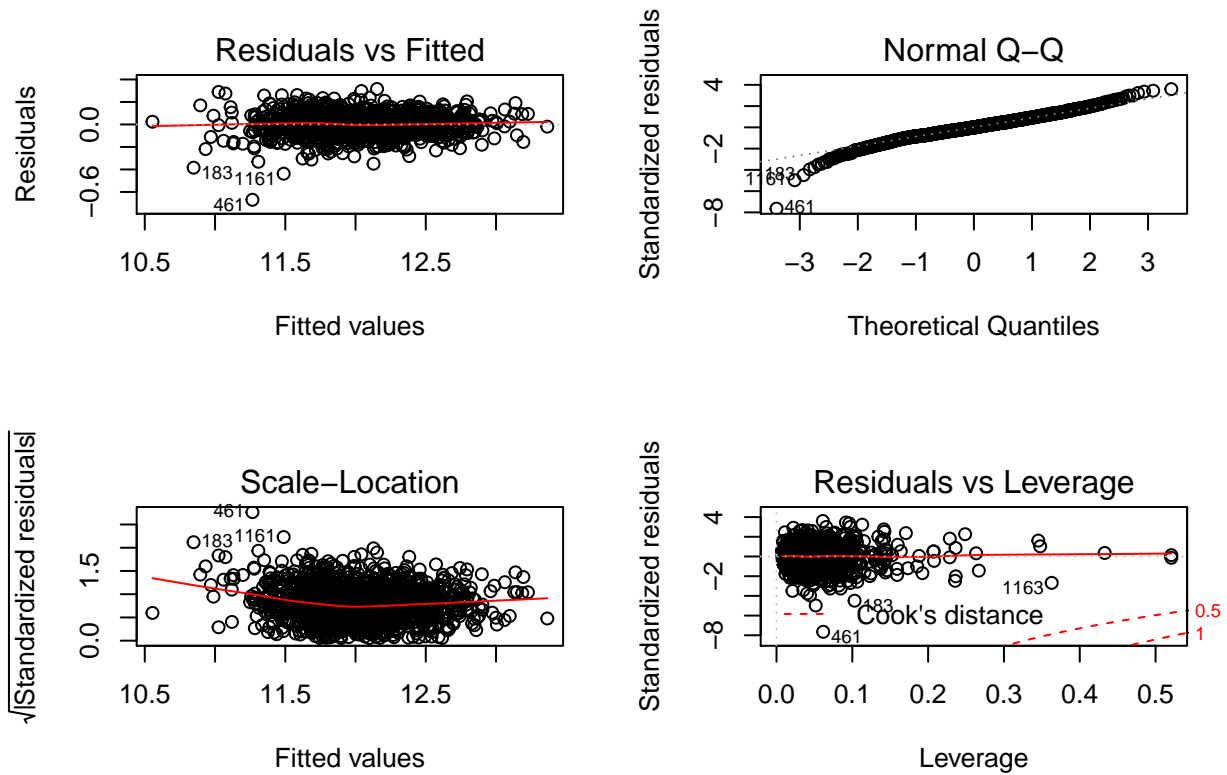
In consequence, the final model contains 22 different variables, two 2-way interactions and one 3-way interaction.

The 18 most important predictors are summarized in the following table.

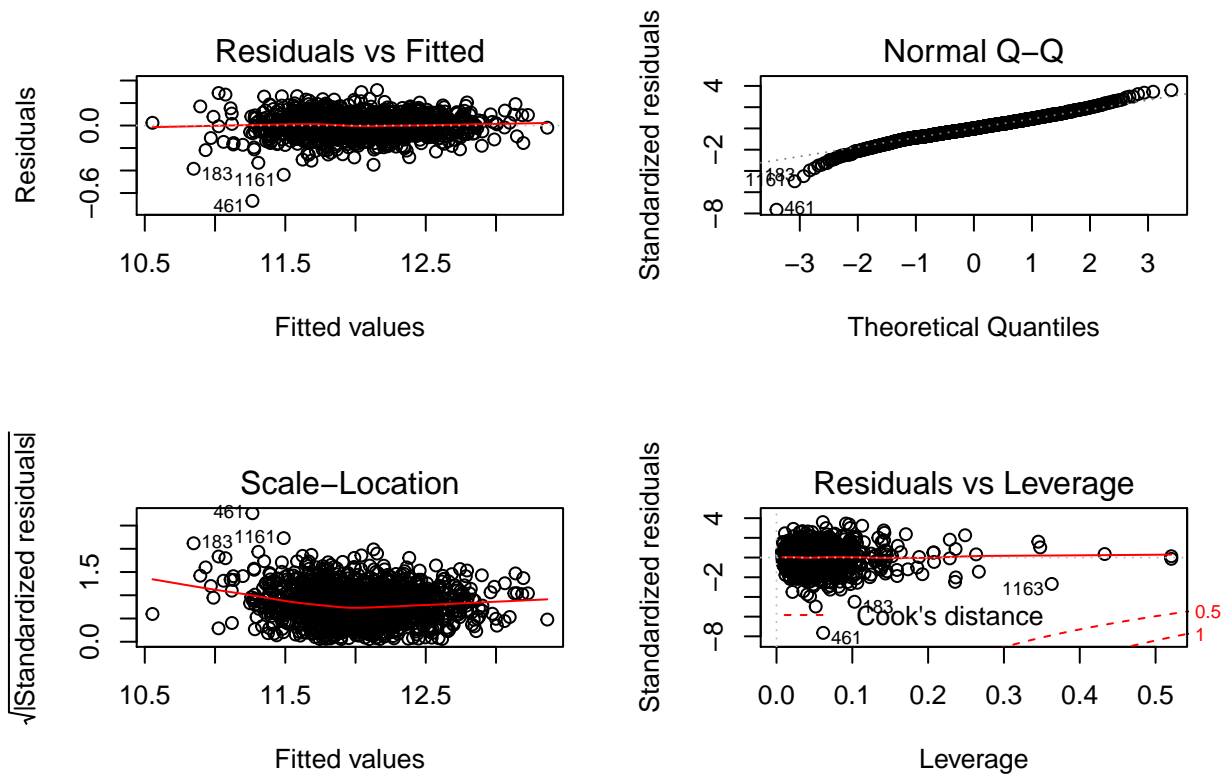
	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	6.5102614	0.4057490	16.045046	0
log(Lot.Area)	0.0987836	0.0075415	13.098737	0
log(Total.Bsmt.SF)	0.1189555	0.0091090	13.059147	0
Overall.Cond	0.0336221	0.0027929	12.038585	0
Baths	0.0468993	0.0045396	10.331210	0
log(TotalSq)	0.5012898	0.0581418	8.621853	0
Bsmt.QualTA	-0.8376454	0.1112764	-7.527611	0
Bsmt.QualGd	-0.8267350	0.1113854	-7.422292	0
Bsmt.QualFa	-0.8262069	0.1114845	-7.410957	0
Bsmt.QualEx	-0.7632015	0.1127644	-6.768105	0
Condition.1Feedr	-0.0745876	0.0112824	-6.610950	0
Kitchen.AbvGr	-0.0910776	0.0142064	-6.411031	0
FoundationPConc	0.0731980	0.0114720	6.380555	0
NeighborhoodGrnHill	0.4438689	0.0705238	6.293886	0
Fireplaces	0.0296062	0.0047109	6.284618	0
Kitchen.Qual	0.0332874	0.0055725	5.973465	0
Bsmt.QualPo	-0.8358514	0.1430638	-5.842509	0
FoundationCBlock	0.0589952	0.0104102	5.667074	0

As one would expect, and similarly to the conclusions from the simple model, we find that the overall condition of the house, lot area, and square footage are critical predictors of the house price. Interestingly, the square footage of the basement is an independent predictor, as is the number of bathrooms (Marc was particularly fascinated with the bathrooms).

Finally, here are the model analytics for the complete model.



4. Assessment of the final model (25 points)




```
##
## Call:
## lm(formula = log(price) ~ log(Lot.Area) + Condition.1 + Overall.Qual +
##     Baths + Neighborhood + Garage.Cars + log(Total.Bsmt.SF) +
##     log(TotalSq) + Overall.Cond + HouseAge + Foundation + Bsmt.Qual +
##     Bsmt.Exposure + Heating.QC + Central.Air + Bedroom.AbvGr +
##     Kitchen.AbvGr + Kitchen.Qual + Functional + Fireplaces +
##     Paved.Drive + Overall.Qual:Garage.Cars + Overall.Qual:log(TotalSq) +
##     Garage.Cars:log(TotalSq) + Overall.Qual:Garage.Cars:log(TotalSq),
##     data = data_train)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.66963 -0.05131  0.00028  0.05415  0.31533
##
## Coefficients:
##                                Estimate Std. Error t value
## (Intercept)                   6.5102614   0.4057490  16.045
## log(Lot.Area)                   0.0987836   0.0075415  13.099
## Condition.1Artery              -0.0736717   0.0150310  -4.901
## Condition.1Feedr              -0.0745876   0.0112824  -6.611
## Condition.1Park               -0.0051531   0.0181031  -0.285
## Condition.1Rail              -0.0456100   0.0148189  -3.078
## Overall.Qual                   0.2068733   0.0752569   2.749
## Baths                          0.0468993   0.0045396  10.331
## NeighborhoodBlueste           0.0203947   0.0416095   0.490
## NeighborhoodBrDale           -0.0722602   0.0345132  -2.094
## NeighborhoodBrkSide          -0.0217135   0.0304547  -0.713
## NeighborhoodClearCr          0.0103092   0.0336539   0.306
## NeighborhoodCollgCr         -0.0159676   0.0273774  -0.583
## NeighborhoodCrawfor          0.0499783   0.0308951   1.618
## NeighborhoodEdwards         -0.1032780   0.0291913  -3.538
## NeighborhoodGilbert         -0.0368491   0.0285011  -1.293
## NeighborhoodGreens           0.0574166   0.0441283   1.301
## NeighborhoodGrnHill          0.4438689   0.0705238   6.294
## NeighborhoodIDOTRR          -0.1470244   0.0322360  -4.561
## NeighborhoodLandmrk         -0.0700365   0.0941311  -0.744
## NeighborhoodMeadowV         -0.1441208   0.0380517  -3.787
## NeighborhoodMitchel         -0.0346045   0.0293105  -1.181
## NeighborhoodNames           -0.0586303   0.0283935  -2.065
## NeighborhoodNoRidge          0.0250677   0.0307515   0.815
## NeighborhoodNPkVill         -0.0218909   0.0381536  -0.574
## NeighborhoodNridgHt          0.0317326   0.0289648   1.096
## NeighborhoodNWAmes          -0.0465829   0.0295620  -1.576
## NeighborhoodOldTown         -0.1206424   0.0292366  -4.126
## NeighborhoodSawyer          -0.0277534   0.0297174  -0.934
## NeighborhoodSawyerW         -0.0496333   0.0287674  -1.725
## NeighborhoodSomerst          0.0576101   0.0274356   2.100
## NeighborhoodStoneBr          0.0743440   0.0318815   2.332
## NeighborhoodSWISU           -0.0658180   0.0339123  -1.941
## NeighborhoodTimber          -0.0245197   0.0315509  -0.777
## NeighborhoodVeenker         -0.0060078   0.0372772  -0.161
## Garage.Cars                   1.2032426   0.2341509   5.139
## log(Total.Bsmt.SF)           0.1189555   0.0091090  13.059
```

## log(TotalSq)	0.5012898	0.0581418	8.622
## Overall.Cond	0.0336221	0.0027929	12.039
## HouseAge	-0.0009157	0.0001870	-4.896
## FoundationCBlock	0.0589952	0.0104102	5.667
## FoundationPConc	0.0731980	0.0114720	6.381
## FoundationSlab	0.0718932	0.0289054	2.487
## FoundationStone	0.0044701	0.0389376	0.115
## FoundationWood	0.0392228	0.0550548	0.712
## Bsmt.QualEx	-0.7632015	0.1127644	-6.768
## Bsmt.QualFa	-0.8262069	0.1114845	-7.411
## Bsmt.QualGd	-0.8267350	0.1113854	-7.422
## Bsmt.QualPo	-0.8358514	0.1430638	-5.843
## Bsmt.QualTA	-0.8376454	0.1112764	-7.528
## Bsmt.ExposureAv	0.1575361	0.0910068	1.731
## Bsmt.ExposureGd	0.1836363	0.0913334	2.011
## Bsmt.ExposureMn	0.1122697	0.0912138	1.231
## Bsmt.ExposureNo	0.1280950	0.0909443	1.408
## Heating.QC	0.0093065	0.0033440	2.783
## Central.AirY	0.0636338	0.0121613	5.232
## Bedroom.AbvGr	-0.0147371	0.0042642	-3.456
## Kitchen.AbvGr	-0.0910776	0.0142064	-6.411
## Kitchen.Qual	0.0332874	0.0055725	5.973
## FunctionalMaj2	-0.1605476	0.0514362	-3.121
## FunctionalMin1	0.0184235	0.0343702	0.536
## FunctionalMin2	0.0407226	0.0339893	1.198
## FunctionalMod	0.0029487	0.0370700	0.080
## FunctionalTyp	0.0888654	0.0310927	2.858
## Fireplaces	0.0296062	0.0047109	6.285
## Paved.DriveP	0.0171341	0.0181069	0.946
## Paved.DriveY	0.0605392	0.0114220	5.300
## Overall.Qual:Garage.Cars	-0.1928486	0.0384268	-5.019
## Overall.Qual:log(TotalSq)	-0.0225185	0.0105342	-2.138
## Garage.Cars:log(TotalSq)	-0.1633273	0.0325288	-5.021
## Overall.Qual:Garage.Cars:log(TotalSq)	0.0269148	0.0052360	5.140
##	Pr(> t)		
## (Intercept)	< 2e-16 ***		
## log(Lot.Area)	< 2e-16 ***		
## Condition.1Artery	1.06e-06 ***		
## Condition.1Feedr	5.39e-11 ***		
## Condition.1Park	0.775950		
## Condition.1Rail	0.002125 **		
## Overall.Qual	0.006055 **		
## Baths	< 2e-16 ***		
## NeighborhoodBlueste	0.624106		
## NeighborhoodBrDale	0.036463 *		
## NeighborhoodBrkSide	0.475977		
## NeighborhoodClearCr	0.759399		
## NeighborhoodCollgCr	0.559823		
## NeighborhoodCrawfor	0.105954		
## NeighborhoodEdwards	0.000416 ***		
## NeighborhoodGilbert	0.196256		
## NeighborhoodGreens	0.193425		
## NeighborhoodGrnHill	4.11e-10 ***		
## NeighborhoodIDOTRR	5.53e-06 ***		

```

## NeighborhoodLandmrk      0.456981
## NeighborhoodMeadowV      0.000159 ***
## NeighborhoodMitchel      0.237952
## NeighborhoodNames        0.039111 *
## NeighborhoodNoRidge      0.415112
## NeighborhoodNPkVill      0.566222
## NeighborhoodNridgHt      0.273459
## NeighborhoodNWames        0.115302
## NeighborhoodOldTown      3.90e-05 ***
## NeighborhoodSawyer        0.350509
## NeighborhoodSawyerW      0.084685 .
## NeighborhoodSomerst      0.035919 *
## NeighborhoodStoneBr      0.019846 *
## NeighborhoodSWISU        0.052476 .
## NeighborhoodTimber        0.437201
## NeighborhoodVeenker       0.871985
## Garage.Cars               3.15e-07 ***
## log(Total.Bsmt.SF)        < 2e-16 ***
## log(TotalSq)              < 2e-16 ***
## Overall.Cond              < 2e-16 ***
## HouseAge                  1.09e-06 ***
## FoundationCBlock          1.76e-08 ***
## FoundationPConc           2.38e-10 ***
## FoundationSlab            0.012990 *
## FoundationStone           0.908619
## FoundationWood            0.476314
## Bsmt.QualEx               1.90e-11 ***
## Bsmt.QualFa               2.14e-13 ***
## Bsmt.QualGd               1.97e-13 ***
## Bsmt.QualPo               6.36e-09 ***
## Bsmt.QualTA               9.13e-14 ***
## Bsmt.ExposureAv           0.083662 .
## Bsmt.ExposureGd           0.044555 *
## Bsmt.ExposureMn           0.218586
## Bsmt.ExposureNo           0.159202
## Heating.QC                0.005457 **
## Central.AirY              1.92e-07 ***
## Bedroom.AbvGr             0.000564 ***
## Kitchen.AbvGr             1.96e-10 ***
## Kitchen.Qual              2.93e-09 ***
## FunctionalMaj2            0.001837 **
## FunctionalMin1            0.592021
## FunctionalMin2            0.231077
## FunctionalMod             0.936611
## FunctionalTyp             0.004324 **
## Fireplaces                4.36e-10 ***
## Paved.DriveP              0.344168
## Paved.DriveY              1.34e-07 ***
## Overall.Qual:Garage.Cars   5.86e-07 ***
## Overall.Qual:log(TotalSq) 0.032715 *
## Garage.Cars:log(TotalSq)   5.79e-07 ***
## Overall.Qual:Garage.Cars:log(TotalSq) 3.12e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

```
##
## Residual standard error: 0.09038 on 1422 degrees of freedom
## Multiple R-squared: 0.9437, Adjusted R-squared: 0.941
## F-statistic: 340.7 on 70 and 1422 DF, p-value: < 2.2e-16
```

- Residual:

Pardoe’s paper “Modeling home prices using realtor data” suggests an issue of heteroscedasticity with variation increasing with sale price. While it caused us concern in the beginning, it was not an issue after we had logged the response variable. The residual plot shows that the variance is relatively constant and no obvious patterns exist, which suggests that the logged response helped account for this non-constant variance. However, there are some outliers in the bottom-left of the plot where properties are overvalued. In addition, the variance in the left part of the plot is slightly higher in predicted prices. We do not throw away outliers in the model.

- Model Evaluation:

From the summary table, our model has a mutiple R-squared of 0.9437, which means that the final model explains about 94 percent of variation in the data. In addition, the residual standard error of 0.09 indicates that the fit on training data is decent. The diagnostic plots suggest that our model actually did a better job in predicting price for those properties in middle and high price ranges than lower price range. The residual plot has higher variance in the lower price range. In addition, the normal QQ plot generally follows a straight line, but with a heavier left tail. In the residual vs leverage plot, it can be observed that there are about 5 high leverage point, but they are not influential because no points have cook’s distance greater than 0.5. In conclusion, our final model did a good job in predicting the prices, especially for the properties in middle and high price ranges.

```
##          bias  max.dev mean.dev    RMES Coverage
## 1 -156.1364 64352.67 11229.45 14983.7      0.966
```

- RMSE discussion and Model Testing:

The RMSE for our final model is 14983.7. After comparing with other groups in the leaderboard, we find that the bias for our model is actually the lowest. However, due to the variance and bias trade-off, our model has a higher deviance or variance for prediction, possibly resulting from outliers in the dataset. Another possible issue would be extrapolation, which affects the prediction accuracy for some properties with extreme prices. The coverage of prediction of 0.966 indicates that our model satisfactorily captures true prices within the prediction interval.

PID pr	ice pri	ce_to_pred_ratio Tot	alSq Ove	rall.Qual Neig	hborhood
535382020	160000	1.334928	2414	5	OldTown
905452050	113000	1.253250	672	4	Edwards
535454070	166000	1.381293	1385	5	NAMES
903400220	214500	1.251595	2134	6	BrkSide
910206010	68104	1.351226	640	2	IDOTRR
916226030	241500	1.217801	1501	7	Timber
910203020	120500	1.208631	778	5	IDOTRR
903429110	179900	1.228883	1944	7	OldTown
905427010	235000	1.243933	2009	6	Edwards
905376090	216000	1.424357	1325	5	Edwards

PID pr	ice pri	ce_to_pred_ratio Tot	alSq Ove	rall.Qual Neig	hborhood
528102010	315000	0.8440589	1980	9	NridgHt
535150070	104900	0.8279967	1738	4	NAMES
902105130	64000	0.7702524	672	5	OldTown

PID pr	ice pri	ce_to_pred_ratio Tot	alSq Ove	rall.Qual Neig	hborhood
534451130	79900	0.7688397	747	4	BrkSide
909101010	110000	0.7873724	1196	6	Edwards
532378110	127000	0.8488299	1040	5	Sawyer
905200160	80000	0.7450666	1006	5	Sawyer
532351140	112000	0.8317835	1902	6	Sawyer
527182020	130000	0.7556536	1204	8	StoneBr
534479120	105000	0.8349797	1376	5	NAmes

- Model result:

The two tables show the top 10 most under- and over-priced properties based on our final model. The real-to-prediction ratio over 1 suggests over-priced and less than 1 under-priced. We also include some other features to compare the under and over-priced properties. One of our teammates comments that, were he to be a property investor in Ames, Iowa, he would keep an eye on the Sawyer neighborhood for buying opportunities, given that there were three here in the undervalued top 10, while considering selling properties that were in the Edwards neighborhood, as there were three here in the overvalued top 10.

5. Conclusion

In essence, the overall house quality, neighborhood, and total square footage were the most important predictors in our model. While this makes sense intuitively, confirming these predictors' relative importance in a quantitative analysis is reassuring.

Interestingly, the number of bathrooms is an independent predictor of house price. Who would have guessed it? Not us!

In reflecting upon the work we have done and our results, we have learned several things about the data analysis process, as well as the pricing of houses.

One of the key aspects that was apparent to us was how important our data-cleaning was to creating a successful model. In the early stages we had fit a model using the data more or less in its original form, with few modifications. After our meticulous data-cleaning session, the improvements were dramatic, as our bias dropped drastically and RMSE considerably.

Another observation we take away from our project is the realization that OLS is often very useful, and that more advanced methods are not always optimal. Over the course of the project we attempted running Ridge, Lasso, Blasso, Trees, Random Forests, and Boosting, with linear regression of the same variables ultimately outperforming them all.

Another thing we took away was the usefulness of tree models, not only in prediction, but also during model selection. We found boosting to be particularly useful in finalizing the variables we wanted to keep in our model through the examining the relative importance plot, and found a decision tree to be useful in indicating to us important interactions to consider in constructing our complex model.

Part IV

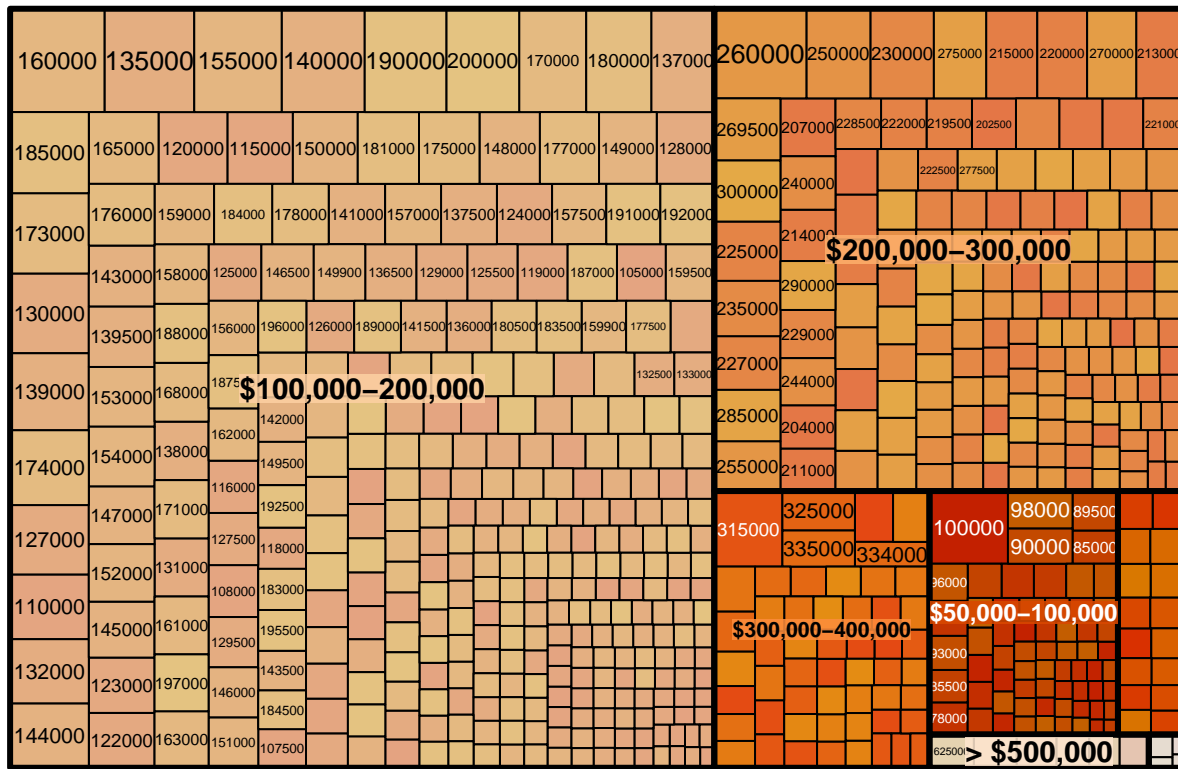
Create predictions for the validation data from your final model and write out to a file `prediction-validation.Rdata`. This should have the same format as the models in Part I and II.

Please see `prediction-validation.Rdata` for predictions.

Appendix

One team member was really into making graphs that were subsequently rejected by the rest of the team. As a token of appreciation to the energetic team member, we include their visualizations as an appendix.

Effect of Total Square Footage on Price Range



Price by Overall Quality

