

Predicting Housing Prices in Ames Iowa through Machine Learning Techniques

Rajesh Arasada, Yung Chou, Nilesh Patel,
Pankaj Sharma, Tim Waterman

Purpose

For Industry: Understand which housing aspects are valued by consumers, in order to accurately price assets, and evaluate undervalued/overvalued assets in the market.

For Consumers: Understand how much they might expect to pay for specific features/aspects of a new home.

Presentation Outline:

- Visualizing Data

- Missing Data

- Feature Selection

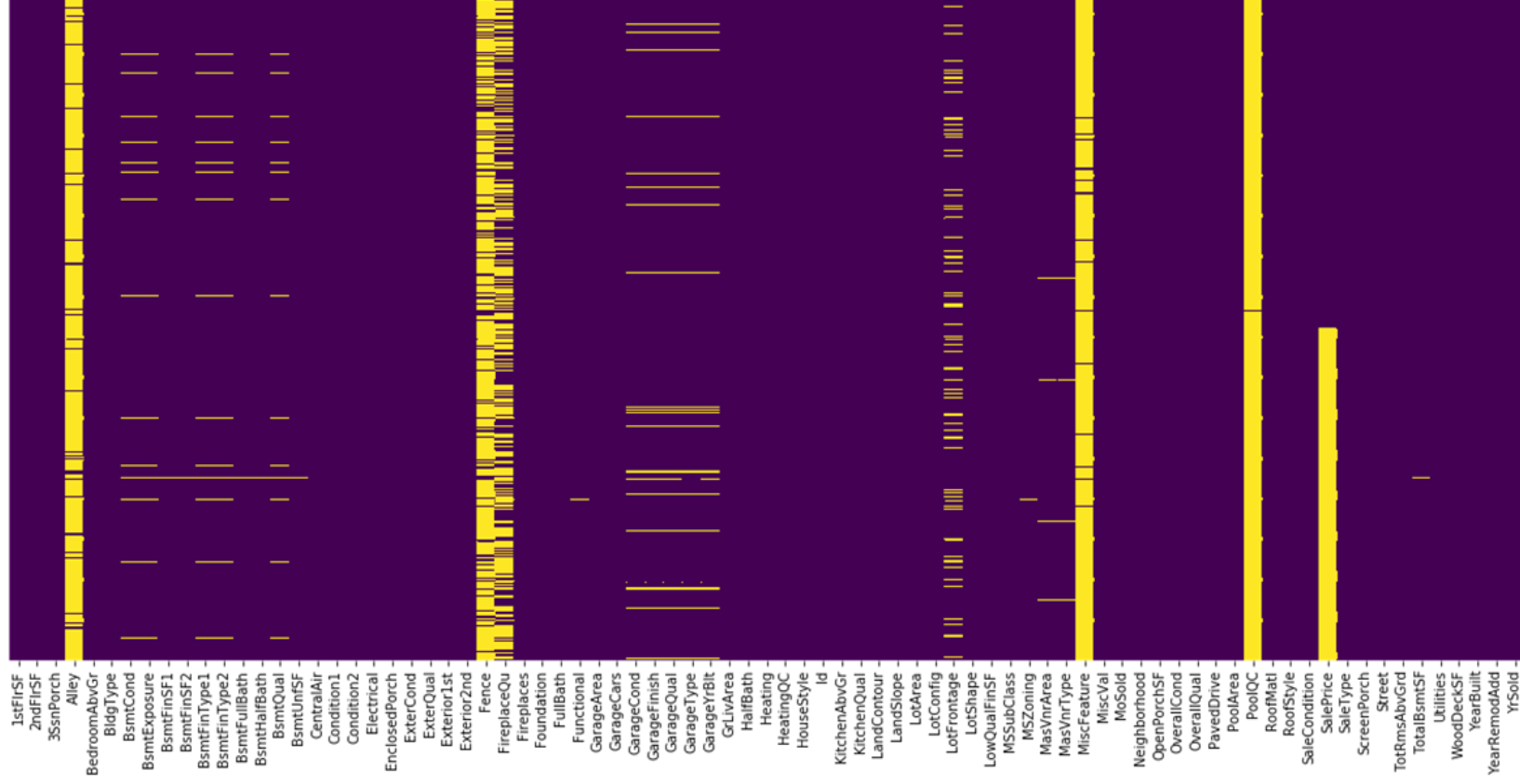
- Data Transformation

- Linear Models

- Tree Models

- KNN

Missingness



Data cleaning

we drop few unwanted columns, which don't have any impact on sale price. Then we change the data type of the 'MSSubClass' features from numeric to string.

- Missing Values
- Flagging as 'No':
BsmtFinType1, BsmtFinType2, GarageType, GarageFinish
Values are Missing, Because there is no Garage and Basement.
- Impute the Mode:
Electrical, Exterior1st, GarageCars, MSZoning, KitchenQual, MasVnrType, SaleType
- Impute Zero:
GarageQual, GarageCond, BsmtCond, BsmtExposure, BsmtQual, FireplaceQu
- Changing categorical ranking into numerical scale values.
- Dummify remaining categorical features.

Features engineering

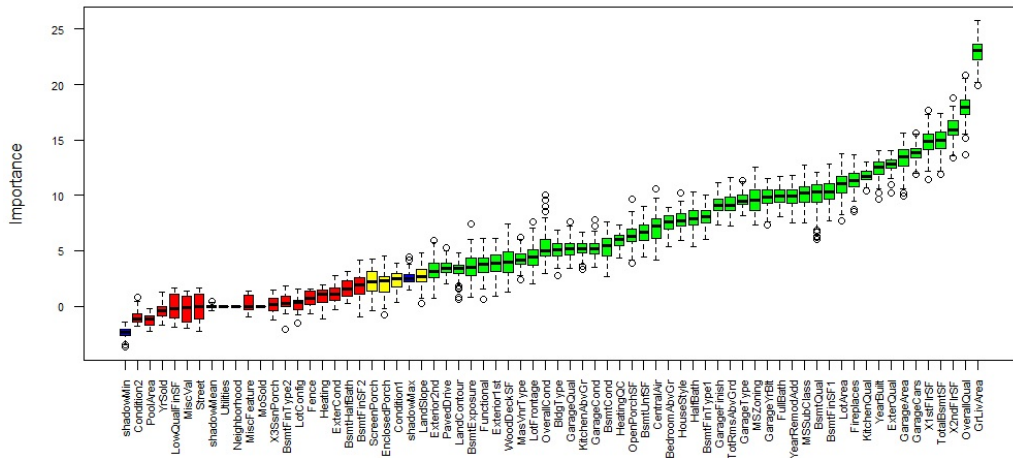
- Derived Features
- Feature:
 - $\text{TotalSF} = \text{TotalBsmtSF} + \text{GrLivArea} + \text{all other areas}$
 - $\text{TotalBaths} = \text{FullBath} + \text{BsmtFullBath} + .5(\text{HalfBath} + \text{BsmtHalfBath})$
 - $\text{YearBuilt_Age} = 2018 - \text{YearBuilt}$
- Dropping columns which are repeated:
 - $\text{TotalBsmtSF} = \text{sum of}(\text{BsmtFinSF1}, \text{BsmtFinSF2}, \text{BsmtUnfSF})$
 - $\text{GrLivAre} = \text{sum of}(\text{1stFlrSF}, \text{2ndFlrSF})$
- Taking log of 'SalePrice'

Understanding and visualizing data

- Missing values
- Boruta as starting point
 - Random Forest

```
> missingness
```

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street
	0.00	0.00	0.00	17.74	0.00	0.00
	Alley	LotShape	LandContour	Utilities	LotConfig	Landslope
	93.77	0.00	0.00	0.00	0.00	0.00
	Neighborhood	Condition1	Condition2	BldgType	HouseStyle	OverallQual
	0.00	0.00	0.00	0.00	0.00	0.00
	OverallCond	YearBuilt	YearRemodAdd	RoofStyle	RoofMatl	Exterior1st
	0.00	0.00	0.00	0.00	0.00	0.00



```
> train.boruta.selected.features.stats
```

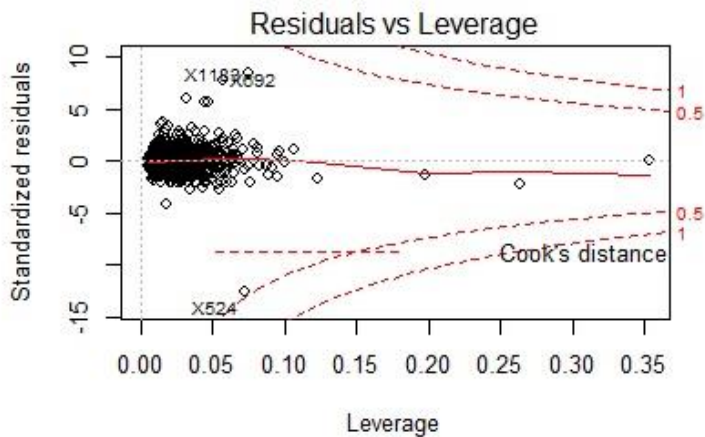
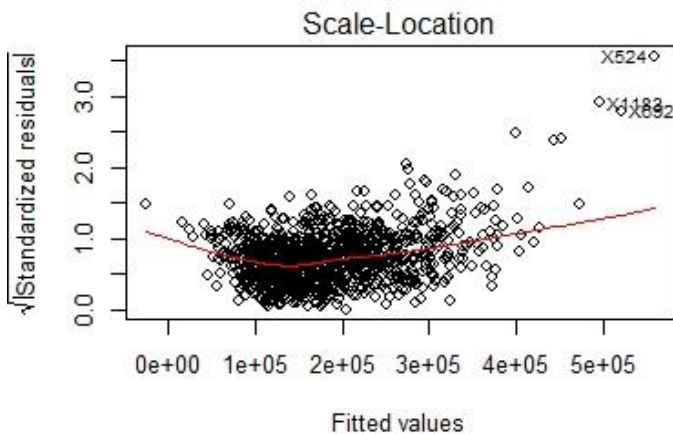
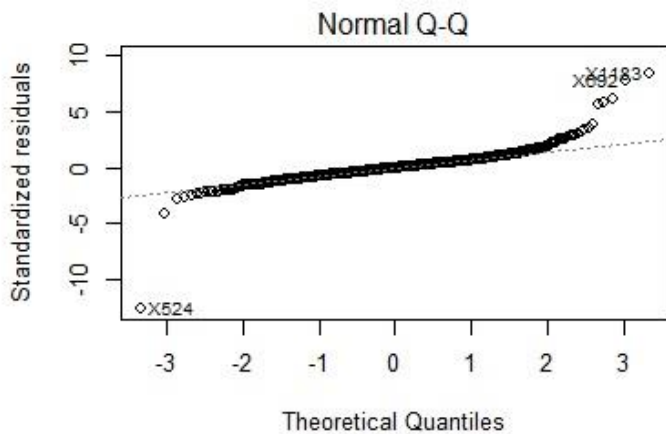
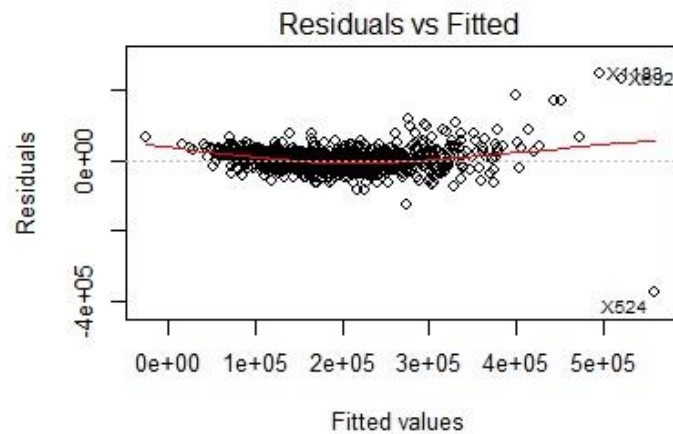
	meanImp	medianImp	minImp	maxImp	normHits	decision
MSSubClass	10.23146182	10.19289121	7.0381690	13.4160559	1.000000000	Confirmed
MSZoning	9.54442241	9.54223872	6.6190851	13.5386744	1.000000000	Confirmed
LotFrontage	4.84000063	4.81041786	1.0591037	8.2433175	0.945891784	Confirmed
LotArea	11.06137068	11.07546492	7.6883059	14.6355573	1.000000000	Confirmed
Street	0.32034836	0.60605652	-2.2139295	1.6233930	0.000000000	Rejected
LandContour	3.29294944	3.33744953	0.8217137	5.0138084	0.825651303	Confirmed
Utilities	0.00000000	0.00000000	0.0000000	0.0000000	0.000000000	Rejected
LotConfig	0.01045448	-0.34702343	-1.6685811	2.2380066	0.000000000	Rejected
Landslope	2.78879262	2.78509174	-0.1224927	6.3589709	0.597194389	Confirmed
Neighborhood	0.00000000	0.00000000	0.0000000	0.0000000	0.000000000	Rejected
Condition1	2.31997212	2.37257459	-0.7026733	4.6660047	0.452905812	Rejected
Condition2	-1.05018471	-1.39913091	-2.5243537	0.5650928	0.000000000	Rejected
BldgType	5.14919430	5.13505760	2.5850106	7.1341928	0.993987976	Confirmed
HouseStyle	7.80085816	7.90203825	4.2218717	10.1629154	1.000000000	Confirmed
OverallQual	17.83794999	17.97926215	13.8490065	20.3389953	1.000000000	Confirmed
OverallCond	5.19890016	4.97132116	2.8419043	9.9911439	0.993987976	Confirmed
YearBuilt	12.45174947	12.49480428	9.2918448	14.1892539	1.000000000	Confirmed
YearRemodAdd	10.02056144	10.08768364	5.0212588	13.4186005	1.000000000	Confirmed

Linear models

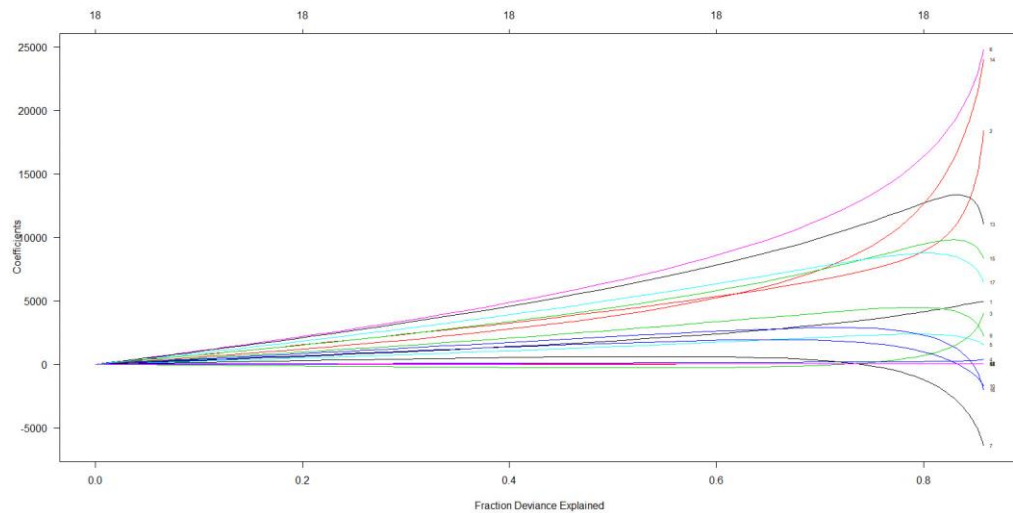
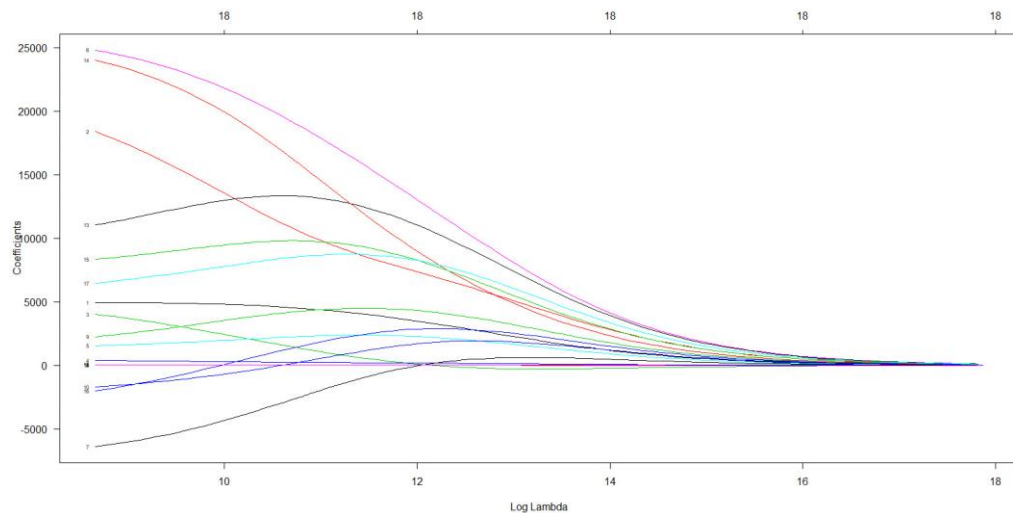
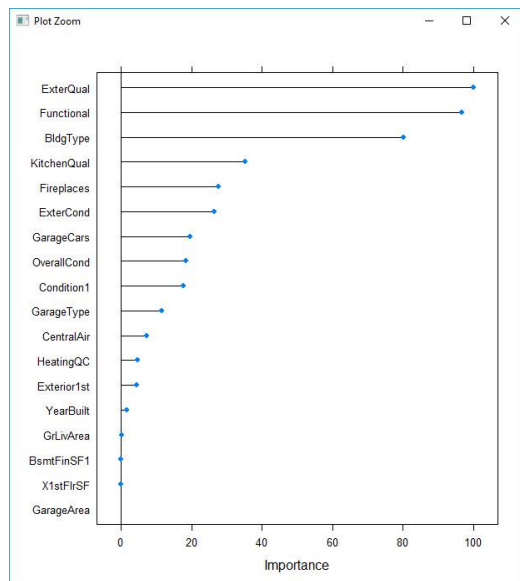
- Individual team models, best “crowdsourced” results selected
- Multivariable linear models
- Ridge
- Lasso
- Elastic

Plot different values for different attempts by different team members for different models

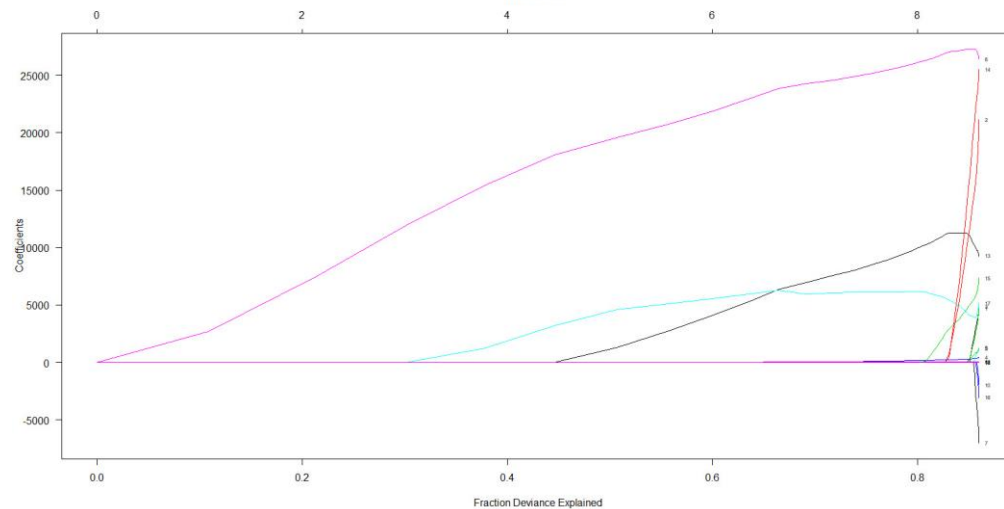
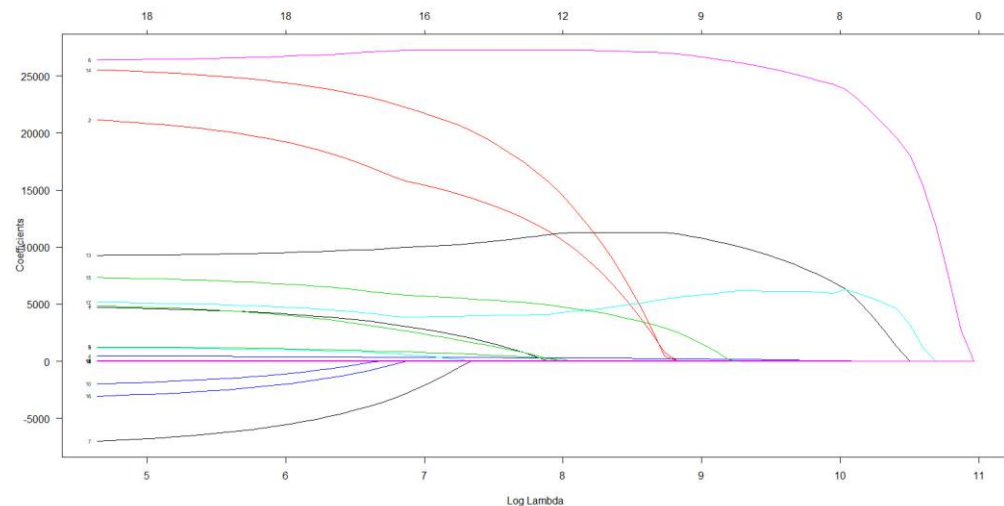
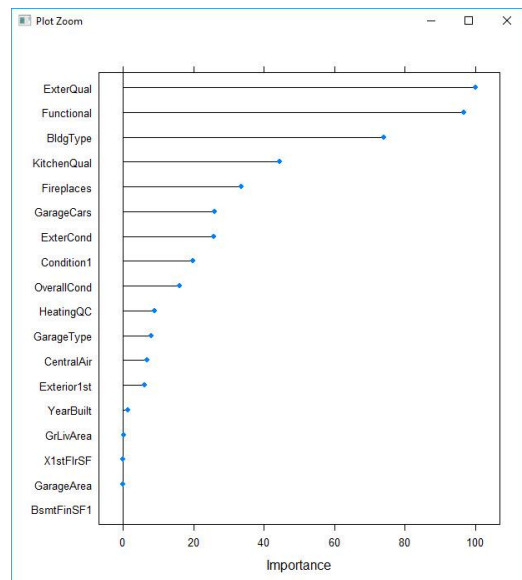
Im



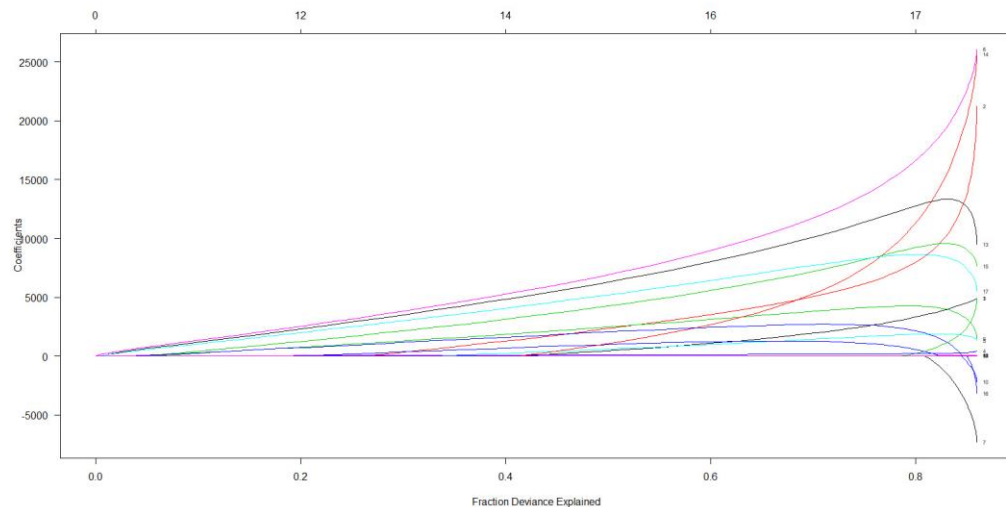
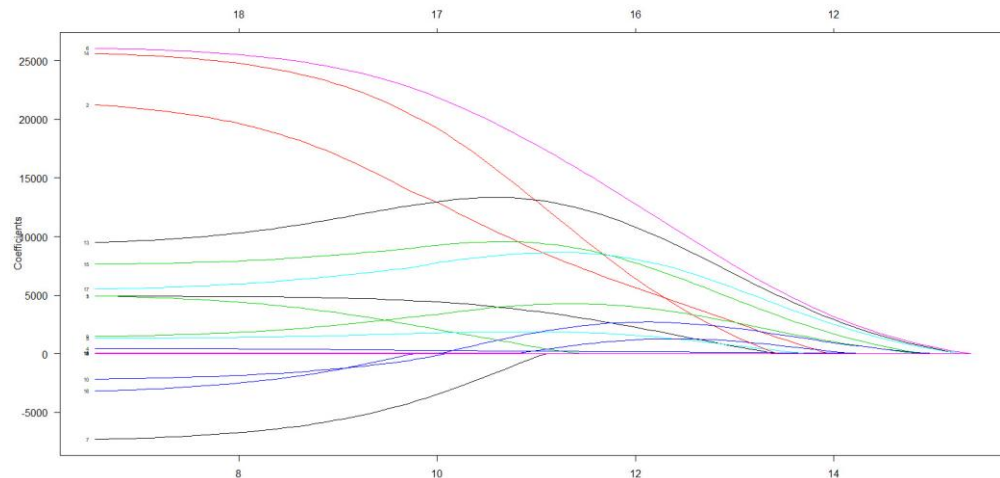
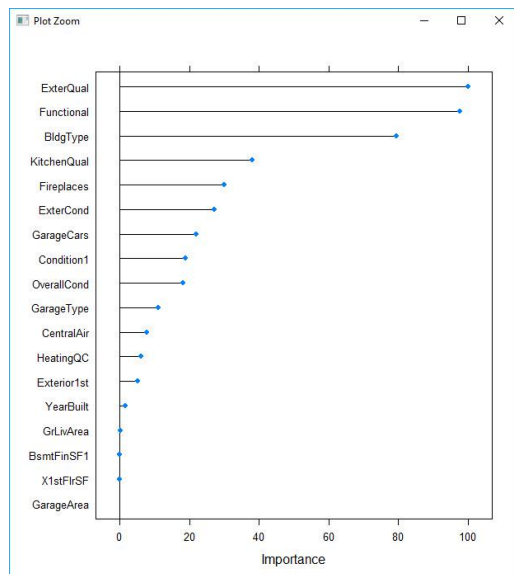
ridge



lasso

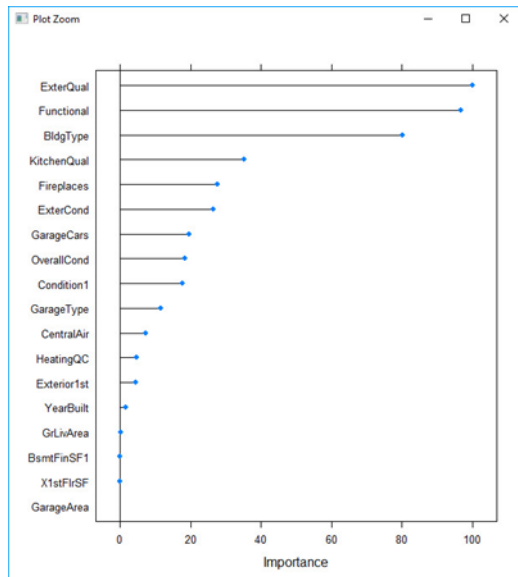


elastic net

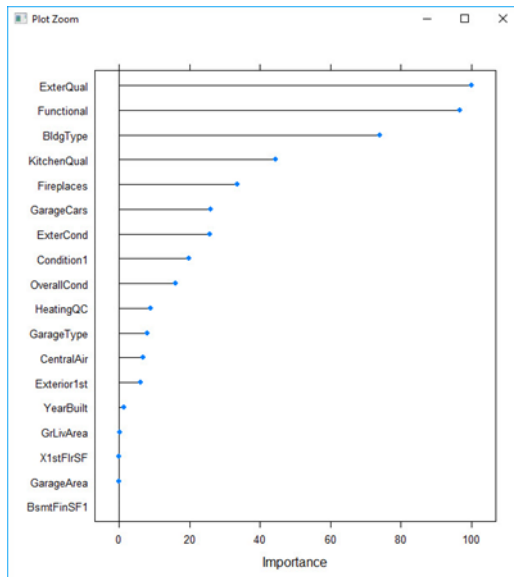


Comparing variables importance

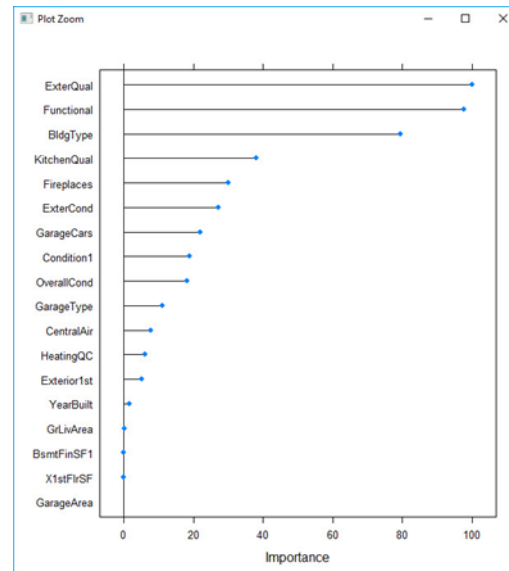
Ridge



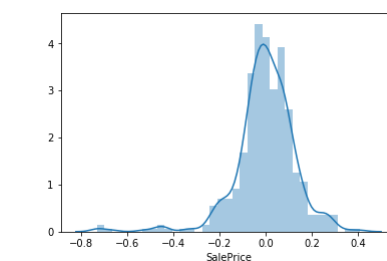
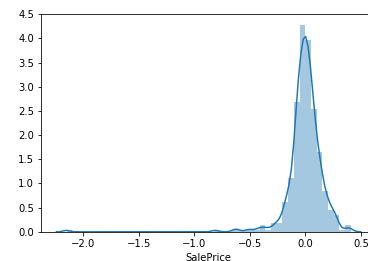
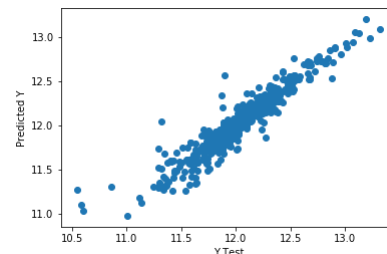
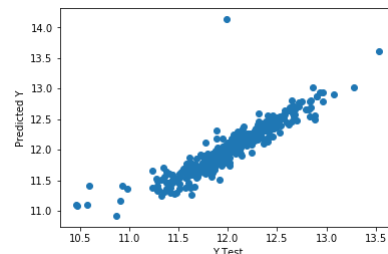
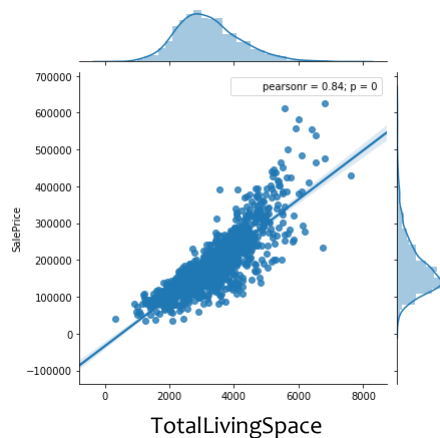
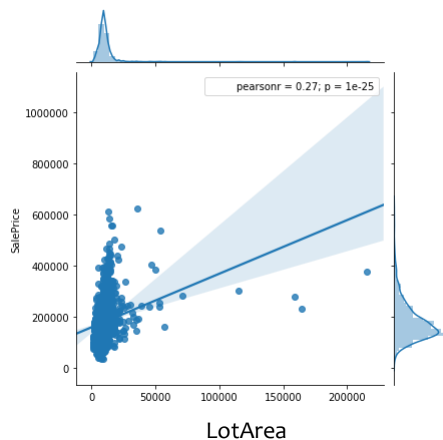
lasso



Elastic



Predicting house prices using linear models



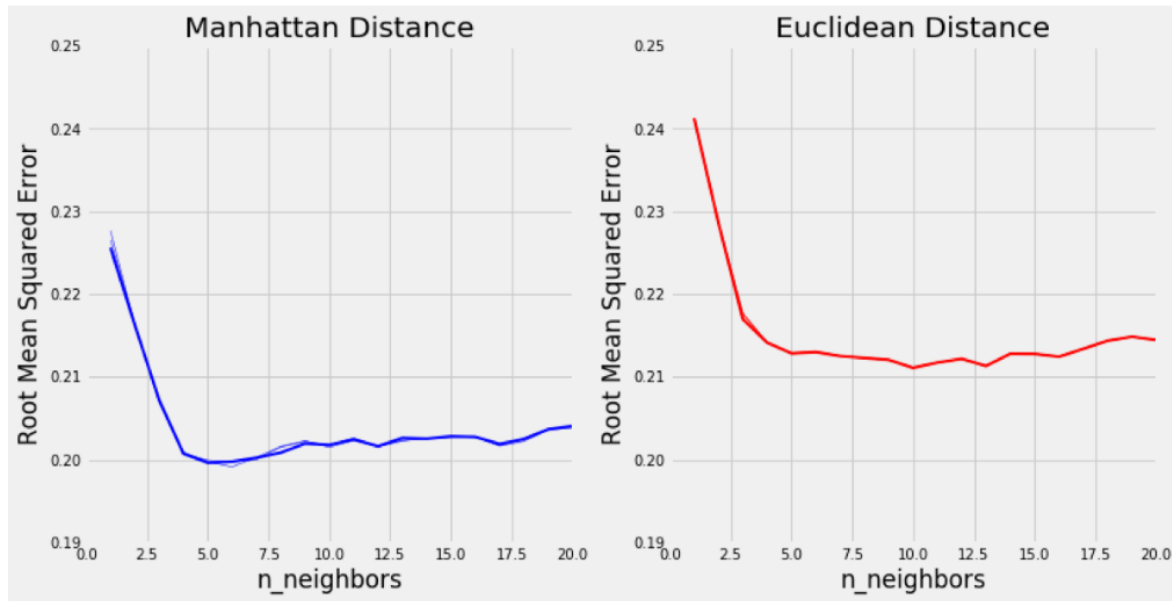
Model	RMSE	Comments
Linear Regression	0.1685	49 features, TotalLivingSpace
Linear Regression	0.1317	Squared # Bedrooms

Predicting house prices using k-nearest neighbors regression

New Features:
TotalLivingSpace instead of **LotArea**
Squared # BedRooms

Normalized the values: 0 - 1

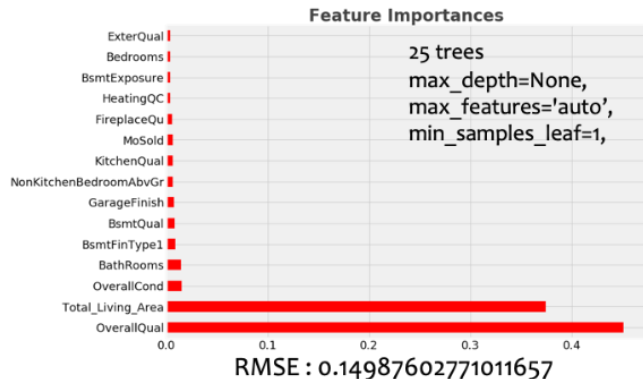
Metrics
Manhattan Distance
Euclidean Distance



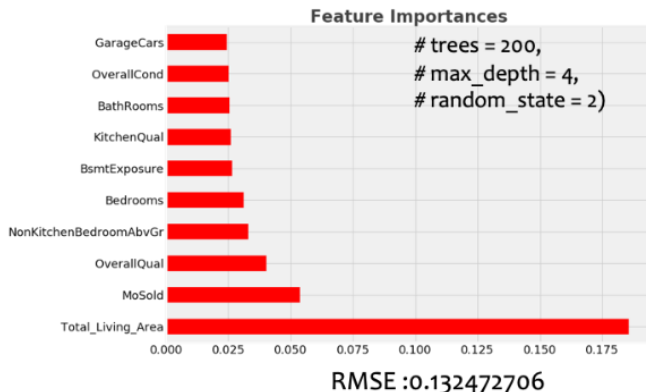
Best **RMSE** obtained when **K = 5** **0.19915480**

Predicting house prices using Tree-based models

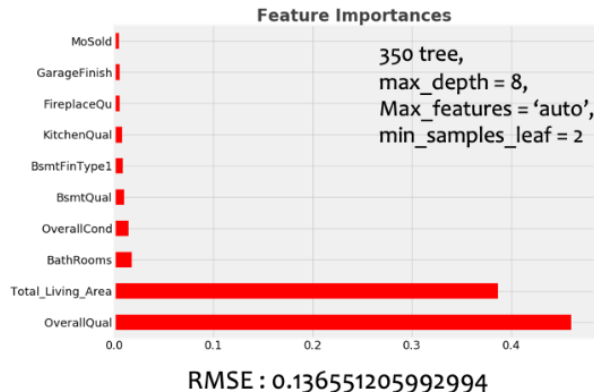
Random Forest



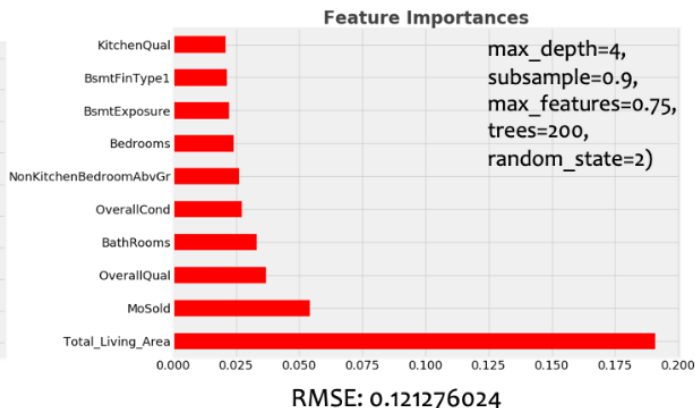
Gradient Boosting Regressor



Optimal Random Forest



Optimal Gradient Boosting Regressor



Decision Tree

0.2383957875036891

Random Forest

0.14987602771011657

0.136551205992994

Boosting Regressor

0.132472706

0.121276024

Stochastic
Boosting Regressor

0.13312094

0.122846123

XGBoost Regressor

1.963953

0.167252

0.1492399

Conclusions and lessons learned

- Feature selection is difficult and subjective!
- Heart on common sense, mind on statistics
- The selected features reveal similar results in all regression models.
- The following features appear to have much influence on a house price:
 - Total Living Area
 - Lot Area
 - Month Sold
 - Overall Quality
- For the tree-based models,
 - Total Living Area
 - Month Sold
 - Number of Bathroom
 - Kitchen Quality