

Machine Learning Project

Ames House Prices Prediction

Yiling Lin

Project Information

Dataset:

Ames Housing dataset

79 explanatory variables describing every aspect of residential homes in Ames, Iowa

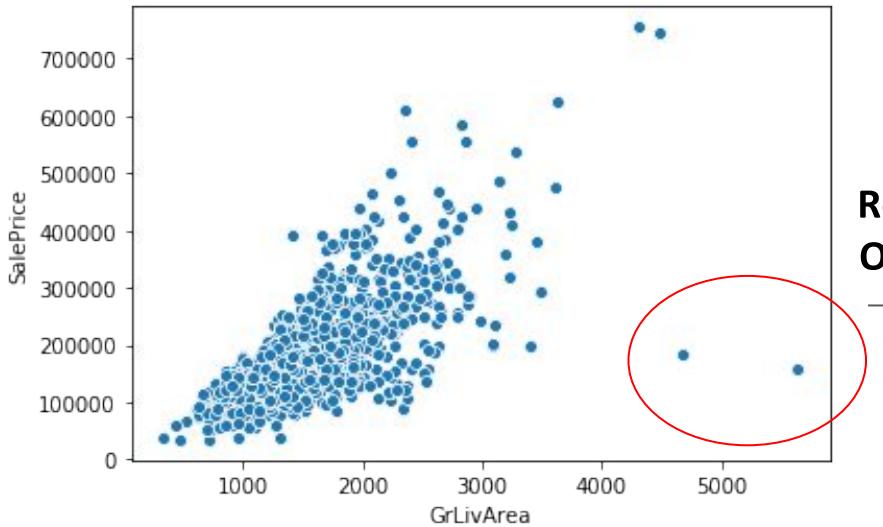
Part 1: Exploratory Data Analysis:

Does Remodeling Increase The Profit of Your Property Sale?

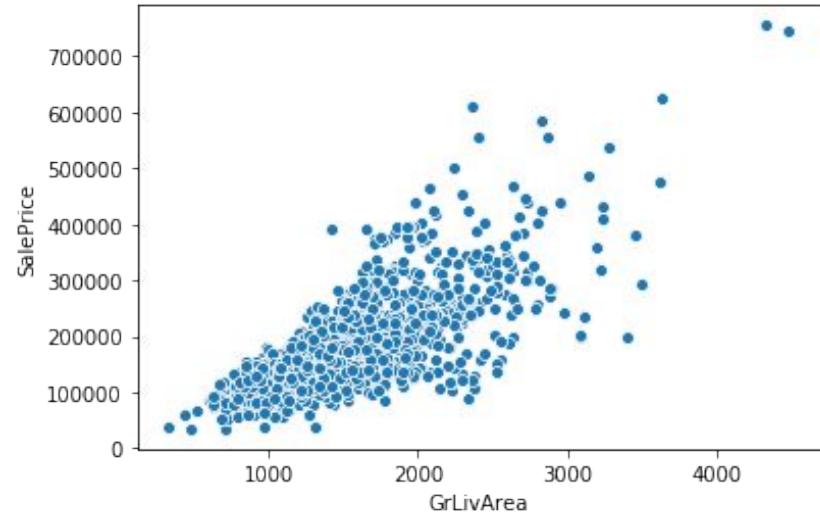
Part 2: Machine Learning:

Predict Sale Prices of houses in Ames, Iowa

Pre-Processing



Remove
Outliers



Feature Engineering

Columns Added

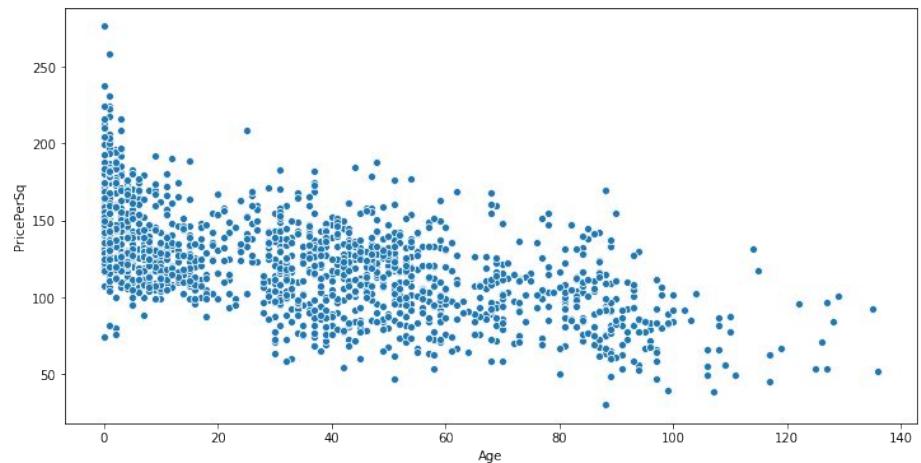
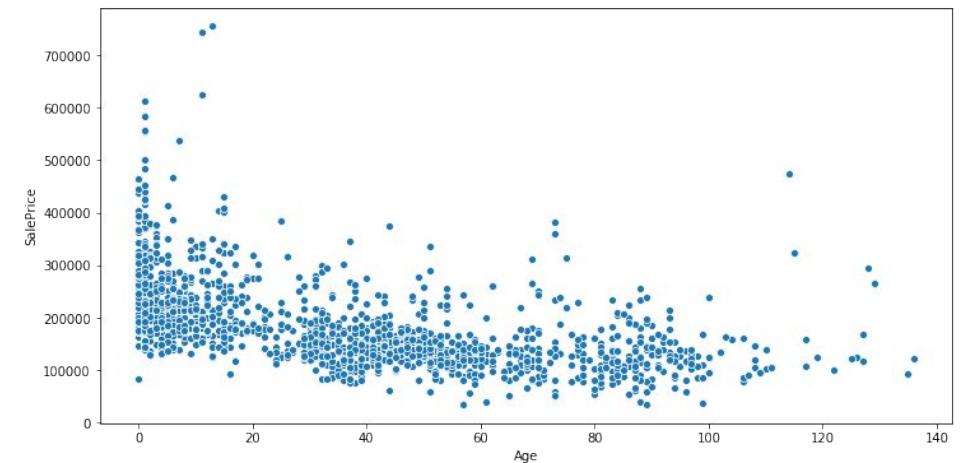
- **Age** = YrSold - YrBuilt
- **Remodel Age** = YrSold - YearRemodAdd
- **Remodeled** - if the house is remodeled
- **HasGarage** - if the house has a garage
- **GarageAge** = YrSold - GarageYrBlt
- **PricePerSq** = Price Per Square Foot
(Sale Prices/Total Ground Living Area)

Correlation Matrix

	LotFrontage	LotArea	OverallQual	OverallCond	YearBuilt	YearRemodAdd	MasVnArea	BsmtFinSF1	BsmtFinSF2	BsmtUnfSF	TotalBsmtSF	1stFlrSF	2ndFlrSF	LwvltQual	GrLivArea	BsmntAbvGr	TotRmsAbvGrd	KitchenAbvGr	Fireplaces	GarageCars	GarageArea	WoodDeckSF	OpenPorchSF	EnclosedPorch	3SsnPorch	ScreenPorch	PoolArea	MscVal	YrsOld	SalePrice	PricePerSq	Age	RemodeAge	GarageAge									
LotFrontage	1	0.39	0.24	-0.056	0.12	0.081	0.16	0.15	0.055	0.14	0.32	0.41	0.066	0.042	0.36	0.077	0.0051	0.2	0.042	0.27	0.039	0.34	0.25	0.061	0.29	0.32	0.082	0.12	0.015	0.075	0.046	0.12	0.005	0.078	0.37	0.15	0.12	-0.081	-0.06				
LotArea	0.39	1	0.091	-0.0027	0.0075	0.0075	0.085	0.18	0.11	-0.0033	0.22	0.27	0.041	0.0055	0.23	0.15	0.05	0.12	0.074	0.12	-0.017	0.18	0.26	-0.032	0.15	0.16	0.17	0.062	-0.017	0.021	0.045	0.036	0.039	-0.014	0.27	0.12	-0.0091	0.0084	0.032				
OverallQual	0.24	0.091	1	-0.09	0.057	0.055	0.4	0.22	-0.056	0.31	0.54	0.47	0.29	-0.03	0.59	0.1	-0.039	0.55	0.27	0.1	-0.18	0.42	0.39	0.55	0.6	0.56	0.24	0.3	-0.11	0.031	0.066	0.045	-0.031	0.027	0.8	0.49	-0.57	-0.55	-0.55				
OverallCond	0.056	-0.0027	-0.09	1	0.38	0.075	-0.13	-0.043	0.04	-0.14	-0.17	0.04	0.031	0.025	0.077	0.0054	0.024	0.054	0.12	-0.19	-0.06	0.013	-0.087	0.005	-0.0022	-0.32	-0.19	-0.15	-0.0027	-0.008	0.007	0.005	0.0028	0.069	0.044	-0.078	0.03	0.38	-0.072	0.33			
YearBuilt	0.12	0.075	0.57	-0.38	1	0.59	0.31	0.25	-0.049	0.15	0.4	0.28	0.007	-0.18	0.19	0.19	-0.038	0.47	0.24	0.071	-0.17	0.091	0.14	0.83	0.54	0.48	0.22	0.18	-0.39	0.032	-0.05	0.0052	0.034	0.013	0.52	0.56	-1	-0.59	-0.82				
YearRemodAdd	0.081	0.075	0.55	0.075	0.59	1	0.18	0.12	-0.067	0.18	0.29	0.24	0.14	-0.062	0.29	0.12	-0.012	0.44	0.18	-0.041	-0.15	0.19	0.11	0.64	0.42	0.37	0.2	0.22	-0.19	0.046	-0.038	-0.0034	-0.01	0.036	0.51	0.46	-0.59	-1	-0.64				
MasVnArea	0.16	0.085	0.4	-0.13	0.31	0.1	1	0.24	-0.072	0.11	0.34	0.32	0.17	-0.069	0.37	0.07	0.027	0.27	0.2	0.01	-0.07	0.27	0.24	0.25	0.36	0.36	0.16	0.1	-0.11	0.02	0.063	-0.021	-0.0071	0.048	0.23	0.31	-0.18	-0.25					
BsmtFinSF1	0.15	0.18	0.22	-0.043	0.25	0.12	0.24	1	-0.05	-0.52	0.47	0.4	-0.16	-0.067	0.14	0.06	0.074	0.048	-0.099	-0.11	-0.083	0.011	0.24	0.15	0.23	0.27	0.2	0.073	0.1	0.029	0.068	0.053	0.0047	0.016	0.41	0.45	-0.25	-0.12	-0.15				
BsmtFinSF2	-0.055	0.11	-0.058	0.04	-0.049	-0.067	-0.072	-0.05	1	-0.21	0.11	0.1	-0.099	0.015	-0.0071	0.16	0.071	-0.076	-0.032	-0.016	-0.041	-0.034	0.048	-0.088	-0.038	-0.017	0.068	0.005	0.0049	0.032	-0.011	0.016	0.05	0.069	0.089								
BsmtUnfSF	0.14	-0.033	0.31	-0.14	0.15	0.18	0.11	-0.52	0.21	1	0.44	0.33	0.034	0.028	0.25	0.42	-0.096	0.29	0.042	0.17	0.03	0.25	0.052	0.19	0.21	0.18	0.018	0.046	0.013	0.024	0.021	0.012	0.035	0.024	-0.041	0.21	0.047	0.15	-0.18	-0.19			
TotalBsmtSF	0.32	0.22	0.54	-0.17	0.04	0.29	0.34	0.47	0.11	0.44	1	0.8	-0.21	-0.034	0.041	0.3	0.025	0.33	-0.066	0.051	-0.07	0.027	0.33	0.33	0.45	0.48	0.23	0.22	-0.097	0.041	0.093	0.038	0.019	-0.015	0.65	0.68	0.4	-0.3	-0.33				
1stFlrSF	0.41	0.27	0.47	-0.14	0.28	0.24	0.32	0.4	0.1	0.33	0.8	1	-0.23	-0.014	0.53	0.29	0.045	0.38	-0.14	0.13	0.075	0.4	0.4	0.23	0.45	0.48	0.24	0.18	-0.064	0.06	0.095	0.063	0.021	0.013	0.63	0.31	0.28	0.24	-0.23				
2ndFlrSF	0.066	0.041	0.29	0.038	0.007	0.014	0.17	0.16	-0.026	0.003	0.21	0.23	1	0.064	0.09	0.18	0.023	0.042	0.02	0.012	0.024	0.076	0.016	0.28	0.32	-0.34	-0.008	0.14	0.069														
LwvltQual	0.042	0.055	-0.03	0.025	-0.18	-0.062	-0.069	-0.067	0.015	0.028	0.034	-0.014	0.064	1	0.14	0.047	0.059	0.0043	0.027	0.11	0.0075	0.13	-0.021	-0.036	-0.094	-0.068	0.025	0.019	0.061	-0.0043	0.027	0.066	0.038	-0.029	-0.06	-0.16	0.18	0.006	0.035				
GrLivArea	0.56	0.23	0.59	-0.077	0.19	0.29	0.37	0.14	-0.001	0.25	0.41	0.53	0.69	0.14	1	0.014	-0.017	0.64	0.42	0.54	0.11	0.83	0.46	0.23	0.48	0.46	0.25	0.3	0.013	0.023	0.11	0.12	-0.016	0.036	0.73	-0.086	0.2	-0.29	-0.23				
BsmntAbvGr	-0.077	0.15	0.01	-0.054	0.19	0.12	0.076	0.66	0.16	0.42	0.42	0.3	0.23	-0.18	-0.047	0.014	1	-0.15	-0.069	0.035	-0.15	-0.041	-0.063	0.13	0.12	0.13	0.17	0.017	-0.057	-0.049	0.0023	0.024	0.045	0.023	0.067	0.23	0.35	0.18	-0.11	-0.12			
BsmtUnfSF	-0.0051	0.05	-0.039	0.12	0.38	-0.012	0.028	0.074	0.071	-0.096	0.025	0.045	0.023	-0.059	0.017	0.15	0.1	-0.054	0.012	0.047	0.038	0.023	0.03	-0.077	0.021	-0.024	0.024	-0.024	0.007	0.045	0.027	0.017	0.009	0.036	0.005	0.075							
FulBlr	0.2	0.12	0.55	-0.19	0.47	0.44	0.27	0.048	-0.076	0.29	0.33	0.38	0.42	-0.0045	0.64	-0.069	-0.054	1	0.13	0.36	0.13	0.55	0.24	0.48	0.47	0.4	0.19	0.25	-0.11	0.036	-0.0075	0.046	0.014	0.019	0.56	0.095	-0.47	-0.44	-0.48				
HalfBlr	0.044	0.074	0.27	-0.06	0.24	0.18	0.02	-0.099	-0.032	0.042	-0.066	0.024	0.011	-0.027	0.024	0.082	0.042	0.003	0.13	0.021	0.023	0.068	0.034	0.2	0.2	0.11	-0.025	0.095	0.024	0.073	0.013	0.0014	0.099	0.28	-0.093	0.24	-0.18	-0.2					
BdrssAbvGr	0.27	0.12	0.01	0.013	0.071	0.021	0.041	0.1	-0.11	0.017	0.051	0.13	0.5	0.11	0.54	0.15	0.047	0.36	0.23	1	0.2	0.06	0.11	-0.065	0.086	0.065	0.047	0.094	0.024	-0.024	0.044	0.073	0.0078	0.036	0.17	-0.36	0.069	0.039	0.063				
KitchenAbvGr	-0.0039	-0.017	-0.18	-0.087	0.17	-0.17	-0.15	-0.037	-0.083	0.041	0.03	0.07	0.073	0.06	0.0075	0.11	-0.041	-0.038	0.13	-0.068	0.2	1	0.26	-0.12	-0.12	-0.05	-0.064	-0.09	-0.07	0.037	-0.025	-0.052	-0.013	0.062	0.032	-0.14	0.3	0.18	0.15	0.12			
TotRmsAbvGrd	0.34	0.18	0.4	-0.056	0.091	0.12	0.27	0.011	-0.034	0.25	0.27	0.4	0.61	0.13	0.83	0.063	-0.023	0.55	0.34	0.68	0.26	1	0.32	0.14	0.36	0.33	0.16	0.22	-0.0057	-0.0062	0.061	0.068	0.025	0.025	0.034	0.54	0.18	-0.059	0.19	-0.15			
Fireplaces	0.25	0.26	0.39	-0.022	0.14	0.11	0.24	0.24	-0.048	0.052	0.33	0.04	0.19	-0.021	0.46	0.13	0.03	0.24	0.2	0.11	0.12	0.32	1	0.043	0.043	1	0.59	0.57	0.22	0.2	0.16	0.024	0.012	0.19	0.074	0.047	0.17	0.03	0.045	0.023	0.015	0.11	-0.045
GarageCars	0.081	-0.032	0.55	-0.32	0.83	0.64	0.25	0.15	-0.088	0.19	0.33	0.23	0.68	-0.036	0.23	0.12	-0.077	0.074	0.48	0.2	0.065	0.12	0.14	0.043	0.57	0.22	0.2	0.3	0.024	0.016	0.069	0.008	0.017	0.024	0.047	0.16	0.15	0.64	-1				
GarageArea	0.19	0.16	0.16	0.15	0.48	0.37	0.36	0.27	-0.017	0.18	0.48	0.48	0.13	-0.066	0.46	0.17	-0.024	0.024	0.032	0.0075	0.073	0.044	-0.052	0.068	0.19	0.075	0.053	0.074	0.077	0.031	1	0.058	0.032	0.011	0.11	0.042	0.05	0.039	0.076				
WoodDeckSF	0.082	0.17	0.24	-0.0027	0.22	0.2	0.16	0.2	0.068	0.0056	0.23	0.24	0.09	0.025	0.25	0.17	0.044	0.19	0.11	0.047	0.059	0.16	0.2	0.22	0.27	0.22	1	0.054	-0.13	-0.033	-0.074	0.065	0.0094	0.023	0.32	0.2	0.22	-0.22					
OpenPorchSF	0.12	0.062	0.3	-0.03	0.18	0.22	0.013	0.005	0.13	0.22	0.18	0.2	0.019	0.03	0.057	0.024	0.025	0.022	0.017	0.022	0.16	0.2	0.21	0.23	0.054	1	0.092	-0.095	0.077	0.028	0.12	0.19	-0.23	-0.23									
EnclosedPorch	0.015	-0.017	0.11	0.07	-0.39	-0.19	-0.11	-0.1	-0.021	0.041	0.063	0.061	0.013	-0.049	0.087	0.011	-0.095	0.042	0.037	0.0057	0.024	0.1	0.13	-0.12	-0.13	-0.092	0.1	0.037	-0.083	0.06	0.018	0.01	0.13	-0.22	0.39	0.19	0.3						
3SsnPorch	0.075	0.021	0.031	0.029	0.034	0.02	0.029	0.02	0.031	0.024	0.016	0.004	0.003	0.029	0.035	0.036	0.036	-0.048	-0.025	-0.0062	0.012	0.024	0.036	0.036	-0.033	0.0052	-0.037	1	0.031	-0.074	0.0034	0.019	0.045	0.039	-0.031	-0.044	-0.023						
ScreenPorch	0.0045	0.066	0.055	0.05	0.038	0.063	0.068	0.012	0.009	0.012	0.027	0.011	0.024	0.032	0.029	0.035	0.036	0.044	0.066	-0.026	0.019	0.027	0.069	0.033	0.0052	0.0074	0.056	1	0.024	0.009	0.011	0.042	0.005	0.039	0.076								
PoolArea	0.12	0.036	0.043	0.025	0.034	0.021	0.053	0.046	0.053	0.074	0.066	0.12	0.045	0.023	0.046	0.013	0.073	0.011	0.03	0.06	0.069	0.027	0.069	0.033	0.011	0.001	0.023	0.007	0.022														
MscVal	0.009	0.039	-0.031	0.034	-0.001	0.004	0.024	0.019	-0.021	0.016																																	

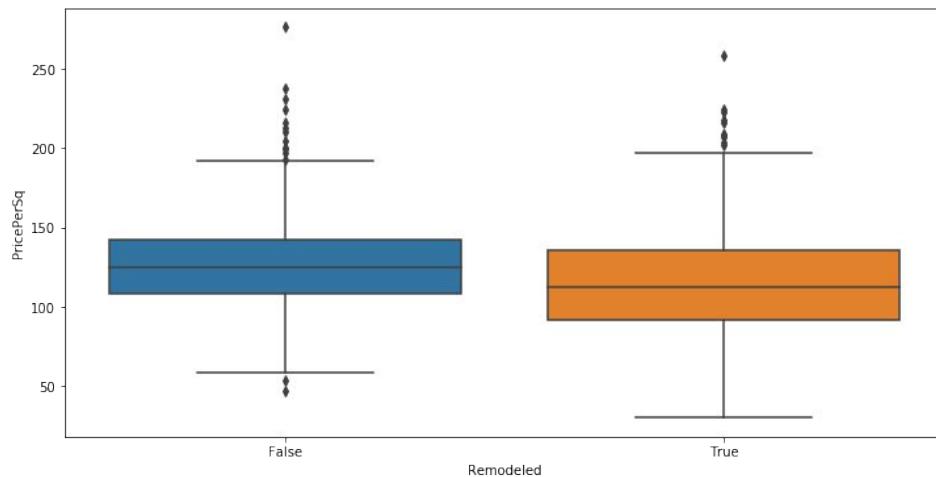
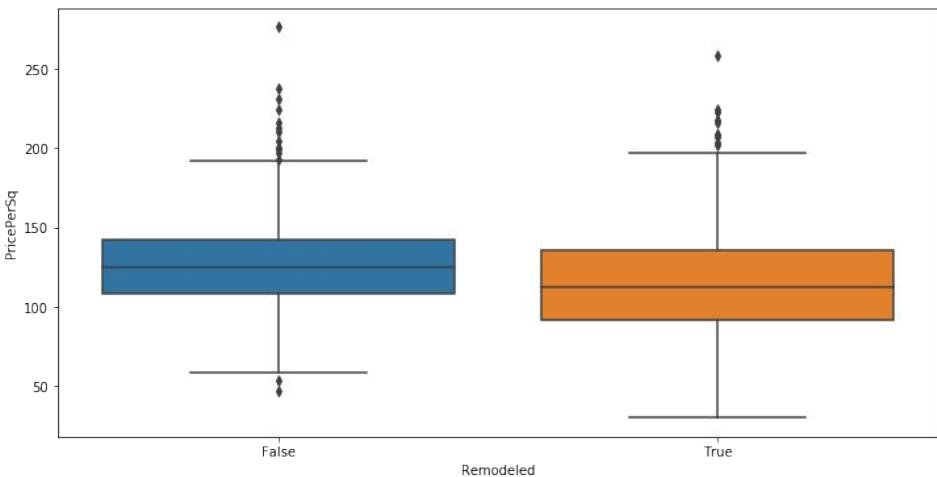
House Age, Remodel, Price

House Age vs Price

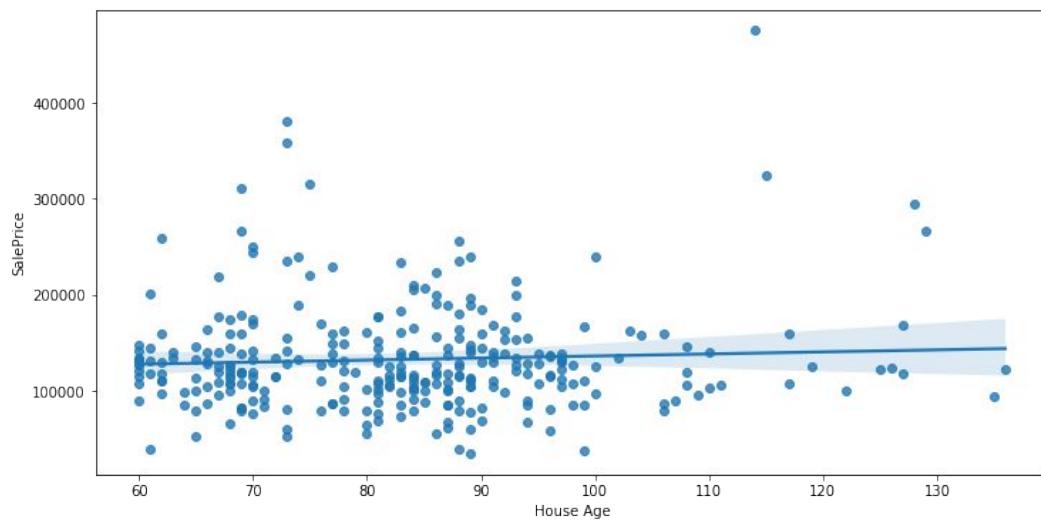
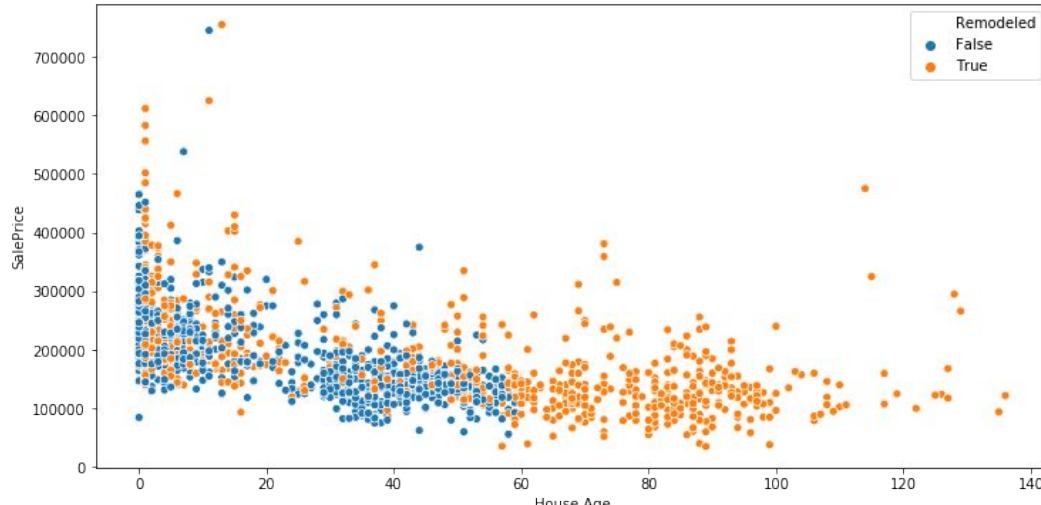


The older, the cheaper

Remodeled Houses vs Sale Prices & Price Per Sqft



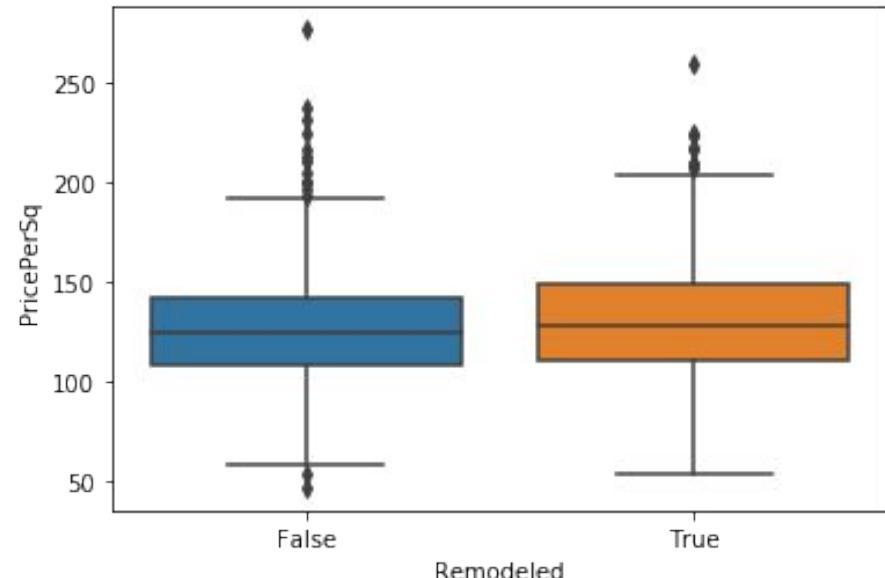
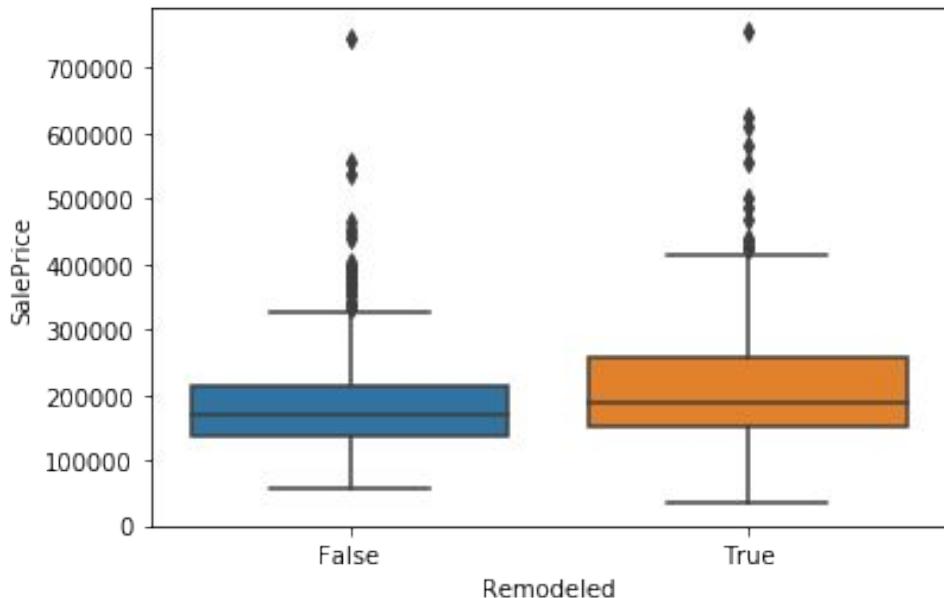
It seems that prices of remodeled houses are lower than non-remodeled houses?
Why?



House Age, Remodel & Sale Price

- Around 44% of remodeled houses in the dataset are houses over 59 years old.
- Houses over 59 years old have no significant price changes.

Remodeled Houses vs Sale Prices & Price Per Sqft 0-59 Year-Old House



Averagey, when the age of a house is below 60, the price of a remodeled house is higher than a non-remodeled one.

House Age vs Avg. Price Per Sqft

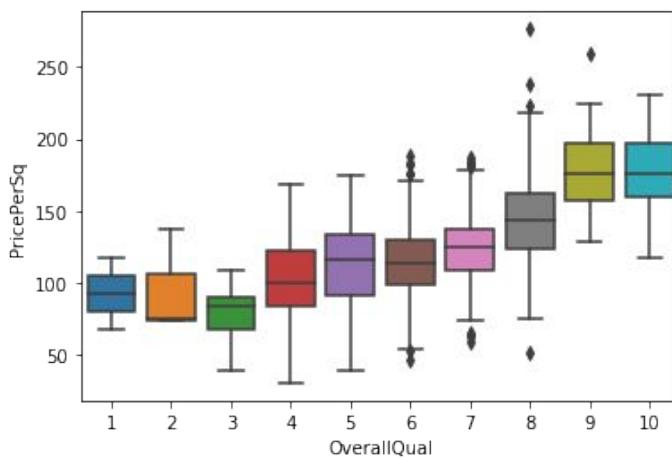
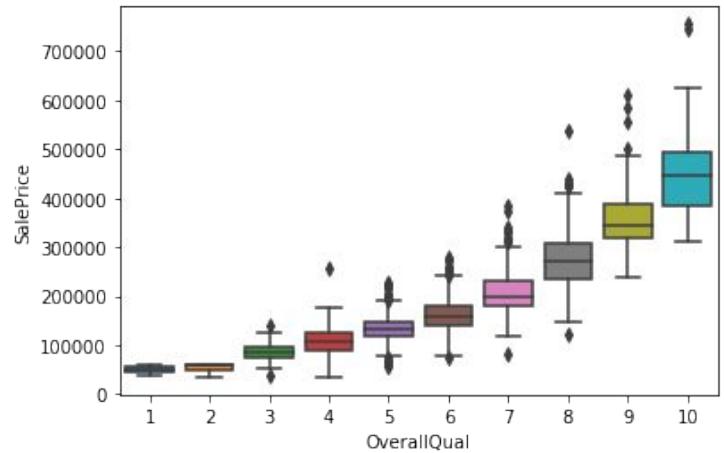
	0-19	20-39	40-59
Remodeled	\$141.14	\$126.20	\$116.54
Non-Remodeled	\$137.14	\$118.43	\$112.93
Increase %	3%	6.6%	3.2%

Remodeling brings the biggest impact on the price when the age of the house is between 20-39 years old.

House Age vs Avg. Sale Price

	0-19	20-39	40-59
Remodeled	\$255613.63	\$190699.10	156695.68
Non-Remodeled	\$227679.55	\$151715.08	139179.70
Increase %	12.3%	25.7%	12.6%

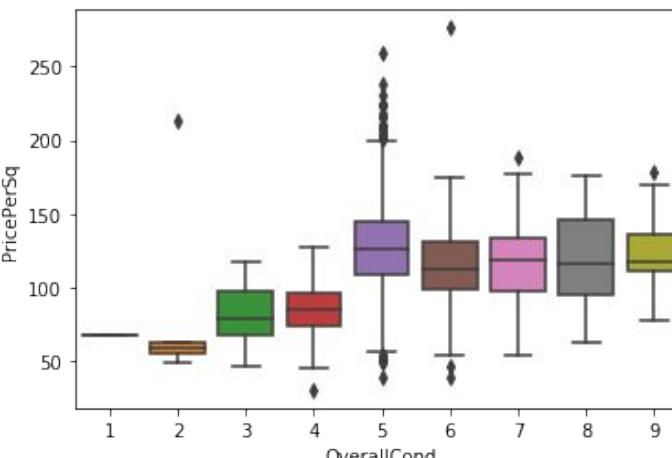
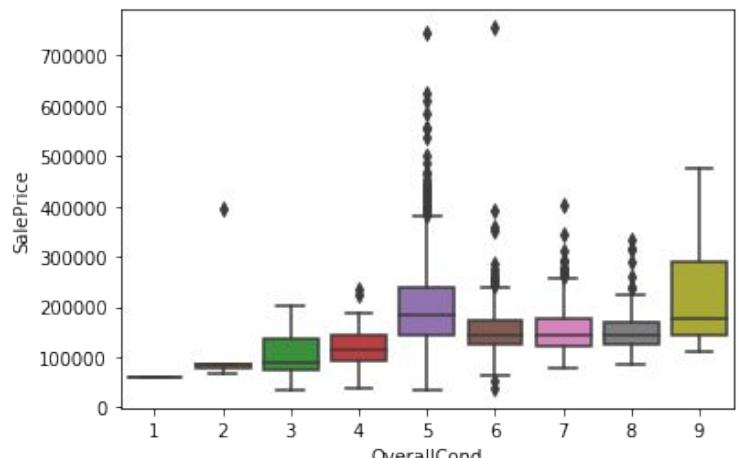
Quality & Condition



Overall Quality

Rates the overall material and finish of the house

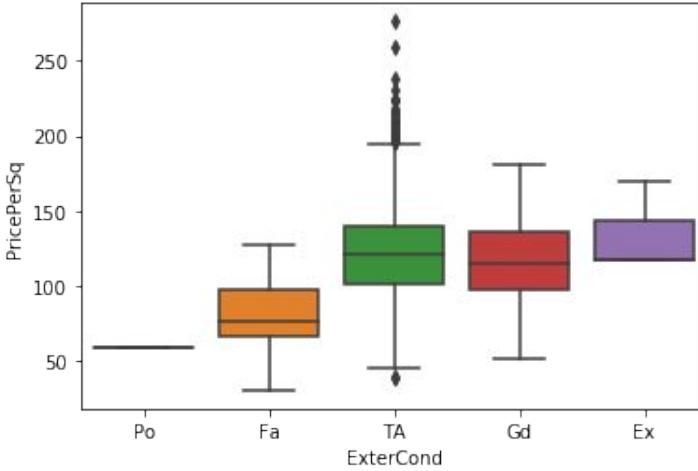
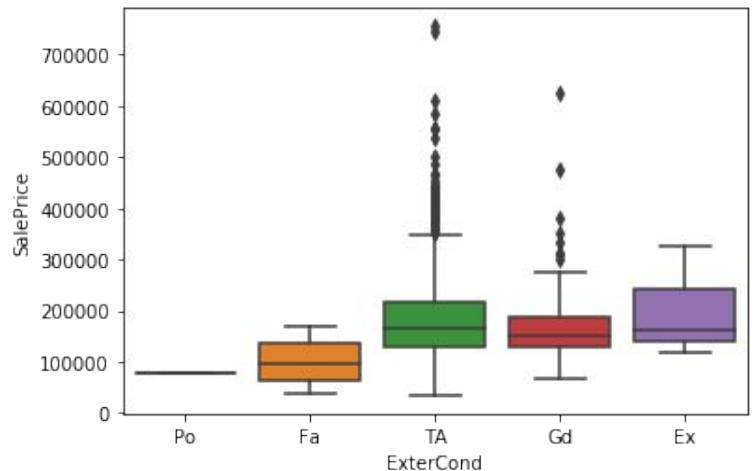
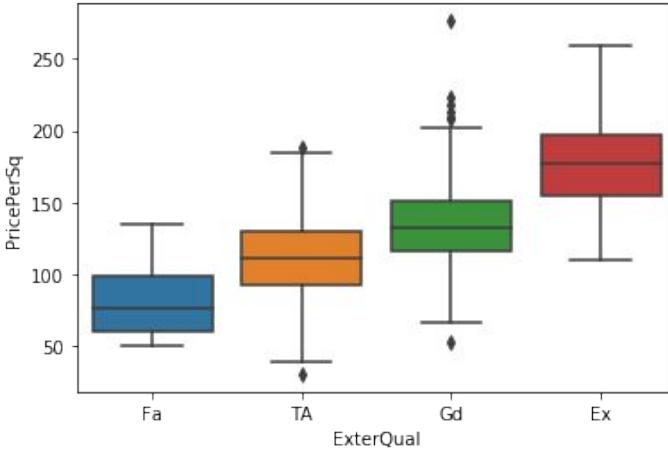
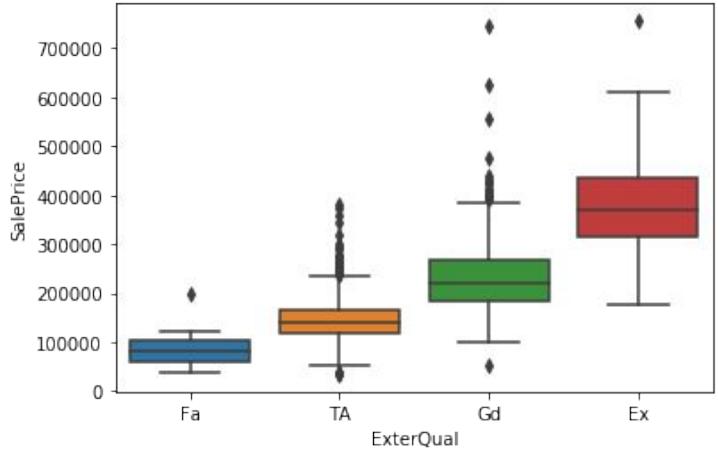
- Price increases when overall quality increase.



Overall Condition

Rates the overall condition of the house

- Better overall condition does not guarantee a better price.



Exterior Quality

Evaluates the quality of the material on the exterior

- Price increases when exterior quality increases.

Exterior Condition

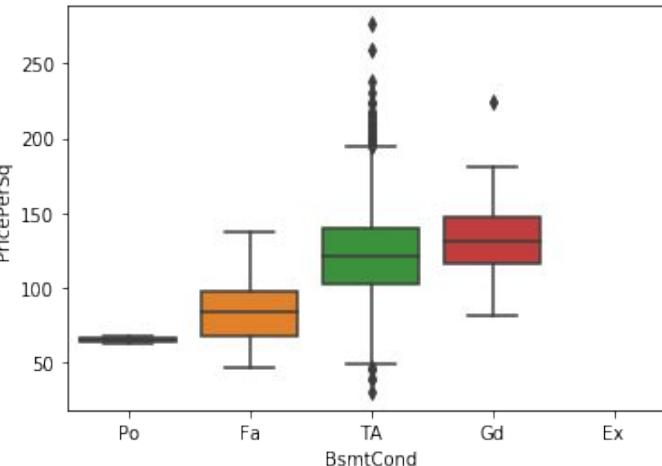
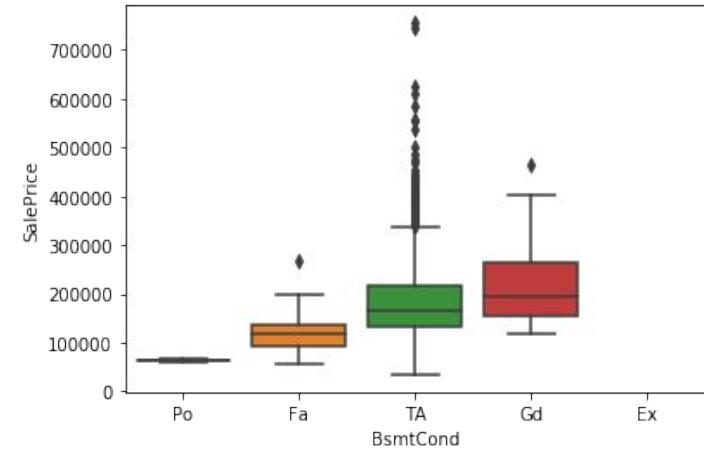
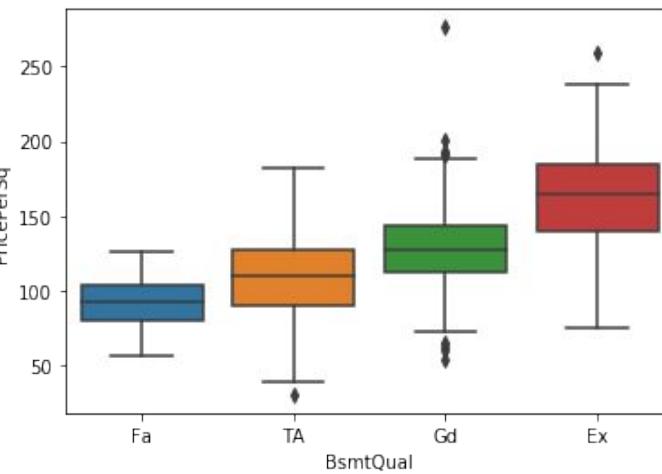
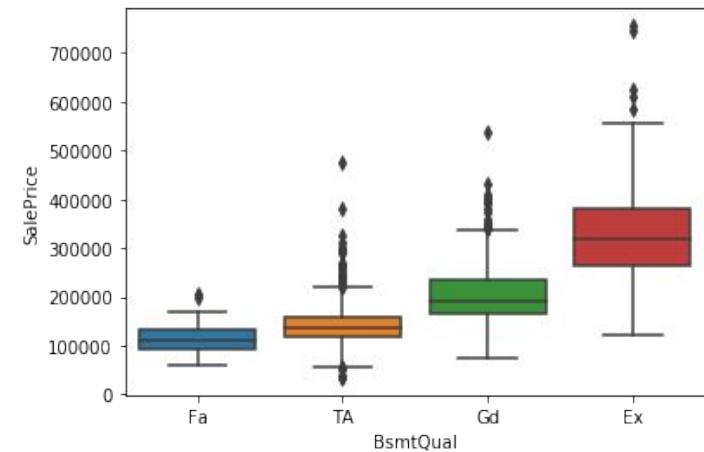
Evaluates the present condition of the material on the exterior

- Better exterior condition does not guarantee a better price.

Ex	Excellent
Gd	Good
TA	Average/Typical
Fa	Fair
Po	Poor

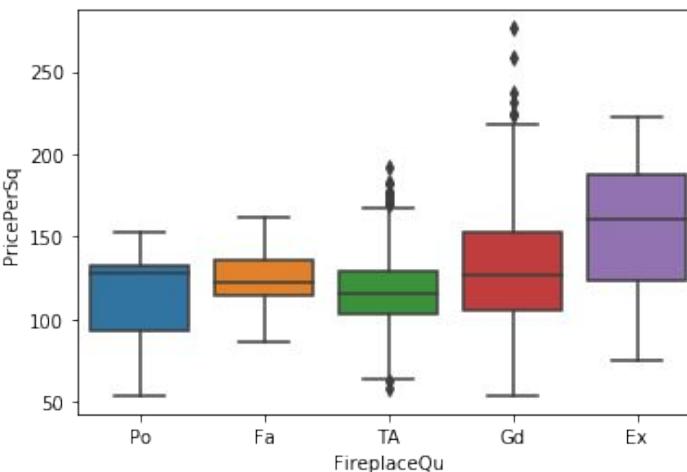
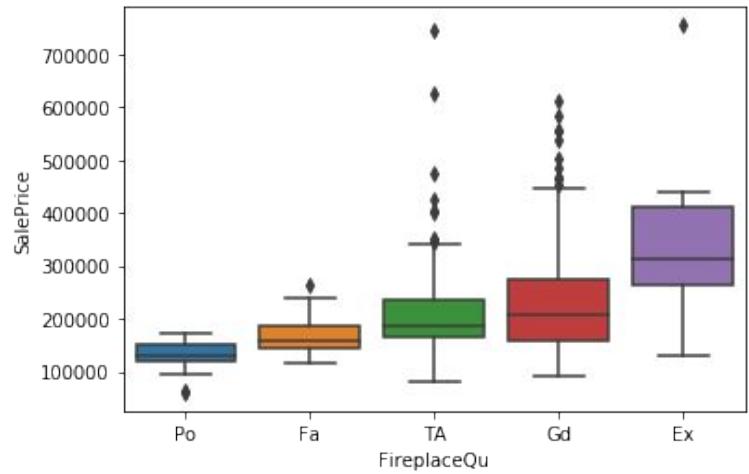
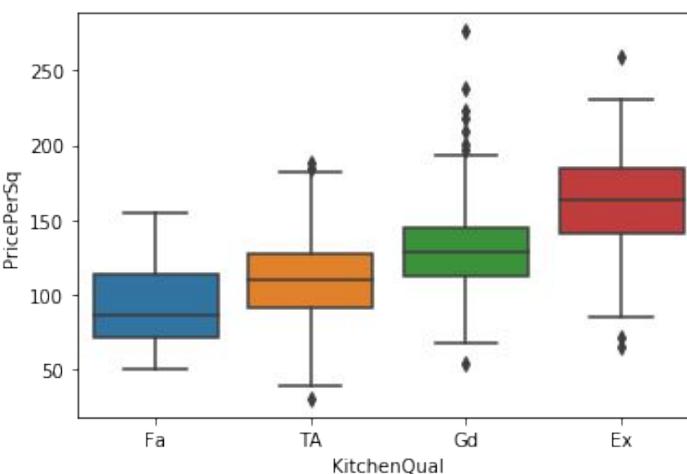
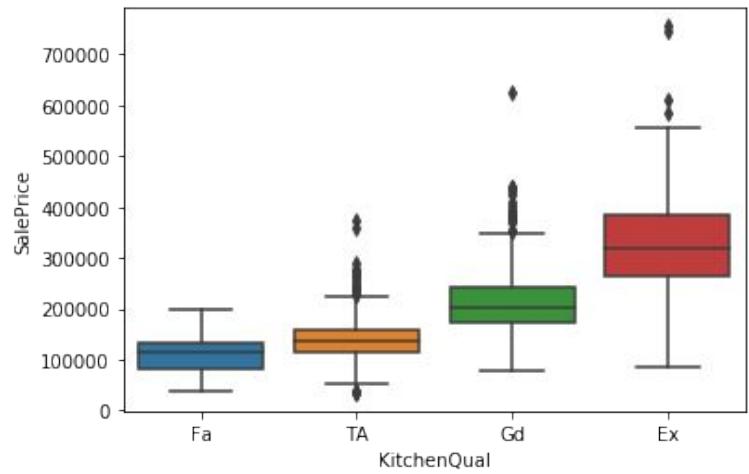
Basement Quality & Basement Condition

- Both the height of a basement and the condition of a basement are positively correlated to house prices.



Basement Quality	
Ex	Excellent (100+ inches)
Gd	Good (90-99 inches)
TA	Typical (80-89 inches)
Fa	Fair (70-79 inches)
Po	Poor (<70 inches)
NA	No Basement

Basement Condition	
Ex	Excellent
Gd	Good
TA	Average/Typical
Fa	Fair
Po	Poor



Kitchen Quality

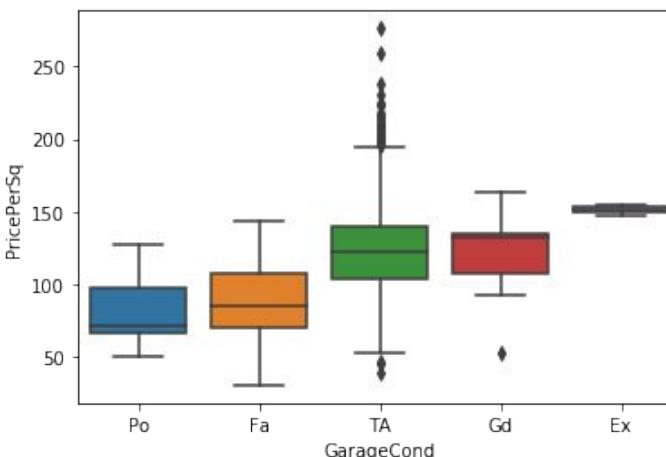
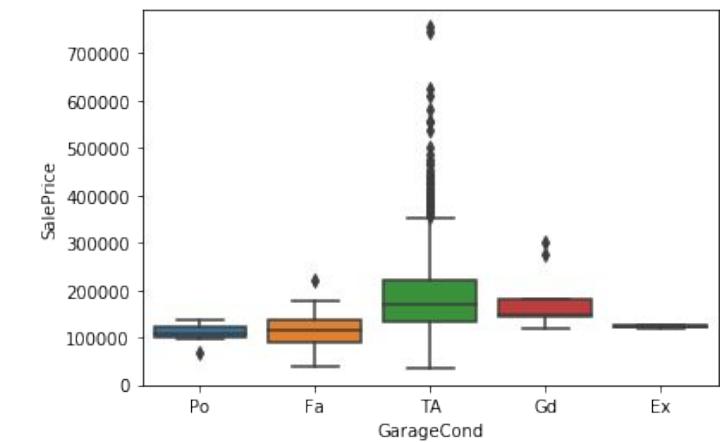
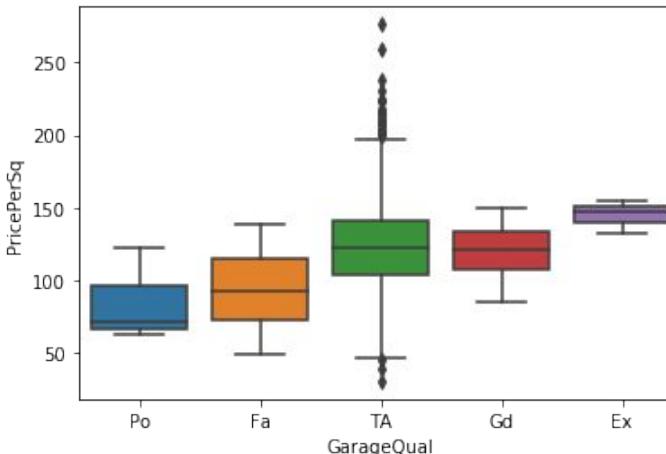
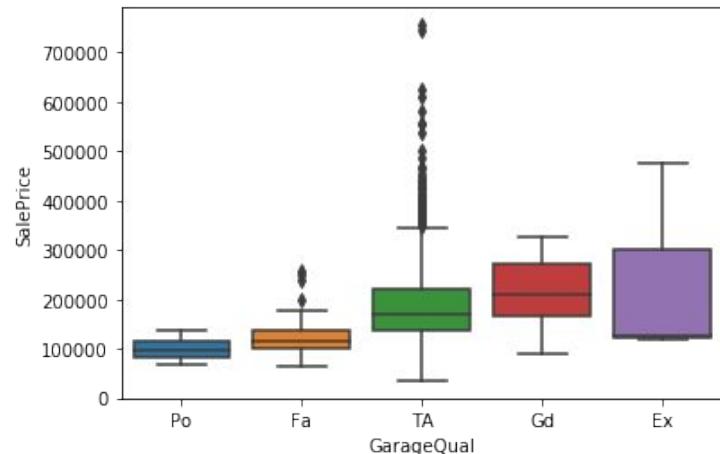
- As kitchen quality increases, sale prices and price per sqft increase.
- Averagely, price per sqft increases by 20% when updating the quality to one level up.

Fireplace Quality

- The distribution of fireplace quality is more skewed than kitchen quality. The avg price per sqft doesn't necessarily increase when quality improves.

Garage Quality & Garage Condition

- If you have a poor or fair quality garage, updating the garage to at least TA(Average/Typical) level will help increase the house price.
- Price Per Sqft doesn't increase from TA Quality to Gd Quality.
- Price Per Sqft increases by 30% from fair quality to typical/average quality. (from poor to fair, price per sqft only increases by 9%).



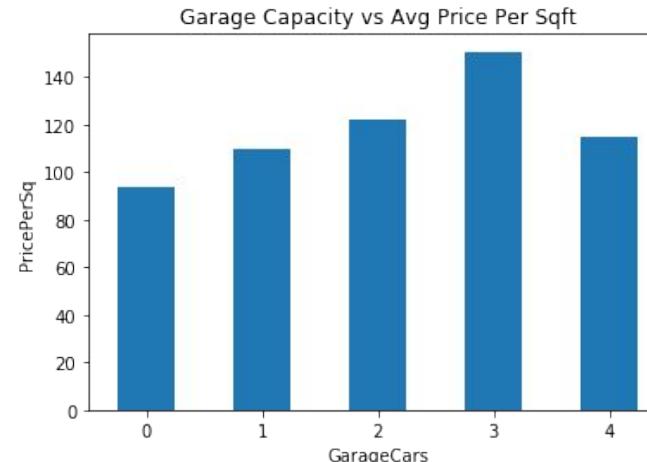
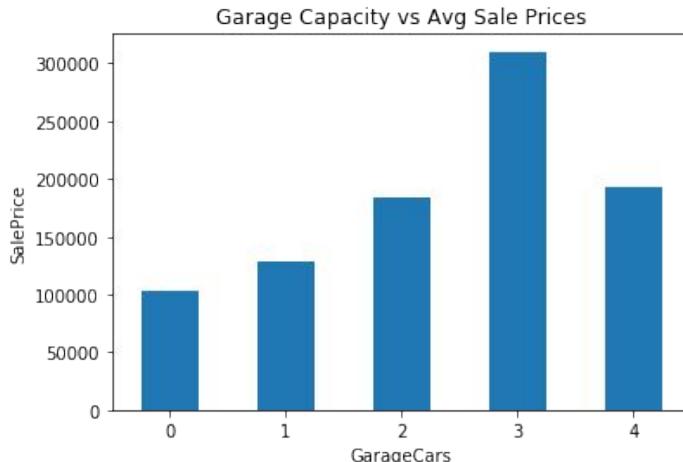
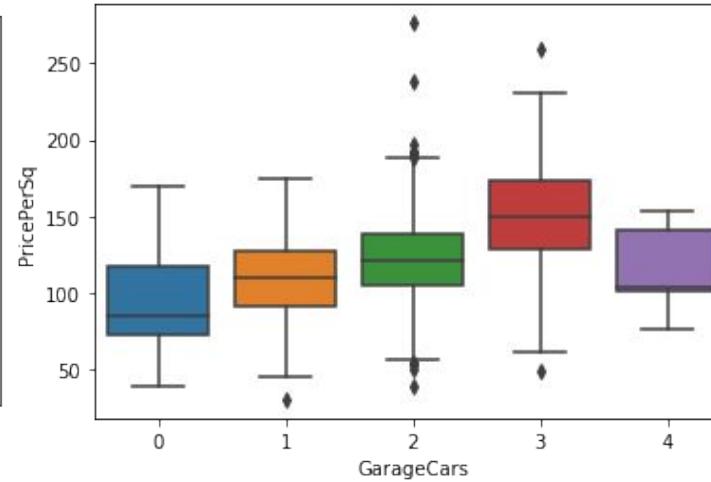
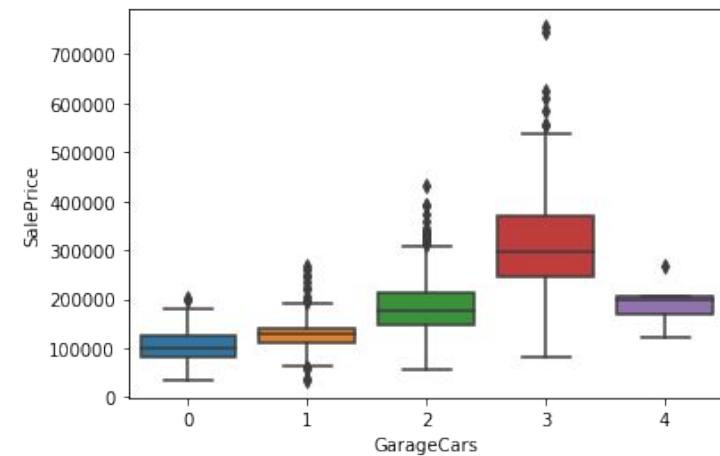
Ex	Excellent
Gd	Good
TA	Average/Typical
Fa	Fair
Po	Poor

- Garage capacity increases prices up to three cars:

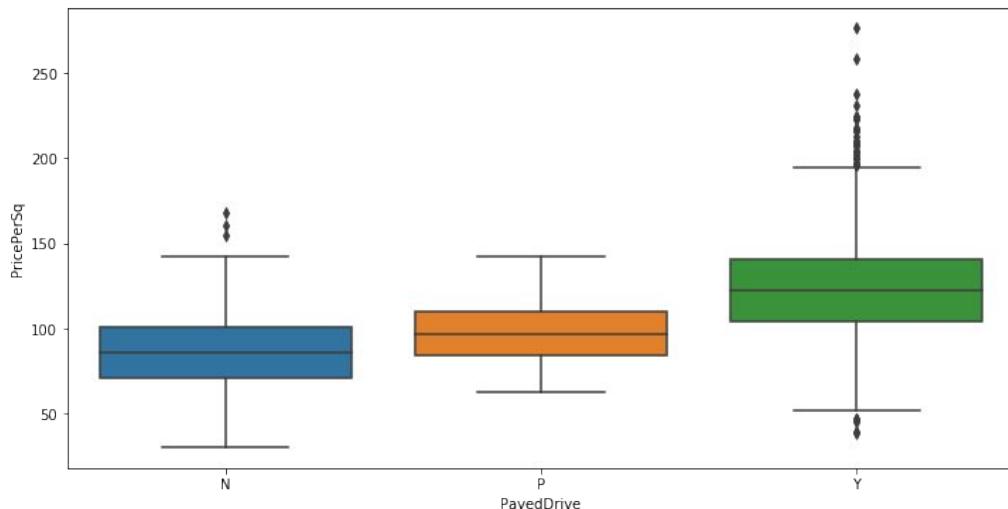
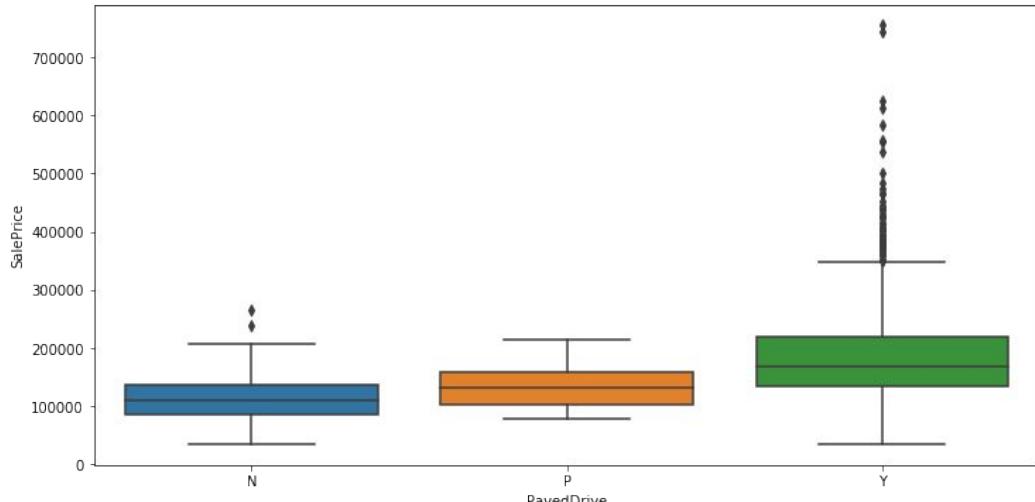
0 - 1 car:
Price per sqft
increases by 16%

1-2 cars:
Price per sqft
increases by 12%.

2-3 cars:
Price per sqft
increases by 23%



Additional Amenities

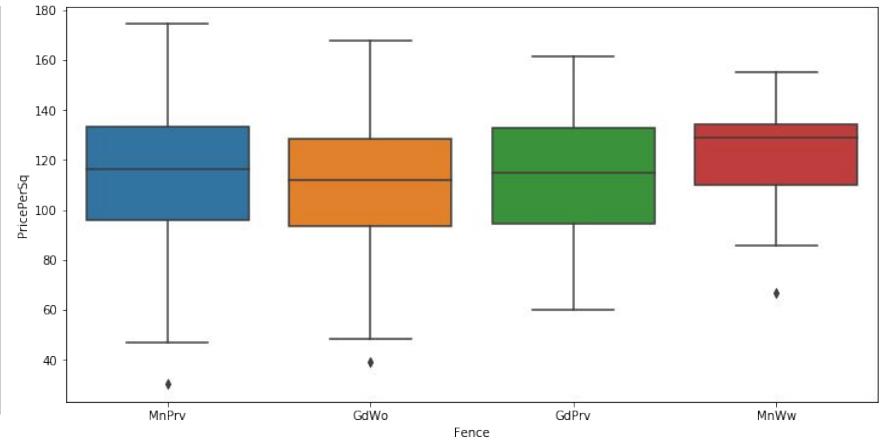
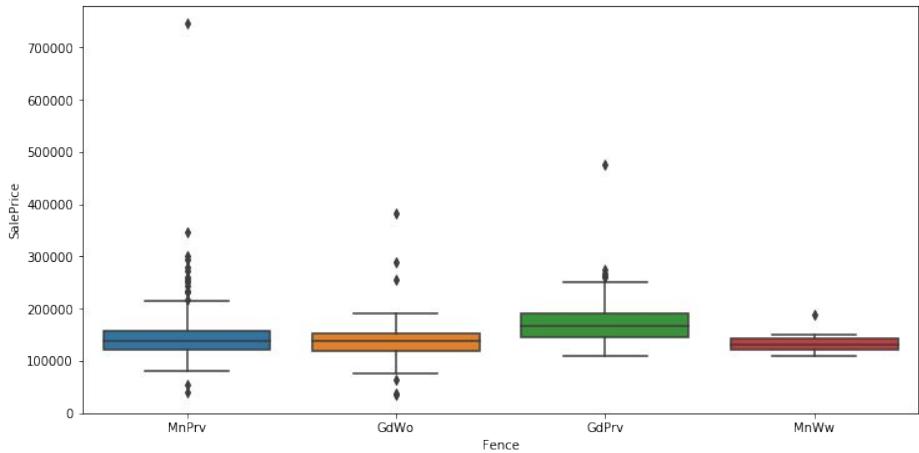


Paved Driveway vs Price

- The price of a house is higher when driveway is paved.
- From no paved driveway to a partial paved driveway, the avg price per sqft increases by 10%.
- From a partial paved driveway to a paved driveway, the avg price per sqft increases by 26%

Y	Paved
P	Partial Pavement
N	Dirt/Gravel

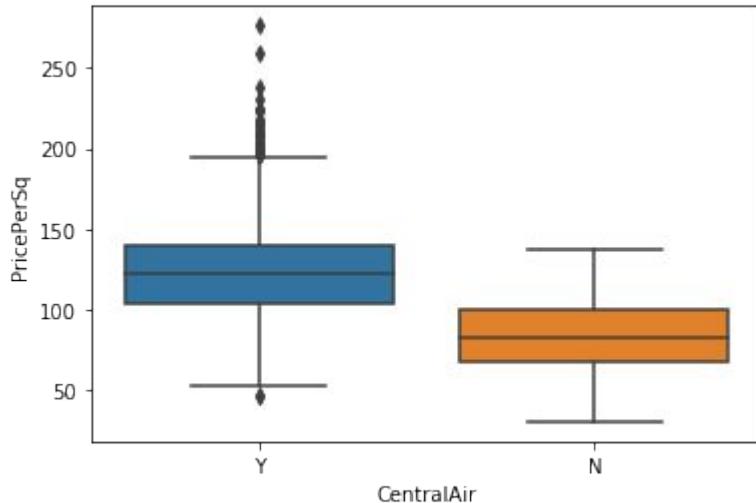
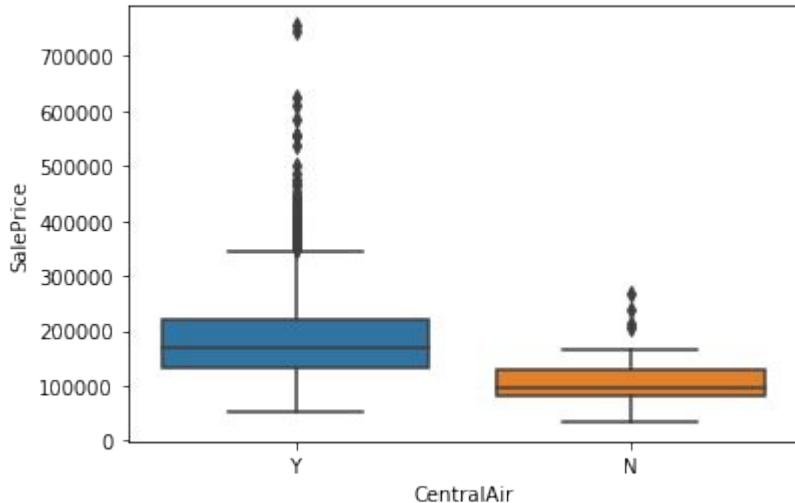
Fence vs Price



GdPrv	Good Privacy
MnPrv	Minimum Privacy
GdWo	Good Wood
MnWw	Minimum Wood/Wire
NA	No Fence

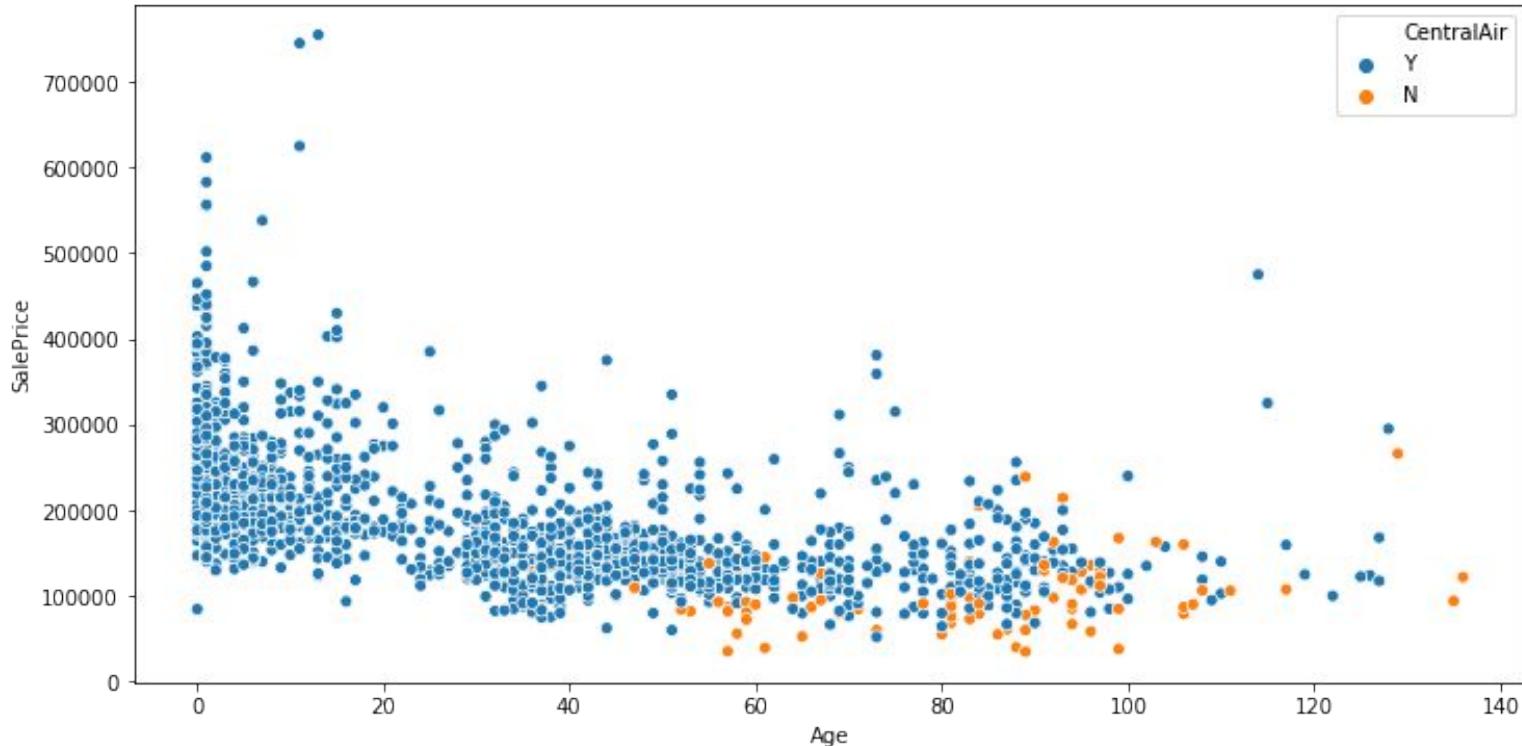
Averagely, the price of a house with a better fence quality is not necessarily higher.

Central Air



Averagey, the sale price of a house featuring central air is 70% higher than the price of a house that has no central air. Price per sqft increase by 46%.

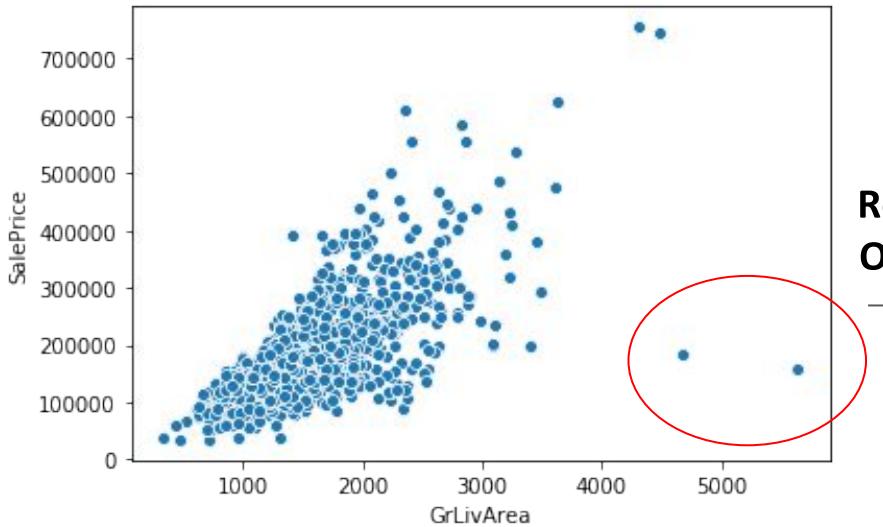
Central Air



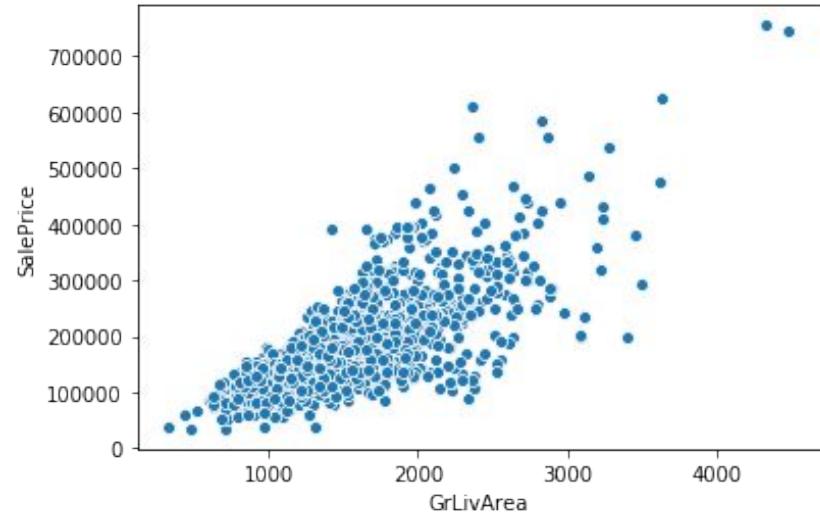
Most of the houses (94%) in the dataset have central air except for houses over 50 years old. However among the same age range, the prices of houses with central air are still higher than houses without central air.

Machine Learning

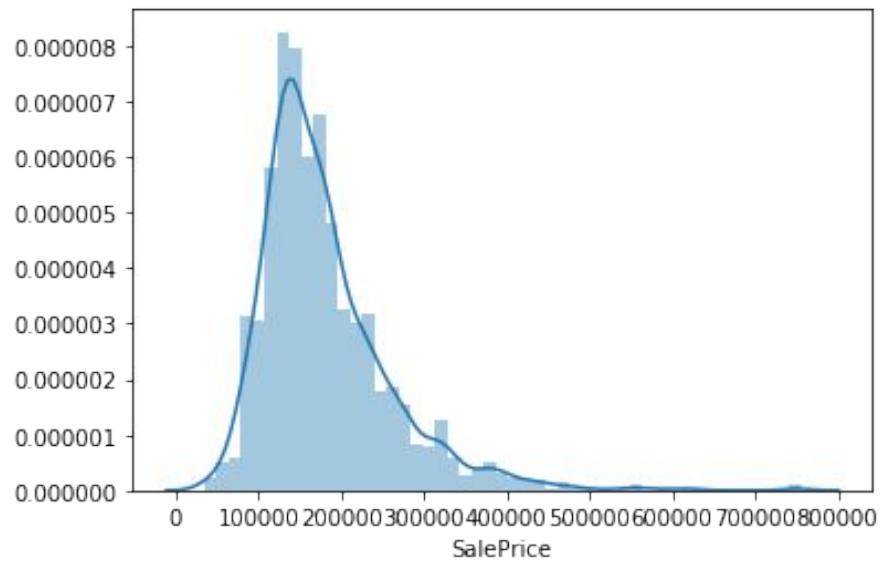
Pre-Processing



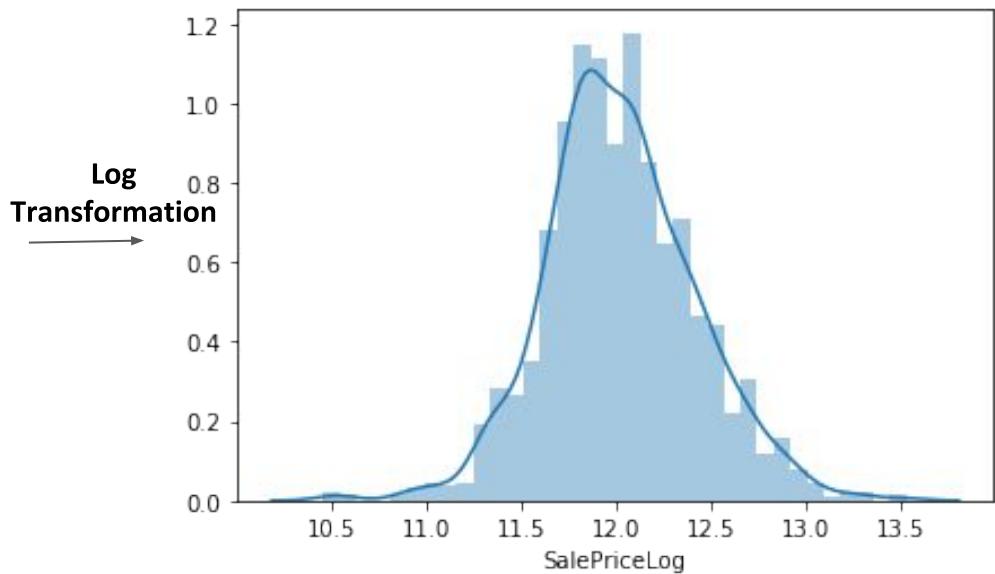
Remove
Outliers



Pre-Processing

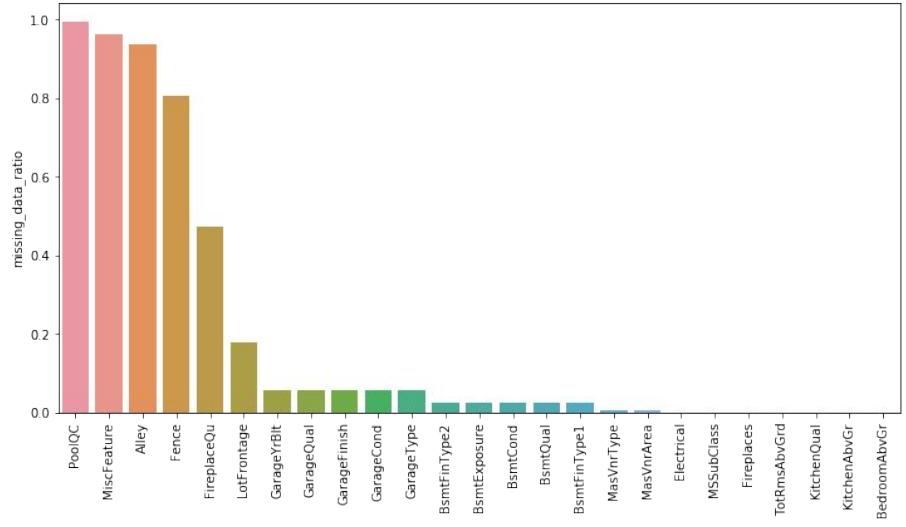


Log
Transformation

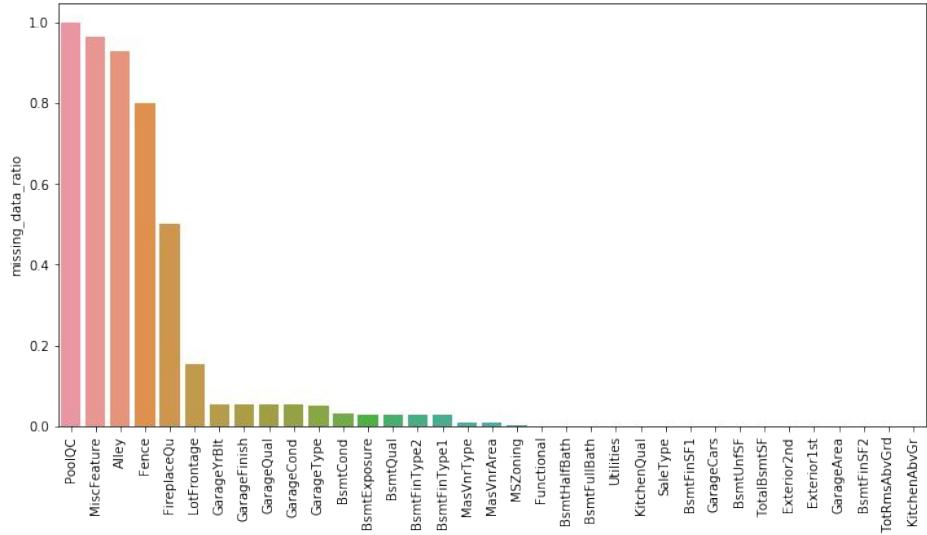


Impute Missing Values

Train



Test



Some of the missing values in these datasets simply mean that a house does not have a specific house feature, such as pool, alley, garage, basement, etc.

Impute Missing Values

No..(NoPool, NoBasement, NoAlley, etc...)	Median	Mode	0
PoolQC MiscFeature Alley Fence FireplaceQu GarageType GarageFinish GarageCond GarageQual BsmtFinType2 BsmtExposure BsmtCond BsmtQual BsmtFinType1 MasVnrType	LotFrontage BsmtFinSF1 BsmtUnfSF TotalBsmtSF	Electrical MSZoning Functional BsmtHalfBath BsmtFullBath Utilities KitchenQual SaleType Exterior1st Exterior2nd	GarageYrBlt MasVnrArea GarageArea BsmtFinSF2 GarageCars

Pre-Processing

Feature Engineering

Columns Added

- **Age** = YrSold - YrBuilt
- **Remodel Age** - YrSold - YearRemodAdd
- **Remodeled** - if the house is remodeled
- **HasGarage** - if the house has a garage
- **GarageAge** = YrSold - GarageYrBlt

Columns Dropped

- YrSold, GarageYrBlt, BsmtFinSF2, BsmtUnfSF, BsmtHalfBath, 2ndFlrSF, MoSold , MiscVal, 3SsnPorch, EnclosedPorch, BedroomAbvGr, LowQualFinSF, MiscFeature, Condition2

Assign numerical numbers to quality variables

- One Hot Encoder does not do a good job
- No -0, Poor - 1, Fair -2,
- TA -3, GD- 4, Ex - 4

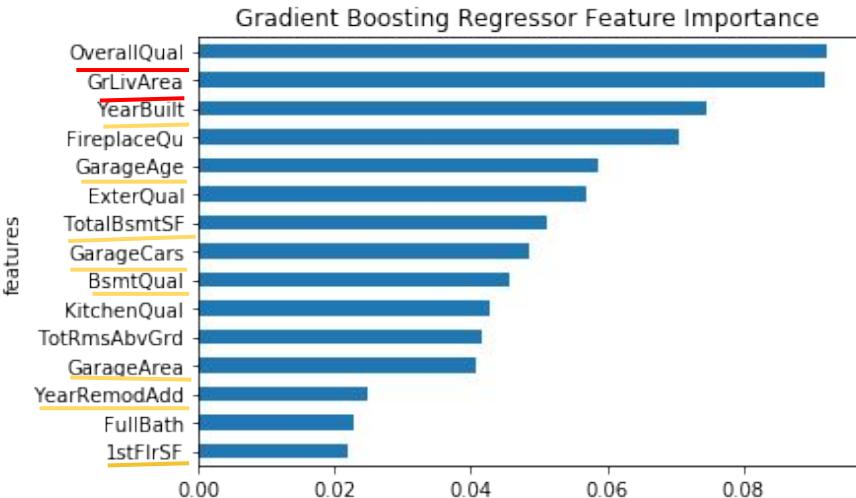
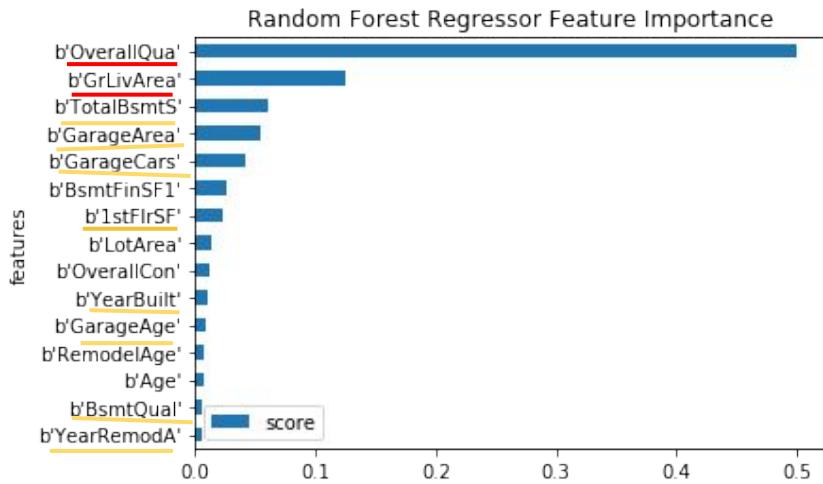
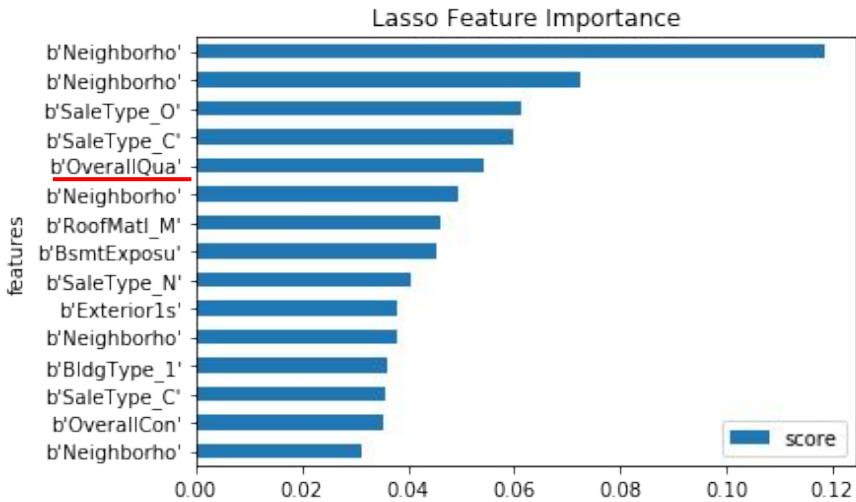
Dummify - get.dummies

- Categorical variables

Machine Learning Models & Scores

	Ridge	Lasso	Random Forest	Gradient Boosting
Train Score	0.9110	0.9205	0.8786	0.9114
Test Score	0.9132	0.9144	0.8682	0.9099
RMSE	0.1141	0.1133	0.1407	0.1163

Gradient Boosting have the best performance on Kaggle.



- **Overall Quality, Above Ground Living Area(GrLivArea)** are the most important indicators when predicting sale prices.
- Important indicators to prediction: Year Built, Garage Area, Garage Cars, Basement Quality, Remodel Year, Total Basement Area, and 1stFlrSF are some of the most important features in both random forest and gradient boosting model.
- Lasso- neighborhood is a good indicator in the lasso model.

Recommendation

- House quality outweighs house condition (quality of house materials is highly important).
- Remodeling a 20 to 39-year-old house makes you have more leverage and possibilities to increase returns.
- It is not necessary to update your garage to a 4-car garage. A 4-car garage does not guarantee a higher house price. Updating a garage to at least the TA(Typical/Average) level will improve the price. Price Per Sqft increases by 30% from fair quality to typical/average quality.
- Central Air is a necessity: Central Air is commonly equipped in the house nowadays. The price per sqft of a house featuring central air is 46% higher than the price per sqft of a house that has no central air.
- Spend money improving the quality in this order: kitchen > driveway > fireplace> fence.