



**THE PRICE IS
(ALWAYS) RIGHT**

WHAT PRICE IS THE FAIR PRICE?

MARK YUNG, ANG JUN SIONG,
JONATHAN SOON, MERIKY LO ALEXANDER



MMJJ REALTY: IOWA CHAPTER



MERIKY

REALTY CHIEF

*Never Been
to Ames nor Iowa*



MERIKY



JONATHAN

INTERN

Lowly Paid



INTERN

Lowly Paid



MARK

TEAM PSYCHIC

Needs More Mana

JUN SIONG



SWMMER

*Perpetually Immersed
(in Work or Water)*



SWMMER

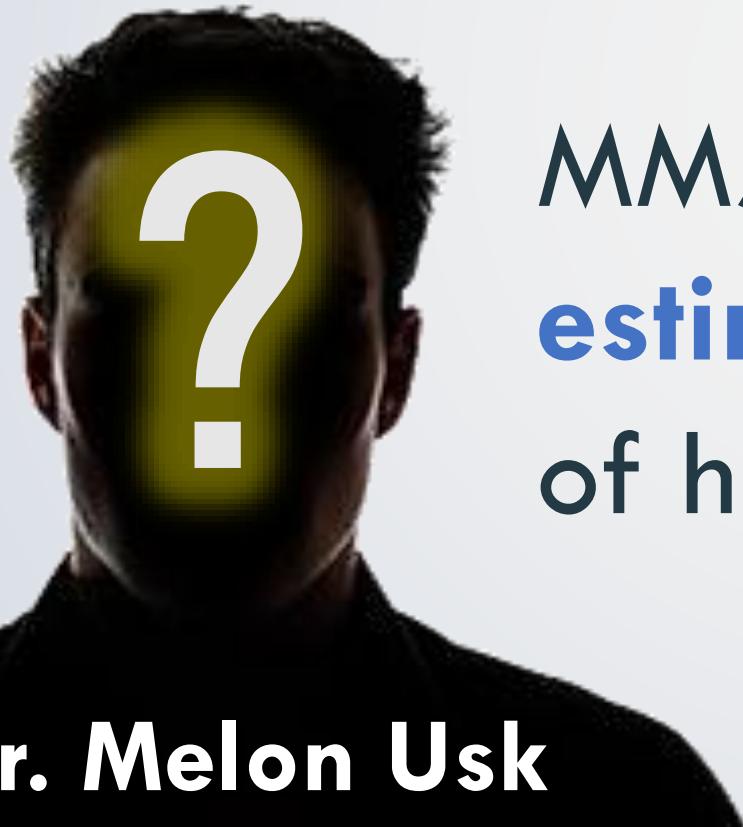
*Perpetually Immersed
(in Work or Water)*



BACKGROUND

WHAT PRICE IS THE FAIRPRICE?
MERIKY, JONATHAN, MARK, JUN SIONG

You - High Profile Property Investors and Home Buyers



MMJJ have been hired to
estimate the price for each
of his **800+ properties**.

Mr. Melon Usk





PROBLEM STATEMENT

What is the Fair Price?

What are the Most
Important Predictors of
Home Prices? How?

Can we Predict
Home Prices
Accurately? How?



WHAT PRICE IS THE FAIRPRICE?
MERIKY, JONATHAN, MARK, JUN SIONG



PROCEDURE & METHODOLOGY

WHAT PRICE IS THE FAIRPRICE?
MERIKY, JONATHAN, MARK, JUN SIONG

Data Source

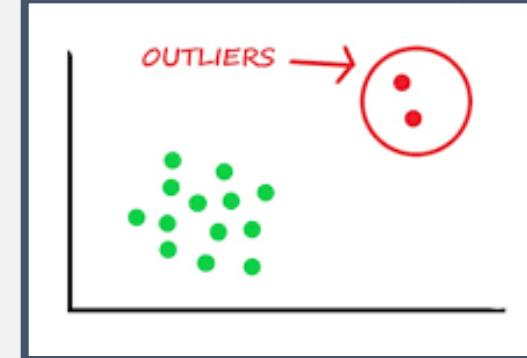


Ames Housing Data Set

- ~2000 Data Points for Home Sales Between 2006 to 2010
- 79 Numerical & Categorical Variables
- Go-to Introductory Regression Dataset



Basic Data Preparation



Remove Outliers

```
0101011011010010110100101101001010110100100  
010101001010110101101001011010010110100101  
0101010010101011010100101101010010110100101  
0101010010101011010100101101010010110100101  
0101001011010100101011010100101101010010100  
0100101010101001010101010010101010100101001  
0101011010101001010101010101010101001010101  
010010001001010101010010010010100101010101  
001001001010101101101011010011010101010101  
010001011010101010101010110101010101010101  
010010110101010101010101010101010101010101
```

Data “Digitization”



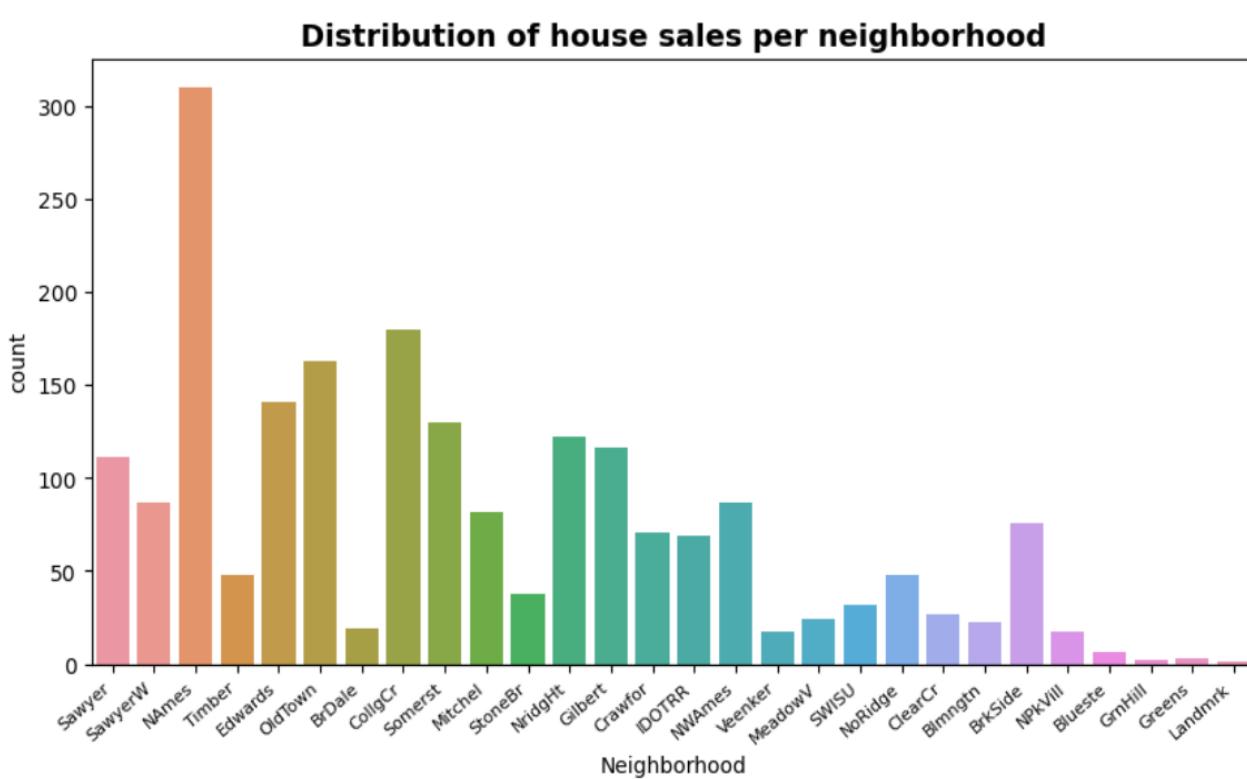
Remove Useless Data



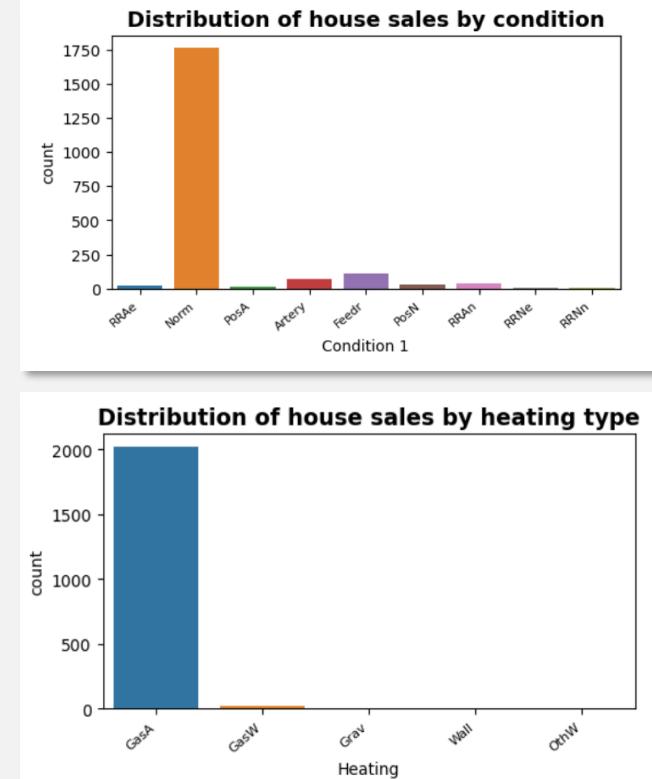
Trim Sparse Variables

Exploratory Data Analysis

Understanding the Dataset



ID Low Variance Predictors





PROCEDURE & METHODOLOGY

WHAT PRICE IS THE FAIRPRICE?
MERIKY, JONATHAN, MARK, JUN SIONG

Feature Engineering

Basement Size

	bsmtfin sf 1	bsmtfin sf 2	bsmt unf sf	total bsmt sf
0	533.0	0.0	192.0	725.0
1	637.0	0.0	276.0	913.0
2	731.0	0.0	326.0	1057.0
3	0.0	0.0	384.0	384.0
4	0.0	0.0	676.0	676.0

Garage Features

	garage score	garage finish	garage cond	garage qual	garage cars
0	144	3	4	4	3
1	144	3	4	4	3
2	64	2	4	4	2
3	192	4	4	4	3
4	96	2	4	4	3

Year Remod Feature

	age sold	year built	year remod/add	yr sold
0	5	1976	2005	2010
1	12	1996	1997	2009
2	3	1953	2007	2010
3	3	2006	2007	2010
4	17	1900	1993	2010

Floor SF Features

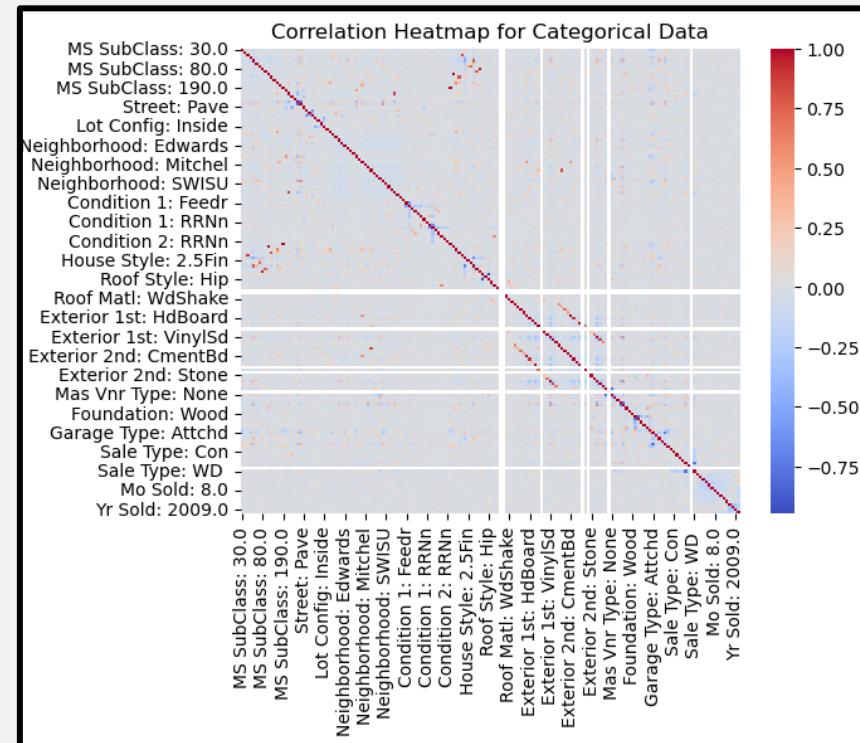
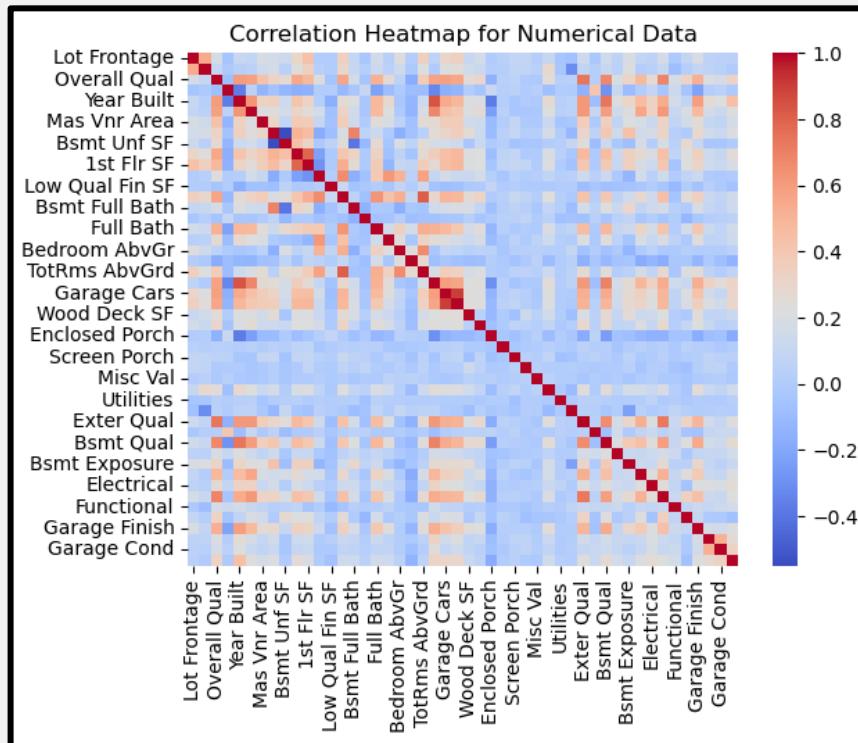
	1st flr sf	2nd flr sf	low qual fin sf	gr liv area
2043	1728	0	0	1728
2044	861	0	0	861
2045	1172	741	0	1913
2046	1200	0	0	1200
2047	1028	776	0	1804



PROCEDURE & METHODOLOGY

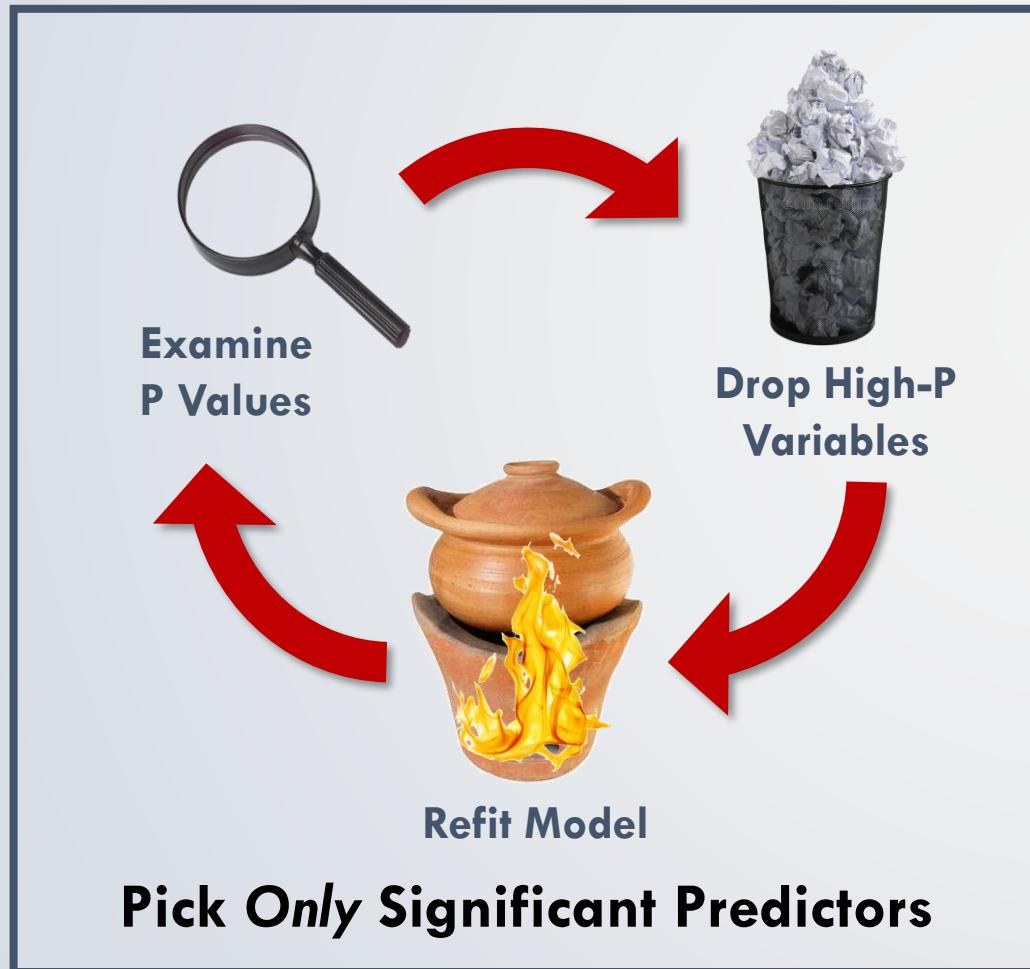
WHAT PRICE IS THE FAIRPRICE?
MERIKY, JONATHAN, MARK, JUN SIONG

Variables Selection Stage 1: (Anti-) Multi-Collinearity

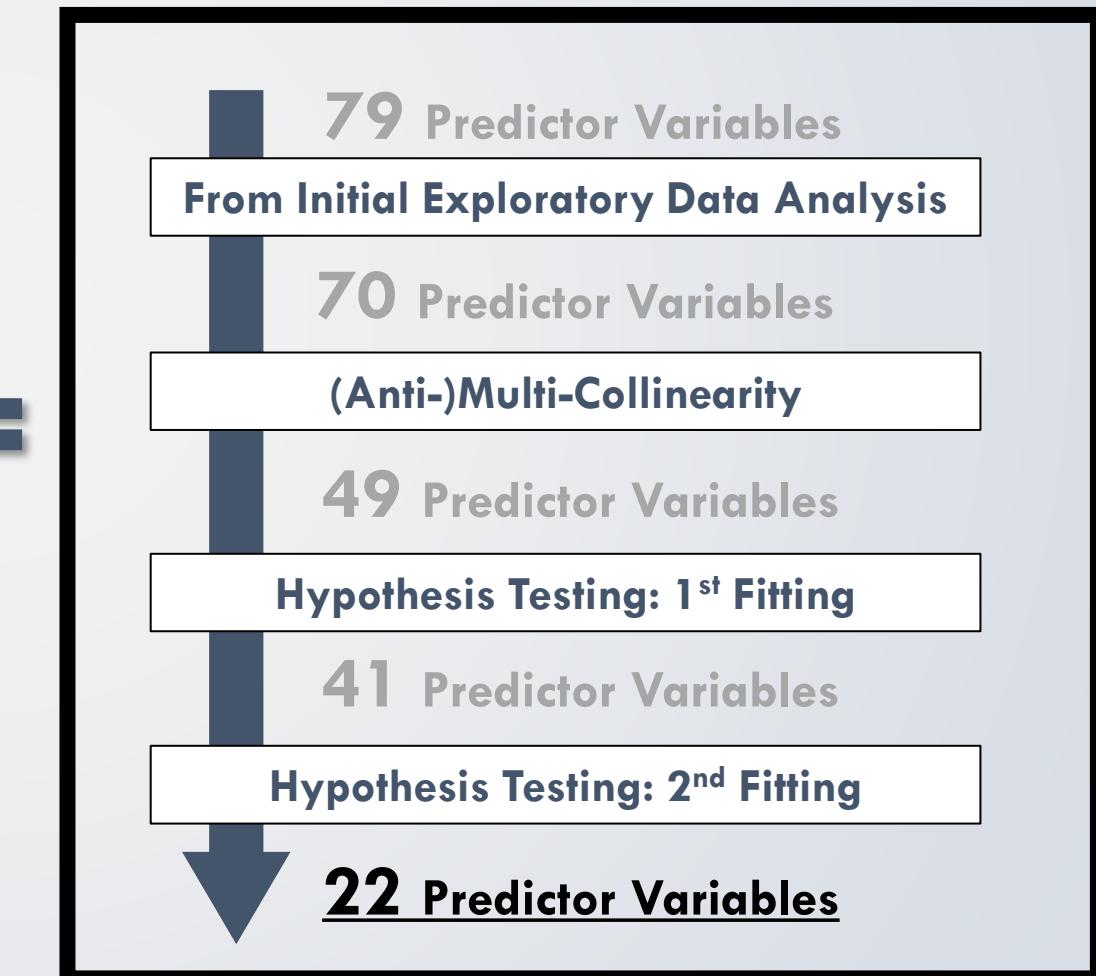


Removed
8 Predictor Variables

Variables Selection Stage 2: Iterative Hypothesis Testing



IN SUMMARY





PRIMARY FINDINGS

2000 Observations
22 Predictors



WHAT PRICE IS THE FAIRPRICE?
MERIKY, JONATHAN, MARK, JUN SIONG





PRIMARY FINDINGS

WHAT PRICE IS THE FAIRPRICE?
MERIKY, JONATHAN, MARK, JUN SIONG

MODELS SUMMARY

M1: Vanilla Model

(Ridge Regression)

Validation R²: 0.91

Kaggle RMSE: \$24,XXX

10% Observations Dropped
for Having Empty Cells

Id	Col 1	...	Col XX	Col XX	Col XX	Col XX
1	Data	...	Data	Data	Data	Data
2	MISSING		Data	Data	Data	Data
3	Data		MISSING	MISSING	Data	Data
4	Data		Data	Data	Data	Data
...
...
YY	Data	...	MISSING	Data	Data	Data
YY	Data	...	Data	Data	Data	Data
YY	Data		Data	Data	Data	Data

Red arrows point to row 2 and row 3, indicating they are dropped due to missing values.





PRIMARY FINDINGS

WHAT PRICE IS THE FAIRPRICE?
MERIKY, JONATHAN, MARK, JUN SIONG

MODELS SUMMARY

M1: Vanilla Model

(Ridge Regression)

Validation R²: 0.91

Kaggle RMSE: \$24,XXX

M2: Drop Columns

(Ridge Regression)

Validation R²: 0.87

Kaggle RMSE: \$29,XXX

4

Variables Dropped for Having Empty Cells

Id	Col 1	...	Col XX	Col XX	Col XX	Col XX
1	Data	...	Data	Data	Data	Data
2	Data		Data	Data	Data	Data
3	Data		MISSING	MISSING	Data	Data
4	Data		Data	Data	Data	Data
...
...
YY	Data	...	MISSING	Data	Data	Data
YY	Data	...	Data	Data	Data	Data
YY	Data		Data	Data	Data	Data





PRIMARY FINDINGS

WHAT PRICE IS THE FAIRPRICE?
MERIKY, JONATHAN, MARK, JUN SIONG

MODELS SUMMARY

M1: Vanilla Model

(Ridge Regression)

Validation R²: 0.91

Kaggle RMSE: \$24,XXX

M2: Drop Columns

(Ridge Regression)

Validation R²: 0.87

Kaggle RMSE: \$29,XXX

M3: Impute Median

(Ridge Regression)

Validation R²: 0.93

Kaggle RMSE: \$21,XXX

Median

**Broad Imputation
Method for Missing Data**

Id	Col 1	...	Col XX	Col XX	Col XX	Col XX
1	Data	...	Data	Data	Data	Data
2	Data		Data	Data	Data	Data
3	MEDIAN		MEDIAN	MEDIAN	Data	Data
4	Data		Data	Data	Data	Data
...
...
YY	Data	...	MEDIAN	Data	Data	Data
YY	Data	...	Data	Data	Data	Data
YY	Data		Data	Data	Data	Data



PRIMARY FINDINGS

WHAT PRICE IS THE FAIRPRICE?
MERIKY, JONATHAN, MARK, JUN SIONG

MODELS SUMMARY

M1: Vanilla Model

(Ridge Regression)

Validation R²: 0.91

Kaggle RMSE: \$24,XXX

M2: Drop Columns

(Ridge Regression)

Validation R²: 0.87

Kaggle RMSE: \$29,XXX

M3: Impute Median

(Ridge Regression)

Validation R²: 0.93

Kaggle RMSE: \$21,XXX

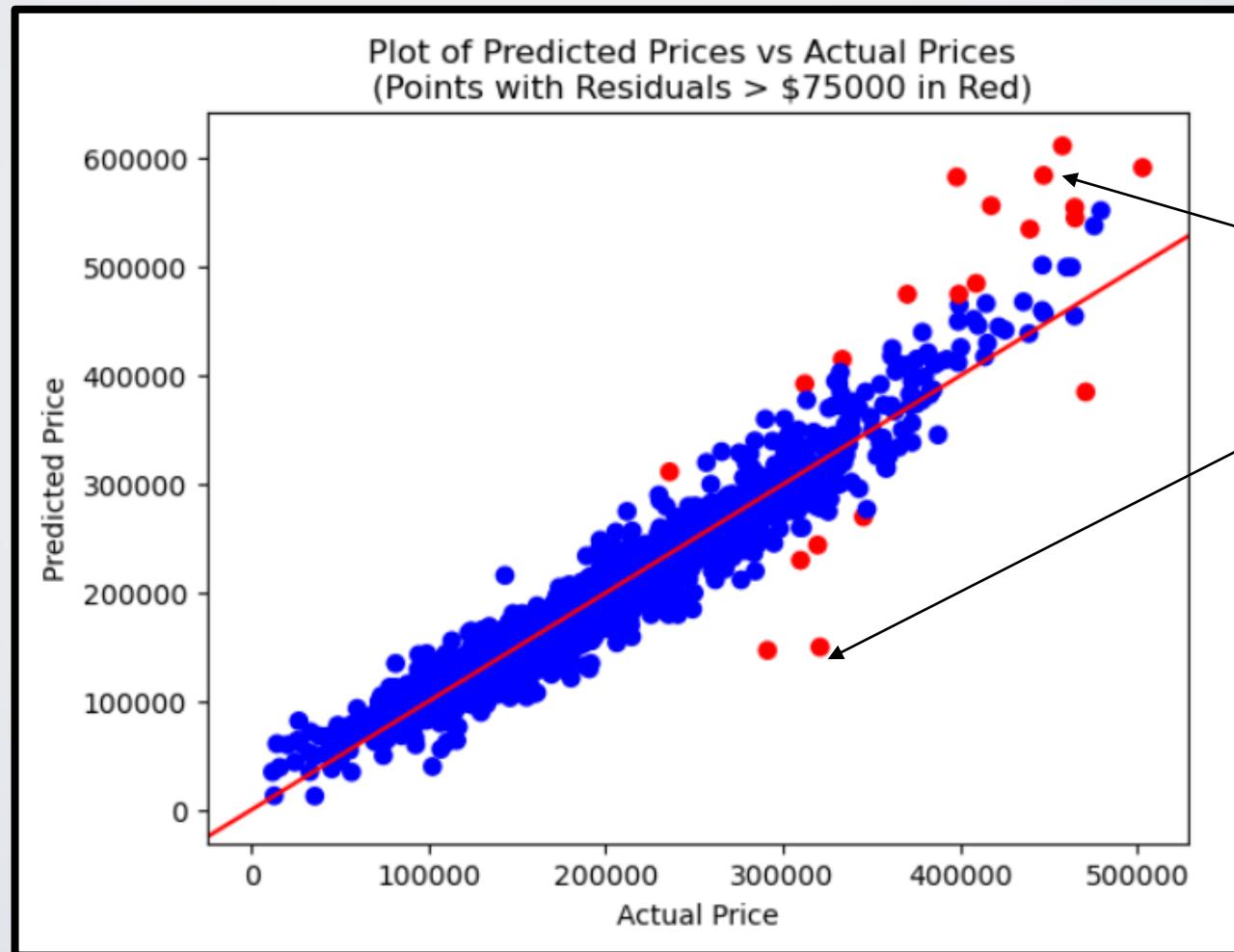
M4: Drop Outliers

(Ridge/ Lasso Regression)

Validation R²: 0.94

Kaggle RMSE: \$20,030

Predicted Prices vs Actual Prices from M3



**Model Refitted
After Removing
Outlying
Observations
(In Red)**





CONCLUSION

WHAT PRICE IS THE FAIRPRICE?
MERIKY, JONATHAN, MARK, JUN SIONG

Problem
Statement
Recap

What is the Fair Price?

What are the Most Important Predictors of Home Prices?
Can we Predict Home Prices Accurately?



Final Model Evaluation

0.94 R^2 + \$20,030 Test RMSE on Kaggle



Most Important Variables in Predicting Home Prices

Neighborhood, Gross Living Area, Condition, Sales Type



Data Preparation Makes the Difference

Same Same (Dataset), But Vastly Different (Results)





RECOMMENDATIONS

WHAT PRICE IS THE FAIRPRICE?
MERIKY, JONATHAN, MARK, JUN SIONG



Employ More Complex Models (GBR, Neural Nets)

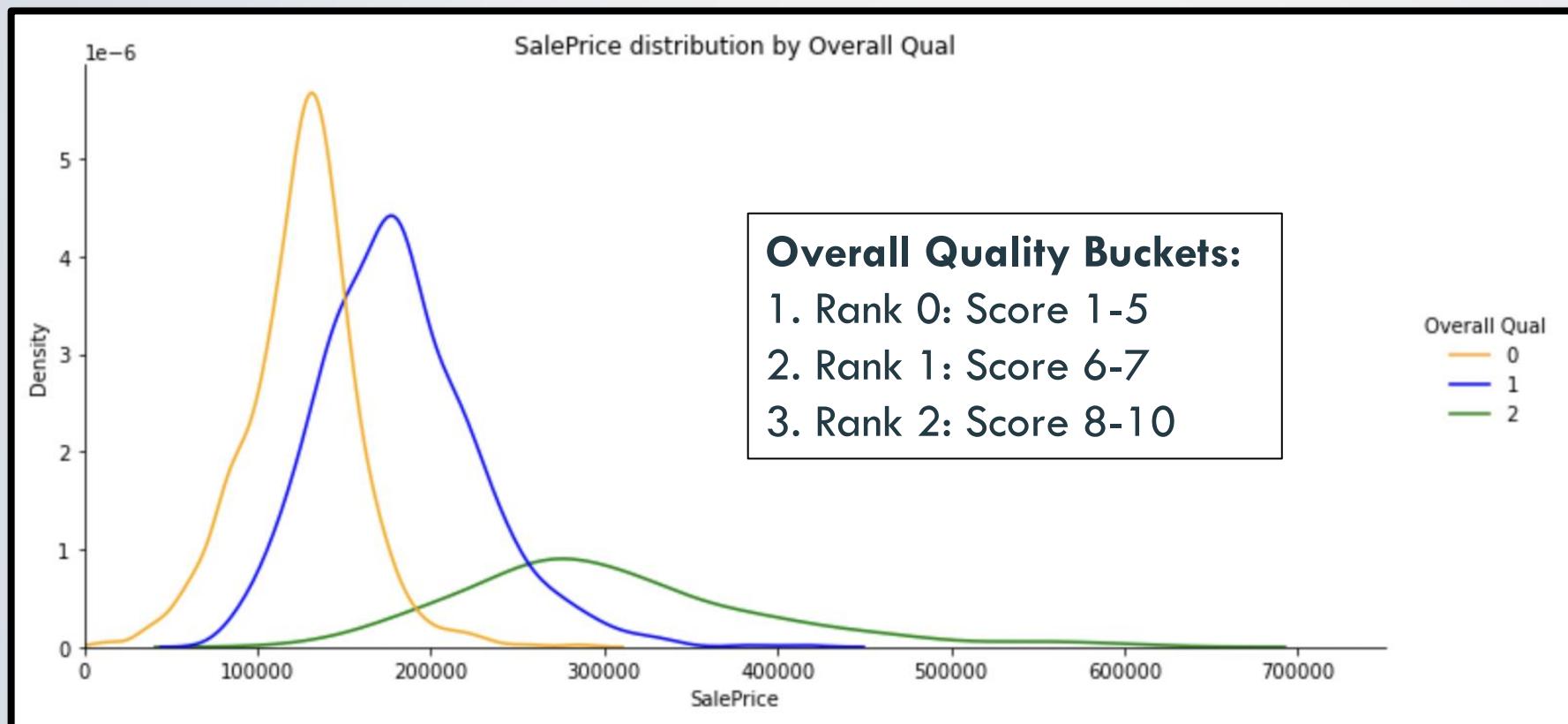
- Unoptimized Gradient Boosting Regression: \$18,XXX Test RMSE on Kaggle
- **BUT**, Involves Tradeoff in Interpretability

Bonus: Implement Market Segmentation...



Market Segments by certain characteristics are affected by **different factors**. Oftentimes, these macro-effects translate into our models as fluctuations in feature dependency.

Translates into different factors and features **affecting model to varying degrees**.





MARKET SEGMENTATION

WHAT PRICE IS THE FAIRPRICE?
MERIKY, JONATHAN, MARK, JUN SIONG

Tailored Models:

Different Market Segments, Different Behaviors

Segments have **different correlation behaviors**



BEST Performance, with Enough Data

	OQ0	OQ0_corr	OQ1	OQ1_corr	OQ2	OQ2_corr
1	1st Flr SF	0.516779	Gr Liv Area	0.658482	Kitchen Qual	0.517902
2	Gr Liv Area	0.487283	Garage Area	0.541589	Gr Liv Area	0.513535
3	Total Bsmt SF	0.465055	Garage Cars	0.526667	Garage Cars	0.512368
4	Central Air_Y	0.439907	Full Bath	0.502532	TotRms AbvGrd	0.501728
5	BsmtFin SF 1	0.437947	1st Flr SF	0.499822	Exter Qual	0.501463
6	Paved Drive	0.410663	Total Bsmt SF	0.492332	Mas Vnr Area	0.496856
7	Year Built	0.408489	TotRms AbvGrd	0.440296	Bsmt Qual	0.495505
8	BsmtFin Type 1	0.402716	Kitchen Qual	0.431487	Garage Area	0.490732
9	Fireplaces	0.390328	Exter Qual	0.420268	Fireplace Qu	0.432934
10	Garage Area	0.385854	Year Remod/Add	0.407365	1st Flr SF	0.410804

RMSE values for each of the subsets:

1. Test RMSE(Overall Qual = 0): 14,550
2. Test RMSE(Overall Qual = 1): 18,713
3. Test RMSE(Overall Qual = 2): 42,278



OUR PRICE IS ALWAYS THE FAIR PRICE

NO ABSD, FOREIGNERS ELIGIBLE,
DM US ON SLACK NOW!!!



RECOMMENDATIONS

WHAT PRICE IS THE FAIRPRICE?
MERIKY, JONATHAN, MARK, JUN SIONG

QUESTIONS?





RECOMMENDATIONS

WHAT PRICE IS THE FAIRPRICE?
MERIKY, JONATHAN, MARK, JUN SIONG

THANKS!

