HOUSE PRICE PREDICTION

Team 5 Site Bai, Smitha Kannanaikkal, Ellie Park, Joy Chen



#### **TABLE OF CONTENTS**



Context, Question and Hypothesis

02

**Dataset** 

Data Gathering and Research Method 03 Models

Linear Regression, Decision Tree, Random Forest & Time Series Models

05
Implications

Suggestions based on models

06 Appendix

Additional Information about Project



#### **CONTEXT**

- Google reported the search "When is the housing market going to crash?" had spiked 2,450% in the past month.
- The real estate market is booming right now and is breaking historic records in April, as home prices rose 21% YoY and the median home-sale price soared nationwide
- With the housing market evolving rapidly, it is important to understand the fluctuation in market trends

## PROJECT QUESTION



#### **KEY QUESTIONS**





- We want to learn different attributes that impact the sale price and be able to boost the price when budget is limited.
- Second: For Buyers
  - By forecasting the count of house sales, we enable customers looking to purchase houses to have an idea when there will be more options listed, and also estimate the price based on supply and demand.
- Our dependent variable is sales price & sales volume for time series models





# DATA GATHERING

#### **DATA GATHERING**

Our Dataset was gathered from Kaggle & was part of the Kaggle contest

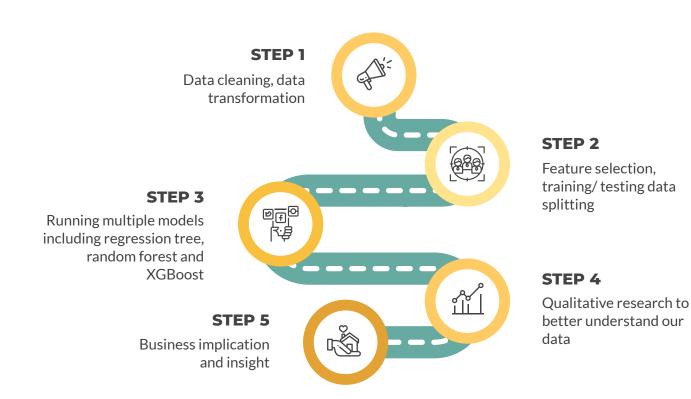
Dataset has 1,457 observations with 80 variables

It includes variables related to features of the houses including home condition, kitchen, porch, fireplace, basement, size of each floor, etc.

It also includes sales related information including date sold, sale type, sale condition, etc.

# **METHODOLOGY**

#### **METHODOLOGY**



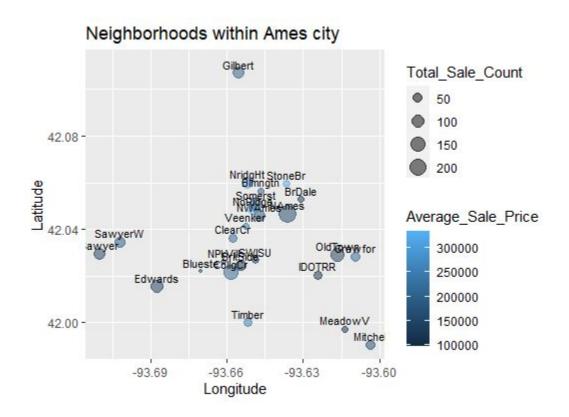


## DATA Visualization

#### **SALES PRICE**



#### **NEIGHBORHOODS**

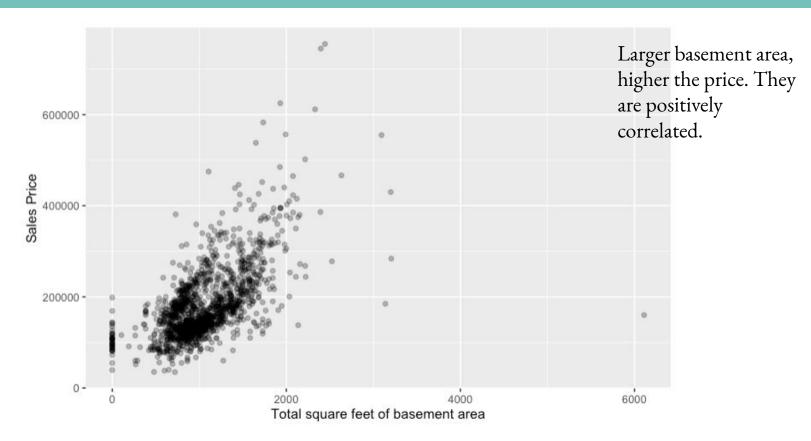


large sale happening where the average sale price is lower

#### **SALES PRICE & MATERIAL**







# **ANALYSIS**For Sellers



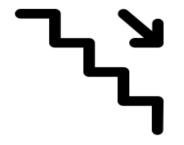
#### REGRESSION

	All in Linear Regression	Backward Feature Selection
Number of Attributes	125	82
R-squared	90.81%	90.68%
Test RMSE	28308.45	28344.37

#### Features impact on the Sales Price

	Neighborhood	Has Basement	Heating Quality
	Stone Brook	Yes	Condition Poor
House Sale Price Prediction	<b>1</b> 20%	<b>12</b> %	<b>↓</b> 29.5%







#### Features impact on the Sales Price

	Neighborhood Stone Brook	Has Basement Yes	Heating Quality Condition Poor
House Sale Price Prediction	<b>1</b> 20%	<b>12</b> %	<b>♦</b> 29.5%
Potential Cause	high rating schools in this area	Means more space, also according to the National Association of Realtors, a finished basement returns about 69 percent of the investment.	lowa is cold in winter

#### **DECISION TREE/ RANDOM FOREST/ XGBOOST**

				smallest 
		Regression tree	Random Forest	XGBoost
RMSE	Training data	38413.77	19754.89	10023.13
	Testing data	45617.69	35995.96	35812.75 (20% of average sales price)

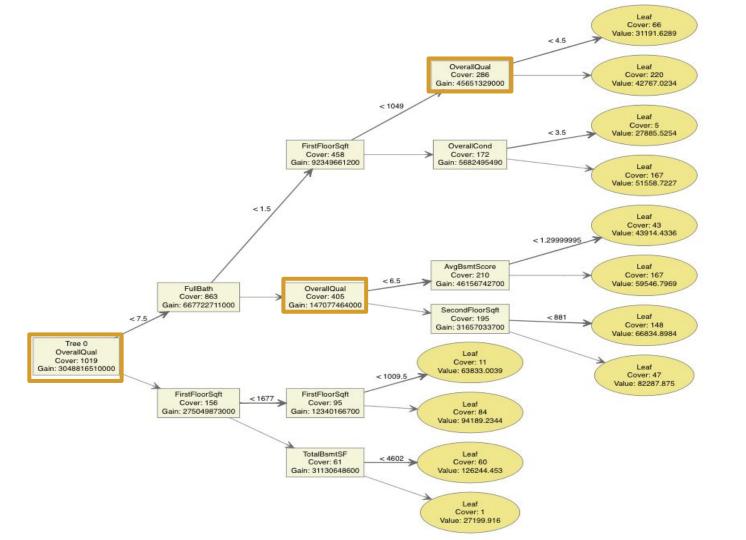
Our regression model has an RMSE of 28K, we think here the regression is better than the trees is probably because our dataset is rather small.

#### **Significant Coefficients From XGBOOST**

OverallQual FirstFloorSaft GarageCars FullBath SecondFloorSaft TotalBsmtSF PctFinBsmt house age OverallCond AvgBsmtScore MSSubClass60 HalfBath BedroomAbvGr ExterQualTA KitchenQualTA Fireplaces MSZoningRM NeighborhoodCrawfor BaseFullBath LotConfigCulDSac NeighborhoodNAmes ExterQualGd MSSubClass90 SaleConditionPartial MSZoningRL ElectricalTypeSBrkr NeighborhoodClearCr LotConfigFR3

Only 10 out of 82 attributes are significant in affecting the sales price, and "Overall Quality" is the most important.



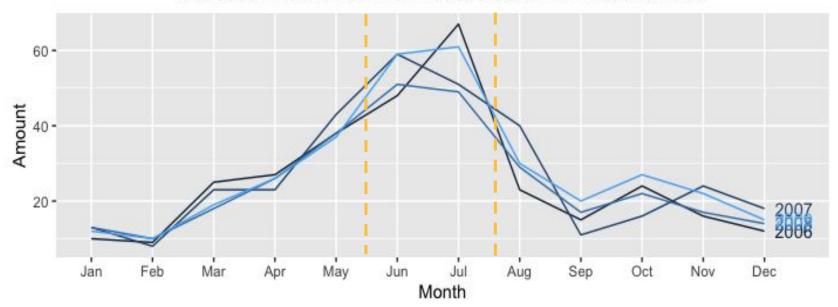


# **ANALYSIS**For Buyers



#### **TIME SERIES**

#### Seasonal Plot: Amount of Houses Sold from 2006 to 2009



Spike between May and July

Potential reason: Moving closer to school before school starts

#### **TIME SERIES**

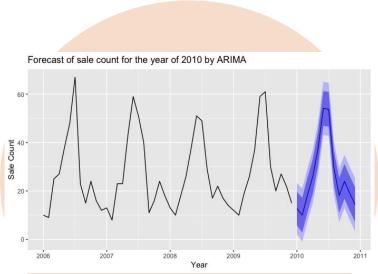
#### Seasonal Plot: Amount of Houses Sold from 2006 to 2009

The transactions probably took place back in Spring because normally it would take 30-60 days to complete the contract if the buyer is taking the loan. As a buyer, you don't want to buy the house during Spring because you would face the most competition despite of the extensive available newly listed houses.

Based on this, we suggest that buying towards the end of summer offers buyers with no need for school district houses both selections and bargains

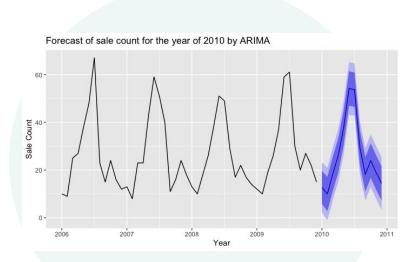


#### **FORECAST**



ETS(M,N,A)
PREDICTION

RMSE: 20.87



#### **ARIMA PREDICTION**

RMSE: 20.42

# IMPLICATIONS



#### **ANALYSIS & INSIGHTS**

#### For sellers:

- Location matters.
- Polish the basement if there is one.
- Install heating system will enable sellers to ask higher price.

#### For buyers:

From late May to mid July, there is a greater demand for houses with sales peaking

 To avoid competition, buyers should look for houses after summer months so they have greater bargain power

#### **LIMITATIONS**

#### **Dataset**

- Our dataset has limited number of observations.
- Is restricted to Iowa, thus, findings may not be applied to other states

### **APPENDIX**

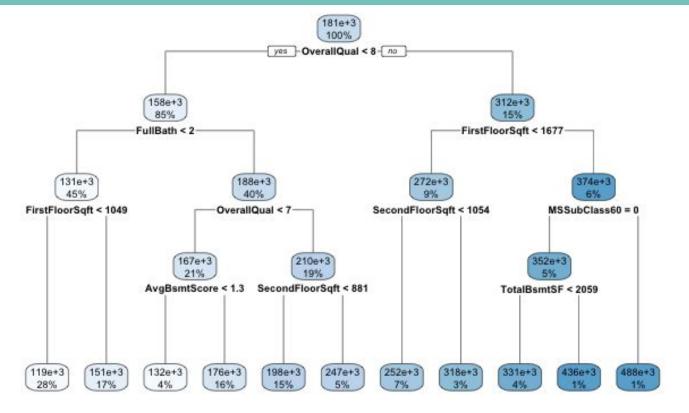


Data cleaning: STEP 1 Manually go over 80 variables such as checking correlation between variables, take out or combined data with very few instances, take log for skewed numeric data. Split data into training and testing (70:30). STEP 2 Run multiple models including regression tree, random forest and XGBoost STEP 3 Time series analysis. STEP 4 Do further qualitative research to better understand our results. STEP 5 STEP 6 Business implication and insight

#### **Decision Tree Result**

#### **Decision Tree**

```
library(rpart)
rpart(
formula = SalePrice ~ .,
data = train,
method = "anova",
control = list(cp = 0,
minsplit = 13
maxdepth = 12,
xval = 10))
```





#### **XGBoosting Result**

#### **XGBOOST**

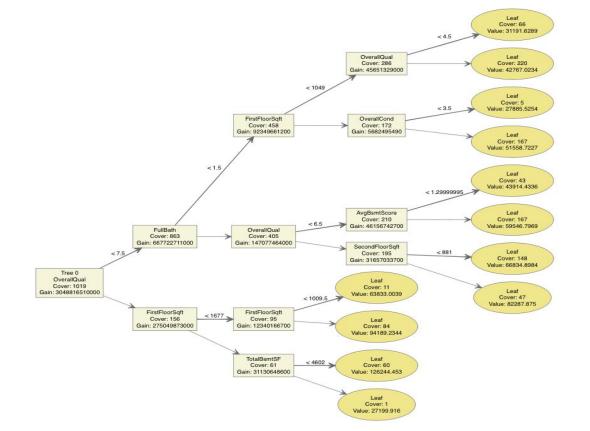
```
xgboost(data = xgb_train,
label = SalePrice
max_depth = 4,
eta = 0.34,
nrounds = 50,
```

nthread = 2,

"reg:squarederror") )

library(xgboost)

objective =



## **THANK YOU**

