# Housing Prices

Elsie He, Owen Ou, Woonsup Kim, Emmelyn Luveta

# TABLE OF CONTENTS

# 01 BACKGROUND

- **Overall material and finish quality**
- **Above ground living area**
- **Basement area**
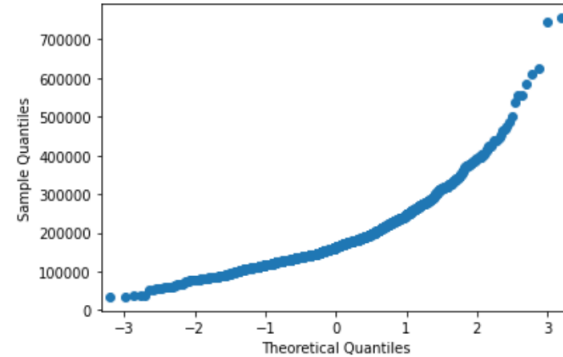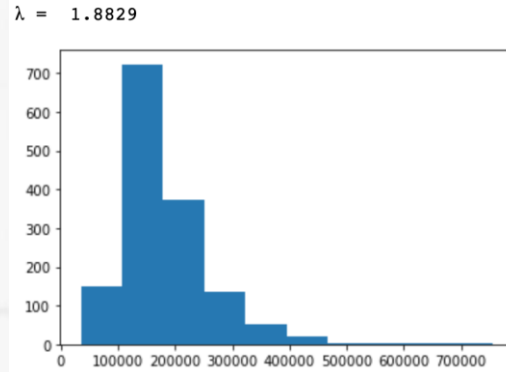- **Size of garage in car capacity**

# 02 DATA

- **Dataset of house prices in Ames, Iowa from 2006-2010.**
- **Dataset was collected by Bart de Cock in 2011.**
- **Raw data includes 1460 data points, 80 features, and 1 response variable, SalePrice.**
- **There are 38 numeric variables.**
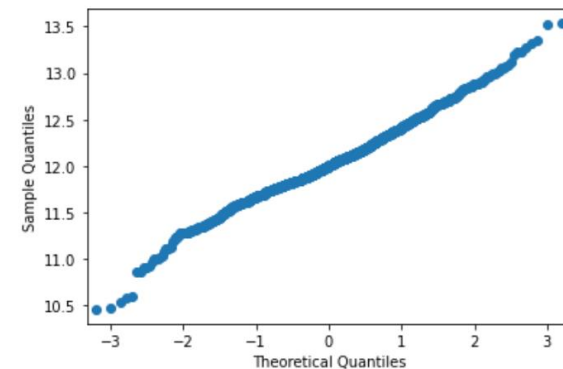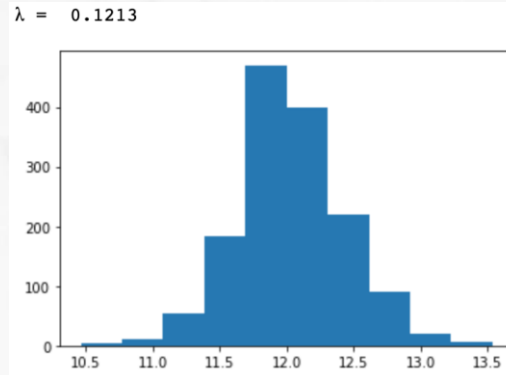- **NA columns.**

# 03 DATA PREPROCESSING
## - Response Variable

**Response**

**Skew: 1.8829**

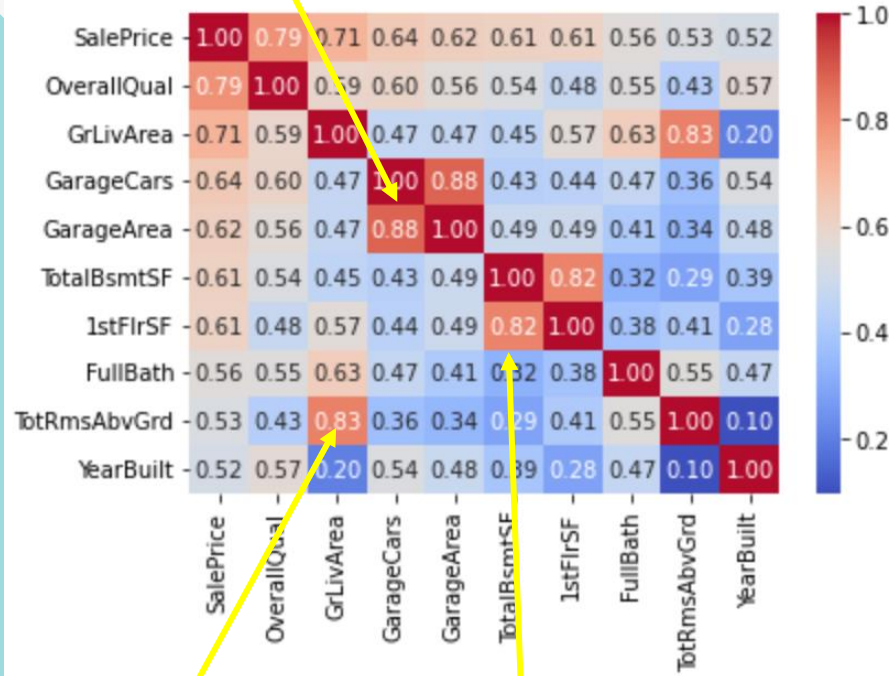**Log(response)**

**Skew: 0.1213**

## - Explanatory Variables

# 03 DATA PREPROCESSING

```
20      1-STORY 1946 & NEWER ALL STYLES
30      1-STORY 1945 & OLDER
40      1-STORY W/FINISHED ATTIC ALL AGES
45      1-1/2 STORY - UNFINISHED ALL AGES
50      1-1/2 STORY FINISHED ALL AGES
60      2-STORY 1946 & NEWER
70      2-STORY 1945 & OLDER
75      2-1/2 STORY ALL AGES
80      SPLIT OR MULTI-LEVEL
85      SPLIT FOYER
90      DUPLEX - ALL STYLES AND AGES
120     1-STORY PUD (Planned Unit Development) - 1946 & NEWER
150     1-1/2 STORY PUD - ALL AGES
160     2-STORY PUD - 1946 & NEWER
180     PUD - MULTILEVEL - INCL SPLIT LEV/FOYER
190     2 FAMILY CONVERSION - ALL STYLES AND AGES
```

```
0       2-story_1946+
1       1-story_1946+
2       2-story_1946+
3       2-story_1945-
4       2-story_1946+
Name: MSSubClass, dtype: object
```

**Dataset observation**

- **Map numerical features that are supposed to be categorical features**
  - **MSSubClass**

**MSSubClass: Identifies the type of dwelling involved in the sale.**

# 03 DATA PREPROCESSING

**BsmtFullBath: Basement full bathrooms**
- **BsmtHalfBath: Basement half bathrooms**
- **FullBath: Full bathrooms above grade**
- **HalfBath: Half baths above grade**

**Dataset observation**

- **Combine 4 bathroom columns into 1**
  - **BsmtFullBath, BsmtHalfBath, HalfBath, FullBath**
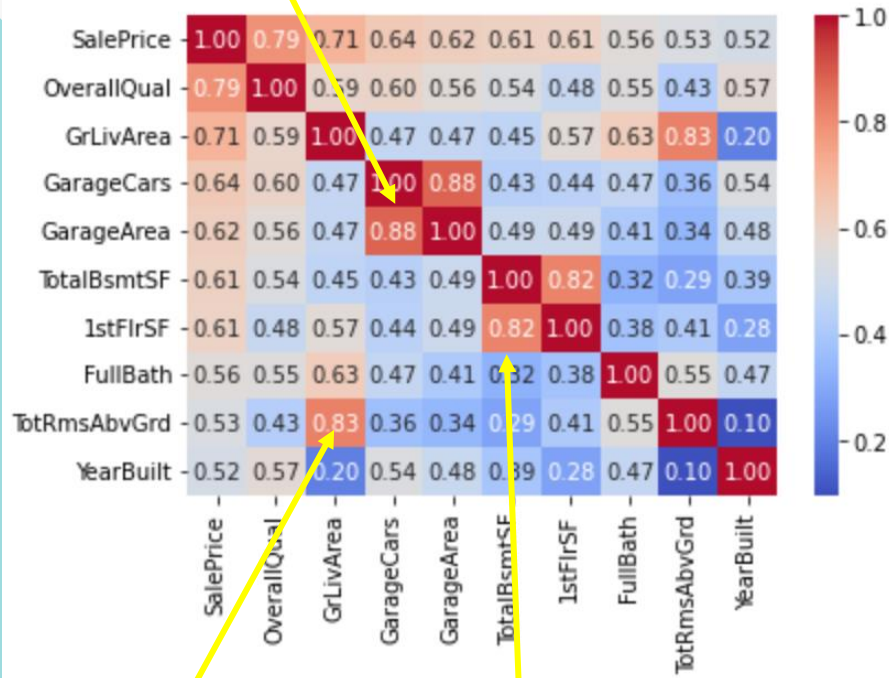
# 03 DATA PREPROCESSING

```
0       2003
1       1976
2       2001
3       1915
4       2000
Name: YearBuilt, dtype: int64
```

- **Remod : indicates whether it has been remodeled**
- **HouseAge: YrSold - YearRemodAdd**
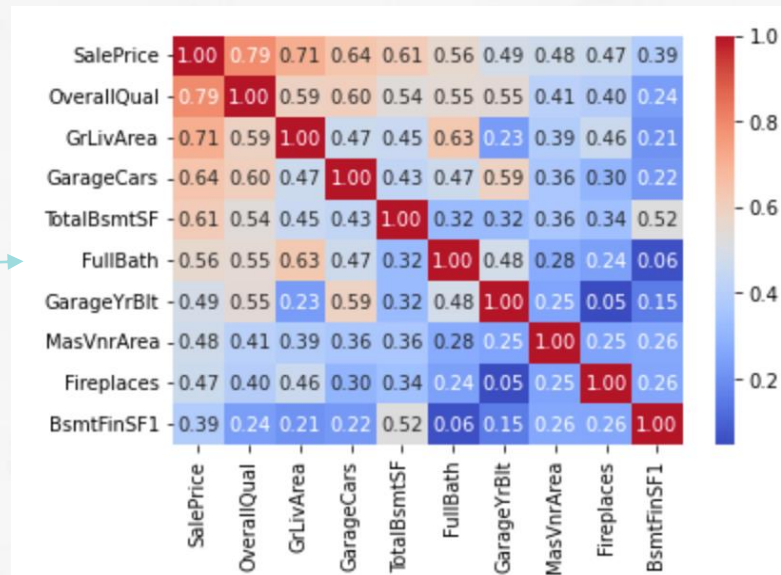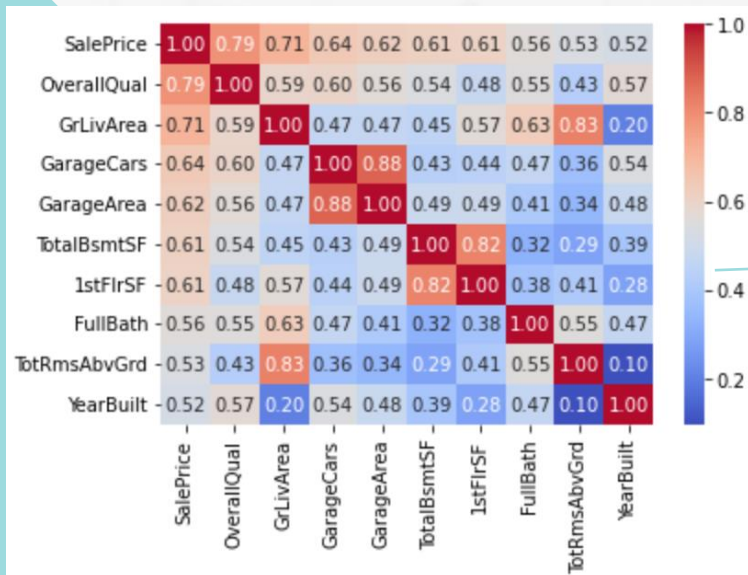
**Dataset observation**

- **YearBuilt: Original construction date**
- **YearRemodAdd: Remodel date (same as construction date if no remodeling or additions)**
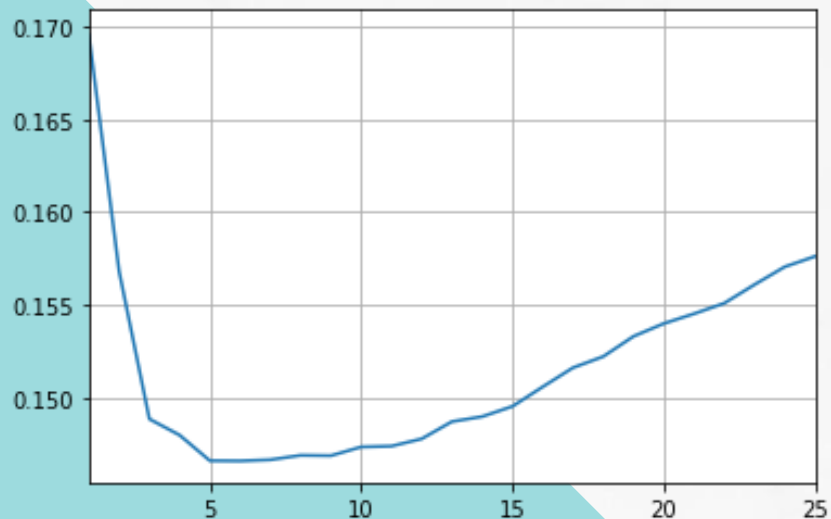- **YrSold: Year Sold**

# 03 DATA PREPROCESSING

# 03 DATA PREPROCESSING

# 04 FEATURE SELECTION

- **Regularization Through Lasso Regression**
- **Removed unrelated features**
- **Reduced the effect of less dependent features**
- **Reduced the number of features from 334 to 117**

# 05 MODELING



**1. Linear Regression**

- ○ Train RMSE = 0.1023
- ○ 10-CV RMSE = 0.1326

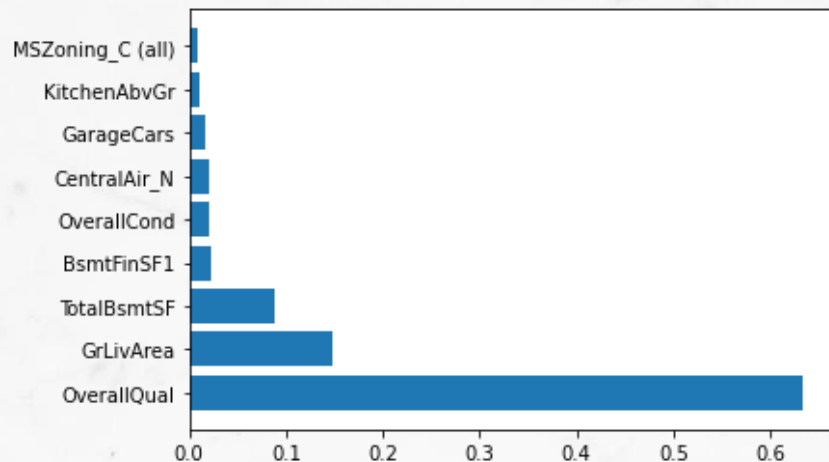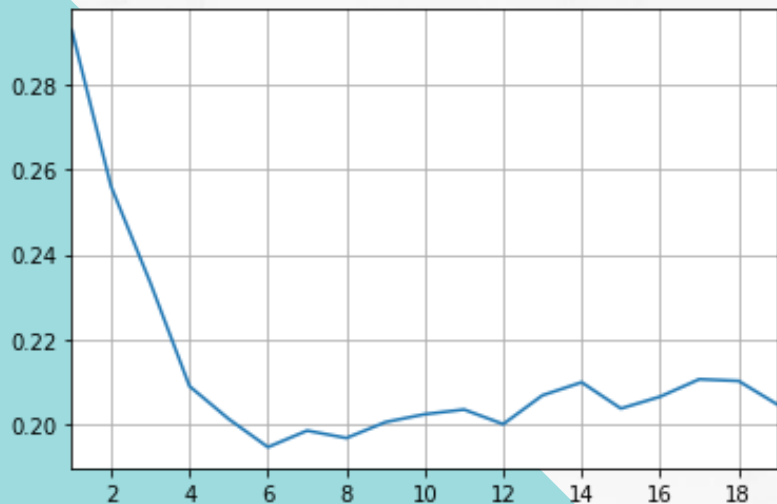**2. KNN**

- ○ nn = 6
- ○ Train RMSE = 0.1334
- ○ 10-CV RMSE = 0.1474

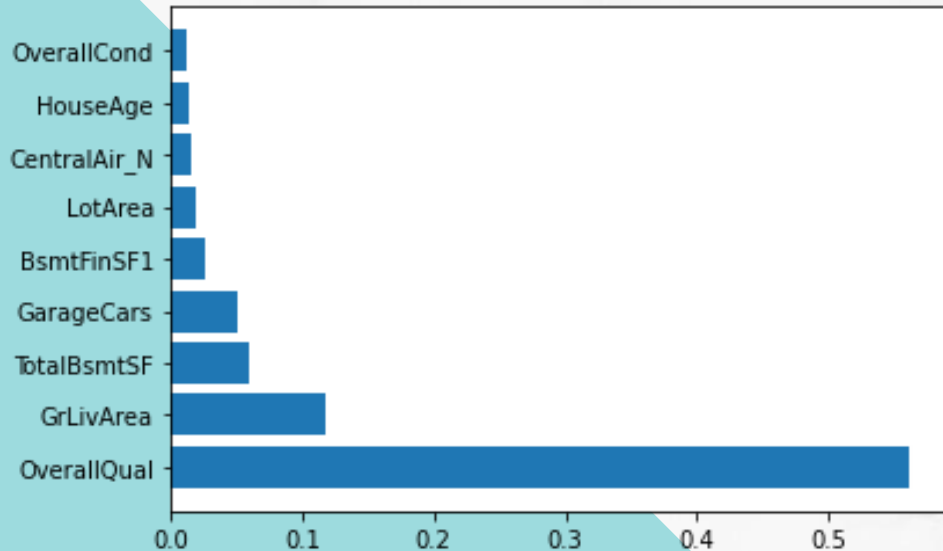# 05 MODELING



## 3. Decision Tree Regression

- Max depth = 6
- Train RMSE = 0.1343
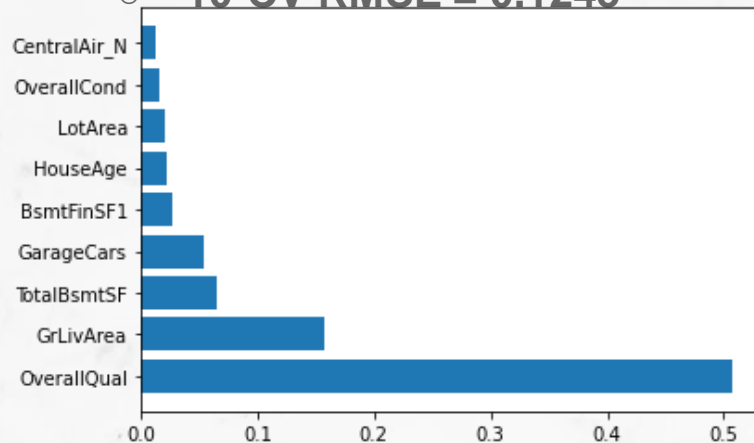- 10-CV RMSE = 0.1924



16

# 05 MODELING



**4. Random Forest**

- ○ **Using 100 trees**
- ○ **No sub-sampling**
- ○ **Train RMSE = 0.05217**
- ○ **10-CV RMSE = 0.1414**

| | i | j | rmse |
|---|---|---|---|
| 0 | 3 | 3 | 0.126778 |
| 1 | 3 | 4 | 0.126703 |
| 2 | 3 | 5 | 0.126928 |
| 3 | 3 | 6 | 0.126386 |
| 4 | 4 | 3 | 0.124517 |
| 5 | 4 | 4 | 0.124473 |
| 6 | 4 | 5 | 0.124713 |
| 7 | 4 | 6 | 0.124377 |
| 8 | 5 | 3 | 0.126461 |
| 9 | 5 | 4 | 0.126232 |
| 10 | 5 | 5 | 0.126498 |
| 11 | 5 | 6 | 0.126630 |
| 12 | 6 | 3 | 0.127997 |
| 13 | 6 | 4 | 0.128322 |
| 14 | 6 | 5 | 0.128316 |
| 15 | 6 | 6 | 0.128188 |

## 5. Gradient Boosting

- ○ **Max depth = 4**
- ○ **Min samples split = 6**
- ○ **Train RMSE = 0.06630**
- ○ **10-CV RMSE = 0.1245**

## Lin Reg

**Train RMSE**
= 0.1023
**10-CV RMSE**
= 0.1326

## KNN Reg

**Train RMSE**
= 0.1334
**10-CV RMSE**
= 0.1474

## Dec Tree

**Train RMSE**
= 0.1343
**10-CV RMSE**
= 0.1924

## Random Forest

**Train RMSE**
**= 0.05217**
**10-CV RMSE**
= 0.1414

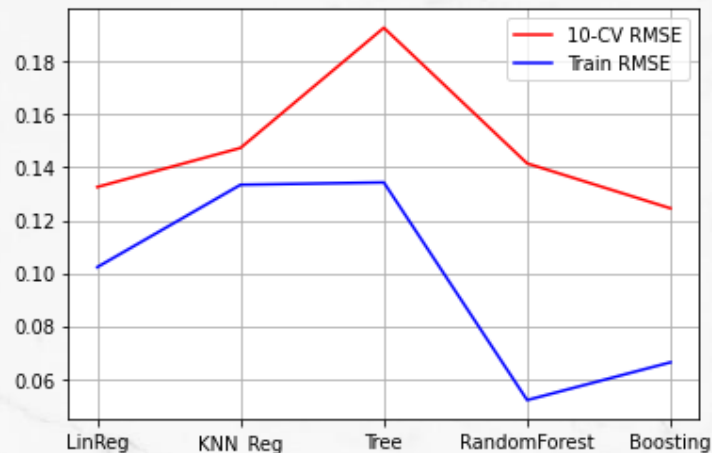## Gradient Boosting

**Train RMSE**
= 0.06630
**10-CV RMSE**
**= 0.1245**

**Random Forest**

- **Lowest Train RMSE of 0.05217**
- **High 10-CV RMSE suggests that there may be an overfit**

**Gradient Boosting**

- **Lowest 10-CV RMSE of 0.1245**
- **Best performance**

# 07 CONCLUSION & RECOMMENDATION

**Based on our dataset, gradient boosting lists the following features to be most influential in housing prices**

## 01 "OverallQual"

Overall material and finish quality

## 02 "GrLivArea"

Above ground living area square feet

## 03 "TotalBsmtSF"

Total square feet of basement area

## 04 "GarageCars"

Size of garage in car capacity

# 07 CONCLUSION & RECOMMENDATION

**Rooms for improvement:**

- PCA analysis when doing feature analysis
- Tune random forest parameters: number of estimators or sub sampling
- Try bagging models for regression

# THANK YOU