# Problem Statement:

To predict the sale price of residential homes in Ames, Iowa

# Data Collection

- Source:
  - Ames, Iowa Assessor's Office
- Period:
  - 2006 to 2010
- Size:
  - 2930 observations, 82 variables

# Data Cleaning

- Missing/Null values
  - Identify
    - Categorical
    - Numerical variables
  - Imputing
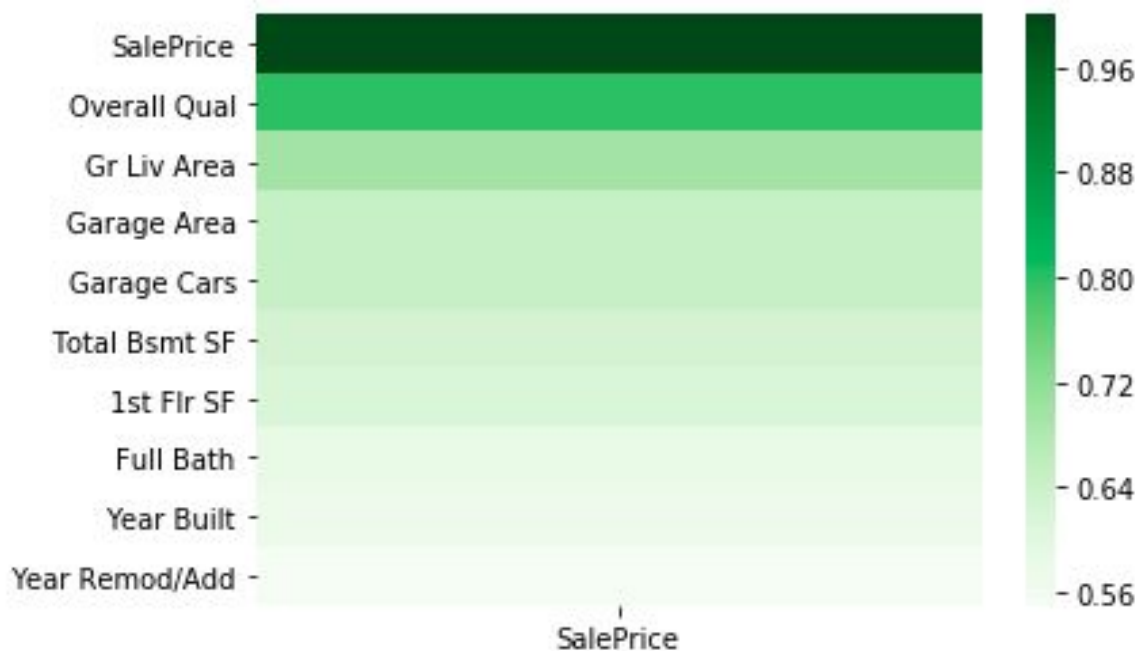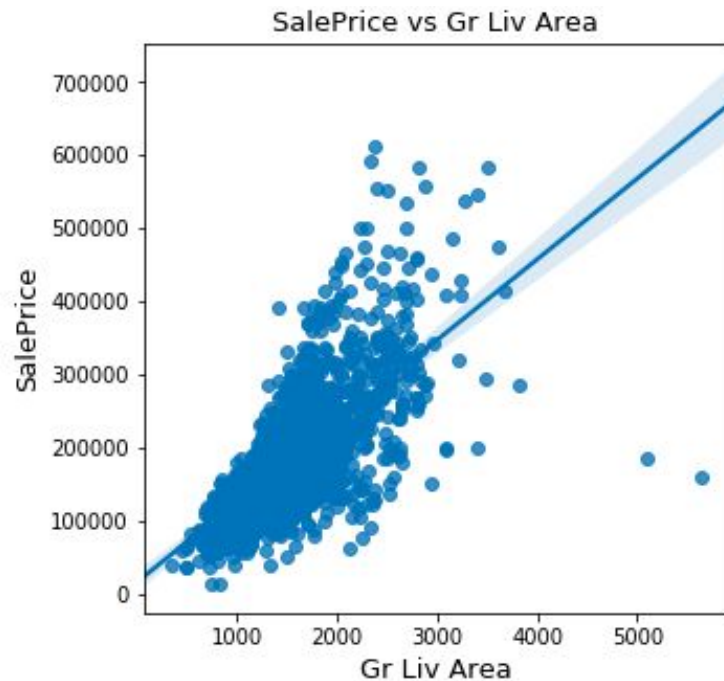  - Combine features with interaction terms

DOES THIS DATA SPARK JOY?

# EDA

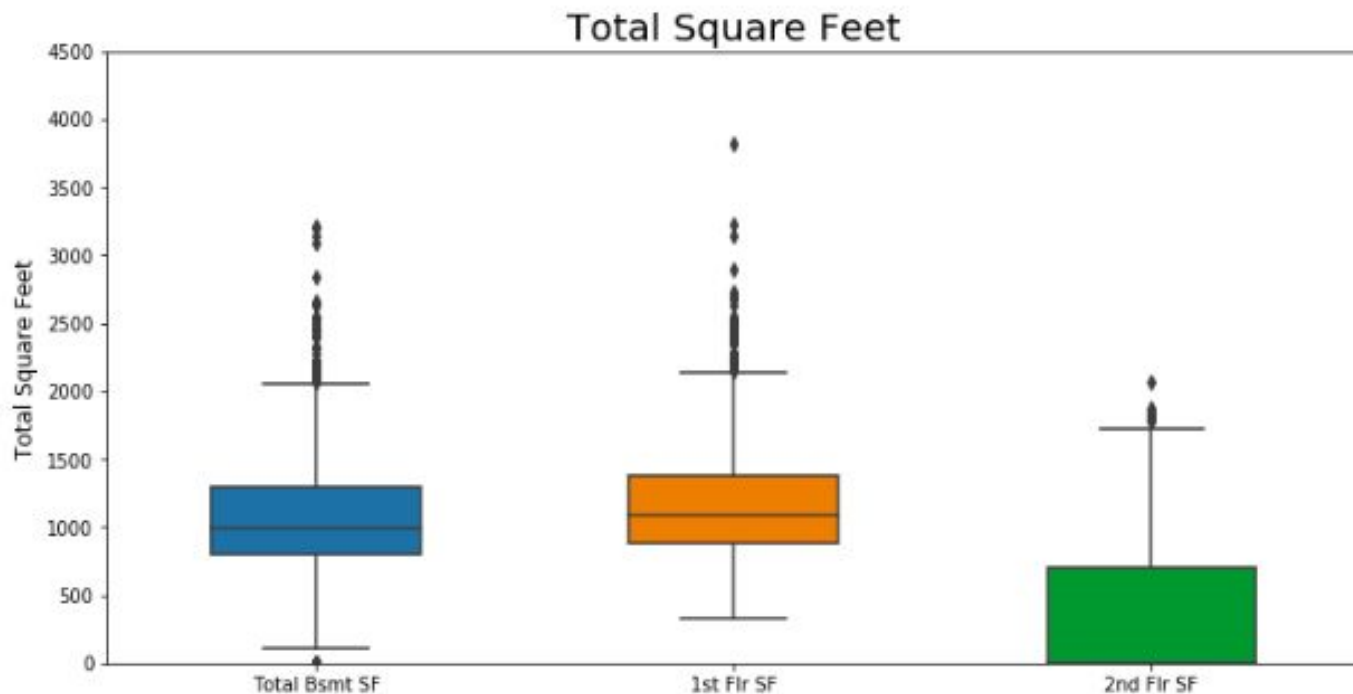- Identify variables with highest correlations with Sale Price

# EDA

- Remove outliers

# EDA

- Remove outliers - cont'
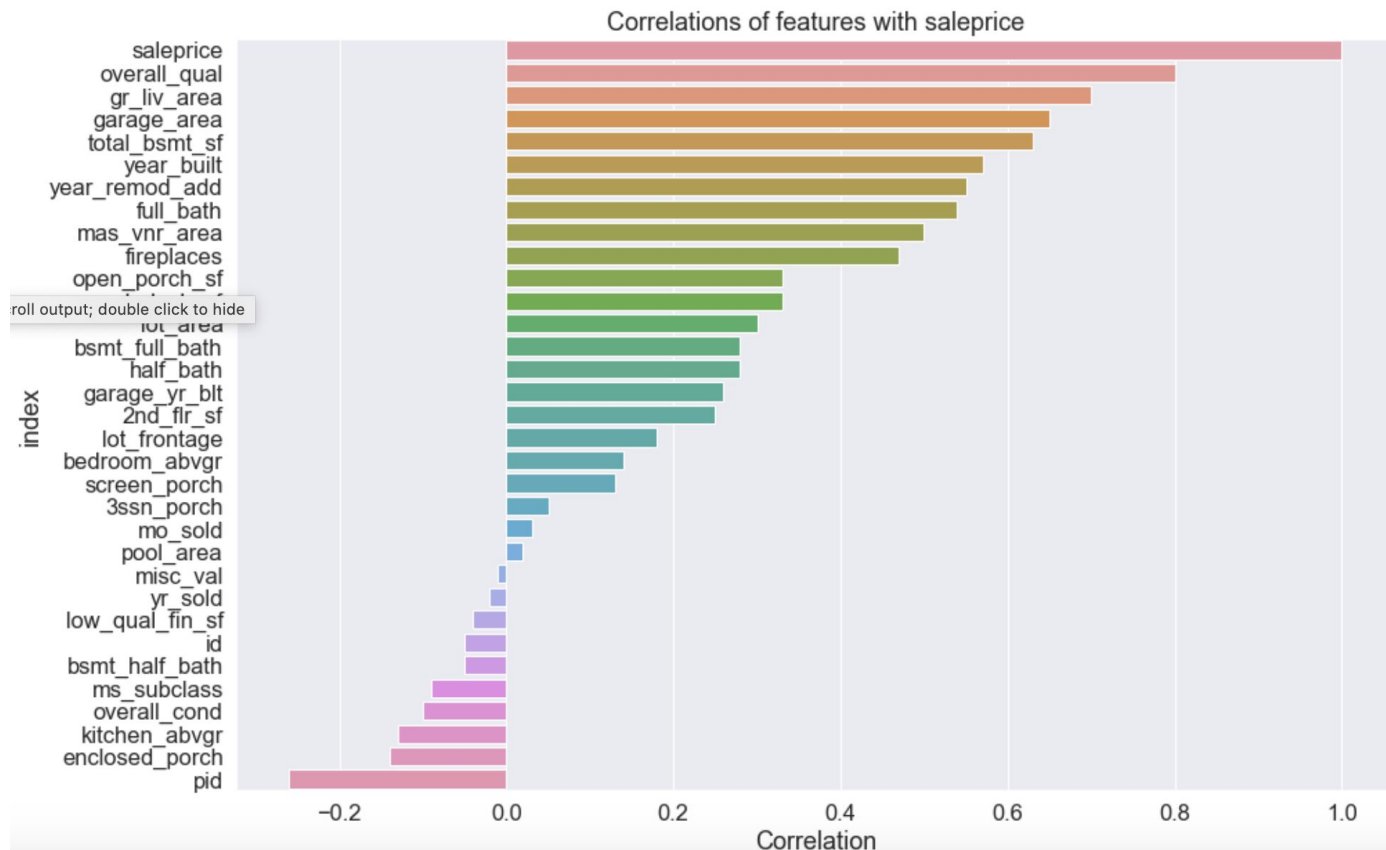
Correlations of features with saleprice

Heatmap & Feature Selection
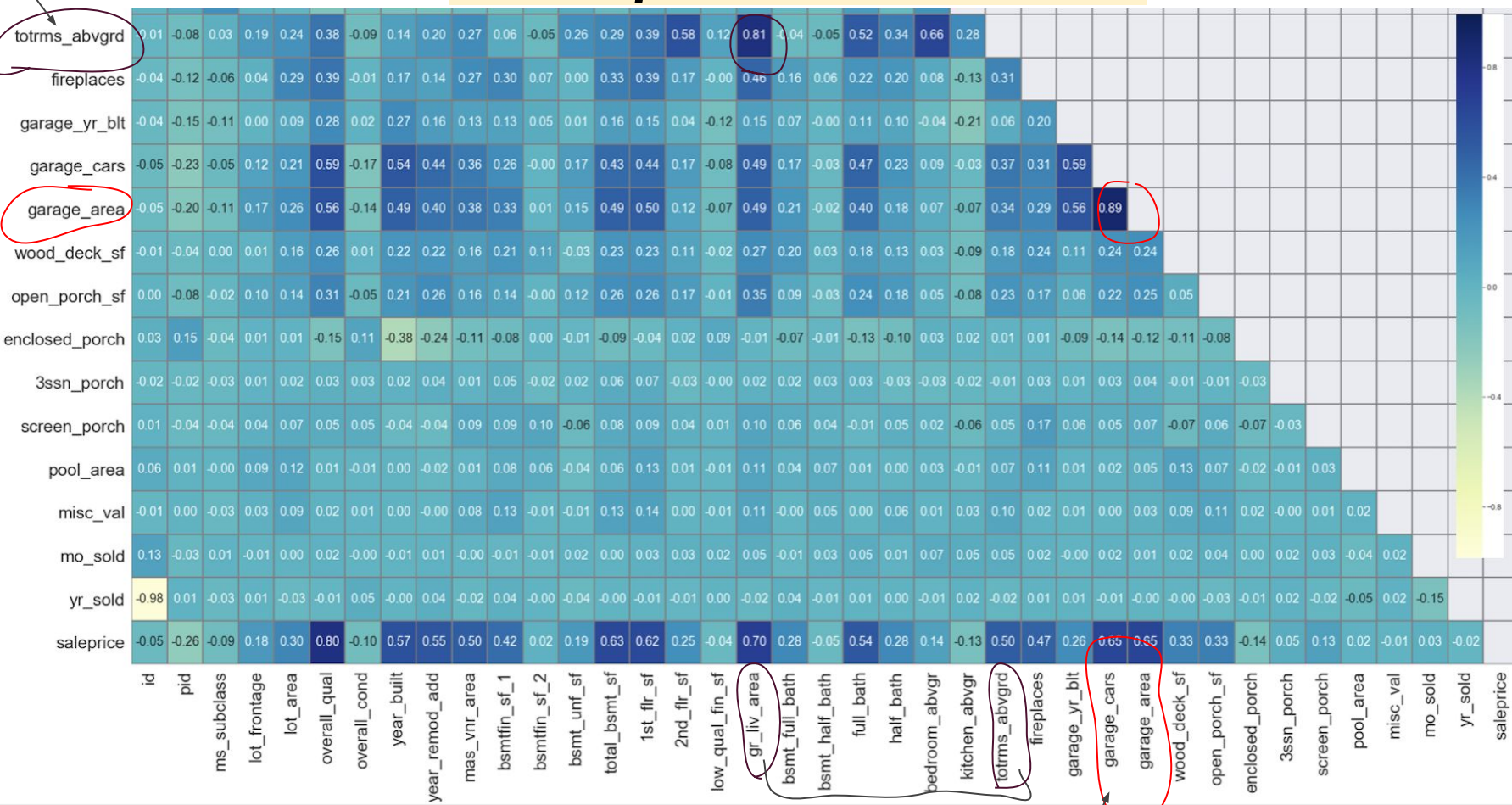
If computers can express?????.....

To avoid this.......

# Ordinal and Categorical features

**Variable**

**Categorical** (qualitative)

**Nominal**
Unordered, categories which are mutually exclusive
e.g. male/female, smoker/non-smoker

**Ordinal**
Ordered, categories which are mutually exclusive
e.g. IOTN 1/2/3/4/5 or minimal/moderate/severe/unberable pain

**Numerical** (quantitative)

**Discrete**
Whole numerical value - typically counts
e.g. number of visits to dentist, DMF

**Continuous**
Can take any value within a range e.g. height in cm, pocket depth in mm

Pen    Pencil    Eraser

Cow    Dog    Cat

Excellent    Good    Bad

Fantastic    Okay    Don't Like

# Pre-processing

## Train-test-split



## Scaling

# Modeling & Evaluation

**What would be a naive guess of the sale price if we did not do any analysis?**

# So how? Was it better or not?

# Train-test split comparison



Test R2 = 0.870
CV R2   = 0.833
Train R2= 0.898

# Better than a naive prediction

| | Root mean squared error |
|---|---|
| **Baseline** | 78,512 |
| **Multiple regression** | 29,714 |
| **Lasso regression** | 28,628 |

# Assumptions of regression

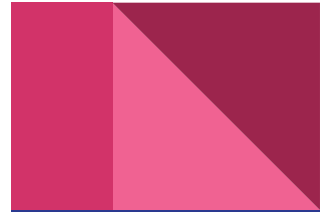1. Predictors must have an approximately linear relationship with the dependent variable
2. Errors are normally distributed
3. Independence of error
4. Constant error variance
5. Independence of Predictors

Histogram plots of unstandardized residuals

QQ plots of unstandardized residuals

Predicted vs Residuals

- Which features appear to add the most value to a home?

Total Living Space (TOTAL BSMT SF + GR LIV AREA) appears to add the most value to a home;

the more space, rooms and land that the house contains, the higher the valuation

Which features hurt the value of a home the most?

By the same token, the factors that add the most value to a home also hurt the value of a home the most.

A smaller Total Living Space and poor Overall quality of the property is detrimental to home value.

- What are things that homeowners could improve in their homes to increase the value?

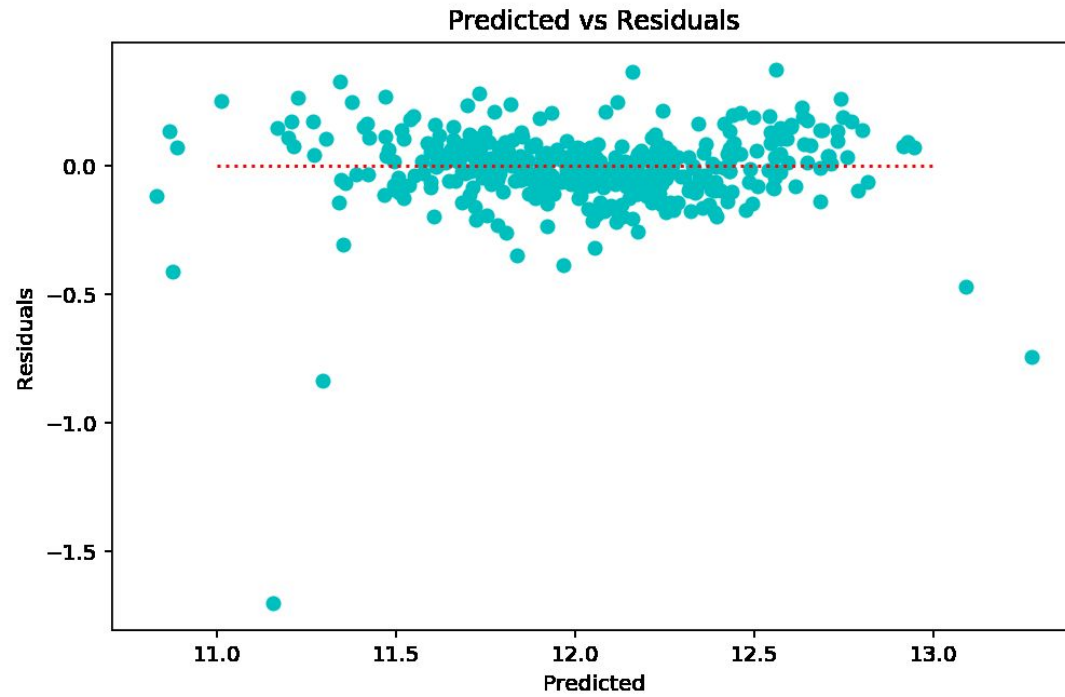Homeowners can improve the value of their homes by renovation using a dual approach to both increase the Total Living Space and Overall quality of the property.

Overall Qual is the best determinant providing the strongest correlations with SalePrice.

Gr Liv Area is the second best determinant providing a strong correlations with SalePrice.

Total Bsmt SF is the third best determinant providing a strong correlations with SalePrice.

# Limitations

- High number of outliers means
- Not all independent variable are normally distributed
- Specific to Ames, Iowa
- 5 years of data from 2006 - 2010

# Conclusion

This project addressed the non-fungibility of real estate property by building a model and make an accurate prediction for property sale price(target variable).

The prospective buyer(s), seller(s) and stakeholder(s) have the need to appraisal a fair market value of a property in order to facilitate transaction, taxation, appraisal, etc.