

An aerial photograph of a suburban neighborhood, showing a grid of streets, houses with various roof colors (red, grey, blue), green lawns, and trees. The image is used as a background for the title slide.

Ames Housing Price Prediction Model

By: Alyse, Sim Yi

Background

Who are we?

HOUSE FLIPPING CONSULTANCY SERVICE

With Core Business in Other Midwestern States



We are expanding in the Neighbouring state of Iowa
with Ames as our first target

Background

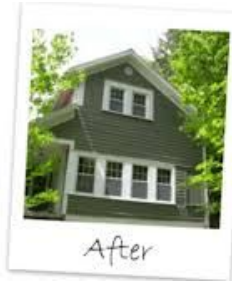
Why Iowa? Why Ames?

- **Seller's Market:**
 - 3.9% Increase in Y-o-Y House Prices (2018)
 - Overall Increase in Number of Houses Sold
- **Booming Home Sales in Ames:**
 - 9.4% Increase in Y-o-Y Number of Houses Sold
- **Highly Profitable:**
 - Cheap renovation and remodelling costs
 - Average flip earns an average of US\$ 54K
 - Top 50 cities in United States for profiting from house flips

Problem Statement

What are we trying to address?

**Identify key driving factors for maximising returns
from house flipping activities**



**WHAT TO
BUY**



**HOW TO
RENOVATE**

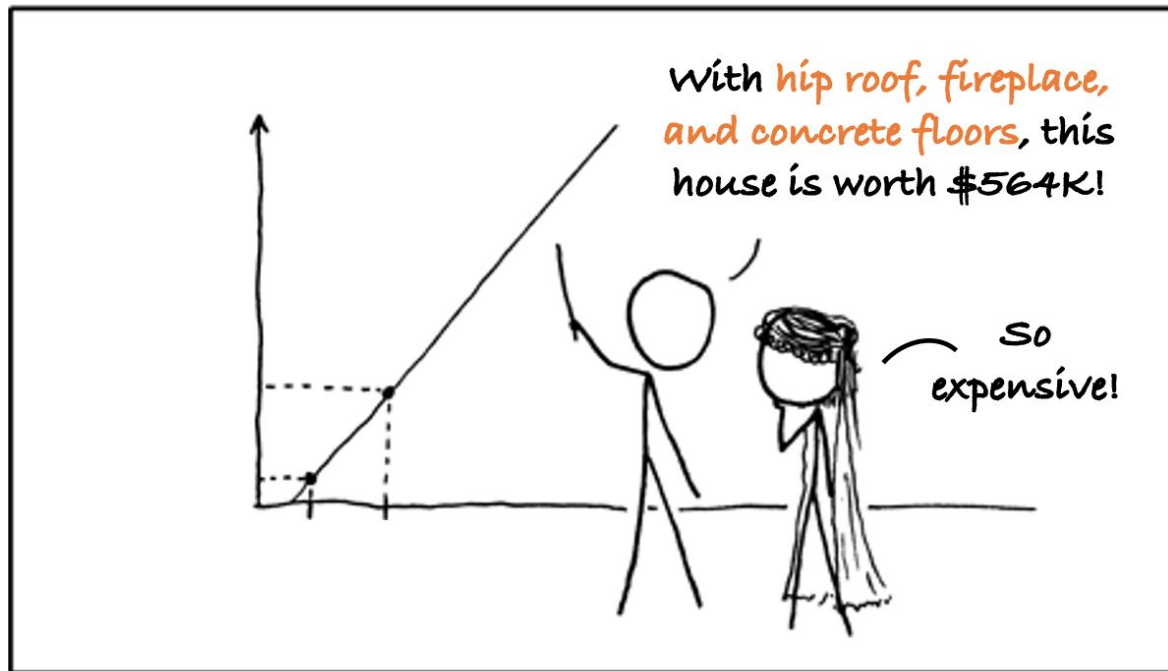


**WHEN TO
SELL**

Problem Statement

How are we going to do it?

Multiple Linear Regression



Project Overview

End-to-end management of project

1. Data Preparation

- a. Data Cleaning
- b. Exploratory Data Analysis

2. Exploratory Visualizations

- a. Correlations with Dependent Variable
- b. Multicollinearity
- c. Skewness

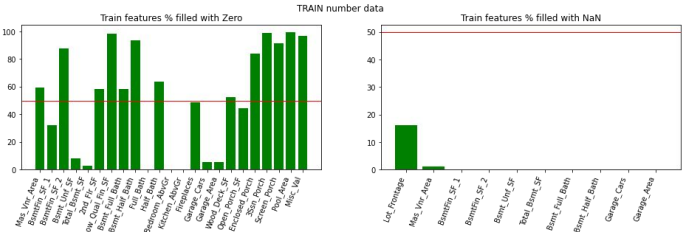
3. Modeling and Testing

4. Inferential Visualizations

5. Conclusion & Recommendations

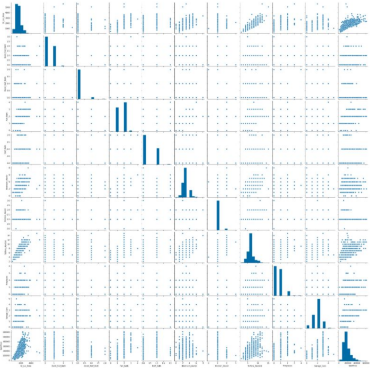
Data Processing

How we classified our data

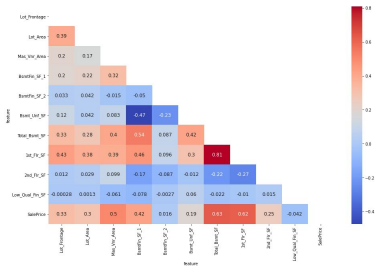


Numerical - Nominal and Ordinal

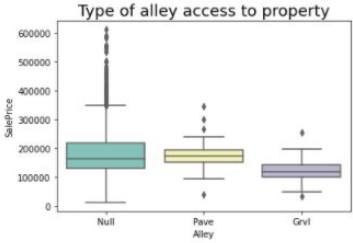
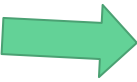
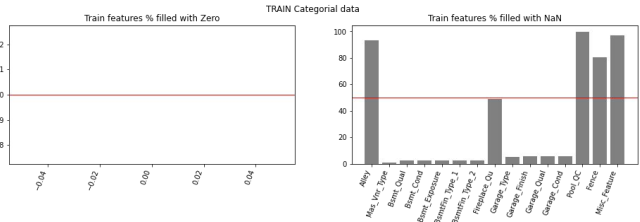
Spread / Scatter



Correlation



Categorical - Nominal and Ordinal



Data Cleaning

How did we prepare our data for modelling?

1. Filling null values with :
 - a. 0 for numerical values
 - b. NA for categorical values
2. Changing Ordinal Categorical data to numerical values

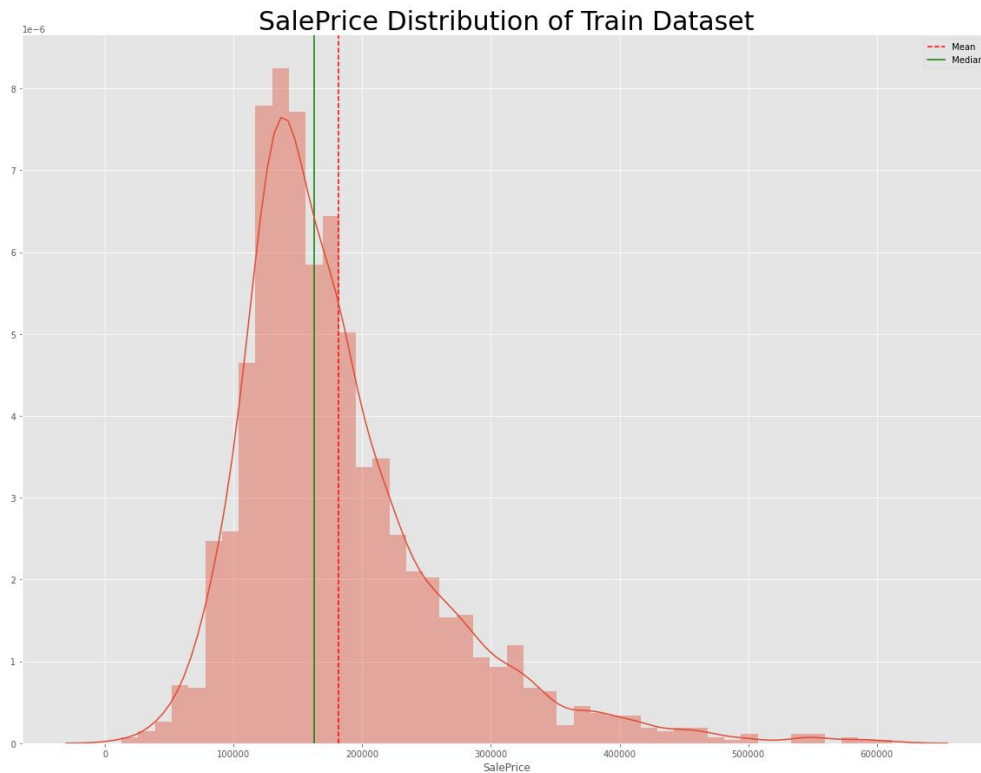
```
1 #filling dtypes object nan values
2 df['Bsmt Qual'].fillna('NA', inplace = True)
3 df['Bsmt Cond'].fillna('NA', inplace = True)
4 df['Bsmt Exposure'].fillna('NA', inplace = True)
5 df['BsmtFin Type 1'].fillna('NA', inplace = True)
6 df['BsmtFin Type 2'].fillna('NA', inplace = True)
7
8 #filling dtypes float64 nan values with number 0
9 df['BsmtFin SF 1'].fillna(0 , inplace = True)
10 df['BsmtFin SF 2'].fillna(0 , inplace = True)
11 df['Bsmt Unf SF'].fillna(0 , inplace = True)
12 df['Total Bsmt SF'].fillna(0 , inplace = True)
13 df['Bsmt Full Bath'].fillna(0 , inplace = True)
14 df['Bsmt Half Bath'].fillna(0 , inplace = True)
```

```
1 #Copied from Data cleaning to apply on Test data
2 labels = {'Exter Qual': {'Ex': 5, 'Gd': 4, 'TA': 3, 'Fa': 2, 'Po': 1} }
3 df_test.replace(labels, inplace=True)
4
5
6 labels = {'Exter Cond': {'Ex': 5, 'Gd': 4, 'TA': 3, 'Fa': 2, 'Po': 1} }
7 df_test.replace(labels, inplace=True)
8
9 labels = {'Bsmt Qual': {'Ex': 5, 'Gd': 4, 'TA': 3, 'Fa': 2, 'Po': 1, 'NA': 0} }
10 df_test.replace(labels, inplace=True)
11
12
13 labels = {'Bsmt Cond': {'Ex': 5, 'Gd': 4, 'TA': 3, 'Fa': 2, 'Po': 1, 'NA': 0} }
14 df_test.replace(labels, inplace=True)
15
16 labels = {'Bsmt Exposure': {'Gd': 4, 'Av': 3, 'Mn': 2, 'No': 1, 'NA': 0} }
17 df_test.replace(labels, inplace=True)
18
```


Exploration of Data

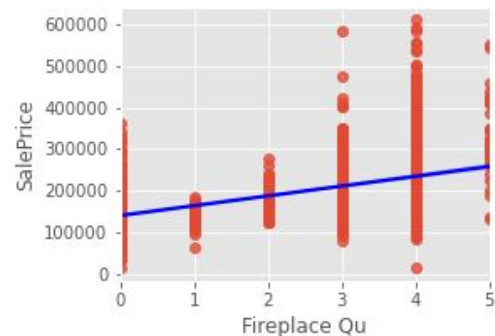
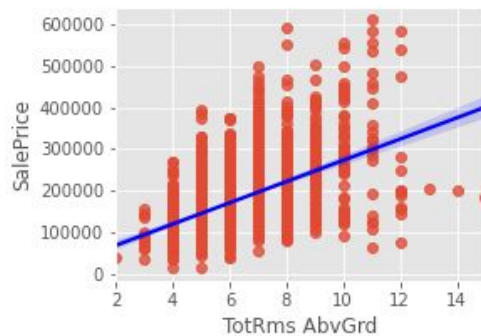
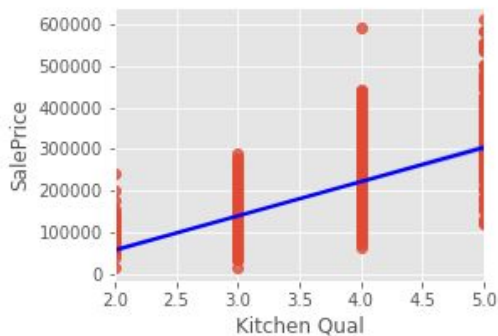
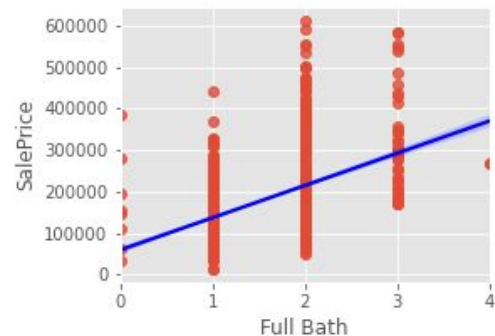
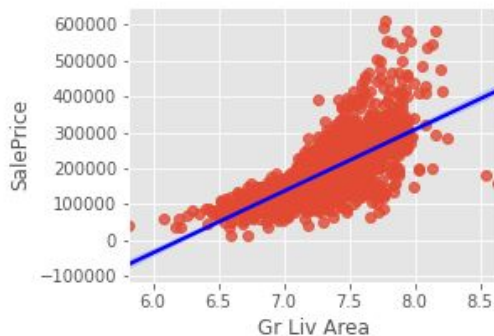
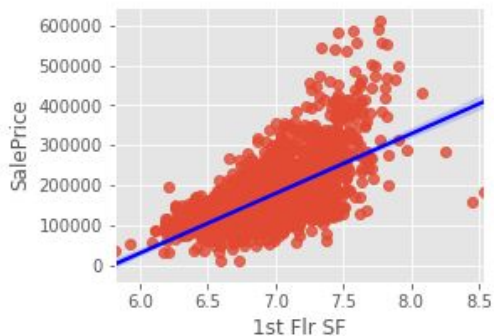
Continuous Data Plotting

1. Saleprice is positively skewed
2. the mean and median are more towards the lower prices



Exploration of Data

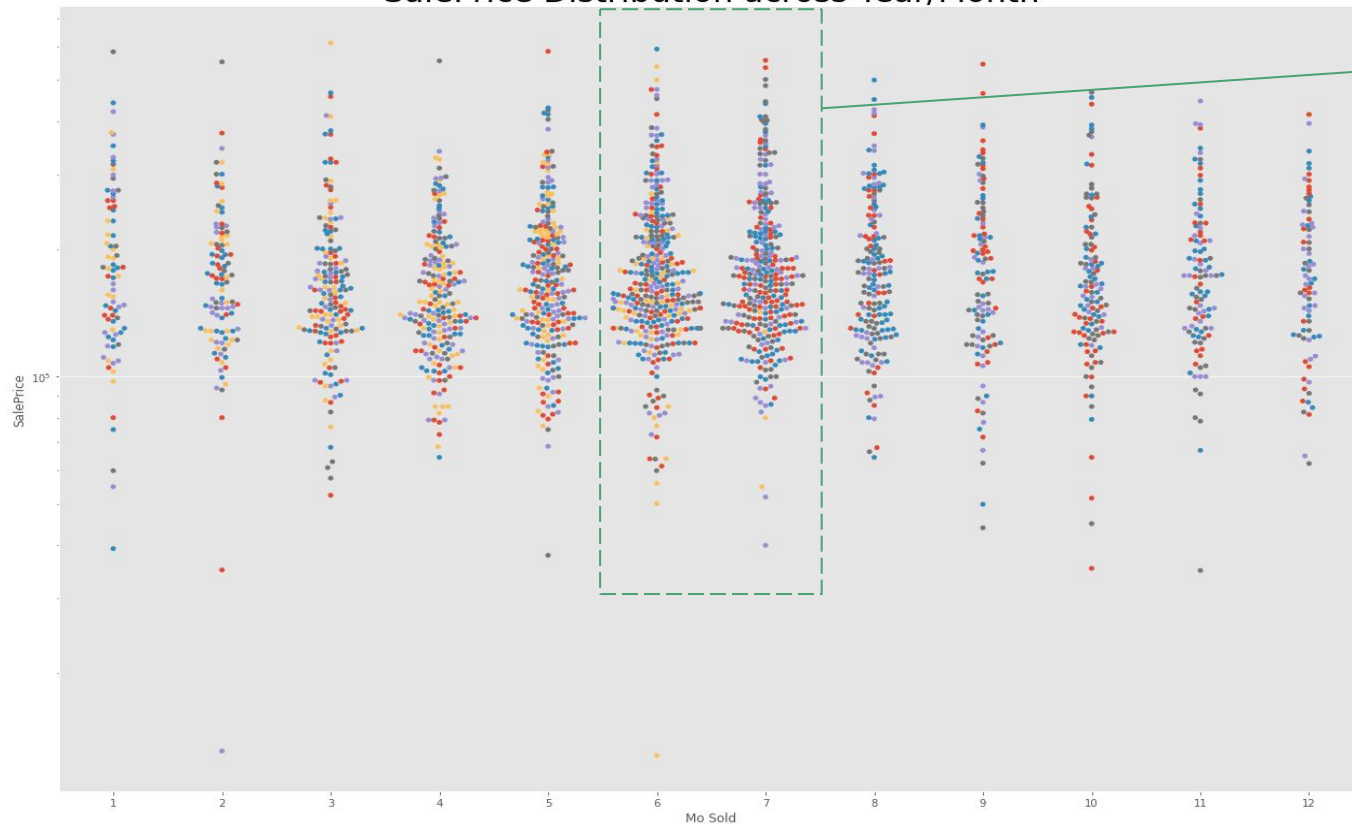
Continuous Data Plotting



Exploration of Data

Time Series Swarmplot

SalePrice Distribution across Year/Month

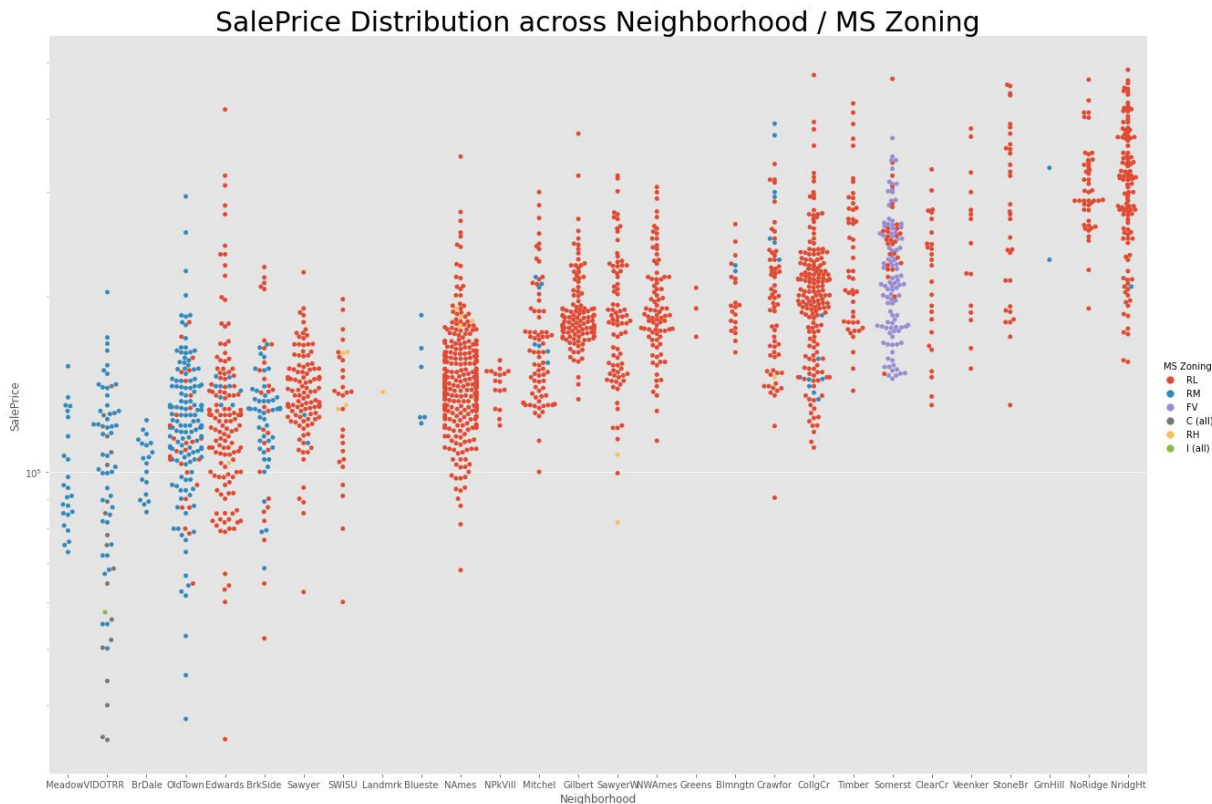


June and July

Highest sales
volume

Exploration of Data

Neighbourhood Distribution Swarmplot



**Neighbourhood factor is
a key determinant of
house prices**



**NOT IMPORTANT
in short-term house
flipping investment**

Legend: MS Zoning

A Agriculture
C Commercial
FV Floating Village Residential
I Industrial
RH Residential High Density
RL Residential Low Density
RP Residential Low Density Park
RM Residential Medium Density

Feature Engineering

Standard Scaler, One_Hot Encoding, Polynomial Features

```
1 num_data = df.select_dtype(['int64', 'float64']).keys()
2 num_data = [x for x in num_data if (x != 'SalePrice') & (x != 'Id')]
3
4 nuns = df[num_data]
5 ss = StandardScaler()
6 ss.fit(nuns)
7 nuns_scaled = ss.transform(nuns)
8
9 nuns_scaled.shape
(2018, 21)
```

```
1 nuns_scaled_pd = pd.DataFrame(nuns_scaled, columns = num_data) #create pd for combining later
2 nuns_scaled_pd.head()
```

Year Built	Year Remodeled	Basement	Central Air	1st Flr SF	Full Bath	Kitchen Qual	Topsoil AreaQual	Garage Finish	Garage Cars	Overall Total Value	Estimate Total Value	Total Flr SF	Fireplace Total Value	
0	0.143092	0.989556	-0.582426	0.268819	-1.511006	0.769545	0.760728	-0.247485	0.312884	0.306899	-1.474812	0.757603	0.164064	-0.968370
1	0.896874	0.914545	0.580207	0.268819	-0.582426	0.769545	0.760728	-0.247485	0.312884	0.306899	-1.474812	0.757603	0.164064	-0.968370
2	-0.632317	1.090609	-0.582426	0.268819	-1.511006	0.769545	0.760728	-0.247485	0.312884	0.306899	-1.474812	0.757603	0.164064	-0.968370
3	1.142166	1.090609	0.580207	0.268819	-1.511006	0.769545	0.760728	-0.247485	0.312884	0.306899	-1.474812	0.757603	0.164064	-0.968370
4	-0.388099	0.424037	-1.686379	0.268819	-0.582426	0.769545	0.760728	-0.247485	0.312884	0.306899	-1.474812	0.757603	0.164064	-0.968370

5 rows x 21 columns

Standard Scaler

Scale the data to the same level across dataset

```
1 obj_processed = pd.get_dummies(df[obj_data], columns = obj_data)
2 obj_processed.head()
```

MS Zoning_C (cat)	MS Zoning_FV	MS Zoning_I (cat)	MS Zoning_RH	MS Zoning_FL	MS Zoning_RM	Street_Grt	Stove_Pave	Contour_Brk	Contour_HLS	Heating_Grwy	Heating_OBHV	H
0	0	0	0	0	1	0	0	1	0	0	0	0
1	0	0	0	0	1	0	0	1	0	0	0	0
2	0	0	0	0	1	0	0	1	0	0	0	0
3	0	0	0	0	1	0	0	1	0	0	0	0
4	0	0	0	0	1	0	0	1	0	0	0	0

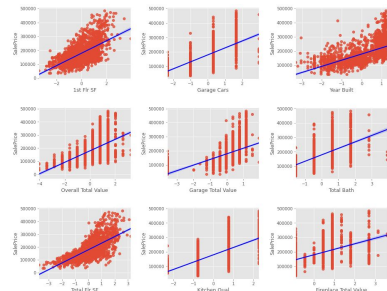
5 rows x 123 columns

```
1 #Remove columns with NA count values
2 na_col = obj_processed.filter(regex = "NA")
3 na_col_keys = na_col.keys()
4 na_col_keys
```

Index(['Neighborhood_Names', 'Mas Vnr Type_NA', 'Garage Type_NA'], dtype='object')

One Hot Encoding

Convert categorical data into numerical data



```
1 print("Garage & Total & Year Features RMSE(L2): " +
2       str(np.sqrt(-cross_val_score(lasso, X=X, y=y, cv=5,
3                                     scoring="neg_mean_squared_error").mean()))
Garage & Total & Year Features RMSE(L2): 38786.89039192311
```

```
1 poly = PolynomialFeatures(degree = 2, interaction_only = True, include_bias=False)
2 X_poly = poly.fit_transform(X)
```

```
1 print("Garage & Total & Year Polynomial Features RMSE(L2): " +
2       str(np.sqrt(-cross_val_score(lasso, X=X_poly, y=y, cv=5,
3                                     scoring="neg_mean_squared_error").mean()))
Garage & Total & Year Polynomial Features RMSE(L2): 23820.97583158794
```

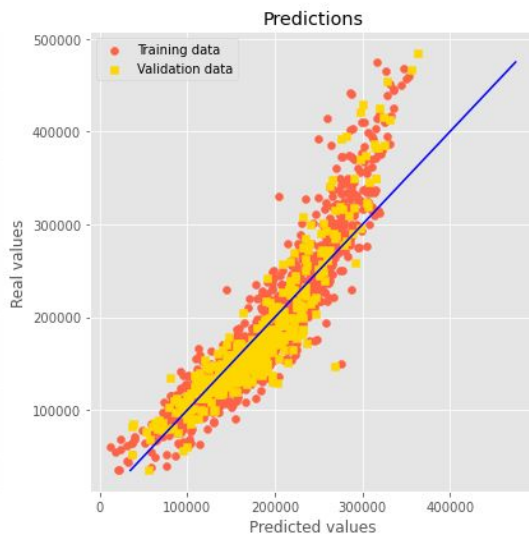
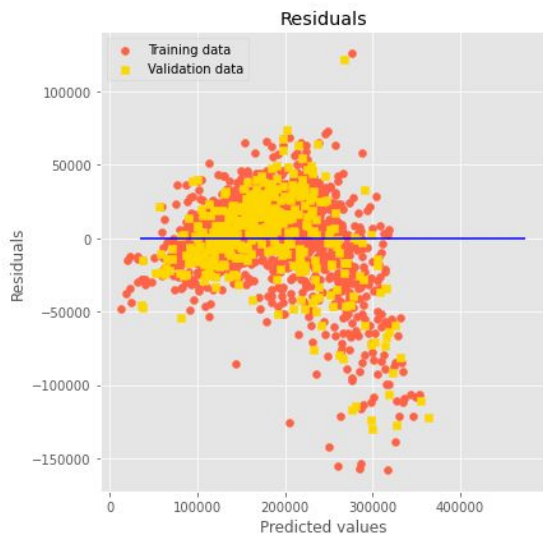
Polynomial Features/ Interaction Terms

Transform data to better fit the best fit line

Model Baseline

Linear Regression

- **Model Baseline**
- Model tuning
- Feature select
- Final model



R2 Score

- Train Score(Lr):
 - 0.898
- Validation Score(Lr):
 - 0.870
- Test Data Estimated score(Lr):
 - $-8.607 \times 10^{+19}$

RMSE Scores

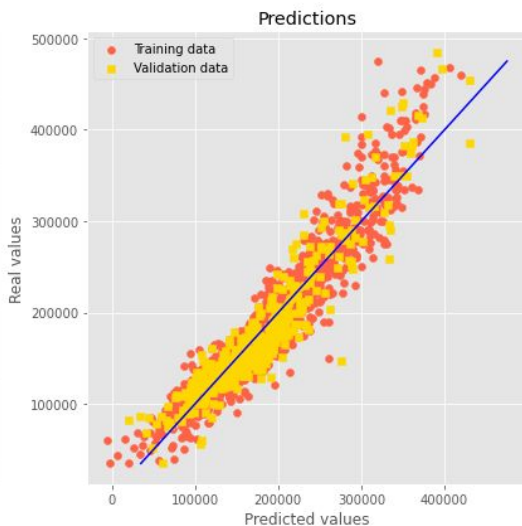
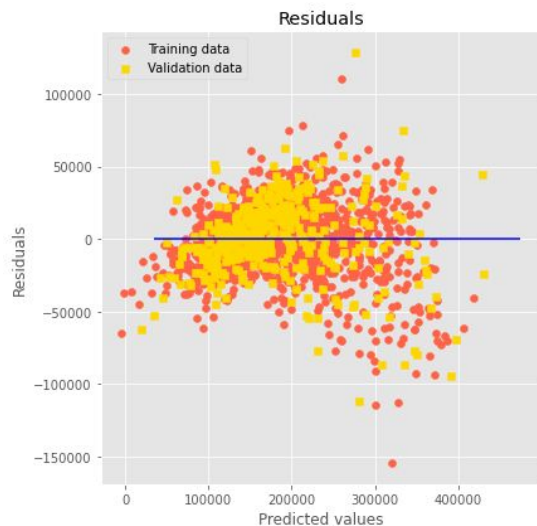
- Train RMSE(Lr):
 - 23224
- Validation RMSE(Lr):
 - 26995
- Test Data Estimated RMSE(Lr):
 - 53515

Model Tuning

Regularization

Using RidgeCV, LassoCV & ElasticNetCV (LassoCV chosen)

- Model Baseline
- **Model tuning**
- Feature select
- Final model



R2 Score

- Train Score(Lr):
 - 0.9029
- Validation Score(Lr):
 - 0.87680
- Test Data Estimated score(Lr):
 - 0.8771

RMSE Scores

- Train RMSE(Lr):
 - 22649
- Validation RMSE(Lr):
 - 26323
- Test Data Estimated RMSE(Lr):
 - 25480

Feature Selection

RFECV

- Model Baseline
- Model tuning
- **Feature select**
- Final model

```
1 rfecv = RFECV(lasso, step = 1, #using lasso model as the results are better among the 3 models
2               min_features_to_select = 30, #minimum of 30 features
3               cv=5,
4               scoring = 'neg_mean_squared_error') #calculating using mse
5
6 rfecv.fit(X_train, y_train)
```

```
RFECV(cv=5, estimator=Lasso(max_iter=2000), min_features_to_select=30,
      scoring='neg_mean_squared_error')
```

```
1 gd_coef = list(compress(X_train.columns, rfecv.support_))
2 print(str(len(gd_coef)) + ' features left after RFECV feature selection')
```

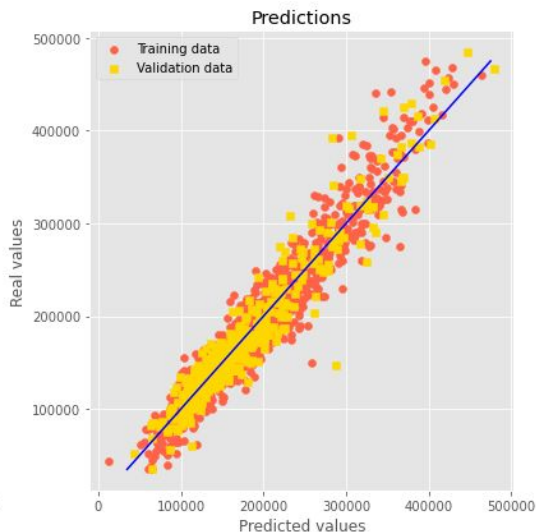
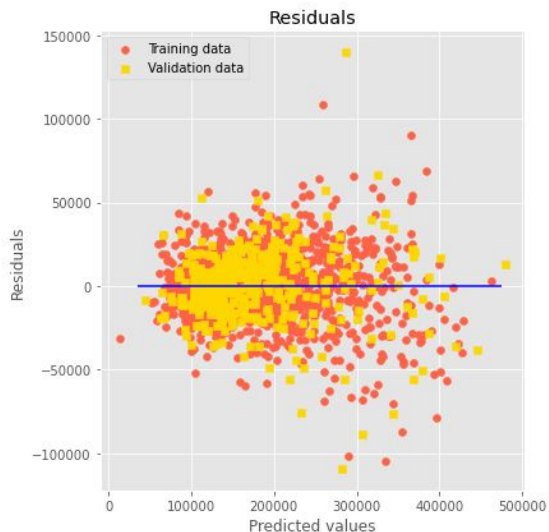
94 features left after RFECV feature selection

Using **RFECV** with the Lasso model to feature select the best coefficient features for further model tuning

Final Model

Final Model (Lasso Regression)

- Model Baseline
- Model tuning
- Feature select
- **Final model**



R2 Score

- Train Score(Lr):
 - 0.9318 (bef: 0.898)
- Validation Score(Lr):
 - 0.9163 (bef: 0.0.870)
- Test Data Estimated score(Lr):
 - 0.9123 (bef: -8.607e+19)

RMSE Scores

- Train RMSE(Lr):
 - 18993 (bef: 23224)
- Validation RMSE(Lr):
 - 21692 (bef: 26995)
- Test Data Estimated RMSE(Lr):
 - 21504 (bef: 53515)

Kaggle Submission

Submission Model (Ridge and Lasso Regression)

Submission and Description	Private Score	Public Score	Use for Final Score
submission_ridge.csv just now by Sim Yi add submission details	28300.91072	28325.32706	<input type="checkbox"/>
submission_lasso.csv a few seconds ago by Sim Yi add submission details	27930.73538	27626.17192	<input type="checkbox"/>

Key challenges

What are some limitations we faced in our modelling exercise?

Data Collection

- Lack of **post-2010 data** that are extremely vital to understanding trend of house prices due to 2008 Global Financial Crisis
- **Demographic & Economic data:** How has the demographic and economic factors over the four years (and after) impacted housing market in Ames?
- **Lacking of certain data** that could be helpful for our recommendations (e.g resale rate of property, key amenities around house, transportation services etc.)

Modelling

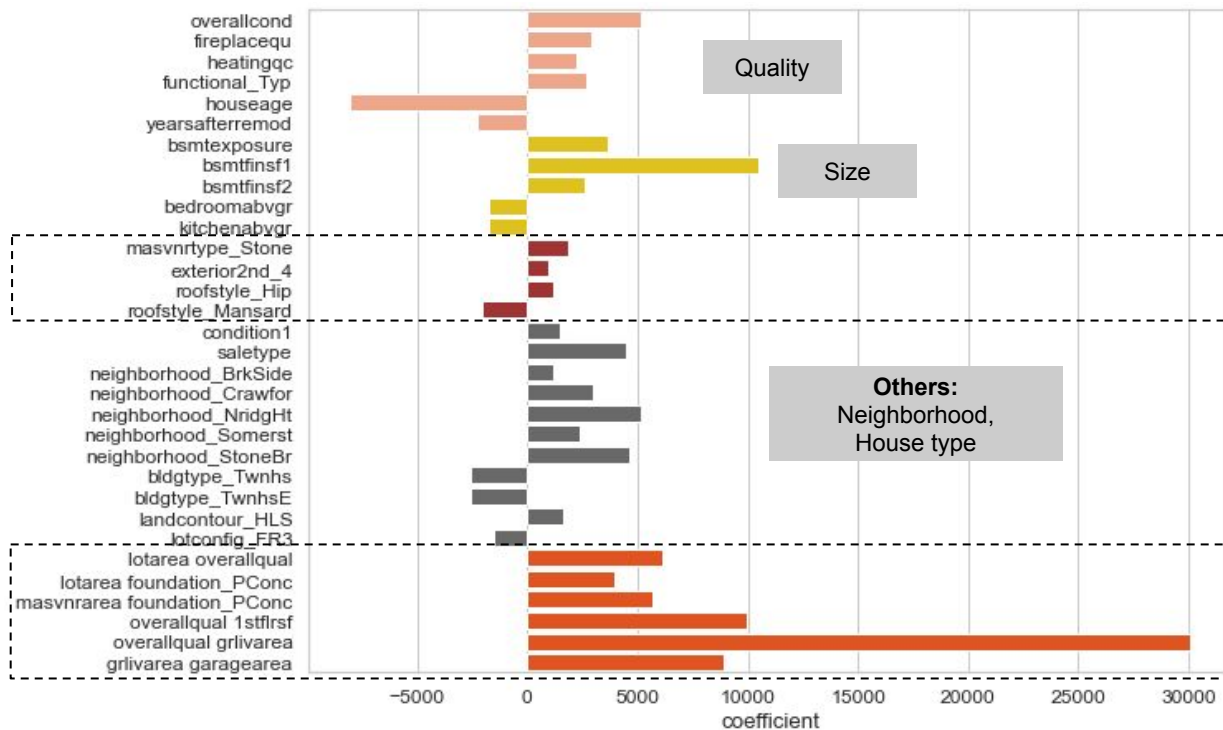
- Many Features were **highly correlated** to each other, making it difficult to choose which features to use (multicollinearity issue)
- Use of other **more advanced modelling techniques** (ie. RandomForest) could be more efficient in our prediction model.

Results from our Model

What can we learn from our regression model?

Renovation
Materials

Quality & Size



Recommendations

How to maximise returns from house flipping?

WHAT TO BUY?

- Poor quality houses + Big area → *Cheapest & Highest Potential*
- Houses with garage and basements → *Lift house prices*

HOW TO RENOVATE?

- Poured Concrete Foundation
- Hip-style roof
- Metal Exterior Covering
- Stone Masonry Veneer

WHEN TO SELL?

- Summer (June-July)

Ensure that renovations lead to much higher quality

Big houses + Excellent quality leads to disproportionately higher sale price.