**Horse price prediction - Kaggel data set**
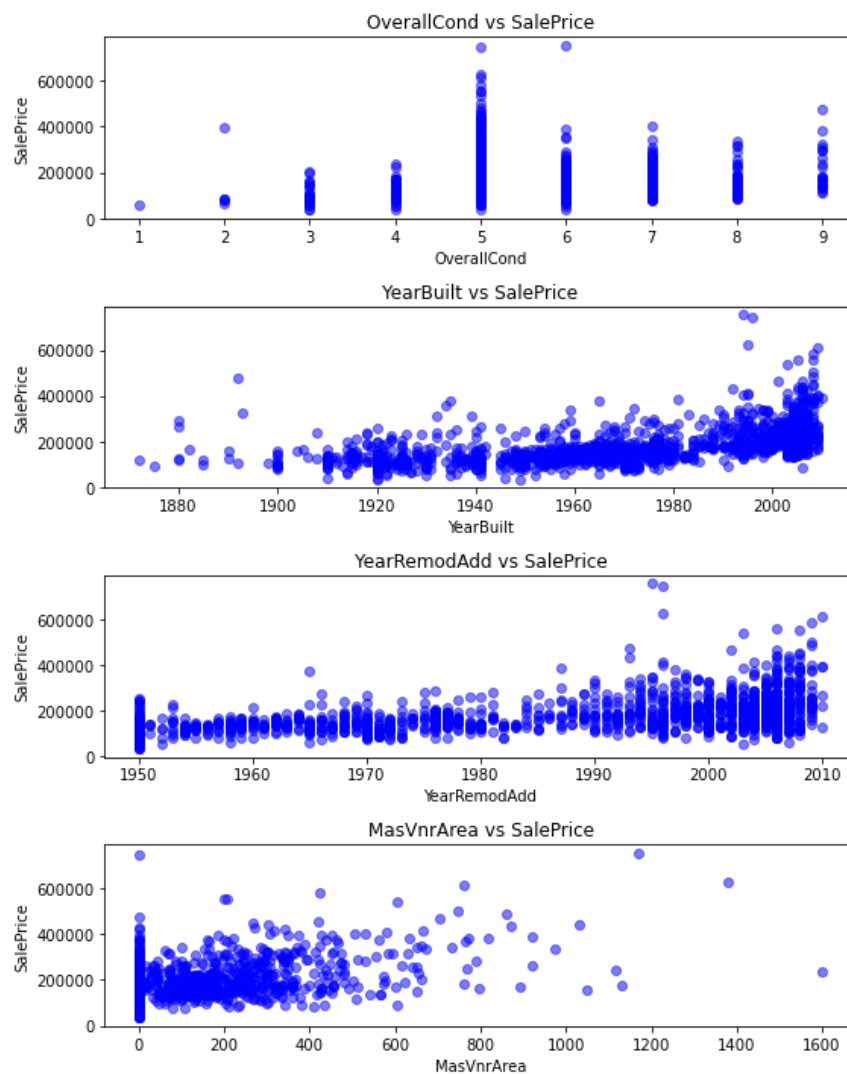
The purpose of this project is to build regression models to predict the price of the house based on the provided data.

**Data Overview:**
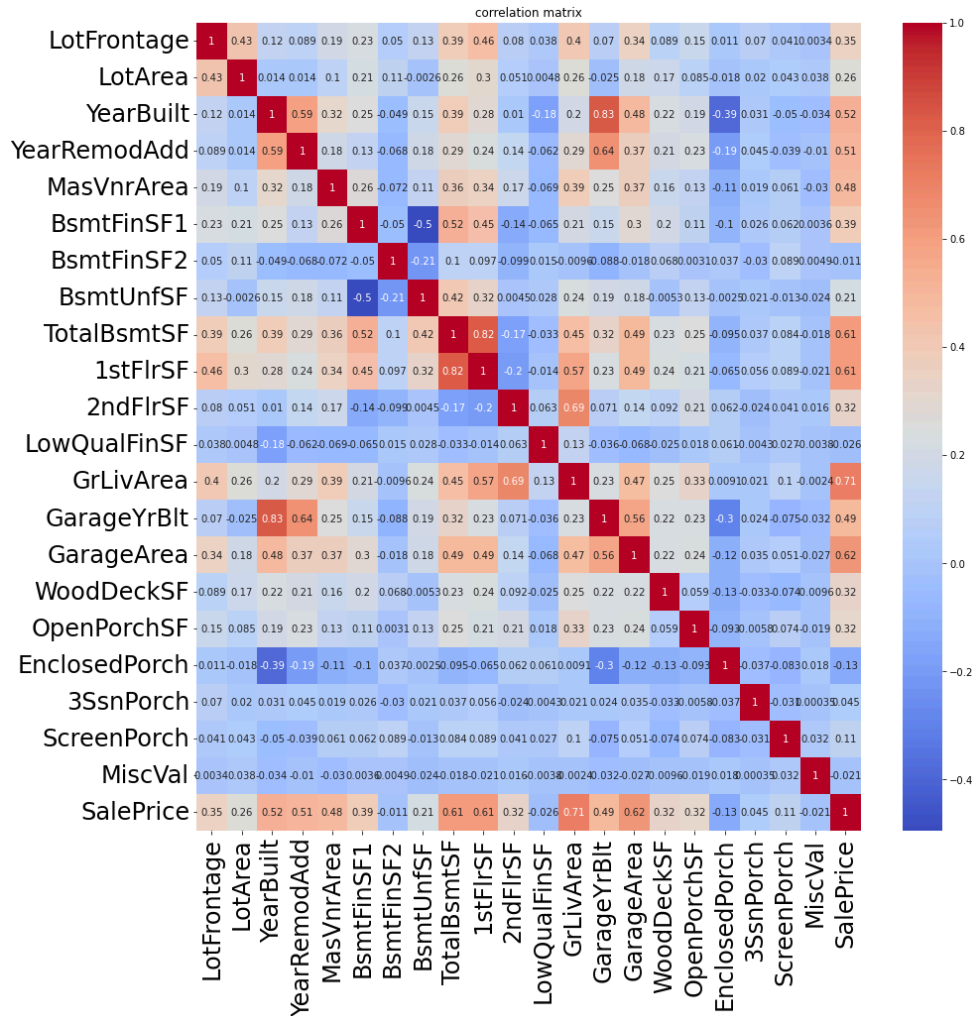Data has 80 columns and 1460 rows. The target variable is SalePrice. 37 columns are numeric and 43 are categorical.

**Scatter plot:**

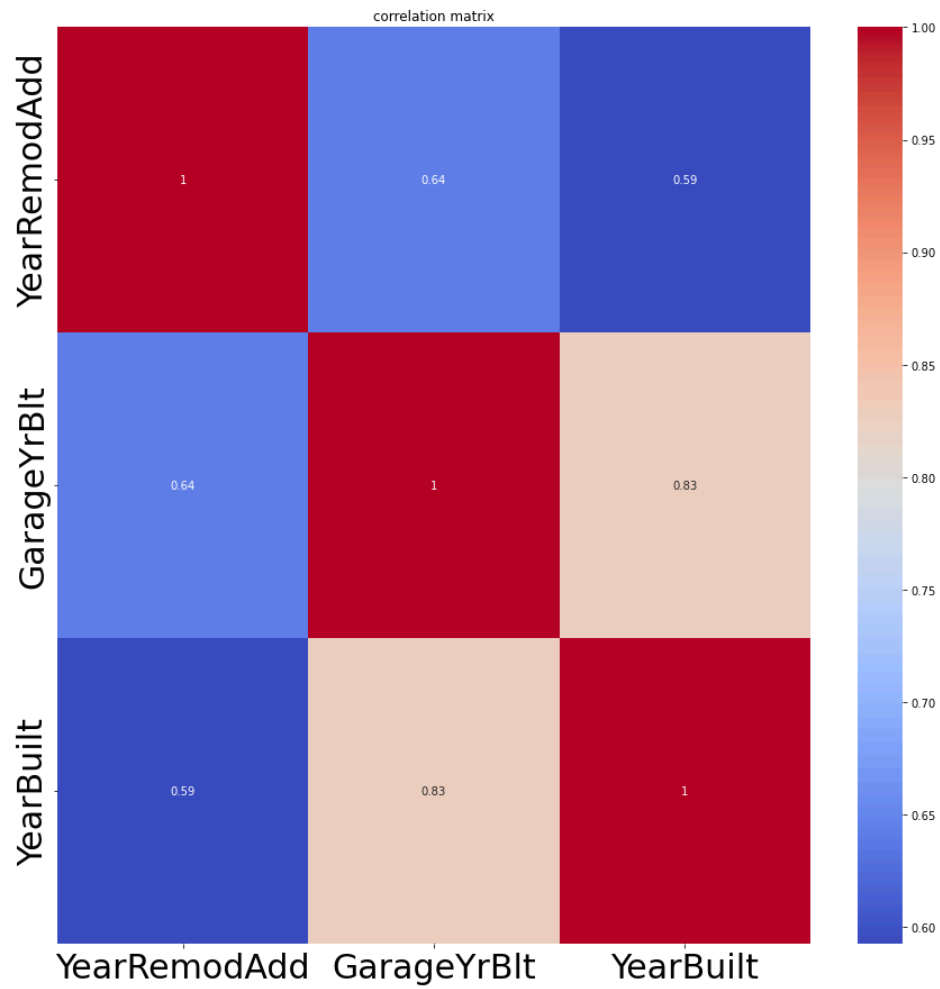Below show the scatter plot of SalePrice plotted as a function of other numeric variables:



**Collinearity:**

Before implementing any model we should examine the collinearity of data. The following is a correlation matrix, showing how correlated different features are:
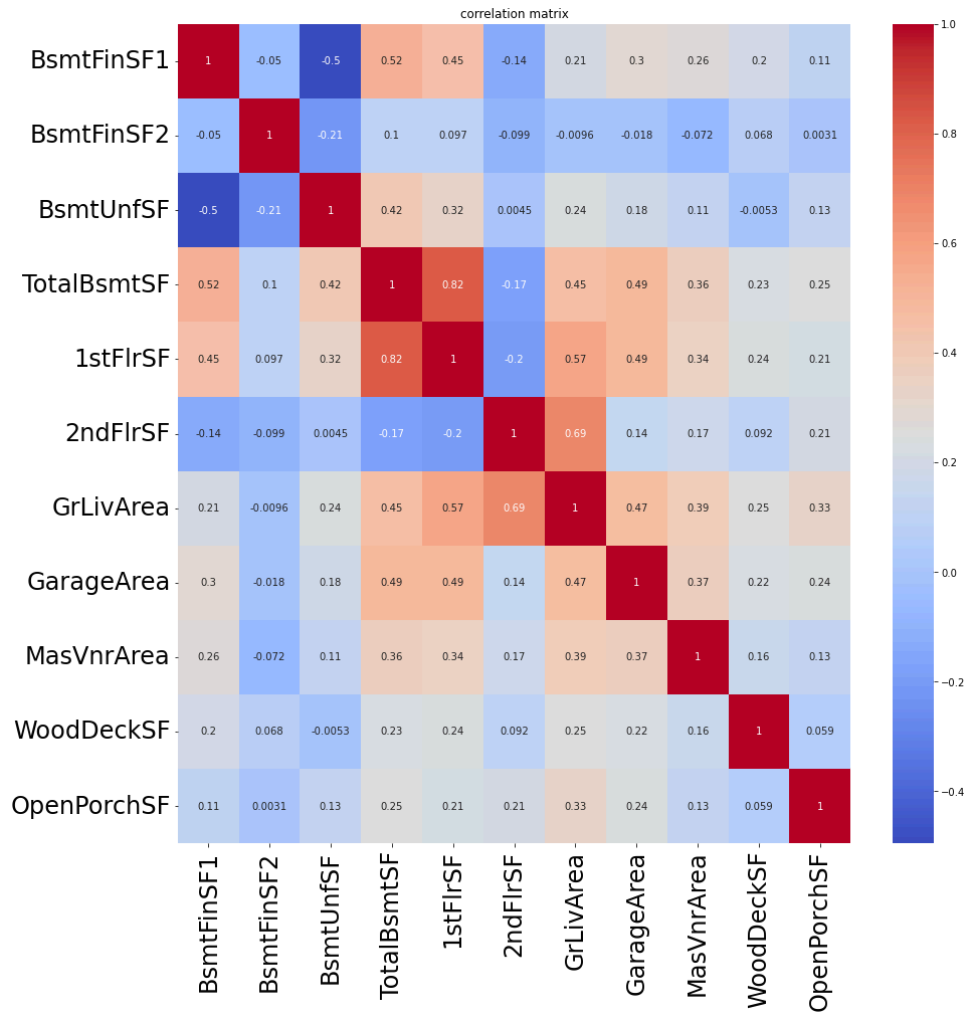
correlation matrix

There are three groups of feature that have high correlation values:
1. Features that include year: ['YearRemodAdd', 'GarageYrBlt', "YearBuilt"]
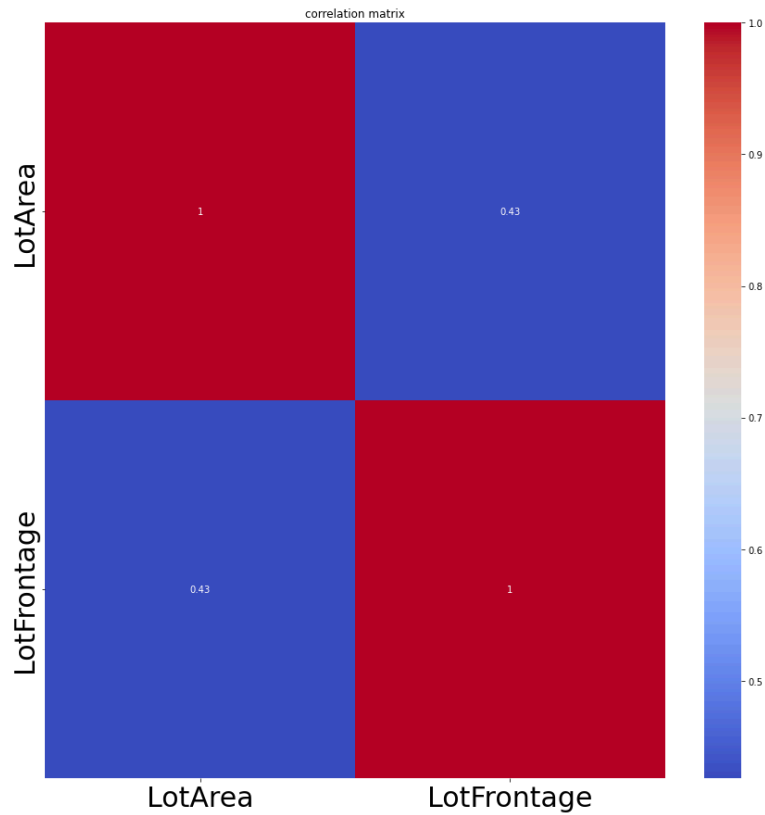
correlation matrix

The above image shows their correlation matrix. The decision was to drop two of these columns and keep only YearBuilt

2. Correlation between features that have units of area. The image below show the correlation matrix for these columns:
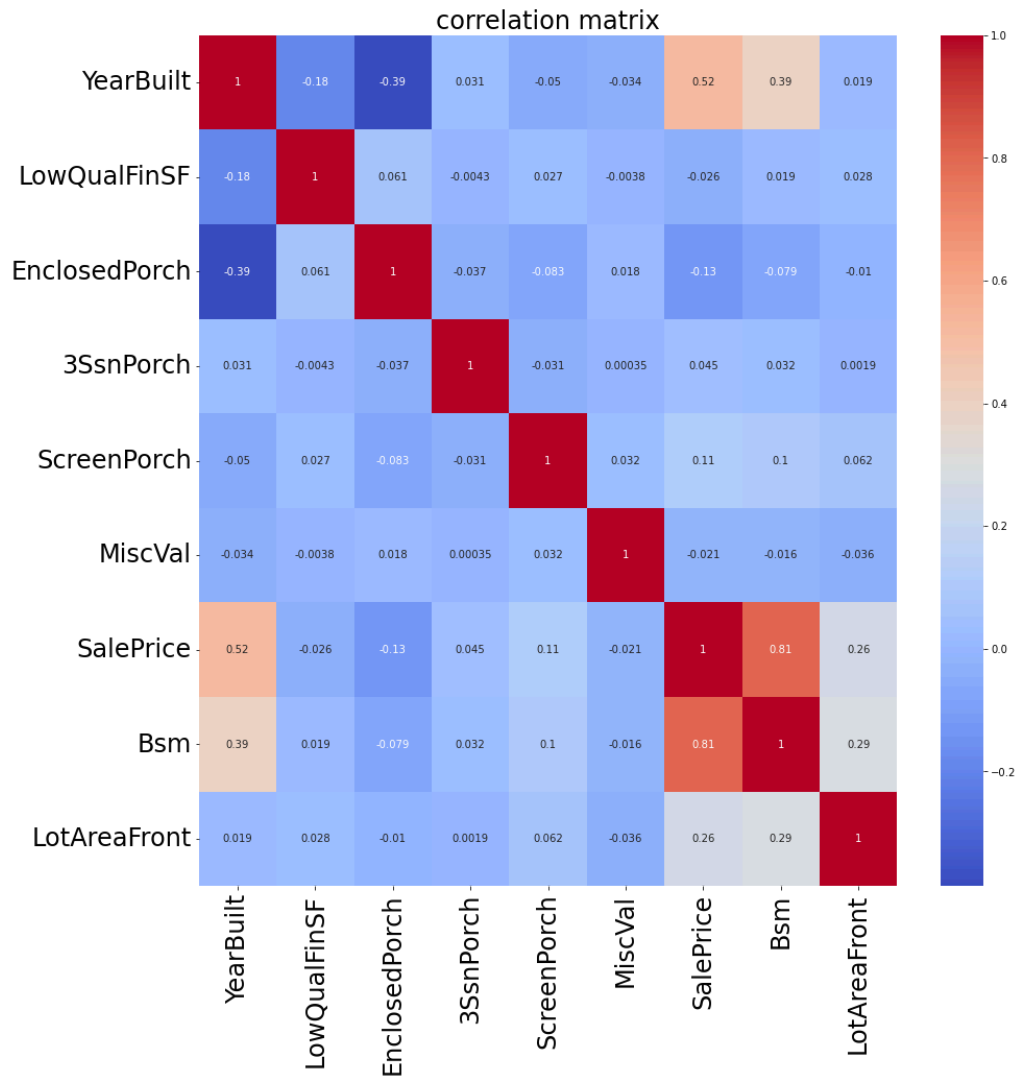
correlation matrix

Since these features all have the same unite, a new feature was created by adding those features and followed by dropping the original columns.

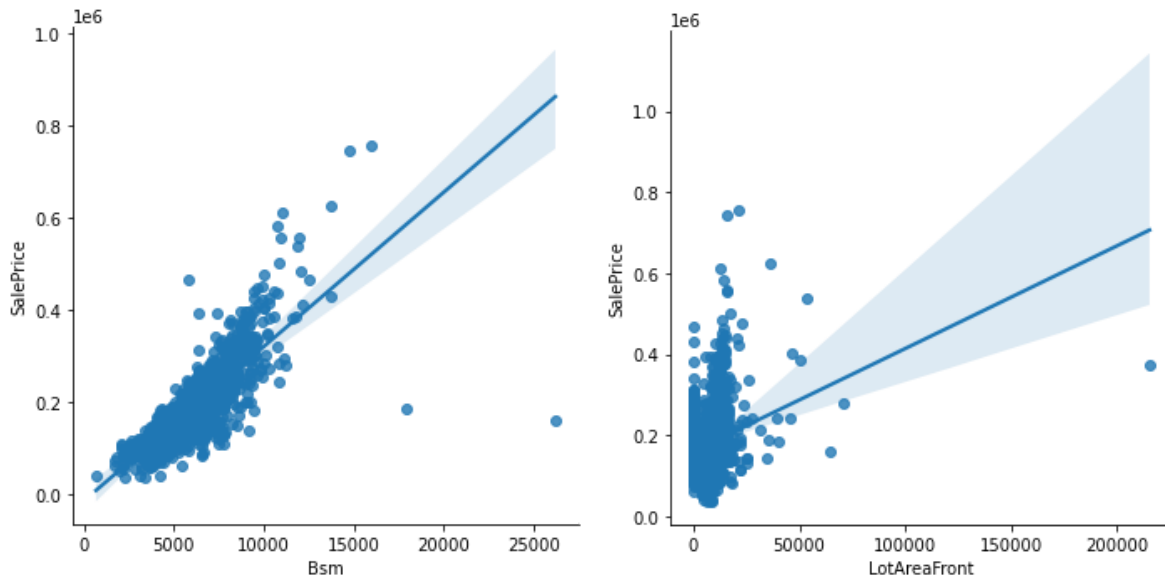3. Correlation between features of type Lot:

correlation matrix

Again the decision was to make a new feature by adding them up followed by dropping the original columns.

Finally the correlation function after feature engineering is plotted in the figure below:

correlation matrix

The following two graphs show how the new feature are correlated with the target variable:

**Dealing with outliers:**

In order the improve the model outliers were removed according to the following criteria:
multiplier = 2.5
Q1 = df[true_numeric_columns].quantile(0.25)
Q3 = df[true_numeric_columns].quantile(0.75)
IQR = Q3 - Q1
lower_bound = Q1 - multiplier * IQR
upper_bound = Q3 + multiplier * IQR

**Data splitting to test and train:**
The data was divided to test and train in order to be able to measure the performance of each model.

**Models:**
Different regression model were implemented and table below compares all of them:

| Model | RMSE |
|---|---|
| Mean | 62834.4 |
| Median | 62113.1 |
| Linear Regression | 10283858416.3 |
| Decision Tree Regression | 31200.8 |
| Decision Tree Regression **Hyperparameter tuning** | 28947.3 |
| **Random Forest Regression (RFR)** | 22059.9 |
| **RFR** hyperparameter tuning | 22245.3 |
| **Gradient Boosting Regression (GBR)** | 20177.0 |
| **GBR** hyperparameter tuning | 19536.5 |
| **XGBoost** hyperparameter tuning | 19971.9 |

According to our results GBR and XGBoost are among the best model.