

# PROJECT 2 - AMES HOUSING DATA AND KAGGLE CHALLENGE

creating a regression model to determine sale price of a house  
in the town of Ames, Iowa

# Background

- Originally published in Journal of Statistics Education, Volume 19, Number 3(2011)<sup>1</sup>
- Ames, Iowa: Alternative to the Boston Housing Data as an End of Semester Regression Project  
Dean De Cock  
Truman State University  
Journal of Statistics Education Volume 19, Number 3(2011), [www.amstat.org/publications/jse/v19n3/decock.pdf](http://www.amstat.org/publications/jse/v19n3/decock.pdf)
- Released as a Kaggle challenge 10 months ago and had participation from 90 teams

# Project objective

- Creating and iteratively refining a regression model

Process:

Regression models covered in class only to be used

Upload prediction values to Kaggle to get a score



# EDA and Data cleaning

- Checking for null values and replacement with valid values if necessary
- Conversion of all ordinal and nominal data types into numerical form
- Dropping of outliers

*df\_train\_alt: Replacing null values*

Column 'Lot Frontage' has a mean 69.06 and median 68 being values close to each other.

This indicates that the distribution is almost Normal and hence the mean value is used for the null values

330 null values in 'Lot Frontage' were filled with mean value 69.055 ¶

```
# df_train_alt: Replacing null values in 'Lot Frontage' with mean value
```

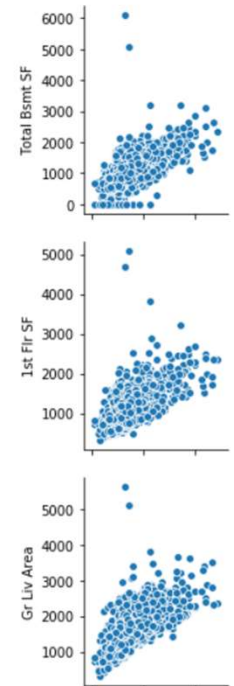
```
df_train_alt['Lot Frontage'] = df_train_alt['Lot Frontage'].fillna(df_train_alt['Lot Frontage'].mean())  
df_train_alt['Lot Frontage'].isnull().sum()
```

```
# changing ordinal ratings to rankings - 'Kitchen Qual'
```

```
mapping = {np.nan: 0, 'Po': 1, 'Fa': 2, 'TA': 3, 'Gd': 4, 'Ex': 5 }  
df_train_alt['Kit_Qual_rk'] = df_train_alt['Kitchen Qual'].apply(lambda x : mapping[x])
```

# Identifying features for the model

- Techniques used include:
  - Viewing heatmap
  - Calculation of Pearson Correlation coefficient
  - Co-linearity between features
  - Filter method
  - Recursive Feature elimination



	Id	PID	MS SubClass	Lot Frontage	Lot Area	Overall Qual	Overall Cond	Year Built	Year Remod/Add	Mas Vnr Area	BsmtFin SF 1	BsmtFin SF 2	Bsmt Unf SF	Total Bsmt SF	1st Flr SF	2nd Flr SF	Low Qual Fin SF	Gr Liv Area	Bsmt Full Bath	Bsmt Half Bath	Full Bath
Mo Sold	0.13	-0.033	0.013	-0.016	0.0032	0.019	-0.0031	-0.0071	0.012	-0.0039	-0.011	-0.014	0.02	0.0039	0.027	0.03	0.024	0.05	-0.0096	0.026	0.049
Yr Sold	-0.98	0.0085	-0.033	0.0084	-0.029	-0.012	0.048	-0.0036	0.043	-0.017	0.038	-1.3e-05	-0.044	-0.0043	-0.0074	-0.012	0.0013	-0.016	0.035	-0.01	0.0071
SalePrice	-0.051	-0.26	-0.087	0.33	0.3	0.8	-0.097	0.57	0.55	0.5	0.42	0.016	0.19	0.63	0.62	0.25	-0.042	0.7	0.28	-0.045	0.54

# Features

## Pearson Correlation coefficients

SalePrice	1.000000
Overall Qual	0.800207
Ex Qual rk	0.712146
Gr Liv Area	0.697038
Kit Qual rk	0.692336
Garage Area	0.649897
Garage Cars	0.647781
Total Bsmt SF	0.629303
1st Flr SF	0.618486
Bsmt_Qual_rk	0.612188
Year Built	0.571849
Gar_Fin_rk	0.557839
Year Remod/Add	0.550370
Firepl_Qu_rk	0.538925
Full Bath	0.537969
FoundationPConc	0.529047
TotRms AbvGrd	0.504014
Mas Vnr Area	0.503579

← **Target variable**

← **Co-linearity between Gr Liv Area and TotRms AbvGrd**  
1 : 0.81

← **Co-linearity between Garage Cars and Garage Area**  
1 : 0.89

# Model designs

## Model 1 - correlation (17 features)

Linear  
Regression

Ridge  
Regression

Lasso  
Regression

Overall Qual  
Year Built  
Year Remod/Add',  
Mas Vnr Area  
Ex\_Qual\_rk  
Bsmt\_Qual\_rk  
Total Bsmt SF  
1st Flr SF  
Gr Liv Area

Full Bath  
Kit\_Qual\_rk  
TotRms AbvGrd  
Firepl\_Qu\_rk  
Garage Yr Blt  
Gar\_Fin\_rk  
Garage Area  
FoundationPConc

## Model 2 - correlation (16 features)

Linear  
Regression

Ridge  
Regression

Lasso  
Regression

Overall Qual  
Year Built  
Year Remod/Add',  
Mas Vnr Area  
Ex\_Qual\_rk  
Bsmt\_Qual\_rk  
Total Bsmt SF  
1st Flr SF  
Gr Liv Area

Full Bath  
Kit\_Qual\_rk  
TotRms AbvGrd  
Firepl\_Qu\_rk  
Garage Yr Blt  
Gar\_Fin\_rk  
Garage Area  
FoundationPConc

## Model 3 - RFE (13 features)

Linear  
Regression

Ridge  
Regression

Lasso  
Regression

MS SubClass60  
MS SubClass75  
MS SubClass120  
MS SubClass150  
Bldg TypeTwnhs  
Roof StyleGable  
Roof StyleGambrel  
Roof StyleMansard

Garage Type2Types  
Garage TypeAttchd  
GarageTypeBuiltIn  
Misc FeatureElev  
Misc FeatureGar2

Polynomial  
(17 features)

# Final model selection

## Model 1 (17 features) Linear Regression model

Kaggle Public score: 33774.73095

Kaggle Private score: 35055.20712

Train R\_square: 0.84745

Test R\_square : 0.85766

RMSE : 29536



## Model 1 (17 features) Ridge Regression model

Kaggle Public score: 33901.40970

Kaggle Private score: 35054.17459

Train R\_square: 0.84754

Test R\_square : 0.85770

RMSE : 29558



## Further refining of the model for better fit

- Consider other features
- Feature Engineering
- Polynomial with degree of 3
- Try other models

