

Project 2  
Ames Iowa  
Housing Data Analysis

# Purpose

Using the Ames Iowa Housing dataset to create a good model for predicting the price of homes at sale

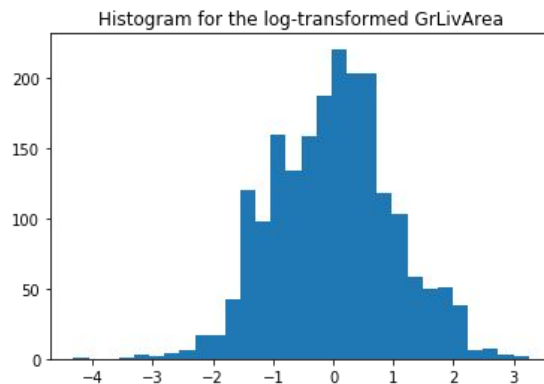
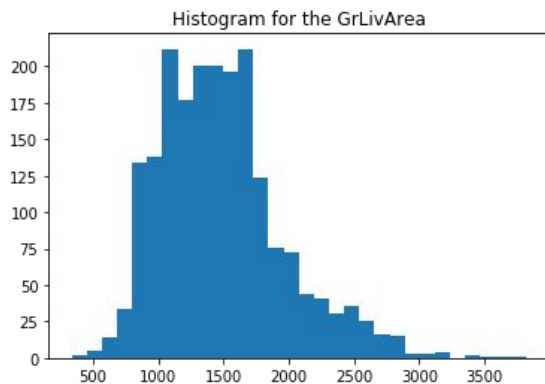
# Feature Engineering

For trying various features, 3 datasets were created.

#1. Training data including some labeled or dummified categorical data

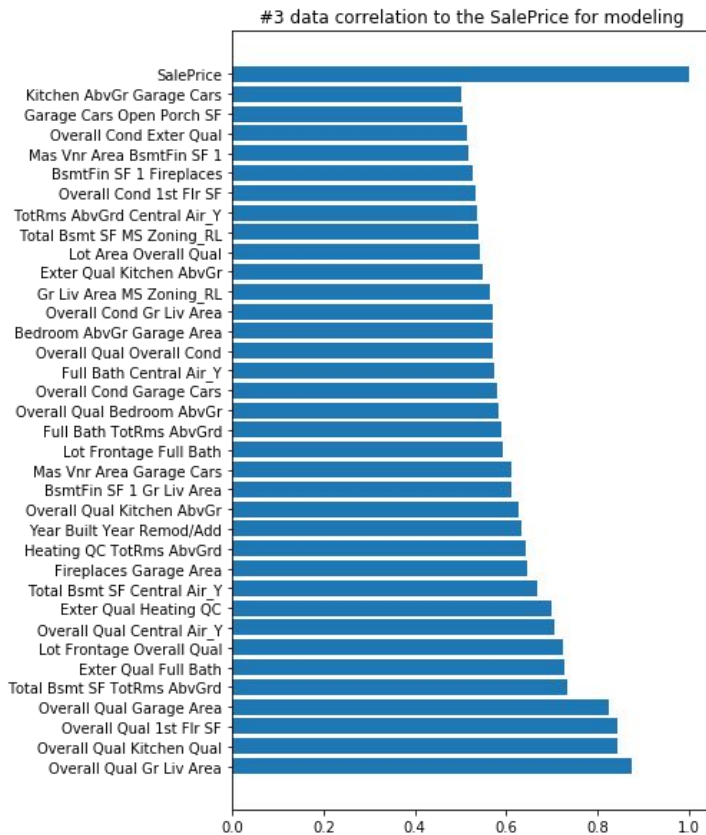
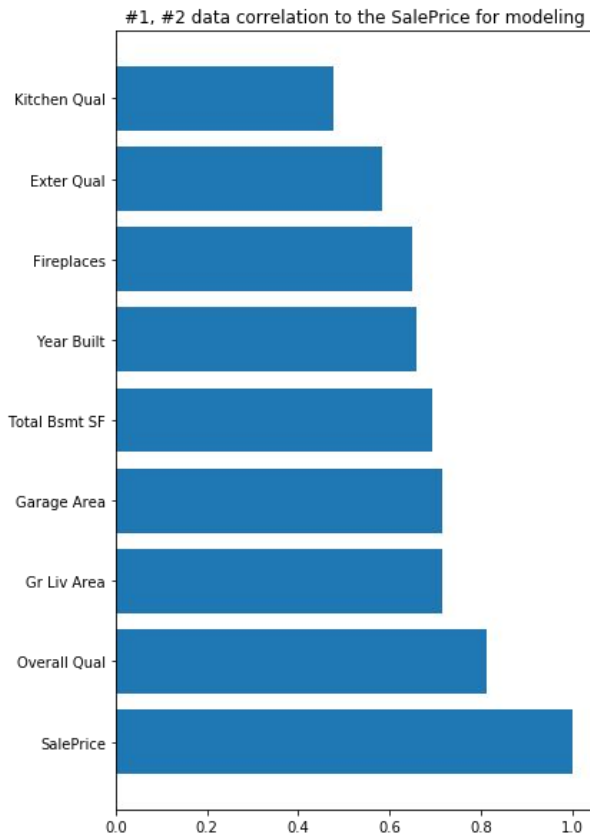
#2. Log-transformed Training data using #1 data

#3. Polynomial or interactive numeric training data using #1 data



# Features for modeling

- These features were used for modeling
- They have more than 0.5 correlation coefficients, but are not highly correlated (less than 0.8) each other



# Modeling

Created 3 models (linear regression, ridge regression, and Lasso regression) for each training dataset

	Linear Regression	Ridge Regression	Lasso Regression
#1 training data	0.8594	0.8593	0.8594
#2 training data	0.8567	0.8572	Failed to make a model
#3 training data	0.9119	0.9128	0.9130

# Modeling

Created 3 models for each training dataset

**Production model = the Lasso Regression using #3 data**

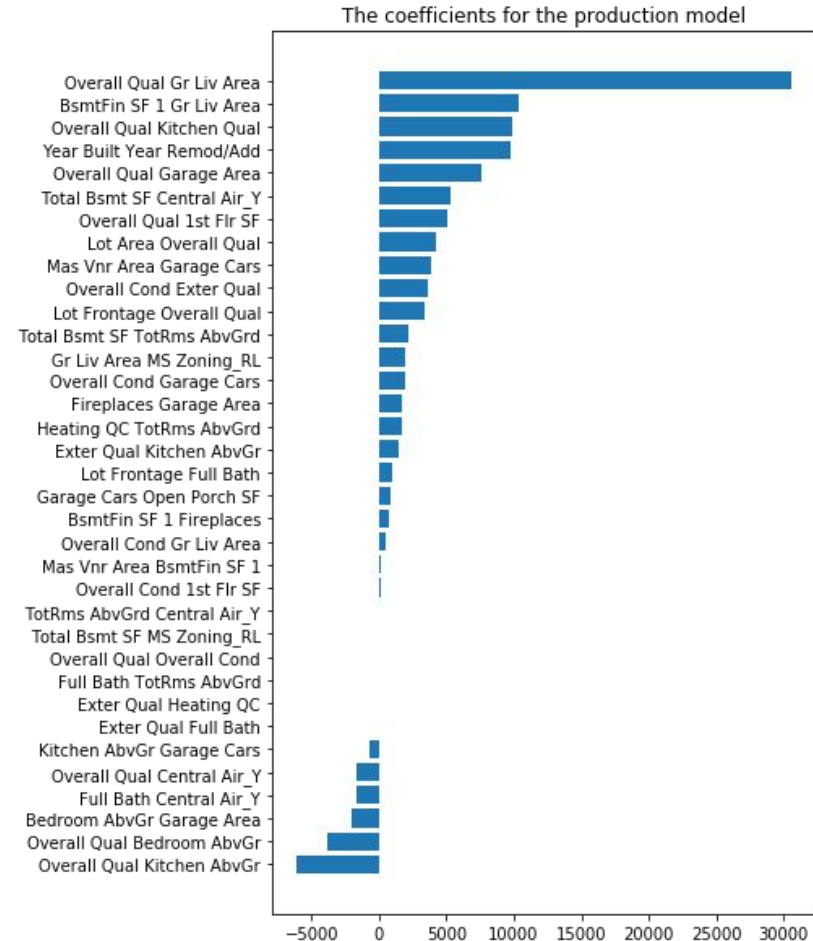
	Linear Regression	Ridge Regression	Lasso Regression
#1 training data	0.8594	0.8593	0.8594
#2 training data	0.8567	0.8572	Failed to make a model
#3 training data	0.9119	0.9128	0.9130



# Production model

- $R^2$  score = 0.91, alpha = 117
- The coefficients show that the 'Overall Qual', the 'Gr Liv Area', the 'Garage Area', etc. have high weights for the sale price

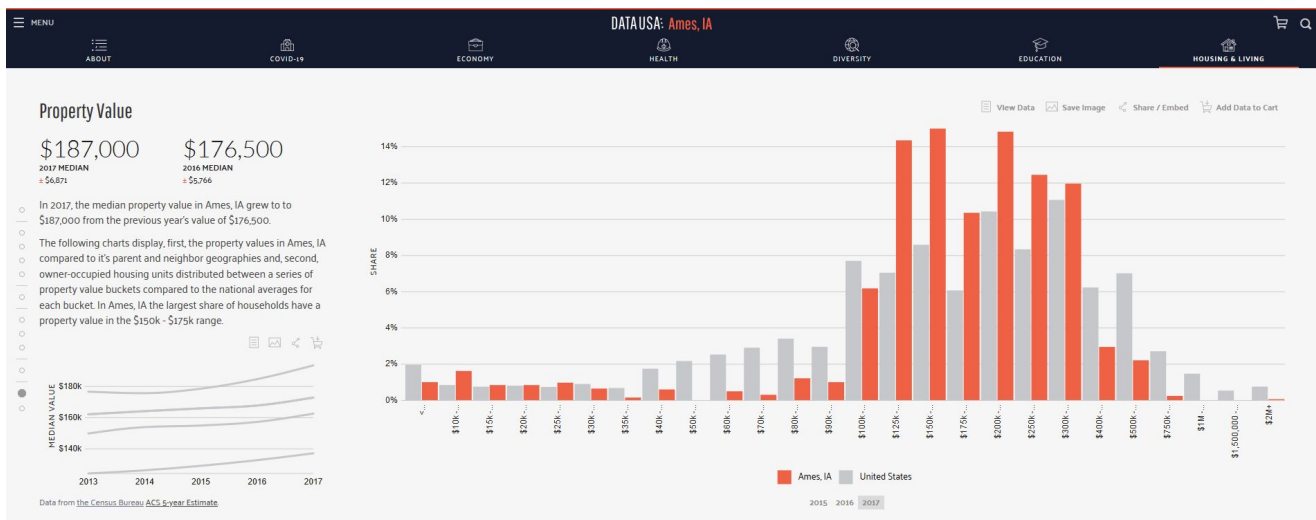
⇒ The quality and size are important to determine the house sale price



# Model constraint

- The production model would be able to be applied to other city as a similar city as Ames
- If the model would be applied to other city like New York, the prediction would be lower

From DATAUSA( \*1)

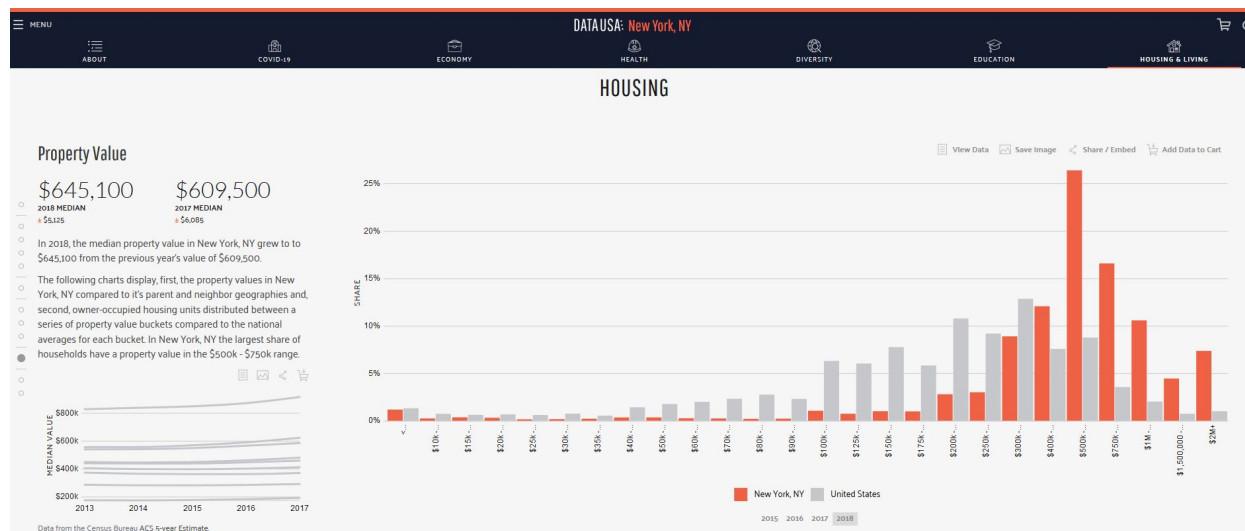




# Model constraint

- The production model would be able to be applied to other city as a similar city as Ames
- If the model would be applied to other city like New York, the prediction would be lower

From DATAUSA( \*2)



# Model constraint

- The production model would be able to be applied to other city similar to Ames
  - If the model would be applied to other city like New York, the prediction would be lower
- ⇒ Models are basing on a training dataset
- ⇒ If we want to make a generalized model, we need data from other cities

# Improve the production model

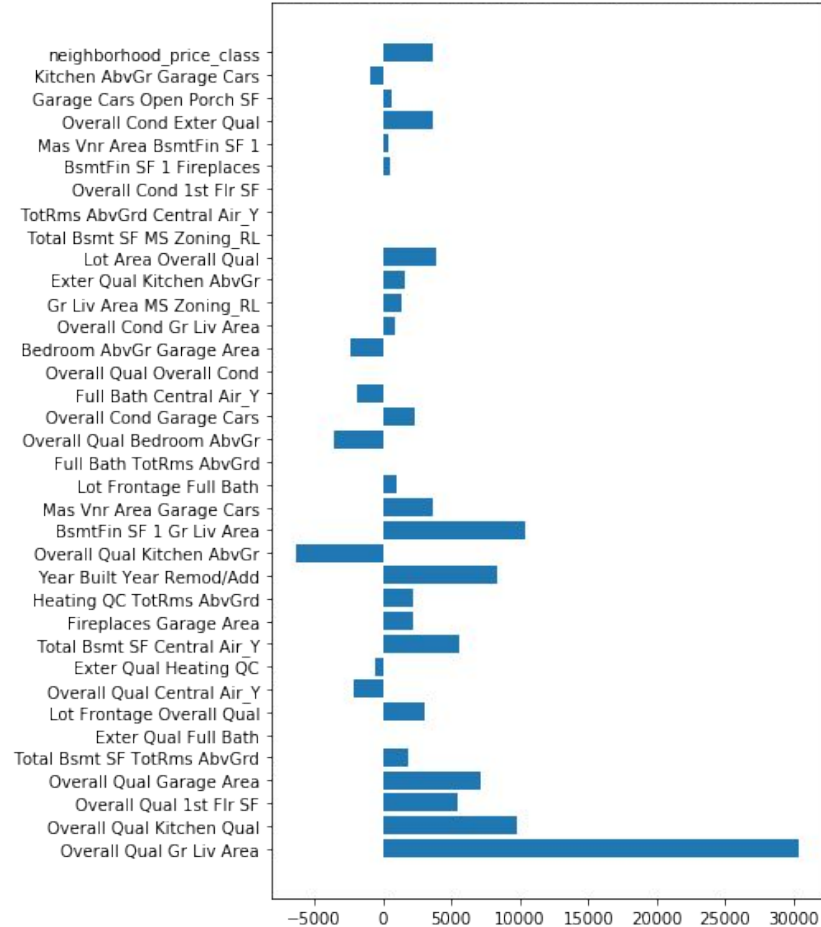
- According to a paper (\* 2), neighborhood information affects house sale price
- The model doesn't have neighborhood information
- Trained the model again adding the neighborhood information labeled into 3 classes basing on the sale price

⇒ The  $R^2$  score = 0.91 the score didn't improve!

⇒ But the coefficient has the power to explain the sale price

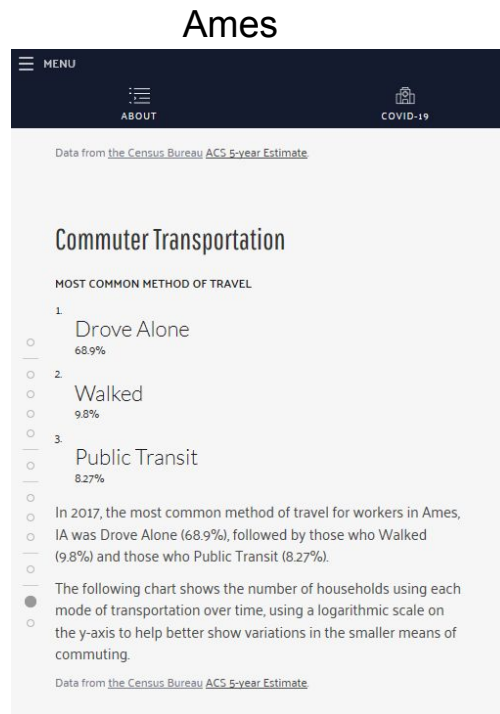
⇒ There might be a better feature for the neighborhood

The coefficients adding the neighborhood information

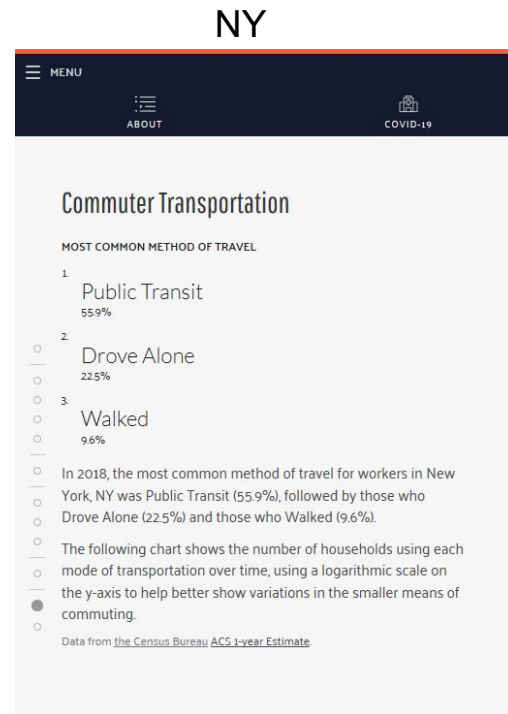


# Suggestion of other neighborhood features

- The distance to a public transportation might affects the sale price



From DATAUSA( \*1)

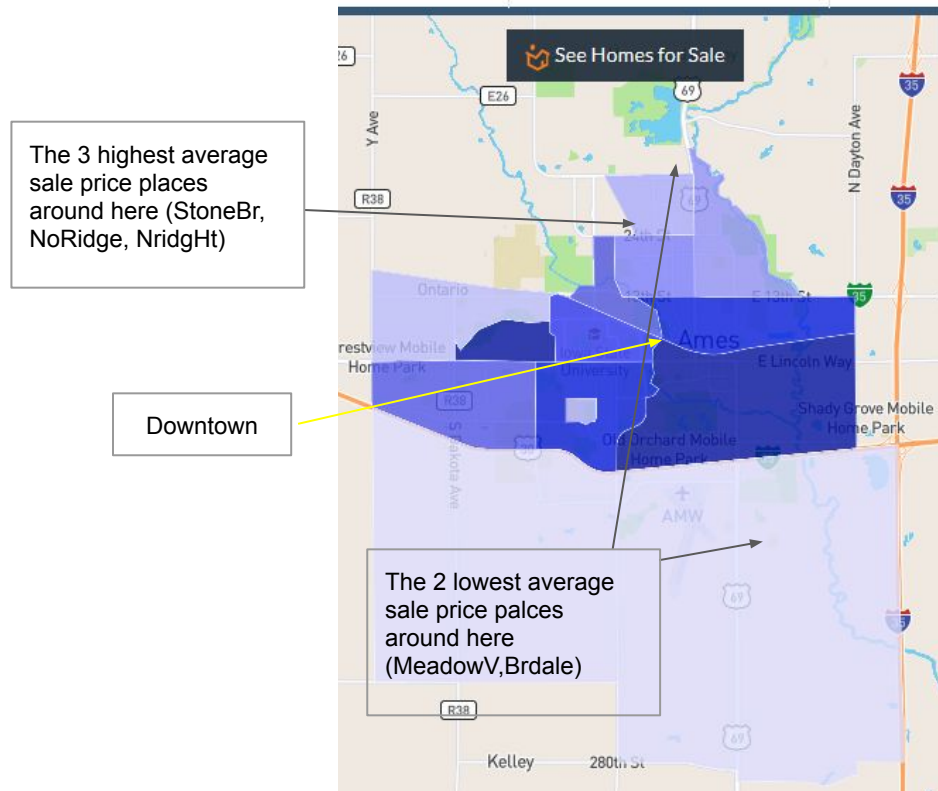


From DATAUSA( \*2)

# Suggestion of other neighborhood features

- The distance to a public transportation might affect the sale price
- Crime rate might affect the sale price
- The distance to downtown might affect the sale price
- Collecting these information might improve the model

Ames crime rate map ( \*4)



# Sammary

- The model is good to predict the sale price in Ames since the  $R^2$  score is about 0.91
- This model can be applied to other cities similar to Ames
- If this model would be applied to other places like New York, the model should be trained for the place's data since their city characteristics are different
- Neighborhood information should be important for a better model to predict house prices
- However neighborhood data this time didn't improve the score but can explain the model
- It would be better to gather more neighborhood data such as distance to a station, crime rates, etc. to create a better model

# Reference

- (\*1) <https://datausa.io/profile/geo/ames-ia/#housing>
- (\*2) <https://datausa.io/profile/geo/new-york-ny/#housing>
- (\*3) [https://www.researchgate.net/publication/304597534\\_THE\\_IMPACT\\_OF\\_NEIGHBORHOOD\\_CHARACTERISTICS\\_ON\\_HOUSING\\_PRICESAN\\_APPLICATION\\_OF\\_HIERARCHICAL\\_LINEAR\\_MODELING](https://www.researchgate.net/publication/304597534_THE_IMPACT_OF_NEIGHBORHOOD_CHARACTERISTICS_ON_HOUSING_PRICESAN_APPLICATION_OF_HIERARCHICAL_LINEAR_MODELING)
- (\*4) <https://www.neighborhoodscout.com/ia/ames/crime>