

---

---

# House Price Prediction

— A Capstone Project —  
By Hope Frost

---

---

# When a house goes on the market, how much will it really sell for?



Photo by melodi2 from <https://freeimages.com> FreImages - Photo by Michael & Christa Richert from <https://freeimages.com> FreImages

# Who might care?

Realtors

Homeowners

Buyers

---

# What factors might determine final sale price?

Realtors consider:

- Year Built
- Overall Condition
- Bedrooms & Bath

But what about:

- Basement Height
- Building Materials
- Heating & cooling

# The Data Set:

## Kaggle: House Price Competition

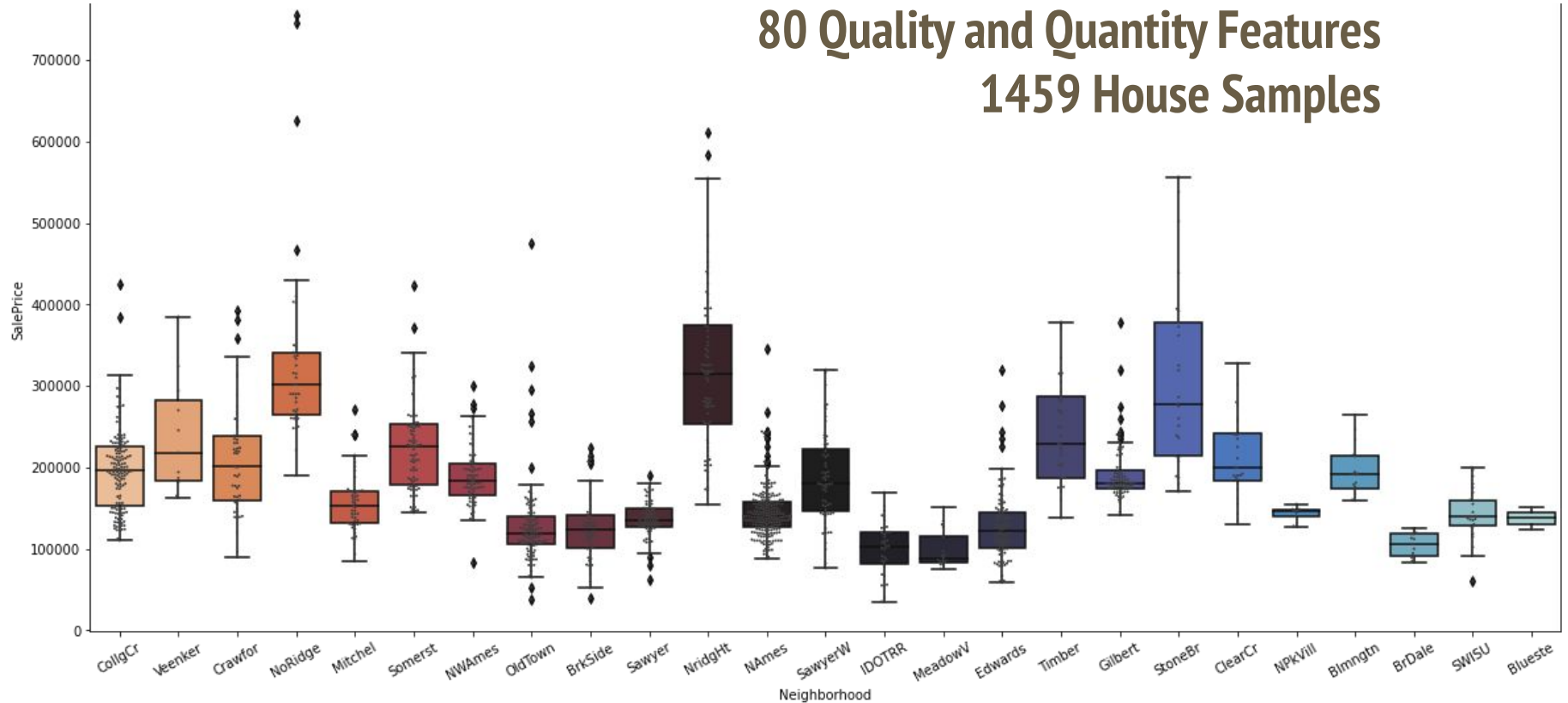
“The Ames Housing dataset was compiled by Dean De Cock for use in data science education”

<https://www.kaggle.com/c/house-prices-advanced-regression-techniques/overview>

# The Competition data contains:

80 Quality and Quantity Features

1459 House Samples



Sales Prices of House in Ames Iowa by Neighborhood

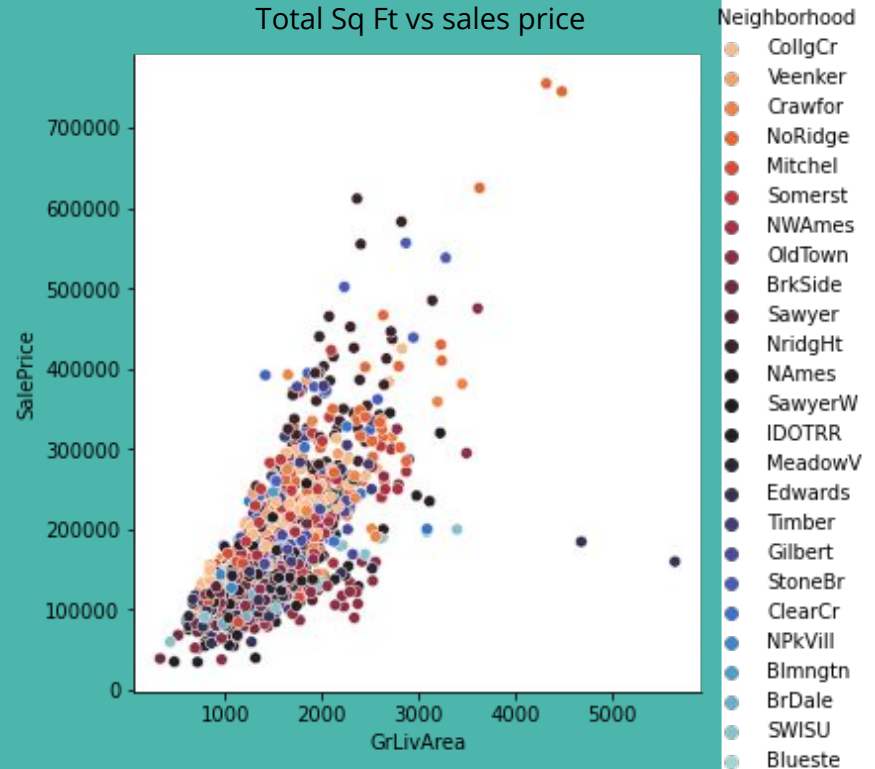
# The Data Process

Impute Missing Data

Convert quality categorical feature to numeric features

Build A few new features

Assess and Drop outliers



# Model Steps

1. One Hot Encode categorical features
2. Split Train and Test sets
3. Scale and Standardize the numeric features

4. Linear Models:  
Linear Regression  
Ridge  
Lasso  
ElasticNet

5. Non-Linear Models:  
Random Forest  
Decision Tree  
Gradient Boost  
SVR  
XGBoost  
LGBMLight



# Model Comparison

	model used	cv_score	MAE	MSE	RMSE	R2 score
<b>model2</b>	Ridge(alpha=1.0, copy_X=True, fit_intercept=Tr...	0.84	19877.8	6.80661e+08	26089.5	0.88
<b>model3</b>	Lasso(alpha=1.0, copy_X=True, fit_intercept=Tr...	0.83	19685	6.65856e+08	25804.2	0.88
<b>model4</b>	ElasticNet(alpha=1.0, copy_X=True, fit_interce...	0.57	34745.7	2.25335e+09	47469.5	0.59
<b>model5</b>	RandomForestRegressor(bootstrap=True, ccp_alph...	0.86	15464.3	5.24048e+08	22892.1	0.9
<b>model6</b>	DecisionTreeRegressor(ccp_alpha=0.0, criterion...	0.77	20760.6	8.88974e+08	29815.7	0.84
<b>model7</b>	GradientBoostingRegressor(alpha=0.9, ccp_alpha...	0.87	13804.9	4.36792e+08	20899.6	0.92
<b>model8</b>	SVR(C=1.0, cache_size=200, coef0=0.0, degree=3...	-0	52181.4	5.29723e+09	72782.1	0.03
<b>model9</b>	GradientBoostingRegressor(alpha=0.9, ccp_alpha...	0.875921	13797.7	4.31391e+08	20770	0.921332
<b>model10</b>	XGBRegressor(base_score=0.5, booster='gbtree',...	0.86	14481.4	4.57588e+08	21391.3	0.92
<b>model11</b>	LGBMRegressor(boosting_type='gbdt', class_weig...	0.87	15216.2	5.09175e+08	22564.9	0.91
<b>model12</b>	<catboost.core.CatBoostRegressor object at 0x7...	0.83	17211.2	5.72056e+08	23917.7	0.9

# Evaluation Metrics

**model9** using a GradientBoostingRegressor is the best **RMSE** score

The final competition score is evaluated on Root-Mean-Squared-Error (RMSE) between the logarithm of the predicted value and the logarithm of the observed sales price.

<b>model9</b>	GradientBoostingRegressor(alpha=0.9, ccp_alpha=0.0, criterion='friedman_mse', init=None, learning_rate=0.01, loss='ls', max_depth=12, max_features='log2', max_leaf_nodes=None, min_impurity_decrease=0.0, min_impurity_split=None, min_samples_leaf=12, min_samples_split=30, min_weight_fraction_leaf=0.0, n_estimators=650, n_iter_no_change=None, presort='deprecated', random_state=648, subsample=1.0, tol=0.0001, validation_fraction=0.1, verbose=0, warm_start=False)	0.875921	13797.7	4.31391e+08	20770	0.921332
---------------	--	----------	---------	-------------	-------	----------

# Predictions and Residuals

## Test Residuals:

STD \$15546 .15

Mean \$13797. 69

50% \$8791.58

## Train Residuals:

STD \$12070.97

Mean \$6824.612

50% \$4281.90

Competition score:

0.14051

Rank at time of submission: 6352



# Next Steps

1. Engineer More Features
2. Correct Skew of Particular Features
3. Explore other models

---

---

# Thank You

To Springboard

And particularly for all the help from:

Silvia Seceleanu

DJ Sarkar

---

---