# AMES

# HOUSING

# PREDICT

DBS-122
Elaine Chen

Purpose: Find the best regression model to predict houses sale price in Ames, IA

Resource:
Data: https://www.kaggle.com/competitions/dsi-122-ames-housing-challenge/data

Google slides template: https://slidesgo.com/

# OVERVIEW
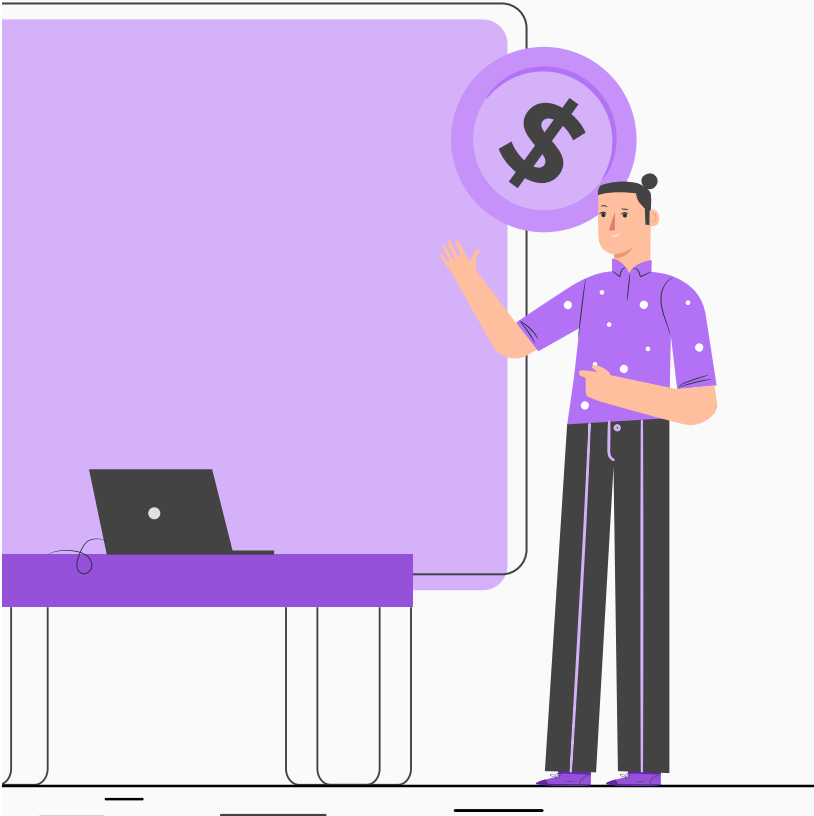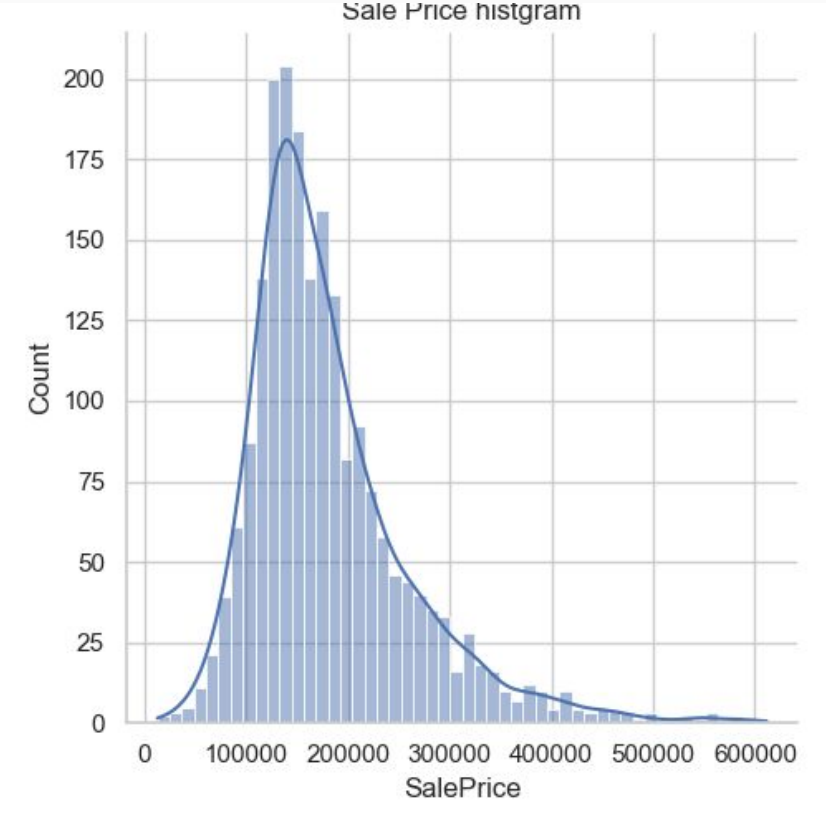
**01**
**INTRODUCTION**

**02**
**DATA CLEANING AND EDA**

**03**
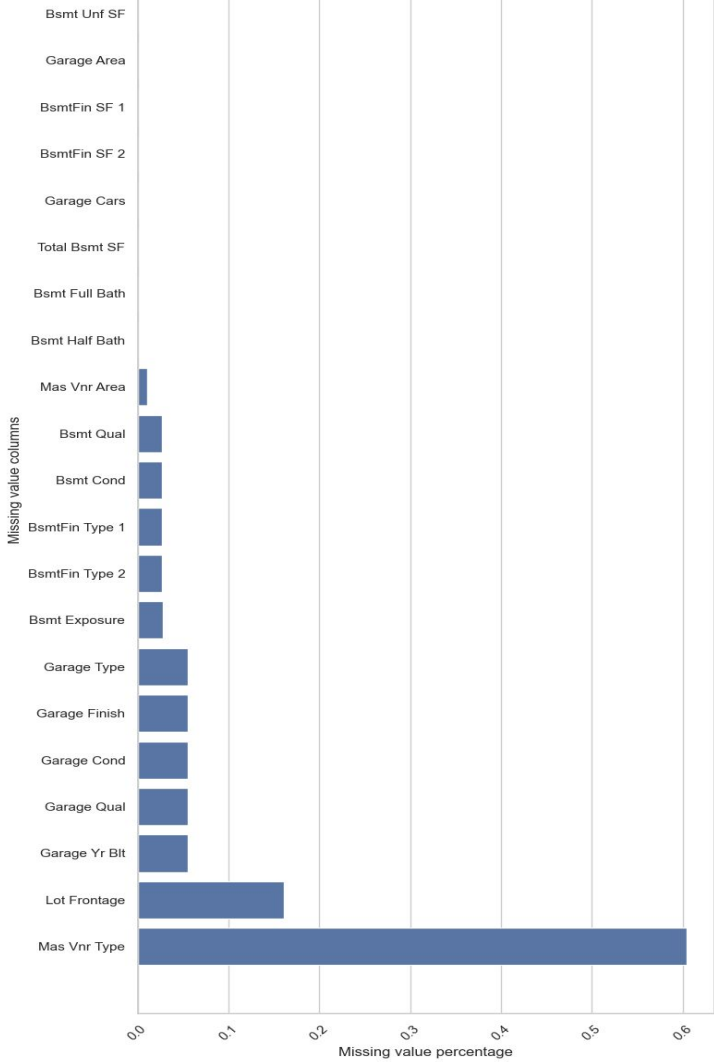**TRAINING TIME AND ASSESSMENT**

**04**
**CONCLUSIONS**

# INTRODUCTION

Training data size: 2,051 rows; 81 columns



Sale Price histgram

Data cleaning:
- Repetitive column: PID
- Missing data
- Outliner

Dropping rows that column missing rate < 0.1%:
- Unfinished basement square feet
- Garage Area
- Type 1 finished square feet
- Quality of second finished area (if present)
- Total basement square feet
- Basement full bathrooms
- Basement half bathrooms

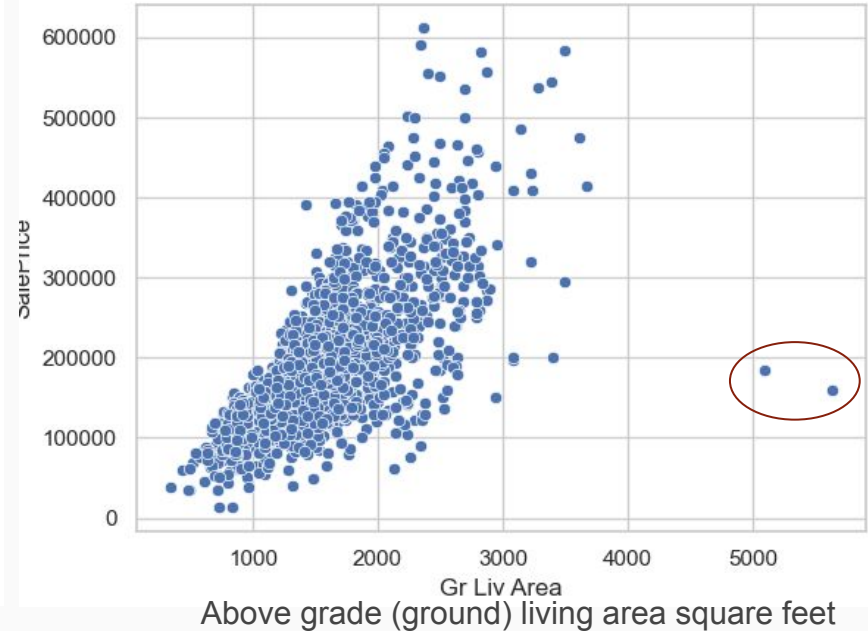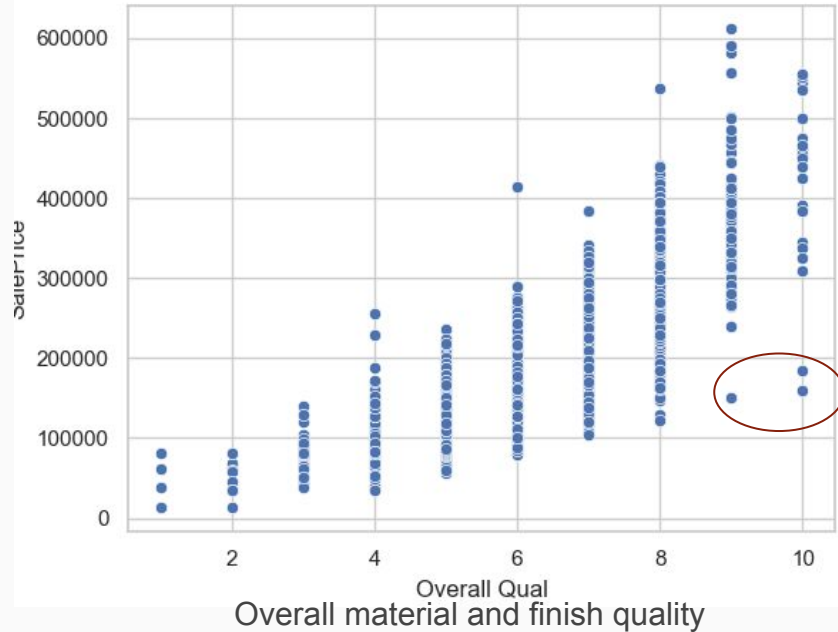Dropping columns which missing rate>80%
- Fence
- Alley
- Miscellaneous feature not covered in other categories
- Pool quality

Filling/Impute missing value columns
- E.g. Using lot frontage by neighborhood mean to impute missing value of lot frontage



Horizontal bar chart — X axis: Missing value percentage (0.0 to 0.6). Y axis: Missing value columns: Bsmt Unf SF, Garage Area, BsmtFin SF 1, BsmtFin SF 2, Garage Cars, Total Bsmt SF, Bsmt Full Bath, Bsmt Half Bath, Mas Vnr Area, Bsmt Qual, Bsmt Cond, BsmtFin Type 1, BsmtFin Type 2, Bsmt Exposure, Garage Type, Garage Finish, Garage Cond, Garage Qual, Garage Yr Blt, Lot Frontage, Mas Vnr Type.

Outliers



Overall material and finish quality



Above grade (ground) living area square feet

# EDA

Select numeric variables that have high correlation with sale price



## Correlation between sale price and other numeric variables

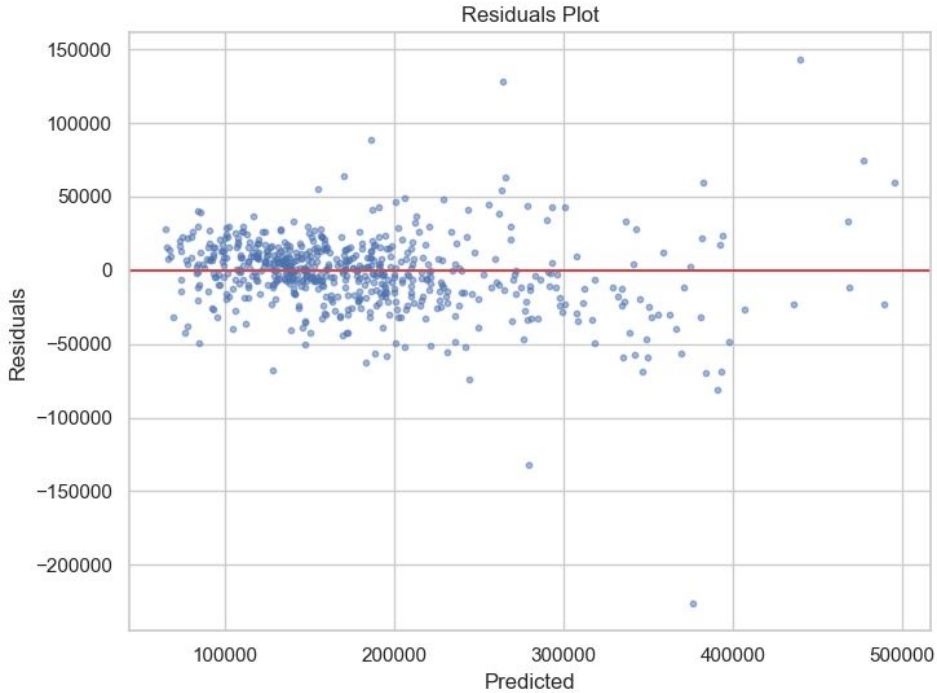| Variable | SalePrice |
|---|---|
| SalePrice | 1 |
| Overall Qual | 0.8 |
| Gr Liv Area | 0.72 |
| Total Bsmt SF | 0.67 |
| Garage Area | 0.66 |
| 1st Flr SF | 0.65 |
| Garage Cars | 0.65 |
| Year Built | 0.57 |
| Year Remod/Add | 0.55 |
| Full Bath | 0.54 |
| Mas Vnr Area | 0.51 |
| TotRms AbvGrd | 0.51 |
| Fireplaces | 0.47 |
| BsmtFin SF 1 | 0.45 |
| Lot Frontage | 0.34 |
| Open Porch SF | 0.34 |
| Wood Deck SF | 0.33 |
| Lot Area | 0.3 |
| Bsmt Full Bath | 0.28 |
| Half Bath | 0.28 |
| Garage Yr Blt | 0.26 |
| 2nd Flr SF | 0.25 |
| Bsmt Unf SF | 0.19 |
| Bedroom AbvGr | 0.14 |
| Screen Porch | 0.13 |
| 3Ssn Porch | 0.049 |
| Mo Sold | 0.032 |
| Pool Area | 0.026 |
| BsmtFin SF 2 | 0.016 |
| Misc Val | -0.01 |
| Yr Sold | -0.015 |
| Low Qual Fin SF | -0.042 |
| Bsmt Half Bath | -0.046 |
| Id | -0.051 |
| MS SubClass | -0.087 |
| Overall Cond | -0.097 |
| Kitchen AbvGr | -0.13 |
| Enclosed Porch | -0.14 |

03.

TRAINING TIME

**NUMERIC VARIABLES**

**CATEGORICAL VARIABLES**

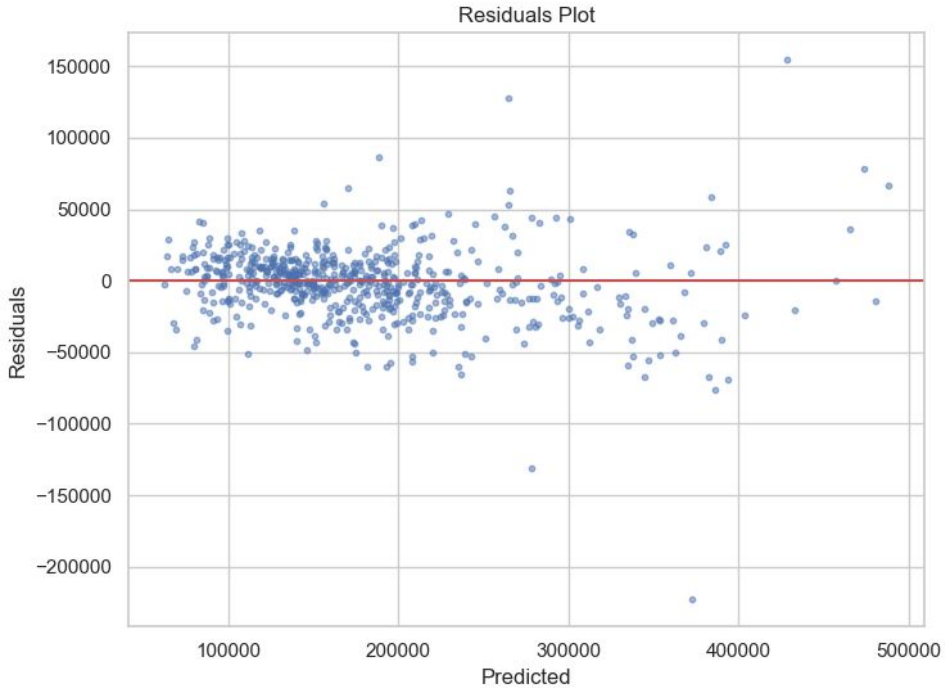1. Selecting numeric variables having >=0.3 (or <= -0.3) correlation with sale price
2. OneHot Encoded categorical features:
   House Style,
   Neighborhood,
   General zoning classification,
   Kitchen quality,
   Roof material,
   Foundation,
3. Add interaction features:
   Garage Cars,
   Overall quality,
   Fireplaces,
   Full bath,
   Total rooms above grade

## OLS COEFFICIENT OF DETERMINATION

### Residuals Plot



| RMSE | 26130.52 |
|---|---|
| Train model R² | 92.99% |
| Test model R² | 88.59% |
| Mean of cross validation score | 86.23% |

Residuals Plot

| RMSE | 26010.26 |
|---|---|
| Train model R² | 92.85% |
| Test model R² | 88.81% |
| Mean of cross validation score | 87.46% |

# CONCLUSIONS

In this Ames house sale price model, LASSO regression has better coefficient of determination over Ordinary Least Squares regression. In training dataset, approximately 92.86% of the variance in the sale price can be explained by the features in the LASSO regression model; 88.81% of sale price can be explained by the features in testing dataset.

# THANKS

Does anyone have any questions?