# Become a Kaggle Master-HW1

**Team ZYMAA**
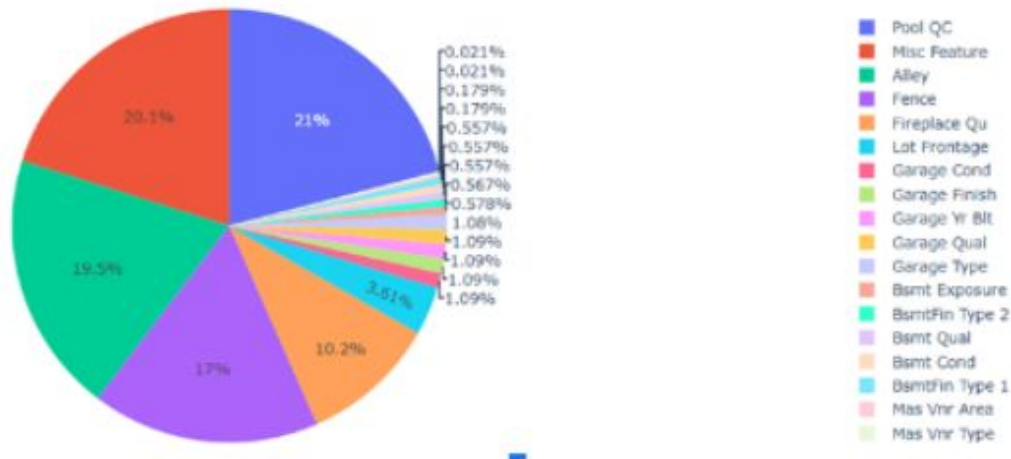
AYARI Mohamed Aziz
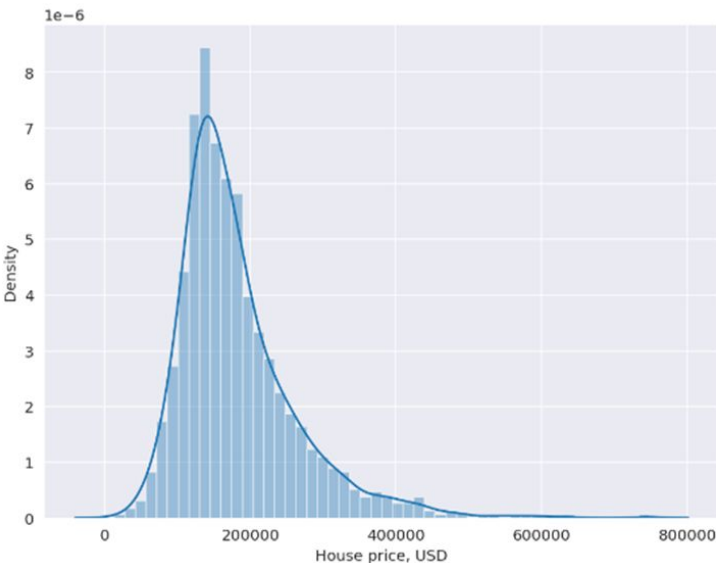YBEGGAZENE Zakaria

# PLAN

# Dataviz

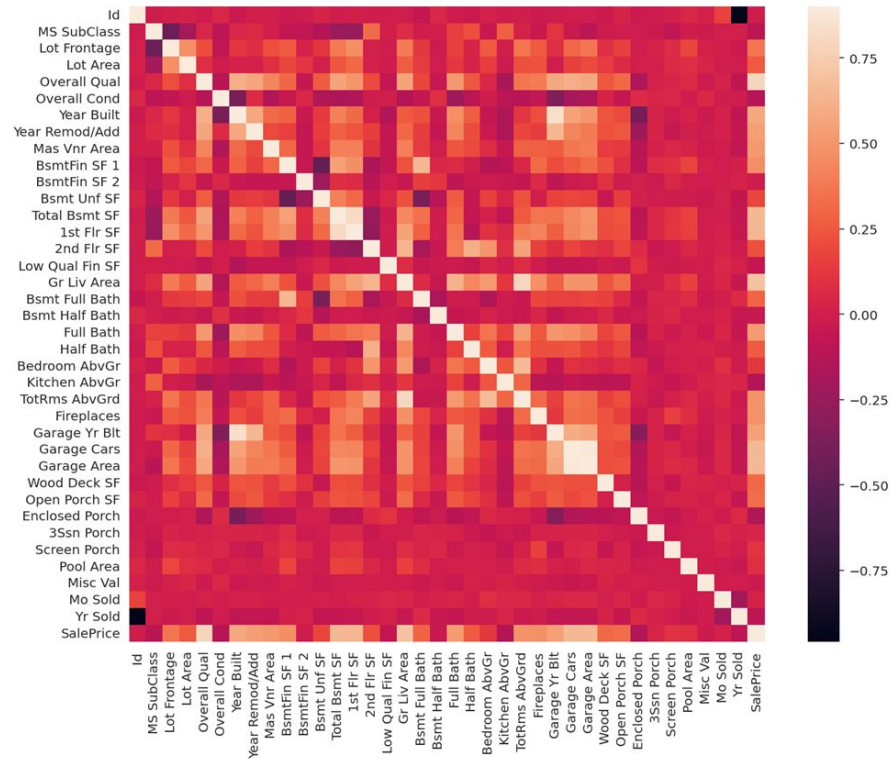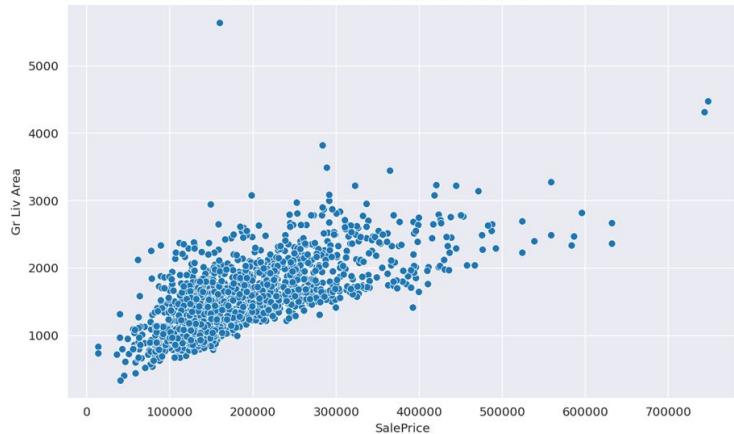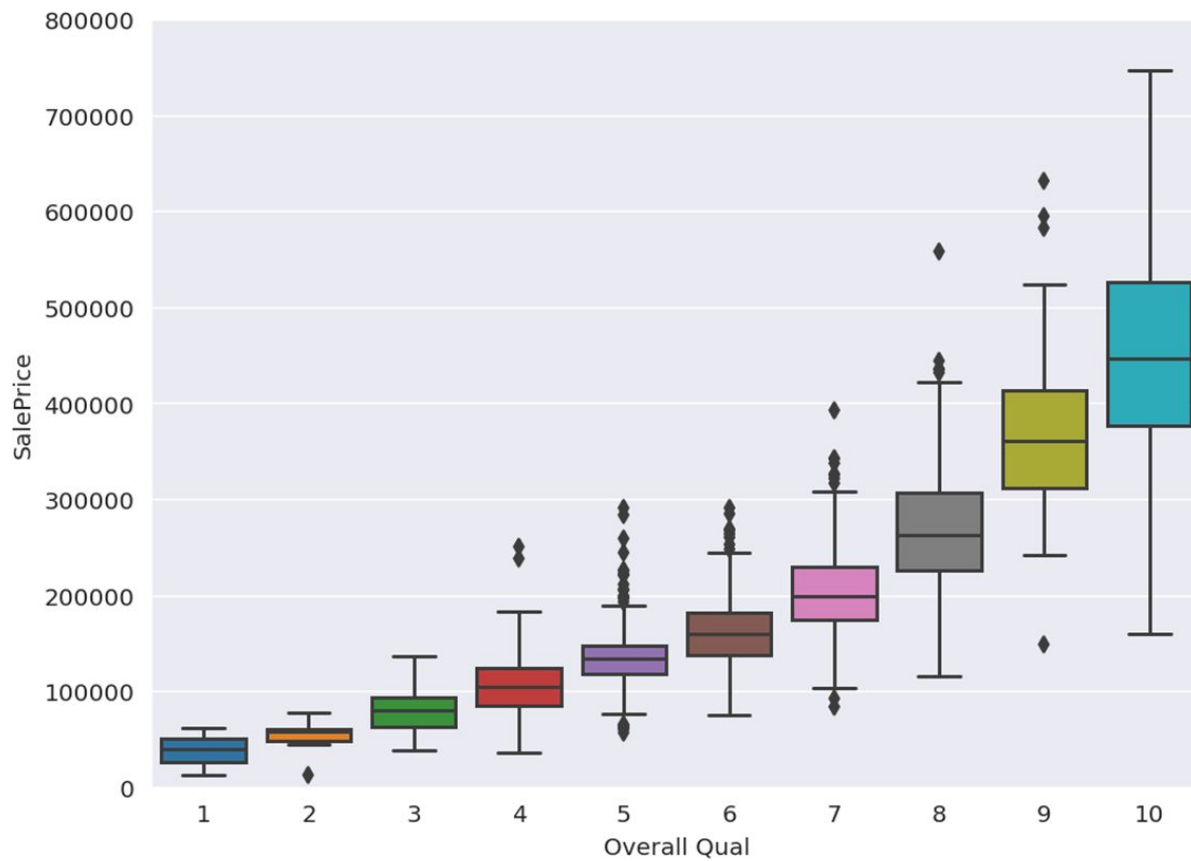# 20 features with highest percentage of missing data



Legend:
- Pool QC
- Misc Feature
- Alley
- Fence
- Fireplace Qu
- Lot Frontage
- Garage Cond
- Garage Finish
- Garage Yr Blt
- Garage Qual
- Garage Type
- Bsmt Exposure
- BsmtFin Type 2
- Bsmt Qual
- Bsmt Cond
- BsmtFin Type 1
- Mas Vnr Area
- Mas Vnr Type

Pie chart values: 21%, 20.1%, 19.5%, 17%, 10.2%, 3.61%, 1.09%, 1.09%, 1.09%, 1.09%, 1.08%, 0.578%, 0.567%, 0.557%, 0.557%, 0.557%, 0.179%, 0.179%, 0.021%, 0.021%

## Price distribution

Correlations between SalePrice and Gr Liv Area

# Data Preprocessing

log transform to the target feature

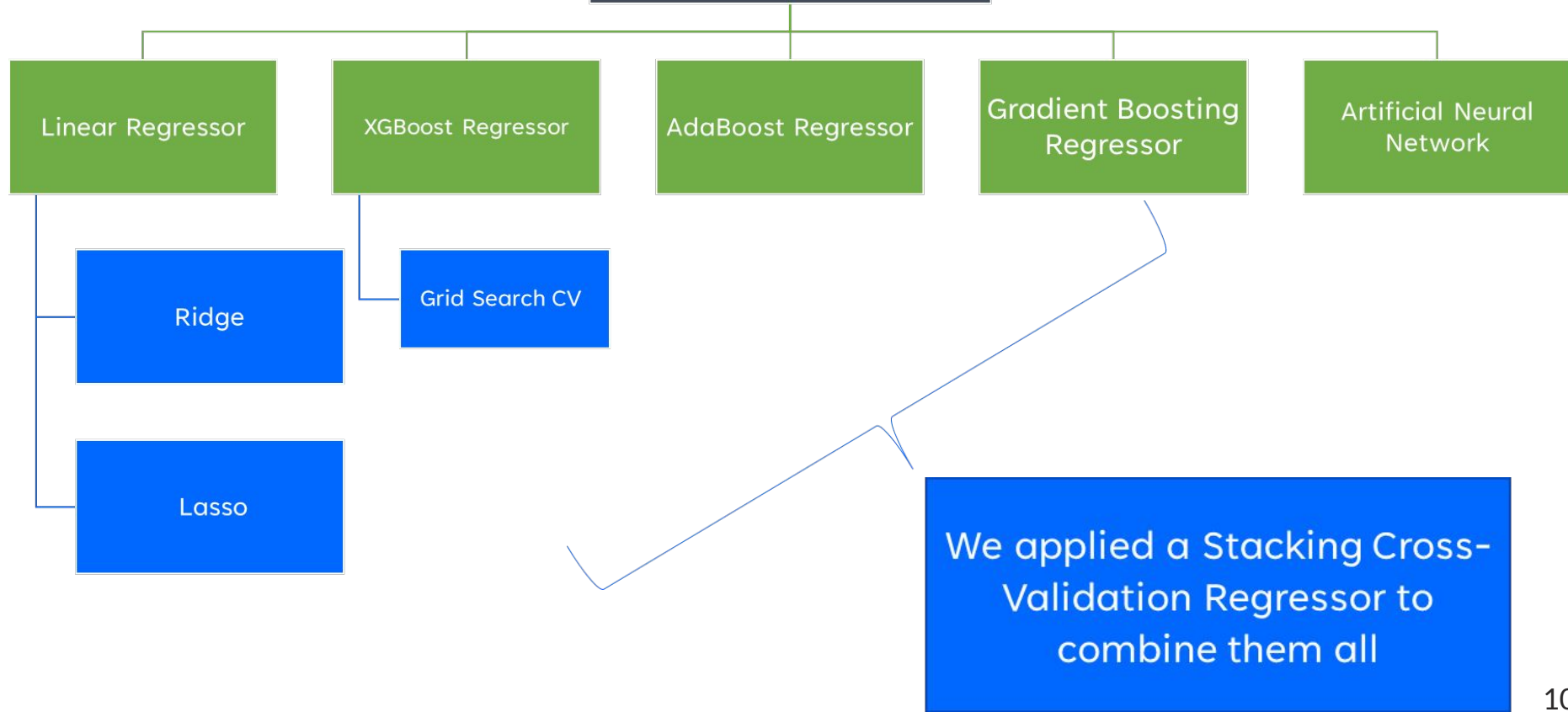Skewness measure to test dataset disparity.

Detect outliers

one-hot encoder get_dummies()

Ordinal Encoder for categorical data
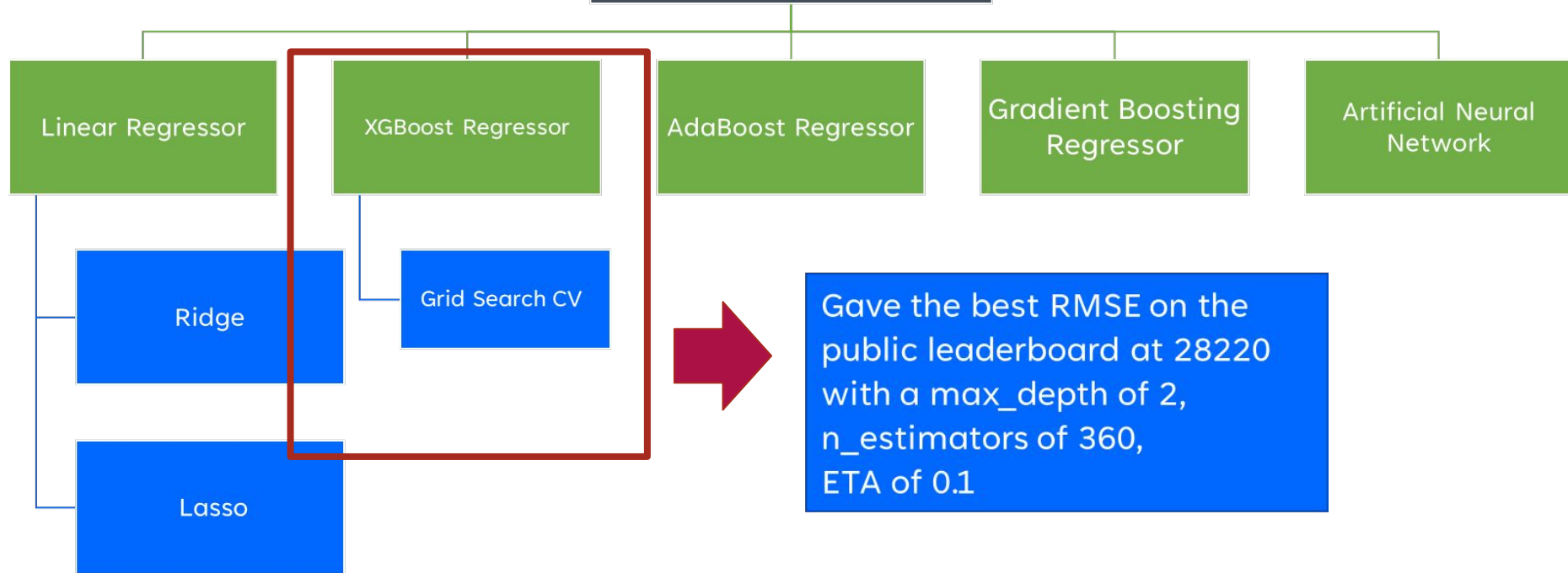
Delete Features with highly missing data

# Models

Tested Models

- Linear Regressor
  - Ridge
  - Lasso
- XGBoost Regressor
  - Grid Search CV
- AdaBoost Regressor
- Gradient Boosting Regressor
- Artificial Neural Network

We applied a Stacking Cross-Validation Regressor to combine them all

# Results

# Conclusion

❖ Exploratory Data Analysis should be given the most time to understand each feature
❖ Data preprocessing plays a substantial role in the performance of any model
❖ Start with simple models and add complexity
❖ Learn to control the bias and variance tradeoff (overfitting vs underfitting)

# Merci pour votre attention

Avez-vous des questions ?