

# *The Principals*

Naren Akurati — [narenakurati@ucla.edu](mailto:narenakurati@ucla.edu)

Ung Lik Teng — [unglikteng@ucla.edu](mailto:unglikteng@ucla.edu)

Alex Tian — [tianxinyuan23@gmail.com](mailto:tianxinyuan23@gmail.com)

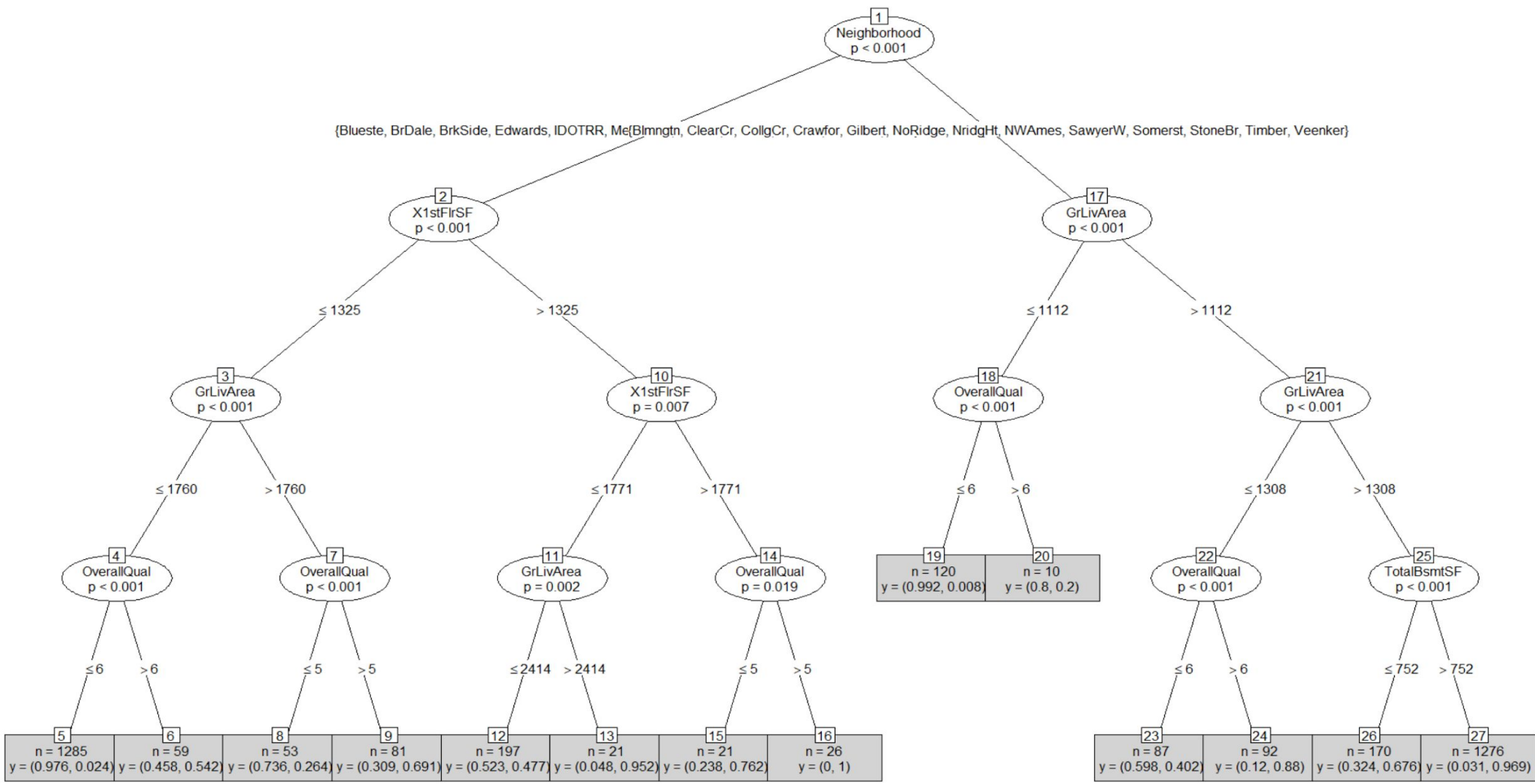
Model

Random Forest

9 features

Public test data accuracy	Private test data accuracy
0.98888	0.97809

Variable Name	Mean Decrease Accuracy
Neighborhood	39.60349
GrLivArea	38.67549
OverallQual	38.07396
TotalBsmtSF	36.62084
X1stFlrSF	35.08231
LotArea	31.96126
BsmtFinSF1	31.47820
X2ndFlrSF	28.62207
BsmtUnfSF	28.22903



# *Procedure*

- Run full model, achieved 97.65% accuracy, but scaled down due to overfitting
- Data exploration
- Remove variables with large number of NAs:
  - *Ob, Street, Alley, Utilities, RoofStyle, BsmtFinSF2, FireplaceQu, PoolQC, Fence*
- Use full model to find top 9 importance features
- Omit remaining NAs in dataset (very few remain)
- Run model with top 9 features
- $mtry=3$  | square root of number of features used

# *Limitations*

- Random forest against logistic, knn, and other methods
  - Random forest robust due to correlation between features
  - Non-parametric: low bias and works well with both linear and non-linear data
- Limitations of random forest
  - Tree correlation potential issue
  - High variance without regularization
  - Easily overfit
  - Need to be aware of highly skewed variables

# *Improvements*

- Further random forest parameter tuning
  - mtry
  - ntree
  - strata
  - nodesize, maxnodes
- Omitting outliers in data
- Use of “CARET” over “randomForest” package
  - Provide grid search function to perform hyperparameter tuning
  - More robust options
- Investigate other methods such as xgboost or using a neural network