

Title

Stuart Miller, Paul Adams, and Chance Robinson

Master of Science in Data Science, Southern Methodist University, USA

1 Introduction

Ramsey and Schafer (2013)

2 Ames, Iowa Data

Kaggle (2016)

3 Analysis Question I

3.1 Question of Interest

Restatement of the problem

3.2 Modeling

TODO: Build and fit the model

We will consider two models: (1) the logarithm of sale price as the response of living room area and (2) the logarithm of sale price as the response of living room area accounting for differences in the three neighborhood of interest (Brookside, Northwest Ames, and Edwards) where Edwards will be used as the reference.

Reduced Model

$$\mu\{\log(\text{SalePrice})\} = \beta_0 + \beta_1(\text{LivingRoomArea}) \quad (1)$$

Full Model

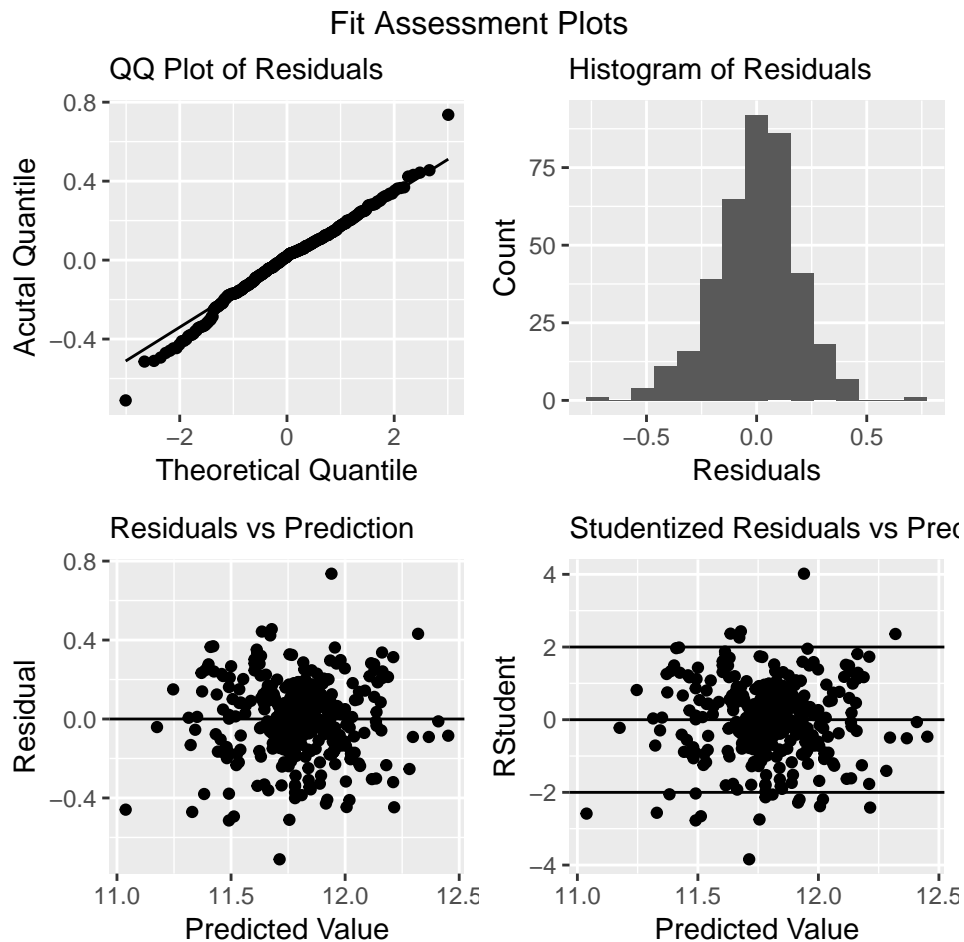
$$\begin{aligned} \mu\{\log(\text{SalePrice})\} = & \beta_0 + \beta_1(\text{LivingRoomArea}) + \beta_2(\text{Brookside}) + \beta_3(\text{NorthwestAmes}) + \\ & \beta_3(\text{Brookside})(\text{LivingRoomArea}) + \beta_4(\text{NorthwestAmes})(\text{LivingRoomArea}) \end{aligned} \quad (2)$$

We will use an extra sums of square test to verify that the interaction terms are useful for the model. The ESS test provides convincing evidence that the interaction terms are useful for the model (p-value < 0.0001); thus, we will continue with the full model.

```
## Analysis of Variance Table
##
## Model 1: log(SalePrice) ~ (GrLivArea) + Neighborhood_BrkSide + Neighborhood_Names
## Model 2: log(SalePrice) ~ (GrLivArea) + Neighborhood_BrkSide + Neighborhood_Names +
##          (GrLivArea) * Neighborhood_BrkSide + (GrLivArea) * Neighborhood_Names
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1      377 14.824
## 2      375 13.441  2    1.3834 19.299 1.053e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

3.3 Model Assumption Assessment

Address each assumption



Model	Adj R^2	CV PRESS
Reduced Model	0.5	2000
Full Model	0.7	1500

3.4 Comparing Competing Models

- Adj R^2
- CV Press

RMSE	CV.Press	Adjusted.R.Squared
0.1910566	12.51675	0.5084024

3.5 Parameters

- Estimates
- Influential points
- Residual plots

3.6 Conclusion

A short summary of the analysis

4 Analysis Question II

4.1 Question of Interest

Restatement of the problem

4.2 Modeling

Type of selection

4.3 Model Assumption Assessment

Address each assumption

4.4 Comparing Competing Models

- Adj R^2
- CV Press
- Kaggle score

4.5 Conclusion

A short summary of the analysis

5 Appendix

Include “well commented” `code` in the appendix!

References

Kaggle (2016). Ames housing dataset. Data retrieved from the Kaggle website, <https://www.kaggle.com/c/house-prices-advanced-regression-techniques/data>.

Ramsey, F. and Schafer, D. (2013). *The Statistical Sleuth: A Course in Methods of Data Analysis*. Brooks/Cole Publishing Company.