

MATH 242 Midterm Project

Jiaxing and Gage

3/02/2024

```
# install.packages("leaps")
# install.packages("glmnet")
# install.packages("DAAG")
# install.packages("neuralnet")
# install.packages("glmnet")
```

```
library(dplyr)
library(readr)
library(ggplot2)
library(neuralnet)
library(glmnet)
library(boot)
library(gridExtra)
library(DAAG)
library(leaps)
```

```
nyc_condos <- read.csv("data/nyc-condos_s24.csv")
nyc_condos_full <- read.csv("data/full_data.csv")
nyc_condos_full <- na.omit(nyc_condos_full)
```

```
# Set the seed
set.seed(123)
# Create an index to randomly sample 70% of the data for training
train_index <- sample(1:nrow(nyc_condos_full), 0.7 * nrow(nyc_condos_full))
# Create the training set
train_data <- nyc_condos_full[train_index, ]
# Create the testing set
test_data <- nyc_condos_full[-train_index, ]
```

```
# summary of dataset
str(nyc_condos)
```

```
## 'data.frame': 200 obs. of 16 variables:
## $ Boro.Block.Lot : chr "1-01613-7501" "1-01171-7501" "3-02237-7519" "4-04955-7512" ...
## $ Condo.Section : chr "0267-R1" "1058-R1" "3457-R1" "0278-R1" ...
## $ Address : chr "1255 5 AVENUE" "200 RIVERSIDE BOULEVARD" "135 MIDDLETON STREET" "1...
## $ Neighborhood : chr "UPPER EAST SIDE (96-110)" "UPPER WEST SIDE (59-79)" "WILLIAMSBURG-...
## $ Building.Classification: chr "R4-CONDOMINIUM" "R4 -ELEVATOR" "R4-ELEVATOR" "R2-CONDOMINIUM" ...
## $ Total.Units : int 59 358 14 4 198 10 60 6 10 20 ...
## $ Year.Built : int 1925 1997 1942 1987 1963 1983 1928 1959 2005 2004 ...
```

```
## $ Gross.SqFt : int 63284 512280 26964 4010 206278 10962 61084 4497 9082 22295 ...
## $ Estimated.Gross.Income : int 1613742 29871047 579187 60391 6266726 392220 742781 92683 242853 72
## $ Gross.Income.per.SqFt : num 25.5 58.3 21.5 15.1 30.4 ...
## $ Estimated.Expense : int 726500 5665817 205466 24782 2044215 162457 417204 37100 73837 21871
## $ Expense.per.SqFt : num 11.48 11.06 7.62 6.18 9.91 ...
## $ Net.Operating.Income : int 887242 24205230 373721 35609 4222511 229763 325577 55583 169016 505
## $ Full.Market.Value : int 6857996 196582995 2914000 239000 32481000 1826000 2048000 437001 13
## $ Market.Value.per.SqFt : num 108.4 383.7 108.1 59.6 157.5 ...
## $ Report.Year : int 2015 2019 2016 2012 2012 2015 2014 2018 2019 2012 ...
```

```
set.seed(250)
## model.matrix() creates our design matrix of predictors
x <- model.matrix(log_Full.Market.Value ~ Total.Units + Estimated.Gross.Income + Gross.Income.per.SqFt)
## select our outcome and convert it into a vector
## instead of a dataframe
y <- train_data %>% select(Full.Market.Value) %>% unlist() %>% as.numeric()
## fit lasso for a range of lambda values (lambda is the tuning parameter
## that controls shrinkage)
cv.out <- cv.glmnet(x, y, alpha = 1)
## pick out the optimal lambda
bestlam <- cv.out$lambda.min
## get coefficients from the Lasso model
lasso <- as.matrix(coef(cv.out, s = bestlam))
t(lasso)
```

```
## (Intercept) Total.Units Estimated.Gross.Income Gross.Income.per.SqFt
## s1 -63438.82 0 7.821129 0
## Estimated.Expense Expense.per.SqFt Net.Operating.Income.per.SqFt Gross.SqFt
## s1 -6.889743 -9544.335 13892.65 -11.65836
```

```
# lm <- lm(Full.Market.Value ~ Year.Built + Gross.SqFt + Gross.Income.per.SqFt + Expense.per.SqFt + log_
# summary(lm)
```

```
lm <- lm(log_Full.Market.Value ~ Gross.SqFt_scaled + Net.Operating.Income.per.SqFt , data = train_data)
summary(lm)
```

```
##
## Call:
## lm(formula = log_Full.Market.Value ~ Gross.SqFt_scaled + Net.Operating.Income.per.SqFt,
## data = train_data)
##
## Residuals:
## Min 1Q Median 3Q Max
## -13.5059 -0.4708 0.0798 0.6007 1.7204
##
## Coefficients:
## Estimate Std. Error t value Pr(>|t|)
## (Intercept) 1.381e+01 1.437e-02 960.7 <2e-16 ***
## Gross.SqFt_scaled 1.839e+01 1.651e-01 111.4 <2e-16 ***
## Net.Operating.Income.per.SqFt 6.592e-02 6.163e-04 107.0 <2e-16 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
```

```
## Residual standard error: 0.8031 on 17732 degrees of freedom
## Multiple R-squared:  0.6141, Adjusted R-squared:  0.6141
## F-statistic: 1.411e+04 on 2 and 17732 DF,  p-value: < 2.2e-16
```

```
predicted_values <- predict(lm, newdata = test_data)

squared_errors <- (predicted_values - test_data$log_Full.Market.Value)^2

MSE <- mean(squared_errors)

MSE
```

```
## [1] 0.6465605
```

Our Model is:

```
Lev <- data.frame(hatvalues(lm)) %>%
ggplot(aes(x = 1:length(hatvalues(lm)), hatvalues(lm))) +
geom_point() +
labs(title = "Leverage",
x = "x", y = "Leverage")

stdresid <- data.frame(rstandard(lm)) %>%
ggplot(aes(x = 1:length(rstandard(lm)), rstandard(lm))) +
geom_point() +
labs(title = "Standardized Residuals",
x = "x", y = "Standardized Residuals")

studresid <- data.frame(rstudent(lm)) %>%
ggplot(aes(x = 1:length(rstudent(lm)), rstudent(lm))) +
geom_point() +
labs(title = "Studentized Residuals",
x = "x", y = "Studentized Residuals")

cooks <- data.frame(cooks.distance(lm)) %>%
ggplot(aes(x = 1:length(cooks.distance(lm)), cooks.distance(lm))) +
geom_point() +
labs(title = "Cooks Distance",
x = "x", y = "Cooks Distance")

p <- length(lm$coeff)
# n <- nrow(train_data)
n <- 22102

Lev <- Lev +
geom_hline(yintercept = 2*p/n, lty = "dashed", col = "orange")
```

```

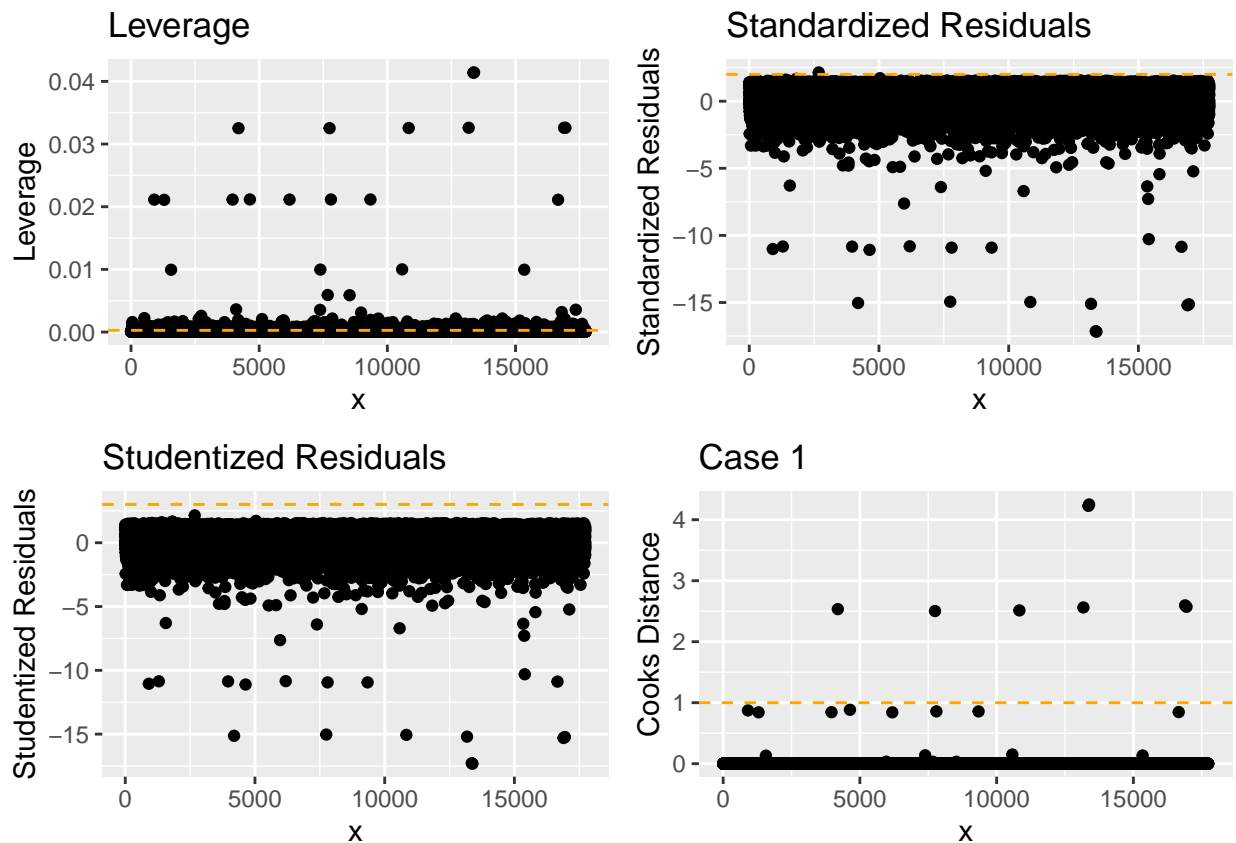
stdresid <- stdresid +
geom_hline(yintercept = 2, lty = "dashed", col = "orange")

studresid <- studresid +
geom_hline(yintercept = 3, lty = "dashed", col = "orange")

cooks <- cooks +
geom_hline(yintercept = 1, lty = "dashed", col = "orange") +
labs(title = "Case 1",
x = "x", y = "Cooks Distance")

grid.arrange(Lev, stdresid, studresid, cooks, ncol = 2)

```



```

# high_leverage <- train_data %>% filter(hatvalues(lm) > 2* p/n)
# high_residuals <- train_data %>% filter(abs(rstandard(lm)) > 2)
# high_influence <- train_data %>% filter(cooks.distance(lm) > 1)
#
#
# print(high_leverage)
# print(high_residuals)
# print(high_influence)]

```

Title:

Abstract

Introduction

Explanation of our variables:

1. **CondoSection:** Identification information for the condominium.
2. **Address:** Street address of the property.
3. **Neighborhood:** Name of the neighborhood where the property is located.
4. **BldgClassification:** Building classification code and description indicating the property's use.
5. **TotalUnits:** Total number of units in the building.
6. **YearBuilt:** Year the building was constructed.
7. **GrossSqFt:** Gross square footage of the building.
8. **EstGrossIncome:** Estimated gross income, calculated as income per square foot multiplied by gross square footage.
9. **GrossIncomePerSqFt:** Estimated gross income per square foot.
10. **EstimatedExpense:** Estimated expense, calculated as expense per square foot multiplied by gross square footage.
11. **ExpensePerSqFt:** Estimated expense per square foot.
12. **NetOperatingIncome:** Net operating income, calculated as estimated gross income minus estimated expense.
13. **FullMarketValue:** Current year's total market value of the property (land and building).
14. **MarketValuePerSqFt:** Market value per square foot, calculated as full market value divided by gross square footage.
15. **ReportYear:** Year of the report.
16. **Boro-Block-Lot:** Borough-Block-Lot location identifier for the property.

Methods

Results

Discussion

Conclusion

References