

HW 4

Evan Cranmer and Xiaoqing Liu

2024-09-25

Examining medv

```
library(MASS)
data("Boston")
?Boston
#Examine medv as a function of crim, zn and indus in a multiple linear regression
medvmodel <- lm(medv ~ crim + zn + indus, data = Boston)
summary(medvmodel)
```

```
##
## Call:
## lm(formula = medv ~ crim + zn + indus, data = Boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -13.070  -4.733  -1.585   2.648  32.423
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  27.39465    0.86484   31.676 < 2e-16 ***
## crim        -0.24863    0.04391   -5.662 2.52e-08 ***
## zn           0.05850    0.01750    3.344 0.000889 ***
## indus       -0.41558    0.06378   -6.515 1.77e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.752 on 502 degrees of freedom
## Multiple R-squared:  0.2937, Adjusted R-squared:  0.2895
## F-statistic: 69.59 on 3 and 502 DF,  p-value: < 2.2e-16
```

1a looking at statistically significant predictors

Crim, zn, and indus are all statistically significant predictors at $\alpha = 0.05$. Their p-values are less than .05.

1b null and alternative hypotheses

Null: The null hypothesis states that the predictors have no effect on variable medv. $H_0: \beta_0 = \beta_1 = \beta_2 = \beta_3 = 0$

Alternative: The alternative hypothesis states that the predictors have a non-zero effect on variable medv.
 $H_a: \beta_0 \neq 0, \beta_1 \neq 0, \beta_2 \neq 0, \beta_3 \neq 0$

1c interpreting the regression coefficients

Crim: For every one unit increase in per capita crime rate, the expected value of median home value decreases by 0.24863 in the \$1000s, holding other variables constant.

zn: For every one unit increase in proportion of residential land zoned for lots over 25,000 sq.ft., the expected value of median home value increases by 0.05850 in the \$1000s, holding other variables constant.

indus: For every one unit increase in proportion of non-retail business acres per town., the expected value of median home value decreases by .41558 in the \$1000s, holding other variables constant.

1d

Interpreting all predictors as the exposure of interest could be problematic because control variables could have interactions with other variables, or just simply not have the equal effects.

1E

```
confint(medvmodel, level = 0.95)
```

```
##                2.5 %       97.5 %  
## (Intercept) 25.6954863 29.09380729  
## crim       -0.3348958 -0.16236084  
## zn         0.0241252  0.09287644  
## indus      -0.5408945 -0.29026116
```

The 95% confidence intervals for the predictors are the following:

crim: (-0.334895, -0.16236084). If we were to take repeated samples of the population and calculate coefficient crim, the true value would be contained in the confidence interval (-0.334895, -0.16236084) 95% of the time.

zn: (0.0241252, 0.09287644). If we were to take repeated samples of the population and calculate coefficient zn, the true value would be contained in the confidence interval (0.0241252, 0.09287644) 95% of the time.

indus: (-0.5408945, -0.29026116). If we were to take repeated samples of the population and calculate coefficient indus, the true value would be contained in the confidence interval (-0.5408945, -0.29026116) 95% of the time.

Since the three variable confidence intervals do not pass through 0, all three are statistically significant and have a non-zero effect. So we would reject the null hypothesis in 1b. In 1a we found that the p-values were less than .05, meaning the variables were statistically significant, and we can reject the null hypothesis that the variables are equal to zero. This corresponded to a confidence interval that did not include zero, which we just found.

1F

$$R^2 = \frac{SS_{Reg}}{SS_Y} = 1 - \frac{RSS}{SS_Y} \text{ and } R^2_{adj} = 1 - \frac{(1 - R^2)(n - 1)}{n - p}$$

```

#storing the observed values for medv
y_medv <- Boston$medv

#storing predicted values from our medv model
y_medv_hat <- predict(medvmodel)

#calculating RSS
rss <- sum((y_medv - y_medv_hat)^2)

#finding SSy
ssy <- sum((y_medv - mean(y_medv))^2)

medv_r2 <- 1 - (rss/ssy)
print(medv_r2)

```

```
## [1] 0.2937136
```

```

##### Now finding R2 adjusted

#storing number of observations
n <- length(y_medv)

#number of predictors not including the intercept
p <- 3

#finding our
medv_r2_adj <- 1 - ((1 - medv_r2)*(n-1) / (n-p-1))
print(medv_r2_adj)

```

```
## [1] 0.2894927
```

Interpreting R squared and Adjusted R squared results

Our R-squared value for this model is approximately .294. Which means that approximately 29.4% of the variance in the dependent variable is explained by the independent variables.

Our adjusted R-squared value for this model is approximately .289. Which means that approximately 28.9% of the variance in the dependent variable is explained by the independent variables, and the value was adjusted for the number of parameters in our model.

2 Fitting a new model

```

#creating new model with zn as only predictor
medvmodel2 <- lm(medv ~ zn, data = Boston)

#performing global f test
anova(medvmodel2, medvmodel)

```

```
## Analysis of Variance Table
##
```

```
## Model 1: medv ~ zn
## Model 2: medv ~ crim + zn + indus
##   Res.Df   RSS Df Sum of Sq      F    Pr(>F)
## 1     504 37167
## 2     502 30170   2    6996.6 58.209 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

After performing a global f test comparing both of the models, we find that the p-value for testing $H_0 : RSS_{Reduced} = RSS_{Full}$ is $< .001$. We can reject H_0 and conclude that the full model is better, because it has a significantly smaller RSS.