

STAT 511: HW #7

Rumil Legaspi

4/19/2021

Contents

Multiple Regression & Brand Preference Dataset	1
1a. Regress degree of brand liking on sweetness only. Write down the estimated regression model.	2
1b. Compute the estimated mean degree of brand liking at each level of sweetness, i.e., what is . .	4
1c. Interpret the intercept coefficient.	4
1d. Interpret the slope coefficient.	4
1e. Is the slope coefficient significant? State the null, alternative, decision rule and conclusion. . .	4
2. Refer to the “Brand Preference” dataset. Code sweetness (X_2) as a dummy variable.	4
2a. Fit a multiple regression model with moisture content, sweetness, and their interaction.	5
2b. Write down the estimated regression equation at each sweetness level	5
***2c. Interpret the slope coefficient in each estimated regression equation in Part (b).	5
2d. Is the interaction coefficient significant at $\alpha = 5\%$? State the null, alternative, decision rule, conclusion.	5
***2e. your answer is NO in Part (d), drop the interaction term and rerun the model. Write down	5
3. Refer to the “Assessed Valuations” dataset (Value.txt)	6
3a. Regress selling price on lot location only. Write down the estimated regression equation. . . .	6
3b. Based on your regression result in Part (a), what is the estimated mean selling price for corner lots? For non-corner lots?	7

Multiple Regression & Brand Preference Dataset

Setting up workspace

```

library(nortest)
library(olsrr)
library(car)
library(lmtest)
library(MASS)
library(tidyverse)

setwd("C:/Users/RUMIL/Desktop/APU/STAT 511 - Millie Mao (Applied Regression Analysis)/Week 10/Week 10/")

brand_data = read.table(file = "Brand.txt", header = FALSE, sep = "")

View(brand_data)

# Adding headers
names(brand_data) <- c("Rating", "Moisture", "Sweetness")

# names(bank_data) <- c("", "")

# Defining dependent and independent vars
Rating = brand_data$Rating #Y
Moisture = brand_data$Moisture #X1
Sweetness = brand_data$Sweetness #X2

```

1a. Regress degree of brand liking on sweetness only. Write down the estimated regression model.

```

# Coding Sweetness as a dummy variable
cat_sweetness <- as.factor(Sweetness)

# Regressing Rating on new Sweetness dummy variable
Sweetness_only <- lm(Rating ~ cat_sweetness, data = brand_data)

summary(Sweetness_only)

##
## Call:
## lm(formula = Rating ~ cat_sweetness, data = brand_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -16.375  -7.312  -0.125   8.688  16.625
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    77.375     3.851  20.094 1.01e-11 ***
## cat_sweetness4    8.750     5.446   1.607   0.13
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 10.89 on 14 degrees of freedom

```

```
## Multiple R-squared:  0.1557, Adjusted R-squared:  0.09539
## F-statistic: 2.582 on 1 and 14 DF,  p-value: 0.1304
```

Estimated Regression model:

$$\hat{Y} = 77.375 + 8.750X$$

```
#Checking how many levels there are in sweetness as a dummy
levels(cat_sweetness)
```

```
## [1] "2" "4"
```

```
typeof(cat_sweetness)
```

```
## [1] "integer"
```

```
attributes(cat_sweetness)
```

```
## $levels
## [1] "2" "4"
##
## $class
## [1] "factor"
```

So in this case since sweetness is a category with only 2 level (sweetness level 2 and sweetness level 4) we can think of our regression model like so:

We will use sweetness level 2 as our reference group

```
#the "2" indicates the level marked as "2" given from our level() function and now setting sweetness level
cat_sweetness_new <- relevel(cat_sweetness, ref = "2")
```

```
#Regressing
```

```
Sweetness_only2 <- lm(Rating ~ cat_sweetness_new, data = brand_data)
```

```
summary(Sweetness_only2)
```

```
##
## Call:
## lm(formula = Rating ~ cat_sweetness_new, data = brand_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -16.375  -7.312  -0.125   8.688  16.625
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      77.375      3.851  20.094 1.01e-11 ***
## cat_sweetness_new4  8.750      5.446   1.607   0.13
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 10.89 on 14 degrees of freedom
## Multiple R-squared:  0.1557, Adjusted R-squared:  0.09539
## F-statistic: 2.582 on 1 and 14 DF,  p-value: 0.1304
```

1b. Compute the estimated mean degree of brand liking at each level of sweetness, i.e., what is

the estimated mean degree of brand liking at sweetness level 2? At sweetness level 4?

After reading the summary we can now say the estimated means when Sweetness level is 2 and 4 to be.

Table 1: Mean Rating at Level 2 or Level 4 Sweetness

$\hat{Y} = B_0 + B_1X_1$	Sweetness = Level 2	Sweetness = Level 4
$\hat{Y} = 77.375 + 8.750X$	$\hat{Y} = 77.375 + 8.750(0)$ $= 77.375$ 77.375	$\hat{Y} = 77.375 + 8.750(1)$ $= 77.375 + 8.750$ 86.125

1c. Interpret the intercept coefficient.

$$B_0 = 77.375$$

The estimated mean Y-value when $X = 0$ which is our reference/baseline group. When put in context, this is the mean Rating when the brand's sweetness level is 2.

1d. Interpret the slope coefficient.

The change in mean Rating for Sweetness level 4 relative to sweetness level 2 is 86.125 ($77.375 + 8.750(1)$).

1e. Is the slope coefficient significant? State the null, alternative, decision rule and conclusion.

Null Hypothesis: $H_0: \beta_j = 0$ (slopes are showing no change), X_j **is not** linearly associated with Y, therefore the partial slope is not significant.

Alternative Hypothesis: $H_1: \beta_j \neq 0$ (slopes are showing change), X_j **is** linearly associated with Y, therefore the partial slope is significant.

Testing the significance of Sweetness Level 4 ($\hat{\beta}_1 = 8.750$) p-value:

Because the p-value for sweetness level 4 is [1] 0.13 and is greater than our alpha (accepted error/significance level) of 0.05 we **fail to reject** our NULL hypothesis and conclude that our partial slope, **Sweetness Level 4** in reference to Sweetness Level of 2, does not show significance (**in our model).

```
####t.test(Sweetness, cat_sweetness_new)
```

2. Refer to the “Brand Preference” dataset. Code sweetness (X_2) as a dummy variable.

2a. Fit a multiple regression model with moisture content, sweetness, and their interaction.

```
#Regressing Rating on Moisture and with Sweetness still as a dummy variable
full_lm <- lm(Rating ~ cat_sweetness + Moisture, data = brand_data)

summary(full_lm)
```

```
##
## Call:
## lm(formula = Rating ~ cat_sweetness + Moisture, data = brand_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.400 -1.762  0.025  1.587  4.200
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    46.4000     2.3129  20.061 3.66e-11 ***
## cat_sweetness4  8.7500     1.3466   6.498 2.01e-05 ***
## Moisture        4.4250     0.3011  14.695 1.78e-09 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.693 on 13 degrees of freedom
## Multiple R-squared:  0.9521, Adjusted R-squared:  0.9447
## F-statistic: 129.1 on 2 and 13 DF,  p-value: 2.658e-09
```

2b. Write down the estimated regression equation at each sweetness level

$$\hat{Y} = 46.40 + 8.75X + 4.425X$$

$$\text{Sweetness level 2: } \hat{Y} = 46.40 + 4.425X$$

$$\text{Sweetness level 4: } \hat{Y} = 46.40 + 8.75X + 4.425X$$

*****2c. Interpret the slope coefficient in each estimated regression equation in Part (b).**

$X_2 = 8.750$: The change in mean Rating for Sweetness level 4 relative to sweetness level 2 while keeping Moisture constant is 55.15.

$X_3 = 4.4250$: The change in mean Rating for a 1 unit increase in Moisture while keeping Sweetness constant is 4.4250

2d. Is the interaction coefficient significant at $\alpha = 5\%$? State the null, alternative, decision rule, conclusion.

*****2e. your answer is NO in Part (d), drop the interaction term and rerun the model. Write down**

the new estimated regression equation at each sweetness level.

3. Refer to the “Assessed Valuations” dataset (Value.txt)

```
value_data = read.table(file = "C:/Users/RUMIL/Desktop/APU/STAT 511 - Millie Mao (Applied Regression Analysis)
View(value_data)

# Adding headers
names(value_data) <- c("Price", "Valuation", "Corner_lot")

# names(bank_data) <- c("", "")

# Defining dependent and independent vars
Price = value_data$Price #Y
Valuation = value_data$Valuation #X1
Corner_lot = value_data$Corner_lot #X2
```

3a. Regress selling price on lot location only. Write down the estimated regression equation.

```
# Coding Corner Location as a dummy variable
cat_corner_lot <- as.factor(Corner_lot)

# For good measure I set 0 as base reference group
cat_corner_lot <- relevel(cat_corner_lot, ref = "0")

location_onlylm <- lm(Price ~ Corner_lot, data = value_data)
summary(location_onlylm)

##
## Call:
## lm(formula = Price ~ Corner_lot, data = value_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -17.854  -5.637   1.157   6.196  16.446
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    81.154      1.198  67.715 < 2e-16 ***
## Corner_lot     -8.523      2.397  -3.556 0.000728 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8.303 on 62 degrees of freedom
## Multiple R-squared:  0.1694, Adjusted R-squared:  0.156
## F-statistic: 12.64 on 1 and 62 DF, p-value: 0.0007283
```

3b. Based on your regression result in Part (a), what is the estimated mean selling price for corner lots? For non-corner lots?

R uses alpha-numeric ordering as the reference group by default so either 'a' or '0' as the base reference group.

Estimated mean selling price for non-corner lots:

The estimated mean selling price for non-corner lots is roughly \$81K

For corner lots:

The estimated mean selling price for corner lots is roughly \$72.5k, which is taken by 81k - 8.5k.

`tinytex::reinstall_tinytex()`