

# STA130 Rstudio Homework

## Problem Set 7

[Xuanqi Wei] ([1009353209]), with Josh Speagle & Scott Schwartz

### Instructions

Complete the exercises in this .Rmd file and submit your .Rmd and knitted .pdf output through [Quercus](#) by 11:59 pm E.T. on Thursday, March 16.

```
library(tidyverse)
```

### Question 1: Multivariate Linear Regression and Mario Kart

In this question, you will revisit the Mario Kart data we looked at in this week's class. This data set contains eBay sales of the game Mario Kart for Nintendo Wii in October 2009 and is available in the `openintro` R package. We have provided a local csv copy, which we will load in below.

```
# load in data
mariokart <- read_csv("mariokart.csv")
glimpse(mariokart)

## Rows: 143
## Columns: 13
## $ ...1      <dbl> 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 1~
## $ id        <dbl> 1.50377e+11, 2.60483e+11, 3.20432e+11, 2.80405e+11, 1.70~
## $ duration  <dbl> 3, 7, 3, 3, 1, 3, 1, 1, 3, 7, 1, 1, 1, 1, 7, 7, 3, 3, 1,~
## $ n_bids    <dbl> 20, 13, 16, 18, 20, 19, 13, 15, 29, 8, 15, 15, 13, 16, 6~
## $ cond      <chr> "new", "used", "new", "new", "new", "new", "new", "used", "new"~
## $ start_pr  <dbl> 0.99, 0.99, 0.99, 0.99, 0.01, 0.99, 0.01, 1.00, 0.99, 19~
## $ ship_pr   <dbl> 4.00, 3.99, 3.50, 0.00, 0.00, 4.00, 0.00, 2.99, 4.00, 4.~
## $ total_pr  <dbl> 51.55, 37.04, 45.50, 44.00, 71.00, 45.00, 37.02, 53.99, ~
## $ ship_sp   <chr> "standard", "firstClass", "firstClass", "standard", "med~
## $ seller_rating <dbl> 1580, 365, 998, 7, 820, 270144, 7284, 4858, 27, 201, 485~
## $ stock_photo <chr> "yes", "yes", "no", "yes", "yes", "yes", "yes", "yes", "~
## $ wheels    <dbl> 1, 1, 1, 1, 2, 0, 0, 2, 1, 1, 2, 2, 2, 1, 0, 1, 1, 2,~
## $ title     <chr> "~~ Wii MARIO KART & WHEEL ~ NINTENDO Wii ~ BRAND NE~
```

Based on documentation in the data set, there are a handful of very high-priced items that were actually bundles of several games/items rather than just Mario Kart. Let's now filter these out.

```
# filter out bundles
mariokart2 <-
  mariokart %>%
  filter(total_pr < 100)
glimpse(mariokart2)
```

```
## Rows: 141
## Columns: 13
```

```
## $ ...1      <dbl> 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 1~
## $ id        <dbl> 1.50377e+11, 2.60483e+11, 3.20432e+11, 2.80405e+11, 1.70~
## $ duration  <dbl> 3, 7, 3, 3, 1, 3, 1, 1, 3, 7, 1, 1, 1, 1, 7, 7, 3, 3, 1,~
## $ n_bids    <dbl> 20, 13, 16, 18, 20, 19, 13, 15, 29, 8, 15, 15, 13, 16, 6~
## $ cond      <chr> "new", "used", "new", "new", "new", "new", "used", "new"~
## $ start_pr  <dbl> 0.99, 0.99, 0.99, 0.99, 0.01, 0.99, 0.01, 1.00, 0.99, 19~
## $ ship_pr   <dbl> 4.00, 3.99, 3.50, 0.00, 0.00, 4.00, 0.00, 2.99, 4.00, 4.~
## $ total_pr  <dbl> 51.55, 37.04, 45.50, 44.00, 71.00, 45.00, 37.02, 53.99, ~
## $ ship_sp   <chr> "standard", "firstClass", "firstClass", "standard", "med~
## $ seller_rating <dbl> 1580, 365, 998, 7, 820, 270144, 7284, 4858, 27, 201, 485~
## $ stock_photo <chr> "yes", "yes", "no", "yes", "yes", "yes", "yes", "yes", "~
## $ wheels    <dbl> 1, 1, 1, 1, 2, 0, 0, 2, 1, 1, 2, 2, 2, 2, 1, 0, 1, 1, 2,~
## $ title     <chr> "-- Wii MARIO KART & WHEEL ~ NINTENDO Wii ~ BRAND NE~
```

(a) Sellers on eBay have the option to include a stock photo as the illustration of the product for sale. Does this choice affect the selling price? Carry out a **univariate (single-variable) linear regression analysis** and predict the mean selling price of the `total_pr` variable for sellers who do and do not use stock photos (`stock_photo`).

*Hint: Your code from Question 4d in HW6 might be helpful here.*

```
mean_p <- lm(total_pr ~ stock_photo, data = mariokart2)
summary(mean_p)$coefficients
```

```
##              Estimate Std. Error  t value    Pr(>|t|)
## (Intercept)  44.327222   1.493540 29.679305 5.092241e-62
## stock_photoyes 4.169159   1.730739  2.408889 1.731116e-02
```

- sellers use stock photos: 48.496381
- sellers don't use stock photos: 44.327222

(b) Sellers are rated by buyers on eBay, captured in the variable `seller_rating`. To simplify our analysis, we will categorize sellers by whether their rating is “low”, “medium”, or “high”. Using `mutate()` and `case_when()`, create a new variable called `seller_rating_tier` that is “low” if `seller_rating`  $\leq$  200, “medium” if  $200 < \text{seller\_rating} \leq 4500$ , and “high” if `seller_rating`  $> 4500$ . Then, carry out a **linear regression analysis** to predict `total_pr` for the “low”, “medium”, and “high” levels of the new `seller_rating_tier` variable.

*Hint: The syntax `lm(y ~ x)` will still work even if `x` is a multi-valued categorical explanatory variable.*

```
mariokart2 <- mariokart2 %>% mutate(seller_rating_tier = ifelse(seller_rating <= 200, "low", ifelse(seller_rating > 200 & seller_rating <= 4500, "medium", "high")))
summary(lm(total_pr ~ seller_rating_tier, data=mariokart2))$coefficients
```

```
##              Estimate Std. Error  t value    Pr(>|t|)
## (Intercept)  49.769792   1.300554 38.268135 2.355060e-75
## seller_rating_tierlow  -4.118396   1.891973 -2.176773 3.119849e-02
## seller_rating_tiermedium -3.050992   1.820776 -1.675654 9.607024e-02
```

- the `total_pr` for the ‘low’ levels of the new ‘`seller_rting_tier`’ is 45.651396
- the `total_pr` for the ‘medium’ levels of the new ‘`seller_rting_tier`’ is 46.72
- the `total_pr` for the ‘high’ levels of the new ‘`seller_rting_tier`’ is 49.769792

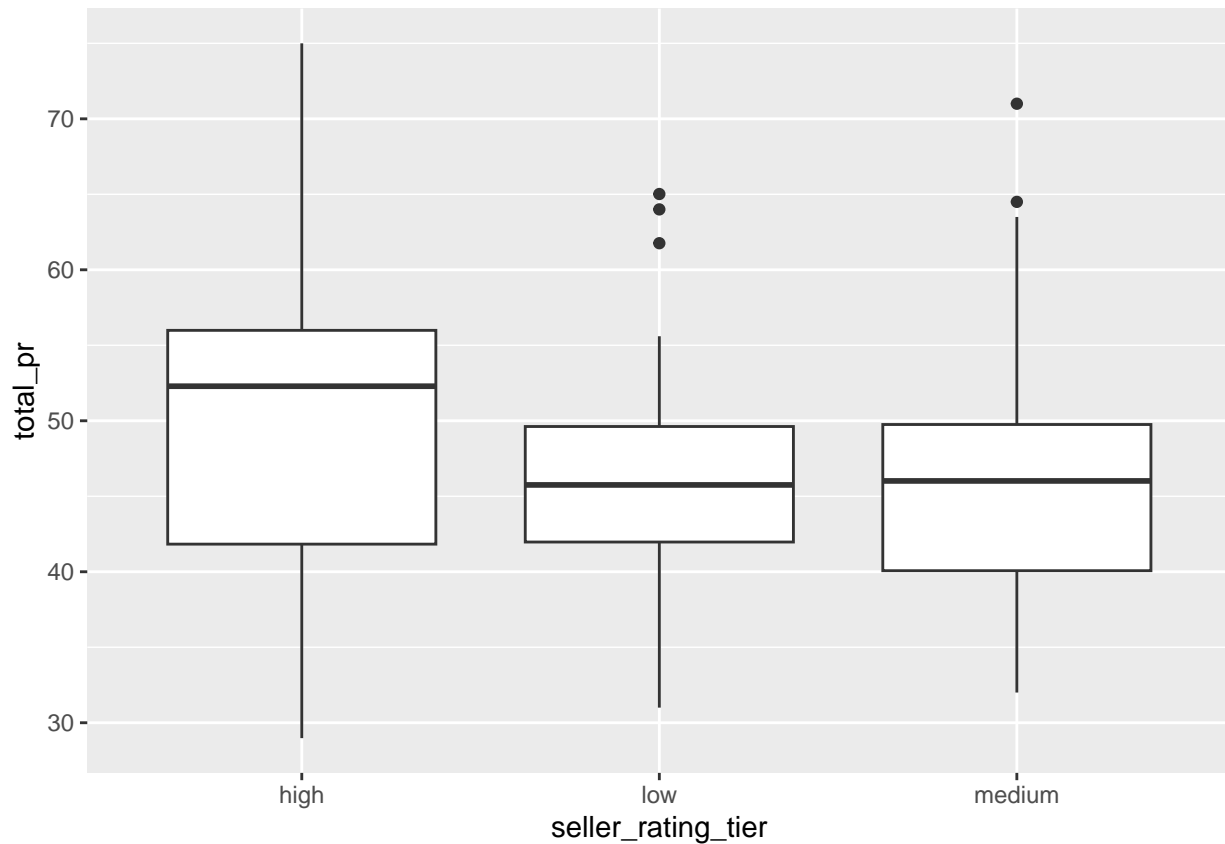
How many indicator variables are in the model? Describe these indicator variables. Which seller rating group is `lm()` treating as the baseline category?

- There are two variables in the model.
- the first: 1 and 0; 1 represent low and 0 represent others

- the second: 1 and 0; 1 represent medium and 0 represent others.
- the group that 'lm()' treating as the baseline category is high

(c) Create **boxplots** of `total_pr` for each category of seller based on `seller_rating_tier`.

```
mariokart2 %>% ggplot(aes(seller_rating_tier, total_pr)) + geom_boxplot()
```



Is this visualization consistent with your estimates from above? Why or why not might this be the case?

- Yes, this visualization is consistent with my estimation from above. Since the boxplot of low is almost symmetric and the mean is around the median. The boxplot of medium is similar to the one of low. The boxplot for high is left skewed, as the mean is smaller than median.

(d) Now, perform an appropriate **multivariate regression analysis** including **interaction terms** to examine whether `seller_rating_tier` has an effect on the relationship between `total_pr` and `duration`.

Note that the full regression model here is:

$$\begin{aligned} \text{total\_pr}_i = & \beta_0 + \beta_1 \times \text{seller\_tier\_low}_i + \beta_2 \times \text{seller\_rating\_tier\_medium}_i + \beta_3 \times \text{duration}_i \\ & + \beta_4 \times \text{seller\_rating\_tier\_low}_i \times \text{duration}_i + \beta_5 \times \text{seller\_rating\_tier\_medium}_i \times \text{duration}_i \\ & + \epsilon_i \end{aligned}$$

*Hint: The syntax for a multivariate interaction model is `lm(y ~ x1 + x2 + x1 * x2)`.*

```
summary(lm(total_pr ~ seller_rating_tier*duration, data=mariokart2))$coefficients
```

```
##              Estimate Std. Error   t value    Pr(>|t|)
## (Intercept)    55.39919    1.9003593  29.1519610 4.101299e-60
## seller_rating_tierlow    -8.185758    3.6192132  -2.2617507 2.531113e-02
## seller_rating_tiermedium  -2.387931    3.0064661  -0.7942651 4.284351e-01
## duration        -2.937082    0.7652626  -3.8380058 1.897989e-04
## seller_rating_tierlow:duration    2.620252    0.9533562   2.7484504 6.807181e-03
## seller_rating_tiermedium:duration  1.538756    0.8856835   1.7373654 8.460333e-02
```

What is the equation of the fitted regression line for sellers with low ratings?

- the equation of the fitted regression line for sellers with low rating is:
- low = 55.399199 - 8.185758 - (2.937082 - 2.620252) x duration

What is the equation of the fitted regression line for sellers with medium ratings?

- the equation of the fitted regression line for sellers with medium rating is:
- medium = 55.399199 - (2.387931 - 1.538756) x duration

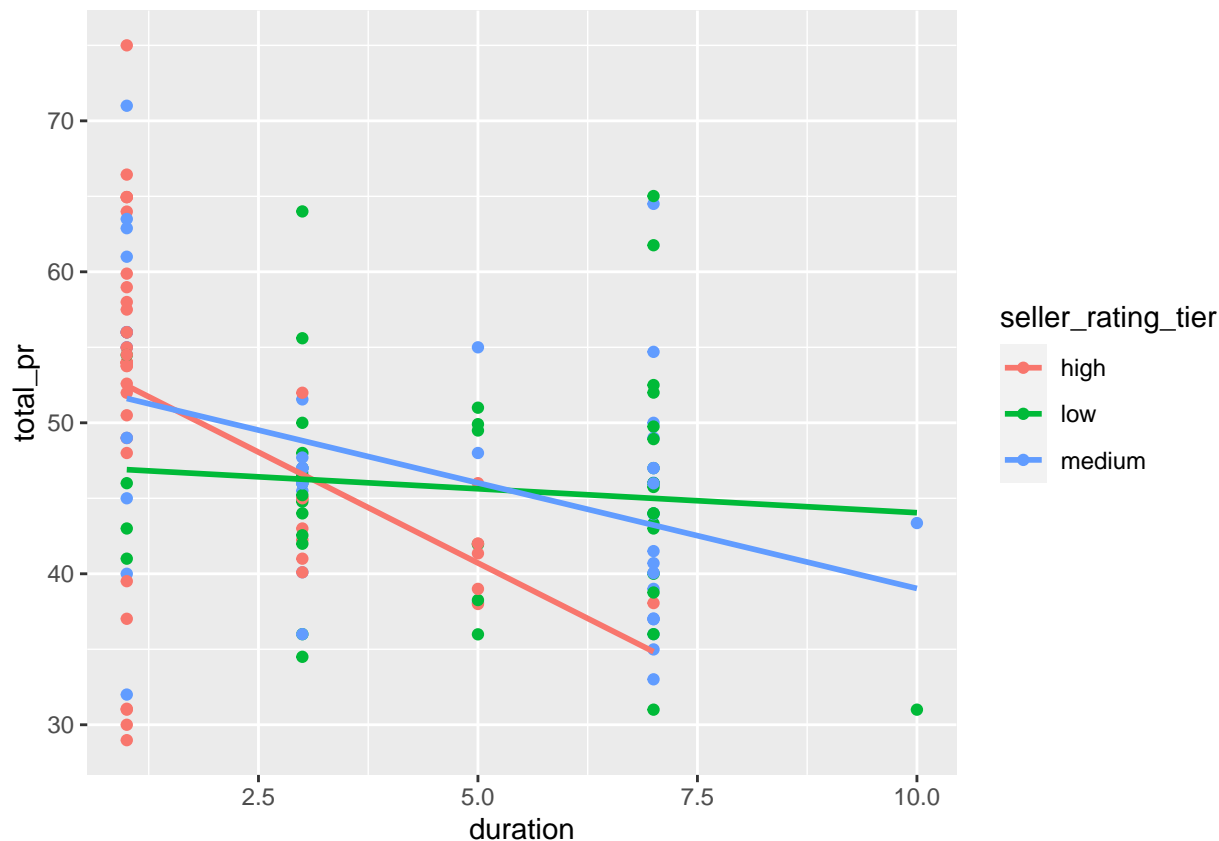
What is the equation of the fitted regression line for sellers with high ratings?

- the equation of the fitted regression line for sellers with medium rating is:
- high = 55.399199 - 2.937082 x duration

(e) Produce an appropriate plot to visualize the fitted relation.

*Hint: Your code from Problem 2d in HW6 might prove useful here.*

```
mariokart2 %>% ggplot(aes(duration, total_pr, color=seller_rating_tier)) +
  geom_point() + geom_smooth(method="lm", se = FALSE)
```



Does the seller rating tier appear to modify the association between duration and total price? Write 1-2 sentences explaining your answer.

- Yes. There is a negative relationship between the duration and the total price.

## Question 2: Predictions and Model Comparison

(a) Divide the data into **testing** and **training** data sets that include 30% and 70% of the data, respectively. Then, fit multivariate linear regression models for the total price `total_pr` using the following combinations of variables (“**features**”) as predictors with training data only:

- `stock_photo`
- `stock_photo`, `duration`, and their interaction
- `seller_rating`
- `stock_photo`, `seller_rating`, and their interaction
- `stock_photo`, `seller_rating`, `duration`, and all interaction terms

*Hint: There are a number of approaches for computing the training/testing splits here. One possibility is you can random sample some fraction  $X$  of the input data using the `sample_frac()` function and then subsequently select the remaining data that has not been sampled using the `anti_join()` function.*

```
set.seed(130)
n <- nrow(mariokart2)
mariokart2 <- mariokart2 %>% rowid_to_column()
```

```
train_ids <- sample(1:n, size=round(0.7*n))
train <- mariokart2 %>% filter(rowid %in% train_ids)
nrow(train)
```

```
## [1] 99
```

```
test <- mariokart2 %>% filter(!(rowid %in% train_ids))
mod1 <- lm(total_pr ~ stock_photo, data=train)
mod2 <- lm(total_pr ~ stock_photo * duration, data=train)
mod3 <- lm(total_pr ~ seller_rating, data=train)
mod4 <- lm(total_pr ~ seller_rating * stock_photo, data=train)
mod5 <- lm(total_pr ~ stock_photo * seller_rating * duration, data=train)
```

(b) Calculate the **root-mean-square-error (RMSE)** for each of the five models from part (a) over both the training and testing datasets (10 values in total) and save the results in a tibble with columns named `model`, `rmse_train`, and `rmse_test`.

As a reminder, for a given response with observed values  $y_1, \dots, y_n$  and corresponding predicted values (from the above models) of  $\hat{y}_1, \dots, \hat{y}_n$ , the RMSE is defined as:

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

*Hint: You can use the syntax `train_linear_model %>% predict(test_data)` to generate predictions to new data values. You can also store your models in a list using the syntax `list(model1, model2, ...)` and access them using the syntax `list[[i]]`.*

```
pred1.test <- predict(mod1, newdata=test)
pred2.test <- predict(mod2, newdata=test)
pred3.test <- predict(mod3, newdata=test)
pred4.test <- predict(mod4, newdata=test)
pred5.test <- predict(mod5, newdata=test)
```

```
pred1.train <- predict(mod1, newdata=train)
pred2.train <- predict(mod2, newdata=train)
pred3.train <- predict(mod3, newdata=train)
pred4.train <- predict(mod4, newdata=train)
pred5.train <- predict(mod5, newdata=train)
```

```
rmse1.train <- sqrt(sum((pred1.train - train$total_pr)^2 / nrow(train)))
rmse2.train <- sqrt(sum((pred2.train - train$total_pr)^2 / nrow(train)))
rmse3.train <- sqrt(sum((pred3.train - train$total_pr)^2 / nrow(train)))
rmse4.train <- sqrt(sum((pred4.train - train$total_pr)^2 / nrow(train)))
rmse5.train <- sqrt(sum((pred5.train - train$total_pr)^2 / nrow(train)))
```

```
rmse1.test <- sqrt(sum((pred1.test - test$total_pr)^2 / nrow(test)))
rmse2.test <- sqrt(sum((pred2.test - test$total_pr)^2 / nrow(test)))
rmse3.test <- sqrt(sum((pred3.test - test$total_pr)^2 / nrow(test)))
rmse4.test <- sqrt(sum((pred4.test - test$total_pr)^2 / nrow(test)))
rmse5.test <- sqrt(sum((pred5.test - test$total_pr)^2 / nrow(test)))
```

```
tibble(model = c(1, 2, 3, 4, 5), rmse_test = c(rmse1.test, rmse2.test, rmse3.test, rmse4.test, rmse5.test),
       rmse_train = c(rmse1.train, rmse2.train, rmse3.train, rmse4.train, rmse5.train))
```

```
## # A tibble: 5 x 3
```

##	model	rmse_test	rmse_train
##	<dbl>	<dbl>	<dbl>
## 1	1	7.71	9.37
## 2	2	7.34	8.66
## 3	3	7.63	9.56
## 4	4	7.66	9.32
## 5	5	6.96	7.95

(c) Based on the results in part (b), write 1-2 sentences discussing which model would you prefer to use for future predictions and why.

- Based on the results in part (b), I prefer the latter one because it has the lower RMSE value.

(d) **(Optional but strongly encouraged)** Make a histogram and boxplot showcasing the distribution of the **effect sizes** over the test data for your preferred model from part (c). As a reminder, the effect size  $e_{ij}$  for object  $i$  with explanatory variable(s)  $x_i \times z_i$  and coefficient  $j > 0$  is defined as

$$e_{ij} = \beta_j \times (x_i \times z_i)$$

such that our linear regression model can be rewritten as

$$y_i = \beta_0 + \sum_{j=1}^m e_{ij} + \epsilon_i$$

*# code you answer here*

Write 1-2 sentences interpreting your results.

*REPLACE THIS TEXT WITH YOUR ANSWER*