

STAT 420: Homework 08

Spring 2019, Guangya Wan

Due: Tuesday, March 26 by 11:30 PM CT

Contents

Assignment	1
Exercise 1 (EPA Emissions Data)	1
Exercise 2 (Hospital SUPPORT Data)	7
Exercise 3 (Fish Data)	10
Exercise 4 (<i>t</i> -test Is a Linear Model)	12

Assignment

Exercise 1 (EPA Emissions Data)

For this exercise we will use the data stored in `epa2015.csv`. It contains detailed descriptions of 4,411 vehicles manufactured in 2015 that were used for fuel economy testing as performed by the Environment Protection Agency. The variables in the dataset are:

- **Make** - manufacturer
- **Model** - model of vehicle
- **ID** - manufacturer defined vehicle identification number within EPA's computer system (not a VIN number)
- **disp** - cubic inch displacement of test vehicle
- **type** - car, truck, or both (for vehicles that meet specifications of both car and truck, like smaller SUVs or crossovers)
- **horse** - rated horsepower, in foot-pounds per second
- **cyl** - number of cylinders
- **lockup** - vehicle has transmission lockup; N or Y
- **drive** - drivetrain system code
 - A = All-wheel drive
 - F = Front-wheel drive
 - P = Part-time 4-wheel drive
 - R = Rear-wheel drive
 - 4 = 4-wheel drive
- **weight** - test weight, in pounds
- **axleratio** - axle ratio
- **nvratio** - n/v ratio (engine speed versus vehicle speed at 50 mph)
- **THC** - total hydrocarbons, in grams per mile (g/mi)
- **CO** - Carbon monoxide (a regulated pollutant), in g/mi
- **CO2** - Carbon dioxide (the primary byproduct of all fossil fuel combustion), in g/mi
- **mpg** - fuel economy, in miles per gallon

We will attempt to model C02 using both `horse` and `type`. In practice we would use many more predictors, but limiting ourselves to these two, one numeric and one factor, will allow us to create a number of plots.

(a) Load the data, and check its structure using `str()`. Verify that `type` is a factor; if not, coerce it to be a factor.

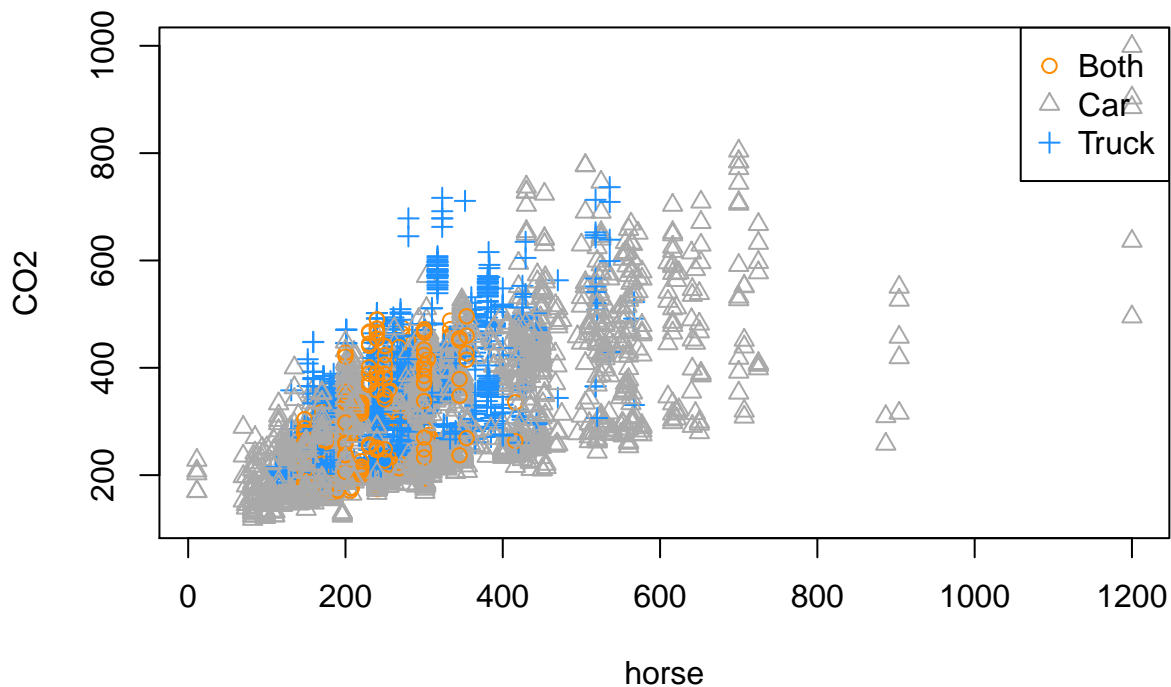
```
epa_data = read.csv('epa2015.csv')
str(epa_data)
```

```
## 'data.frame': 4411 obs. of 16 variables:
## $ Make : Factor w/ 30 levels "aston martin",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ Model : Factor w/ 635 levels "1500 2WD","1500 4X4",...: 189 189 479 479 579 579 582 582 583 583 ...
## $ ID : Factor w/ 872 levels "08-UF2H","0C00007",...: 82 82 180 180 149 149 136 136 148 148 ...
## $ disp : num 5.9 5.9 6 6 6 6 4.7 4.7 4.7 4.7 ...
## $ type : Factor w/ 3 levels "Both","Car","Truck": 2 2 2 2 2 2 2 2 2 2 ...
## $ horse : int 510 510 552 552 565 565 420 420 430 430 ...
## $ cyl : int 12 12 12 12 12 12 8 8 8 8 ...
## $ lockup : Factor w/ 2 levels "N","Y": 2 2 2 2 2 2 1 1 2 2 ...
## $ drive : Factor w/ 5 levels "4","A","F","P",...: 5 5 5 5 5 5 5 5 5 5 ...
## $ weight : int 4500 4500 4750 4750 4250 4250 4000 4000 4000 4000 ...
## $ axleratio: num 3.46 3.46 2.73 2.73 3.73 3.73 3.91 3.91 4.18 4.18 ...
## $ nvratio : num 31 31 22.4 22.4 33.6 33.6 38.6 38.6 36.2 36.2 ...
## $ THC : num 0.0251 0.0022 0.0269 0.0008 0.0248 ...
## $ CO : num 0.12 0.0118 0.5 0.06 0.61 ...
## $ C02 : num 550 344 512 297 603 ...
## $ mpg : num 16.1 25.8 17.3 29.9 14.8 25.3 16 26.5 16.7 28.8 ...
```

Type is a factor variable here.

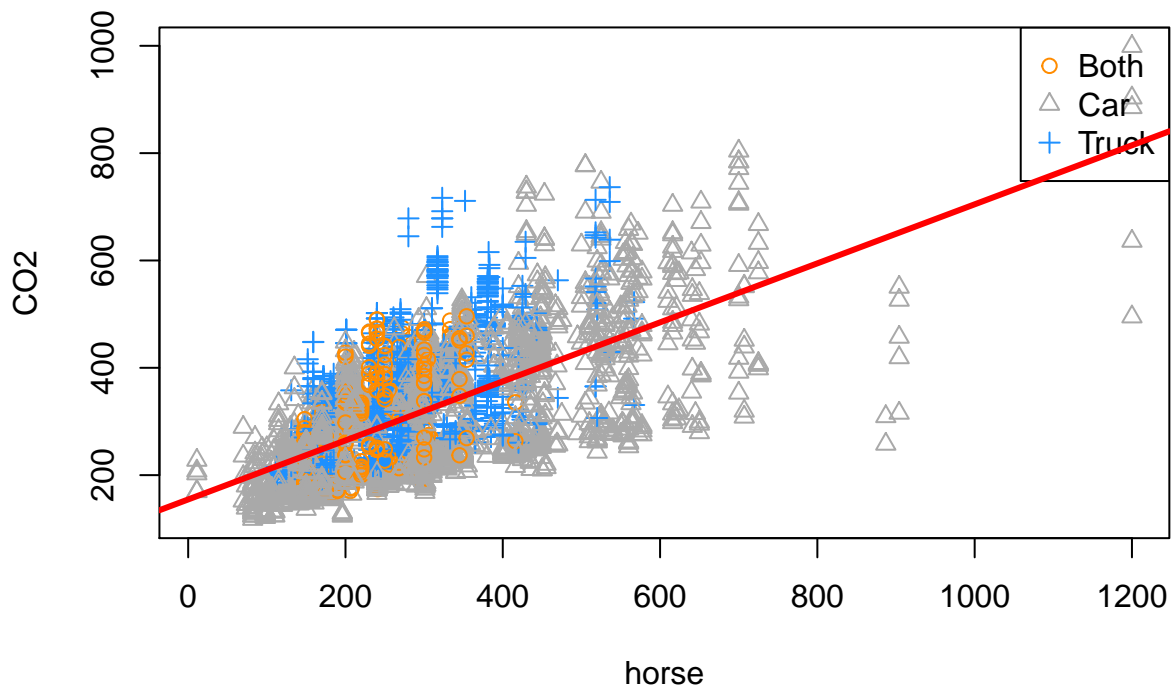
(b) Make a scatterplot of C02 versus `horse`. Use a different color point for each vehicle `type`. Which color is which `type`?

```
plot_colors = c("Darkorange", "Darkgrey", "Dodgerblue")
plot(C02 ~ horse, data = epa_data, col = plot_colors[type], pch = as.numeric(type))
legend("topright", c("Both", "Car", "Truck"), col = plot_colors, pch = c(1, 2, 3))
```



(c) Fit a SLR model with `CO2` as the response and only `horse` as the predictor. Recreate your plot and add the fitted regression line. Comment on how well this line models the data. Give an estimate for the average change in `CO2` for a one foot-pound per second increase in `horse` for a vehicle of type `truck`. Give a 95% prediction interval using this model for the `CO2` of a Subaru Impreza Wagon, which is a vehicle with 148 horsepower and is considered type `Both`. (Interestingly, the dataset gives the wrong drivetrain for most Subarus in this dataset, as they are almost all listed as `F`, when they are in fact all-wheel drive.)

```
epa_slr = lm(CO2~horse,data = epa_data)
plot(CO2 ~ horse, data = epa_data, col = plot_colors[type], pch = as.numeric(type))
legend("topright", c("Both", "Car","Truck"), col = plot_colors, pch = c(1, 2,3))
abline(epa_slr, lwd = 3, col = "red")
```



```
epa_slr$coefficients['horse']
```

```
##      horse
## 0.5498996
```

```
predict(epa_slr,data.frame(horse = 148),interval = 'predict',level = 0.95)
```

```
##      fit      lwr      upr
## 1 236.103 61.35505 410.8509
```

The average change in CO2 for one hp per second increase in horse for type truck is 0.5498996, and a 95% prediction interval using this model for the CO2 of a Subaru Impreza Wagon is (61.35505,410.8509)

(d) Fit an additive multiple regression model with CO2 as the response and **horse** and **type** as the predictors. Recreate your plot and add the fitted regression “lines” with the same colors as their respective points. Comment on how well these lines model the data. Give an estimate for the average change in CO2 for a one foot-pound per second increase in **horse** for a vehicle of type **truck**. Give a 95% prediction interval using this model for the CO2 of a Subaru Impreza Wagon, which is a vehicle with 148 horsepower and is considered type **Both**.

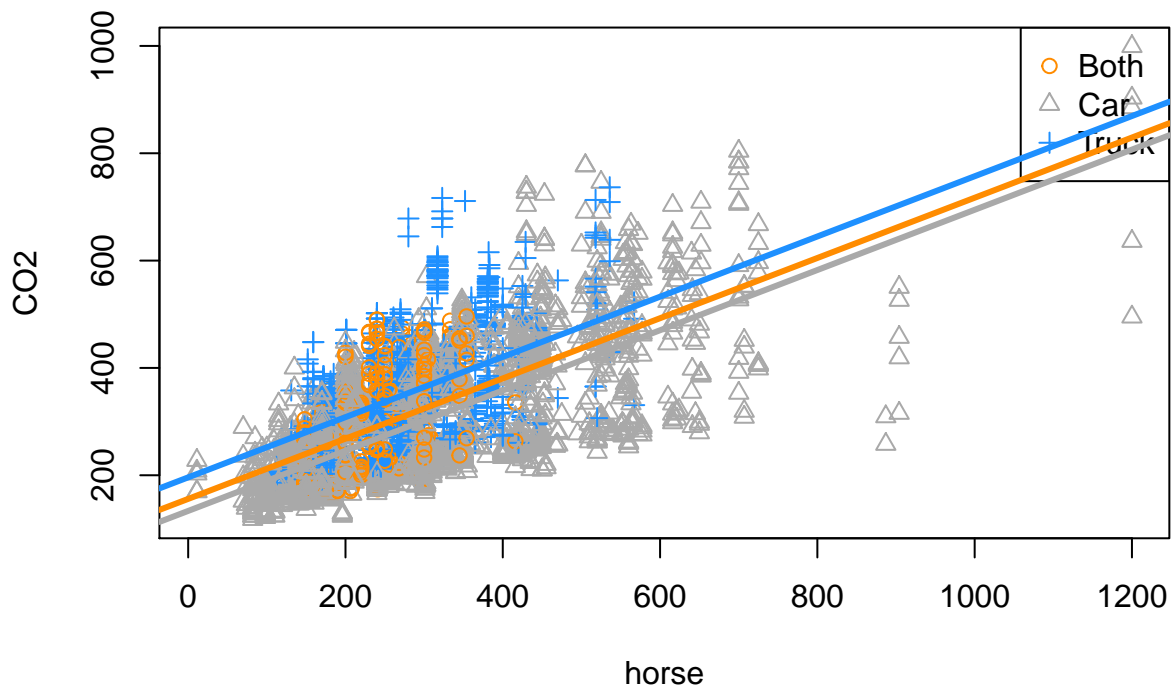
```
epa_add = lm(CO2~horse+type,data = epa_data)
int_both = coef(epa_add)[1]
int_car= coef(epa_add)[1] + coef(epa_add)[3]
int_tru = coef(epa_add)[1] + coef(epa_add)[4]
```

```

slope_all = coef(epa_add)[2]

plot(CO2 ~ horse, data = epa_data, col = plot_colors[type], pch = as.numeric(type))
legend("topright", c("Both", "Car", "Truck"), col = plot_colors, pch = c(1, 2, 3))
abline(int_both, slope_all, lwd = 3, col = "Darkorange")
abline(int_car, slope_all, lwd = 3, col = "Darkgrey")
abline(int_tru, slope_all, lwd = 3, col = "Dodgerblue")

```



```

slope_all

##      horse
## 0.5611008

predict(epa_add, data.frame(horse = 148, type = 'Both'), interval = 'predict', level = 0.95)

##      fit      lwr      upr
## 1 239.0251 71.67975 406.3704

```

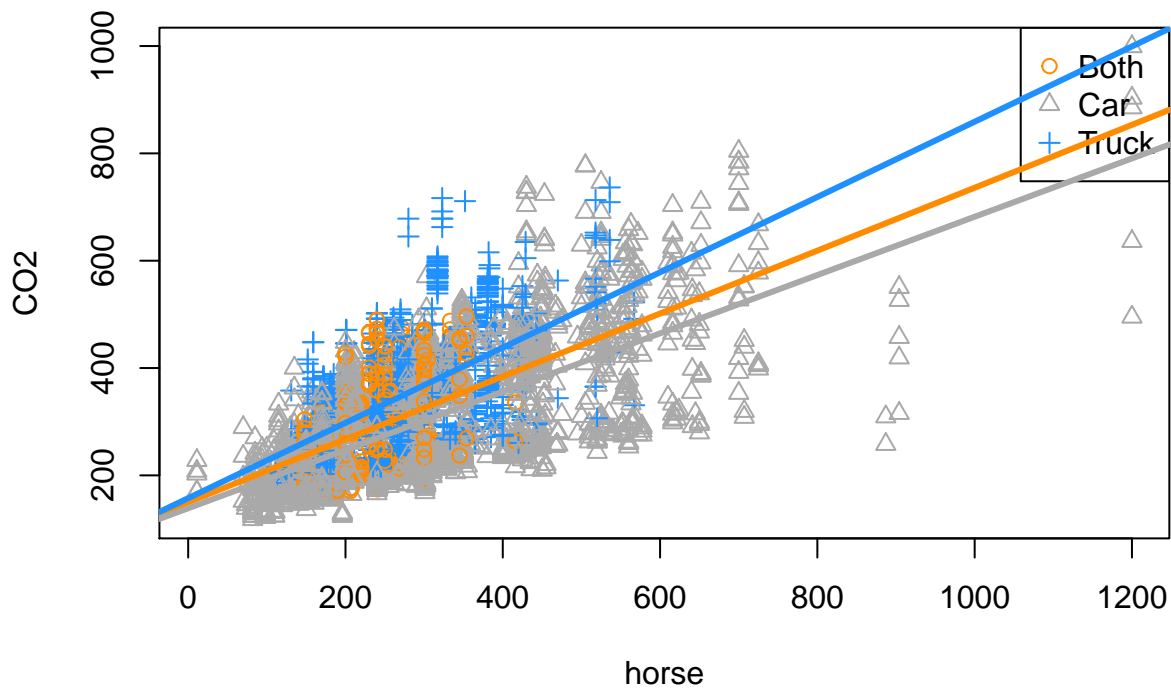
The average change in CO2 for one hp per second increase in horse for type truck is 0.5611008, and a 95% prediction interval using this model for the CO2 of a Subaru Impreza Wagon is (71.67975, 406.3704)

(e) Fit an interaction multiple regression model with CO2 as the response and horse and type as the predictors. Recreate your plot and add the fitted regression “lines” with the same colors as their respective points. Comment on how well these lines model the data. Give an estimate for the average change in CO2 for

a one foot-pound per second increase in `horse` for a vehicle of type `truck`. Give a 95% prediction interval using this model for the CO2 of a Subaru Impreza Wagon, which is a vehicle with 148 horsepower and is considered type `Both`.

```
epa_int = lm(CO2~horse*type,data = epa_data)
int_both = coef(epa_int)[1]
int_car= coef(epa_int)[1] + coef(epa_int)[3]
int_tru = coef(epa_int)[1] + coef(epa_int)[4]
slope_both = coef(epa_int)[2]
slope_car = coef(epa_int)[2] + coef(epa_int)[5]
slope_tru = coef(epa_int)[2] + coef(epa_int)[6]

plot(CO2 ~ horse, data = epa_data, col = plot_colors[type], pch = as.numeric(type))
legend("topright", c("Both", "Car", "Truck"), col = plot_colors, pch = c(1, 2, 3))
abline(int_both,slope_both, lwd = 3, col = "Darkorange")
abline(int_car,slope_car, lwd = 3, col = "Darkgrey")
abline(int_tru,slope_tru, lwd = 3, col = "Dodgerblue")
```



```
slope_tru
```

```
##      horse
## 0.7013883
```

```
predict(epa_int,data.frame(horse = 148,type = 'Both'),interval = 'predict',level = 0.95)
```

```
##          fit          lwr          upr
## 1 236.6339 69.29175 403.9761
```

The average change in CO2 for one fp per second increase in horse for type truck is 0.7013883, and a 95% prediction interval using this model for the CO2 of a Subaru Impreza Wagon is (69.29175,403.9761)

(f) You will perform F -tests later in the exercise, but for now, based solely on the three previous plots, which model is preferred: SLR, additive, or interaction?

I think interaction is the best model here because based on the graph, it is the only one that is able to show the higher average increase trends of truck compared to other two types.

(g) Use an ANOVA F -test to compare the SLR and additive models. Based on this test and a significance level of $\alpha = 0.01$, which model is preferred?

```
anova(epa_slr,epa_add)
```

```
## Analysis of Variance Table
##
## Model 1: CO2 ~ horse
## Model 2: CO2 ~ horse + type
##   Res.Df      RSS Df Sum of Sq    F    Pr(>F)
## 1    4409 35012540
## 2    4407 32054899  2    2957641 203.31 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

According to the anova table results, the p value is smaller than 0.01, so additive model is preferred

(h) Use an ANOVA F -test to compare the additive and interaction models. Based on this test and a significance level of $\alpha = 0.01$, which model is preferred?

```
anova(epa_add,epa_int)
```

```
## Analysis of Variance Table
##
## Model 1: CO2 ~ horse + type
## Model 2: CO2 ~ horse * type
##   Res.Df      RSS Df Sum of Sq    F    Pr(>F)
## 1    4407 32054899
## 2    4405 31894278  2    160621 11.092 1.567e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

According to the anova table results, the p value is smaller than 0.01, so interaction model is preferred

Exercise 2 (Hospital SUPPORT Data)

For this exercise we will use the data stored in `hospital.csv`. It contains a random sample of 580 seriously ill hospitalized patients from a famous study called “SUPPORT” (Study to Understand Prognoses Preferences Outcomes and Risks of Treatment). As the name suggests, the purpose of the study was to determine what factors affected or predicted outcomes, such as how long a patient remained in the hospital. The variables in the dataset are:

- **Days** - Days to death or hospital discharge
- **Age** - Age on day of hospital admission
- **Sex** - female or male
- **Comorbidity** - Patient diagnosed with more than one chronic disease
- **EdYears** - Years of education
- **Education** - Education level; high or low
- **Income** - Income level; high or low
- **Charges** - Hospital charges, in dollars
- **Care** - Level of care required; high or low
- **Race** - Non-white or white
- **Pressure** - Blood pressure, in mmHg
- **Blood** - White blood cell count, in gm/dL
- **Rate** - Heart rate, in bpm

For this exercise, we will use **Charges**, **Pressure**, **Care**, and **Race** to model **Days**.

(a) Load the data, and check its structure using `str()`. Verify that **Care** and **Race** are factors; if not, coerce them to be factors. What are the levels of **Care** and **Race**?

```
hos_data = read.csv('hospital.csv')
str(hos_data)

## 'data.frame':    580 obs. of  13 variables:
## $ Days      : int  8 14 21 4 11 9 25 26 9 16 ...
## $ Age       : num  42.3 63.7 41.5 42 52.1 ...
## $ Sex       : Factor w/ 2 levels "female","male": 1 1 2 2 2 2 1 1 2 2 ...
## $ Comorbidity: Factor w/ 2 levels "no","yes": 1 1 2 2 2 2 2 1 2 2 ...
## $ EdYears   : int   11 22 18 16 8 12 12 13 16 30 ...
## $ Education : Factor w/ 2 levels "high","low": 2 1 1 1 2 2 2 1 1 1 ...
## $ Income    : Factor w/ 2 levels "high","low": 1 1 1 1 1 1 2 1 1 1 ...
## $ Charges   : num   9914 283303 320843 4173 13414 ...
## $ Care      : Factor w/ 2 levels "high","low": 2 1 1 2 2 2 2 1 2 2 ...
## $ Race      : Factor w/ 2 levels "non-white","white": 1 2 2 2 2 2 2 2 2 2 ...
## $ Pressure  : int   84 69 66 97 89 57 99 115 93 102 ...
## $ Blood     : num   11.3 30.1 0.2 10.8 6.4 ...
## $ Rate      : int   94 108 130 88 92 114 150 132 86 90 ...
```

Both **Care** and **Race** are factor variables

(b) Fit an additive multiple regression model with **Days** as the response using **Charges**, **Pressure**, **Care**, and **Race** as predictors. What does R choose as the reference level for **Care** and **Race**?

```
hos_add = lm(Days~Charges+Pressure+Care+Race,data = hos_data)
```

R chooses reference level based on alphabetic order, so for **Care**, high is the reference level, for **Race**, non-white is the one

(c) Fit a multiple regression model with **Days** as the response. Use the main effects of **Charges**, **Pressure**, **Care**, and **Race**, as well as the interaction of **Care** with each of the numeric predictors as predictors. (that is, the interaction of **Care** with **Charges** and the interaction of **Care** with **Pressure**). Use a statistical test to compare this model to the additive model using a significance level of $\alpha = 0.01$. Which do you prefer?


```
hos_int_care = lm(Days~Charges * Care +Pressure * Care +Race , data = hos_data)
anova(hos_add,hos_int_care)
```

```
## Analysis of Variance Table
##
## Model 1: Days ~ Charges + Pressure + Care + Race
## Model 2: Days ~ Charges * Care + Pressure * Care + Race
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1      575 155596
## 2      573 152996   2    2600.2 4.8692 0.008001 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

I prefer the interaction model based on the results of anova table, which suggests that the p value is smaller than 0.01

(d) Fit a multiple regression model with **Days** as the response. Use the predictors from the model in (c) as well as the interaction of **Race** with each of the numeric predictors. (that is, the interaction of **Race** with **Charges** and the interaction of **Race** with **Pressure**). Use a statistical test to compare this model to the additive model using a significance level of $\alpha = 0.01$. Which do you prefer?

```
hos_int_care_race = lm(Days~Charges * Race +Pressure * Race +Charges * Care +Pressure * Care , data = hos_data)
anova(hos_add,hos_int_care_race)
```

```
## Analysis of Variance Table
##
## Model 1: Days ~ Charges + Pressure + Care + Race
## Model 2: Days ~ Charges * Race + Pressure * Race + Charges * Care + Pressure *
##           Care
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1      575 155596
## 2      571 147730   4    7866.6 7.6014 5.699e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

I prefer the interaction model based on the results of anova table, which suggests that the p value is smaller than 0.01

(e) Using the model in (d), give an estimate of the change in average **Days** for a one-unit increase in **Pressure** for a "white" patient that required a high level of care.

```
coefficients(hos_int_care_race)[4] +coefficients(hos_int_care_race)[7]
```

```
## Pressure
## 0.1423975
```

The estimate of the change in average **Days** for a one-unit increase in **Pressure** for a "white" patient that required a high level of care is 0.1423975

(f) Find a model using the four predictors that we have been considering that is more flexible than the model in (d) and that is also statistically significant as compared to the model in (d) at a significance level of $\alpha = 0.01$.

```
three_int_model = lm(Days ~ Charges * Race * Pressure + Charges * Pressure * Care, data = hos_data)
anova(hos_int_care_race, three_int_model)
```

```
## Analysis of Variance Table
##
## Model 1: Days ~ Charges * Race + Pressure * Race + Charges * Care + Pressure *
##      Care
## Model 2: Days ~ Charges * Race * Pressure + Charges * Pressure * Care
##   Res.Df    RSS Df Sum of Sq      F    Pr(>F)
## 1      571 147730
## 2      568 130055   3    17675 25.731 1.28e-15 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Exercise 3 (Fish Data)

For this exercise we will use the data stored in `fish.csv`. It contains data for 158 fish of 7 different species all gathered from the same lake in one season. The variables in the dataset are:

- **Species** - Common name (*Latin name*)
 - 1 = Bream (*Abramis brama*)
 - 2 = Whitewish (*Leuciscus idus*)
 - 3 = Roach (*Leuciscus rutilus*)
 - 4 = (*Abramis bjoerkna*)
 - 5 = Smelt (*Osmerus eperlanus*)
 - 6 = Pike (*Esox Lucius*)
 - 7 = Perch (*Perca fluviatilis*)
- **Weight** - Weight of the fish, in grams
- **Length1** - Length from the nose to the beginning of the tail, in cm
- **Length2** - Length from the nose to the notch of the tail, in cm
- **Length3** - Length from the nose to the end of the tail, in cm
- **HeightPct** - Maximal height as % of Length3
- **WidthPct** - Maximal width as % of Length3
- **Sex** - 0 = female, 1 = male

We will attempt to predict **Weight** using **Length1**, **HeightPct**, and **WidthPct**.

(a) Use R to fit the model

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_1 x_2 + \beta_5 x_1 x_3 + \beta_6 x_2 x_3 + \beta_7 x_1 x_2 x_3 + \epsilon,$$

where

- Y is **Weight**
- x_1 is **Length1**
- x_2 is **HeightPct**
- x_3 is **WidthPct**.

Report the estimated coefficients of the model.

```
fish_data = read.csv('fish.csv')
threeway_fish_model = lm(Weight~Length1 * HeightPct * WidthPct,data = fish_data)
coef(threeway_fish_model)
```

```
##              (Intercept)              Length1
##          -2070.0789433          101.8805290
##              HeightPct              WidthPct
##          99.6568303          137.1369880
##      Length1:HeightPct      Length1:WidthPct
##          -4.2017983          -6.3483423
##      HeightPct:WidthPct Length1:HeightPct:WidthPct
##          -7.9976087          0.3512455
```

(b) Consider fitting a smaller model in R.

```
fish_smaller = lm(Weight ~ Length1 + HeightPct * WidthPct, data = fish_data)
```

Use a statistical test to compare this model with the previous. Report the following:

- The null and alternative hypotheses in terms of the model given in (a)
- The value of the test statistic
- The p-value of the test
- A statistical decision using a significance level of $\alpha = 0.05$
- Which model you prefer

```
anova(fish_smaller,threeway_fish_model)
```

```
## Analysis of Variance Table
##
## Model 1: Weight ~ Length1 + HeightPct * WidthPct
## Model 2: Weight ~ Length1 * HeightPct * WidthPct
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1     153 2480494
## 2     150 1868778   3     611716 16.367 2.972e-09 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The null hypothesis is that the model is

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_1 x_2 + \beta_5 x_1 x_3 + \beta_6 x_2 x_3 + \beta_7 x_1 x_2 x_3 + \epsilon,$$

, the alternative is that the model is

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_6 x_2 x_3 + \epsilon,$$

. The value of F statistic is 16.367, and p value is 2.972e-09, which is smaller than 0.05. Thus, we reject the null model and therefore the alternative model is preferred

(c) Give an expression based on the model in (a) for the true change in average weight for a 1 cm increase in Length1 for a fish with a HeightPct of 20 and a WidthPct of 10.

```
coef(threeway_fish_model)[2] + 20 * coef(threeway_fish_model)[5] + 10 * coef(threeway_fish_model)[6] + 10 * coef(threeway_fish_model)[7]

## Length1
## 24.61024
```

(d) Give an expression based on the smaller model in (b) for the true change in average weight for a 1 cm increase in `Length1` for a fish with a `HeightPct` of 20 and a `WidthPct` of 10.

```
coefficients(fish_smaller)[2]

## Length1
## 31.48522
```

Exercise 4 (*t*-test Is a Linear Model)

In this exercise, we will try to convince ourselves that a two-sample *t*-test assuming equal variance is the same as a *t*-test for the coefficient in front of a single factor variable in a linear model.

First we setup the data frame that we will use throughout.

```
n = 16

ex4 = data.frame(
  groups = c(rep("A", n / 2), rep("B", n / 2)),
  values = rep(0, n))
str(ex4)

## 'data.frame': 16 obs. of 2 variables:
## $ groups: Factor w/ 2 levels "A","B": 1 1 1 1 1 1 1 1 2 2 ...
## $ values: num 0 0 0 0 0 0 0 0 0 0 ...
```

We will use a total sample size of 16, 8 for each group. The `groups` variable splits the data into two groups, A and B, which will be the grouping variable for the *t*-test and a factor variable in a regression. The `values` variable will store simulated data.

We will repeat the following process a number of times.

```
ex4$values = rnorm(n, mean = 10, sd = 3) # simulate data
summary(lm(values ~ groups, data = ex4))

##
## Call:
## lm(formula = values ~ groups, data = ex4)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.5343 -1.6173 -0.7131  1.9922  5.0602
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  10.5686      0.9188  11.503 1.61e-08 ***
```

```
## groupsB      -0.3142      1.2993  -0.242      0.812
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.599 on 14 degrees of freedom
## Multiple R-squared:  0.00416,    Adjusted R-squared:  -0.06697
## F-statistic: 0.05849 on 1 and 14 DF,  p-value: 0.8124
```

```
t.test(values ~ groups, data = ex4, var.equal = TRUE)
```

```
##
## Two Sample t-test
##
## data:  values by groups
## t = 0.24185, df = 14, p-value = 0.8124
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  -2.472566  3.101047
## sample estimates:
## mean in group A mean in group B
##      10.56859      10.25435
```

We use `lm()` to test

$$H_0 : \beta_1 = 0$$

for the model

$$Y = \beta_0 + \beta_1 x_1 + \epsilon$$

where Y are the values of interest, and x_1 is a dummy variable that splits the data in two. We will let R take care of the dummy variable.

We use `t.test()` to test

$$H_0 : \mu_A = \mu_B$$

where μ_A is the mean for the A group, and μ_B is the mean for the B group.

The following code sets up some variables for storage.

```
num_sims = 100
lm_t = rep(0, num_sims)
lm_p = rep(0, num_sims)
tt_t = rep(0, num_sims)
tt_p = rep(0, num_sims)
```

- `lm_t` will store the test statistic for the test $H_0 : \beta_1 = 0$.
- `lm_p` will store the p-value for the test $H_0 : \beta_1 = 0$.
- `tt_t` will store the test statistic for the test $H_0 : \mu_A = \mu_B$.
- `tt_p` will store the p-value for the test $H_0 : \mu_A = \mu_B$.

The variable `num_sims` controls how many times we will repeat this process, which we have chosen to be 100.

(a) Set a seed equal to your UIN. Then write code that repeats the above process 100 times. Each time, store the appropriate values in `lm_t`, `lm_p`, `tt_t`, and `tt_p`. Specifically, each time you should use `ex4$values = rnorm(n, mean = 10, sd = 3)` to update the data. The grouping will always stay the same.

```
seed = 672086209
for (i in 1:100) {
  ex4$values = rnorm(n, mean = 10, sd = 3) # simulate data
  lm_t[i] = as.numeric(coef(summary(lm(values ~ groups, data = ex4)))[2,3])
  lm_p[i] = as.numeric(coef(summary(lm(values ~ groups, data = ex4)))[2,4])
  tt_t[i] = as.numeric(t.test(values ~ groups, data = ex4, var.equal = TRUE)$statistic)
  tt_p[i] = as.numeric(t.test(values ~ groups, data = ex4, var.equal = TRUE)$p.value)
}
```

(b) Report the value obtained by running `mean(lm_t == tt_t)`, which tells us what proportion of the test statistics are equal. The result may be extremely surprising!

```
mean(lm_t == tt_t)
```

```
## [1] 0
```

(c) Report the value obtained by running `mean(lm_p == tt_p)`, which tells us what proportion of the p-values are equal. The result may be extremely surprising!

```
mean(lm_p == tt_p)
```

```
## [1] 0.18
```

(d) If you have done everything correctly so far, your answers to the last two parts won't indicate the equivalence we want to show! What the heck is going on here? The first issue is one of using a computer to do calculations. When a computer checks for equality, it demands **equality**; nothing can be different. However, when a computer performs calculations, it can only do so with a certain level of precision. So if we calculate two quantities we know to be analytically equal, they can differ numerically. Instead of `mean(lm_p == tt_p)` run `all.equal(lm_p, tt_p)`. This will perform a similar calculation, but with a very small error tolerance for each equality. What is the result of running this code? What does it mean?

```
all.equal(lm_p, tt_p)
```

```
## [1] TRUE
```

This means that the pvalue obtained from two sample t test and t-test for single factor variable are the exactly the same

(e) Your answer in (d) should now make much more sense. Then what is going on with the test statistics? Take a look at the values stored in `lm_t` and `tt_t`. What do you notice? Is there a relationship between the two? Can you explain why this is happening?

```
lm_t
```

```
## [1] -0.705858532 -3.121343364 0.688460472 -0.055788060 1.273530642
## [6] 1.338612922 1.278047281 -0.094952651 -1.031571620 -0.057072219
## [11] 0.129812240 0.973981045 -1.010490866 0.183142970 -0.779916663
## [16] -2.074420925 -0.778379368 0.008597025 -0.110781910 1.653819617
## [21] -0.540436185 -0.800086695 -2.264738944 -0.565841430 2.623466111
## [26] -0.706312634 -0.829484574 0.402541223 -1.032285998 -1.034118611
## [31] -0.542773354 -1.652379084 -1.119390959 -1.115645118 1.367073981
## [36] -1.319911304 -2.630280362 0.966416921 0.296761800 0.514085965
## [41] -0.552609341 0.659033506 -0.259747871 -0.741098588 0.603352279
## [46] -0.976513065 -0.891429008 2.426706458 0.177186513 0.434440129
## [51] -0.476338734 0.636277324 -0.972862861 0.919837638 1.236928961
## [56] 1.113681280 -1.138851054 0.572014560 0.237042800 -1.736135523
## [61] 1.301990029 0.717132804 1.357501081 0.368448498 -0.298432812
## [66] 0.654476873 -1.441069173 1.162517351 -0.159033196 -0.021164239
## [71] -3.607870727 -0.337486840 -1.056564679 0.678503557 -0.700810586
## [76] 2.444756219 1.090692411 -1.089014902 -0.443973090 1.069079836
## [81] -0.895417462 1.231184767 -0.187957265 1.431017305 0.063916197
## [86] 0.326317488 -2.638057583 -0.476790909 1.528349427 0.780489794
## [91] -0.355117293 -0.895062342 -0.820646986 0.694186074 0.847067759
## [96] -0.209992627 0.250841547 0.641075778 0.110503606 0.466334278
```

tt_t

```
## [1] 0.705858532 3.121343364 -0.688460472 0.055788060 -1.273530642
## [6] -1.338612922 -1.278047281 0.094952651 1.031571620 0.057072219
## [11] -0.129812240 -0.973981045 1.010490866 -0.183142970 0.779916663
## [16] 2.074420925 0.778379368 -0.008597025 0.110781910 -1.653819617
## [21] 0.540436185 0.800086695 2.264738944 0.565841430 -2.623466111
## [26] 0.706312634 0.829484574 -0.402541223 1.032285998 1.034118611
## [31] 0.542773354 1.652379084 1.119390959 1.115645118 -1.367073981
## [36] 1.319911304 2.630280362 -0.966416921 -0.296761800 -0.514085965
## [41] 0.552609341 -0.659033506 0.259747871 0.741098588 -0.603352279
## [46] 0.976513065 0.891429008 -2.426706458 -0.177186513 -0.434440129
## [51] 0.476338734 -0.636277324 0.972862861 -0.919837638 -1.236928961
## [56] -1.113681280 1.138851054 -0.572014560 -0.237042800 1.736135523
## [61] -1.301990029 -0.717132804 -1.357501081 -0.368448498 0.298432812
## [66] -0.654476873 1.441069173 -1.162517351 0.159033196 0.021164239
## [71] 3.607870727 0.337486840 1.056564679 -0.678503557 0.700810586
## [76] -2.444756219 -1.090692411 1.089014902 0.443973090 -1.069079836
## [81] 0.895417462 -1.231184767 0.187957265 -1.431017305 -0.063916197
## [86] -0.326317488 2.638057583 0.476790909 -1.528349427 -0.780489794
## [91] 0.355117293 0.895062342 0.820646986 -0.694186074 -0.847067759
## [96] 0.209992627 -0.250841547 -0.641075778 -0.110503606 -0.466334278
```

The t statistic obtained from the two tests are exactly opposite to each other. This is because for instance, if in one test the mean values in level B is greater than mean value in level A, since in t test for one single variable is testing whether Beta is 0, which essentially mean(B) - mean(A) because A is the reference level, it will give positive t value in this case; For two sample t test, by alphabetic order, it is comparing whether mean(A) - mean(B) is 0 or not, which is exactly the opposite of the former case. So if mean values in B > mean values in A, it will output negative t value, and vice versa. Thus, these two test would output exactly the opposite t value, although the p value will be the same