

## Week 7 - Homework

STAT 420, Summer 2022, Ilya Andreev, [iandre3@illinois.edu](mailto:iandre3@illinois.edu)

### Exercise 1 (EPA Emissions Data)

For this exercise, we will use the data stored in [epa2017.csv](#). It contains detailed descriptions of vehicles manufactured in 2017 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)
- **C0** - Carbon monoxide (a regulated pollutant), in g/mi
- **C02** - 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.

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

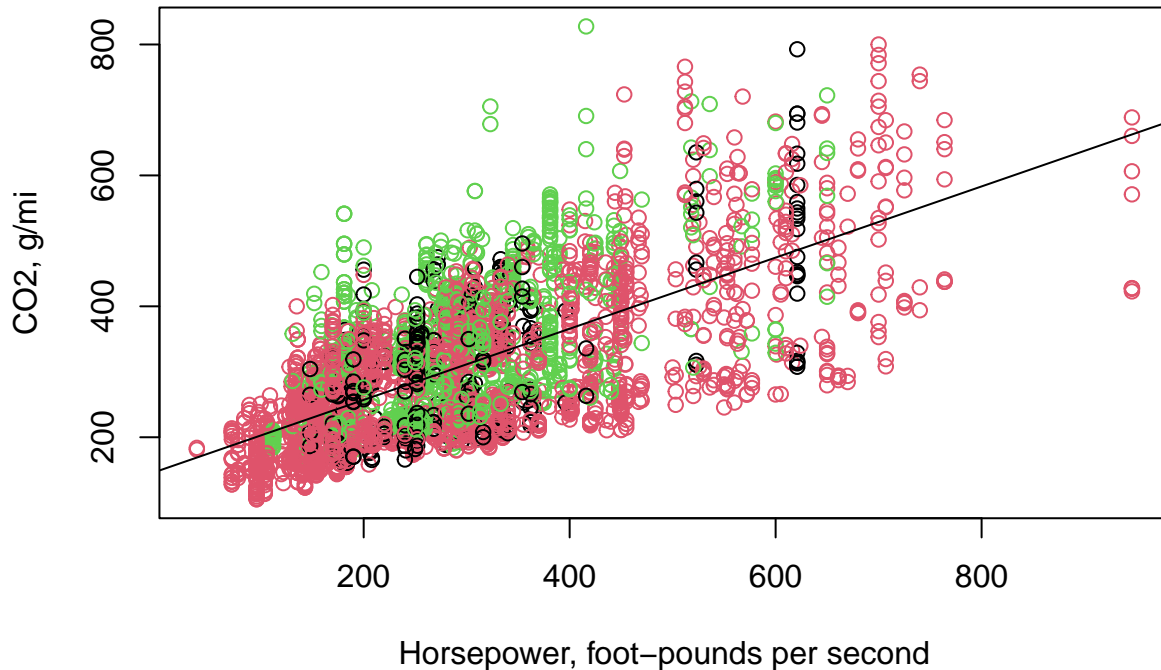
```
data = read.csv("epa2017.csv")
data$type = as.factor(data$type)
```

(a) Do the following:

- Make a scatterplot of **C02** versus **horse**. Use a different color point for each vehicle **type**.
- Fit a simple linear regression model with **C02** as the response and only **horse** as the predictor.
- Add the fitted regression line to the scatterplot. Comment on how well this line models the data.
- Give an estimate for the average change in **C02** for a one foot-pound per second increase in **horse** for a vehicle of type **car**.

- Give a 90% prediction interval using this model for the C02 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.)

```
plot(data$C02 ~ data$horse,
     col=data$type,
     ylab="C02, g/mi",
     xlab="Horsepower, foot-pounds per second")
m = lm(C02 ~ horse, data)
abline(m)
```



While our naive regression seems to fit the data relatively well, any careful reader will surely notice that different types of cars should have different regressions fit to them.

The estimate of the average change in C02 for a one foot-pound per second increase in **horse** for a vehicle of type **car** is 0.5436.

The 90% prediction interval for the C02 of a Subaru Impreza Wagon is (91.5033, 366.0446).

(b) Do the following:

- Make a scatterplot of C02 versus **horse**. Use a different color point for each vehicle **type**.
- Fit an additive multiple regression model with C02 as the response and **horse** and **type** as the predictors.
- Add the fitted regression “lines” to the scatterplot with the same colors as their respective points (one line for each vehicle type). Comment on how well this line models the data.
- Give an estimate for the average change in C02 for a one foot-pound per second increase in **horse** for a vehicle of type **car**.
- Give a 90% prediction interval using this model for the C02 of a Subaru Impreza Wagon, which is a vehicle with 148 horsepower and is considered type **Both**.

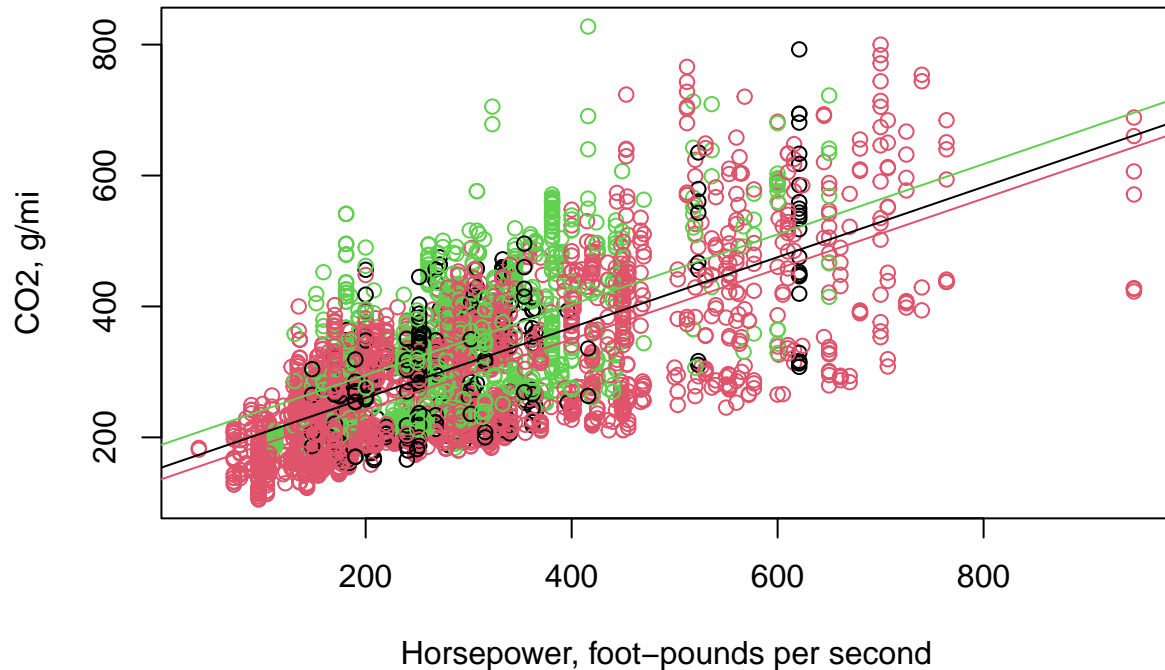
```
plot(data$C02 ~ data$horse,
     col=data$type,
     ylab="C02, g/mi",
     xlab="Horsepower, foot-pounds per second")
m = lm(C02 ~ horse + type, data)
```

```

int_truck = coef(m)["(Intercept)"] + coef(m)["typeTruck"]
slope_truck = coef(m)["horse"]
int_car = coef(m)["(Intercept)"] + coef(m)["typeCar"]
slope_car = coef(m)["horse"]
int_both = coef(m)["(Intercept)"]
slope_both = coef(m)["horse"]

abline(int_truck, slope_truck, col=data$type[data$type == "Truck"])
abline(int_car, slope_car, col=data$type[data$type == "Car"])
abline(int_both, slope_both, col=data$type[data$type == "Both"])

```



Clearly an additive model does not yield much improvement over an SLR. The reason for this is that all 3 fits are constrained to have the same slope, whereas it is visible that the difference in distributions manifests itself through the difference in slopes.

The estimate of the average change in CO2 for a one foot-pound per second increase in `horse` for a vehicle of type `car` is 0.5372.

The 90% prediction interval for the CO2 of a Subaru Impreza Wagon is (100.0012, 364.8952).

(c) Do the following:

- Make a scatterplot of CO2 versus `horse`. Use a different color point for each vehicle `type`.
- Fit an interaction multiple regression model with CO2 as the response and `horse` and `type` as the predictors.
- Add the fitted regression “lines” to the scatterplot with the same colors as their respective points (one line for each vehicle type). 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 `car`.
- Give a 90% 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`.

```

plot(data$CO2 ~ data$horse,
     col=data$type,

```

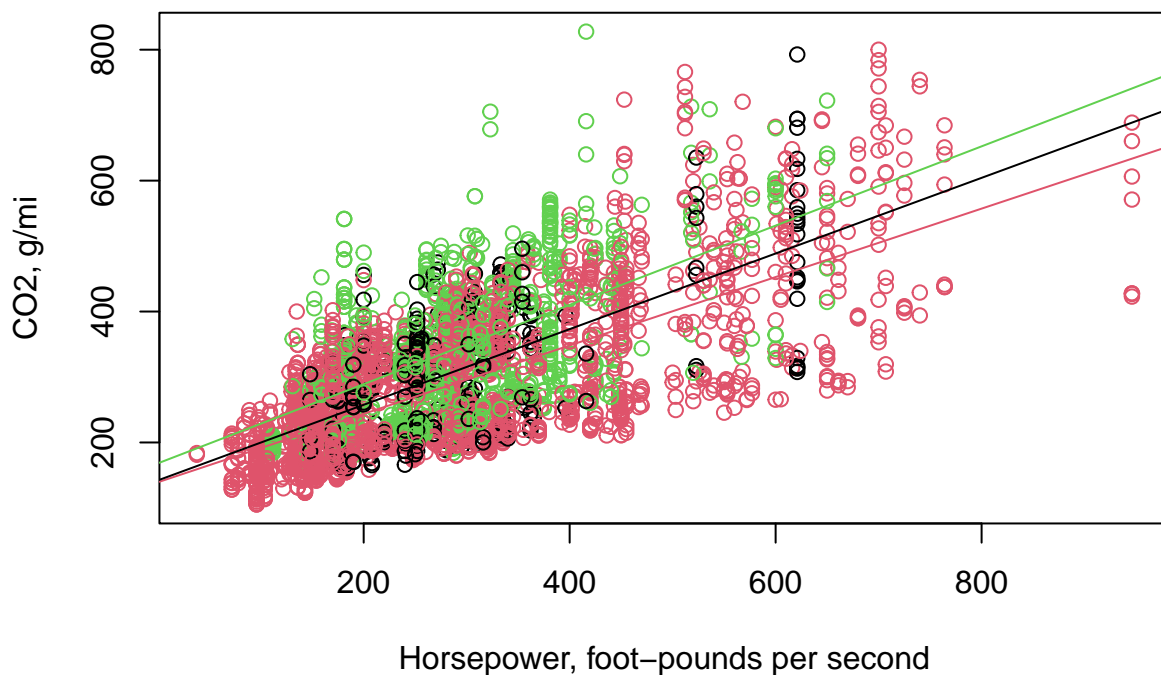
```

ylab="CO2, g/mi",
xlab="Horsepower, foot-pounds per second")
m = lm(CO2 ~ horse * type, data)

int_truck = coef(m)["(Intercept)"] + coef(m)["typeTruck"]
slope_truck = coef(m)["horse"] + coef(m)["horse:typeTruck"]
int_car = coef(m)["(Intercept)"] + coef(m)["typeCar"]
slope_car = coef(m)["horse"] + coef(m)["horse:typeCar"]
int_both = coef(m)["(Intercept)"]
slope_both = coef(m)["horse"]

abline(int_truck, slope_truck, col=data$type[data$type == "Truck"])
abline(int_car, slope_car, col=data$type[data$type == "Car"])
abline(int_both, slope_both, col=data$type[data$type == "Both"])

```



The interactive model explains variation between types of vehicles much better, as can be seen from the plot.

The estimate of the average change in CO2 for a one foot-pound per second increase in `horse` for a vehicle of type `car` is 0.5226.

The 90% prediction interval for the CO2 of a Subaru Impreza Wagon is (95.0055, 359.9619).

(d) Based on the previous plots, you probably already have an opinion on the best model. Now use an ANOVA  $F$ -test to compare the additive and interaction models. Based on this test and a significance level of  $\alpha = 0.10$ , which model is preferred?

```

additive = lm(CO2 ~ horse + type, data)
interactive = lm(CO2 ~ horse * type, data)
anova(additive, interactive)

```

```

## 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   4033 26073761
## 2   4031 26007441   2     66320 5.14 0.0059 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

At  $\alpha = 0.10$ , we prefer the interactive model.

---

## Exercise 2 (Hospital SUPPORT Data, White Blood Cells)

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 **Age**, **Education**, **Income**, and **Sex** in an attempt to model **Blood**. Essentially, we are attempting to model white blood cell count using only demographic information.

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

```
data = read.csv("hospital.csv")
data$Education = as.factor(data$Education)
data$Income = as.factor(data$Income)
data$Sex = as.factor(data$Sex)
```

```
levels(data$Education)
```

```
## [1] "high" "low"
```

```
levels(data$Income)
```

```
## [1] "high" "low"
```

```
levels(data$Sex)
```

```
## [1] "female" "male"
```

(b) Fit an additive multiple regression model with **Blood** as the response using **Age**, **Education**, **Income**, and **Sex** as predictors. What does R choose as the reference level for **Education**, **Income**, and **Sex**?

```
additive = lm(Blood ~ Age + Education + Income + Sex, data)
summary(additive)
```

```
##
## Call:
## lm(formula = Blood ~ Age + Education + Income + Sex, data = data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -12.79  -5.13  -1.58   3.07  47.37
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   10.8662     1.4284   7.61 1.1e-13 ***
## Age           0.0283     0.0207   1.37  0.1725
## Educationlow  0.5967     0.7557   0.79  0.4301
## Incomelow     0.1867     0.7139   0.26  0.7938
## Sexmale      -1.8714     0.6613  -2.83  0.0048 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.88 on 575 degrees of freedom
## Multiple R-squared:  0.0197, Adjusted R-squared:  0.0129
## F-statistic: 2.9 on 4 and 575 DF, p-value: 0.0216
```

R used `high` as the reference level for `Education`, `high` as the reference level for `Income`, and `female` as the reference level for `Sex`.

(c) Fit a multiple regression model with `Blood` as the response. Use the main effects of `Age`, `Education`, `Income`, and `Sex`, as well as the interaction of `Sex` with `Age` and the interaction of `Sex` and `Income`. Use a statistical test to compare this model to the additive model using a significance level of  $\alpha = 0.10$ . Which do you prefer?

```
interactive_small = lm(Blood ~ Age + Education + Income + Sex + Sex * Age + Sex * Income, data)
anova(additive, interactive_small)
```

```
## Analysis of Variance Table
##
## Model 1: Blood ~ Age + Education + Income + Sex
## Model 2: Blood ~ Age + Education + Income + Sex + Sex * Age + Sex * Income
##   Res.Df  RSS Df Sum of Sq    F Pr(>F)
## 1      575 35694
## 2      573 35423   2      271 2.19  0.11
```

Based on the ANOVA test, we prefer the additive model at  $\alpha = 0.10$ .

(d) Fit a model similar to that in (c), but additionally add the interaction between `Income` and `Age` as well as a three-way interaction between `Age`, `Income`, and `Sex`. Use a statistical test to compare this model to the preferred model from (c) using a significance level of  $\alpha = 0.10$ . Which do you prefer?

```
interactive_large = lm(Blood ~ Age + Education + Income + Sex + Sex * Age + Sex * Income + Income * Age +
  Age * Income * Sex, data)
anova(additive, interactive_large)
```

```
## Analysis of Variance Table
##
## Model 1: Blood ~ Age + Education + Income + Sex
## Model 2: Blood ~ Age + Education + Income + Sex + Sex * Age + Sex * Income +
##           Income * Age + Age * Income * Sex
##   Res.Df  RSS Df Sum of Sq    F Pr(>F)
## 1      575 35694
```

```
## 2      571 35166  4      528 2.14  0.074 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

This time, at  $\alpha = 0.10$  we prefer the large interactive model over the additive model.

(e) Using the model in (d), give an estimate of the change in average **Blood** for a one-unit increase in **Age** for a highly educated, low income, male patient.

An estimate of the change in average **Blood** for a one-unit increase in **Age** for a highly educated, low income, male patient is 0.0053.

### Exercise 3 (Hospital SUPPORT Data, Stay Duration)

For this exercise, we will again 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 **Blood**, **Pressure**, and **Rate** in an attempt to model **Days**. Essentially, we are attempting to model the time spent in the hospital using only health metrics measured at the hospital.

Consider 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 **Days**
- $x_1$  is **Blood**
- $x_2$  is **Pressure**
- $x_3$  is **Rate**.

(a) Fit the model above. Also fit a smaller model using the provided R code.

```
hospital = read.csv("hospital.csv")
days_add = lm(Days ~ Pressure + Blood + Rate, hospital)
days_int = lm(Days ~ Blood * Pressure * Rate, hospital)
```

Use a statistical test to compare the two models. Report the following:

- The null and alternative hypotheses in terms of the model given in the exercise description

- The value of the test statistic
- The p-value of the test
- A statistical decision using a significance level of  $\alpha = 0.10$
- Which model you prefer

```
test = anova(days_add, days_int)
```

$H_0 : \beta_4 = \beta_5 = \beta_6 = \beta_7 = 0$  vs  $H_1$  : at least one of the following:  $\beta_4 \neq 0, \beta_5 \neq 0, \beta_6 \neq 0, \beta_7 \neq 0$ .

The test statistic value is 2.0426.

The p-value of the test is 0.087.

The statistical decision at  $\alpha = 0.10$  is to reject the null hypothesis.

We prefer the interaction model.

(b) Give an expression based on the model in the exercise description for the true change in length of hospital stay in days for a 1 bpm increase in **Rate** for a patient with a **Pressure** of 139 mmHg and a **Blood** of 10 gm/dL. Your answer should be a linear function of the  $\beta$ s.

$$\Delta Days = \beta_3 + 10\beta_5 + 139\beta_6 + 1390\beta_7$$

(c) Give an expression based on the additive model in part (a) for the true change in length of hospital stay in days for a 1 bpm increase in **Rate** for a patient with a **Pressure** of 139 mmHg and a **Blood** of 10 gm/dL. Your answer should be a linear function of the  $\beta$ s.

$$\Delta Days = \beta_3$$

\*\*\*

## 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 two-level factor variable (dummy variable) in a linear model.

First, we set up the data frame that we will use throughout.

```
n = 30

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

```
## 'data.frame':   30 obs. of  2 variables:
## $ groups: chr  "A" "A" "A" "A" ...
## $ values: num  0 0 0 0 0 0 0 0 0 0 ...
```

We will use a total sample size of 30, 15 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.

```
set.seed(20)
sim_data$values = rnorm(n, mean = 42, sd = 3.5) # simulate response data
summary(lm(values ~ groups, data = sim_data))
```



```
##
## Call:
## lm(formula = values ~ groups, data = sim_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -9.04  -1.11  -0.14   2.23   7.33
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   40.922      0.950   43.07  <2e-16 ***
## groupsB        0.029      1.344    0.02    0.98
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.68 on 28 degrees of freedom
## Multiple R-squared:  1.66e-05, Adjusted R-squared: -0.0357
## F-statistic: 0.000465 on 1 and 28 DF, p-value: 0.983
t.test(values ~ groups, data = sim_data, var.equal = TRUE)

##
## Two Sample t-test
##
## data:  values by groups
## t = -0.022, df = 28, p-value = 1
## alternative hypothesis: true difference in means between group A and group B is not equal to 0
## 95 percent confidence interval:
##  -2.781  2.723
## sample estimates:
## mean in group A mean in group B
##           40.92           40.95
```

We use `lm()` to test

$$H_0 : \beta_1 = 0$$

for the model

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

where  $Y$  is 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  $\mu_A$  is the mean for the A group, and  $\mu_B$  is the mean for the B group.

The following code sets up some variables for storage.

```
num_sims = 300
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 300.

(a) Set a seed equal to your birthday. Then write code that repeats the above process 300 times. Each time, store the appropriate values in `lm_t`, `lm_p`, `tt_t`, and `tt_p`. Specifically, each time you should use `sim_data$values = rnorm(n, mean = 42, sd = 3.5)` to update the data. The grouping will always stay the same.

```
set.seed(24111999)

for (i in 1:num_sims) {
  sim_data$values = rnorm(n, mean = 42, sd = 3.5)
  m = lm(values ~ groups, data = sim_data)
  lm_t[i] = summary(m)$coefficients[2, 3]
  lm_p[i] = summary(m)$coefficients[2, 4]

  test = t.test(values ~ groups, data = sim_data, var.equal = TRUE)
  tt_t[i] = test$statistic
  tt_p[i] = test$p.value
}
```

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

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

```
## [1] 0.02667
```

(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
```

(e) Your answer in (d) should now make much more sense. Then what is going on with the test statistics? 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?

The test statistic is also equal up to some precision if we ignore the sign. One might notice that the signs of the test statistics for the linear regression test and for the two-sample *t*-test are always the exact opposite.