

Introduction to linear regression

John Cruz

The Human Freedom Index is a report that attempts to summarize the idea of “freedom” through a bunch of different variables for many countries around the globe. It serves as a rough objective measure for the relationships between the different types of freedom - whether it’s political, religious, economical or personal freedom - and other social and economic circumstances. The Human Freedom Index is an annually co-published report by the Cato Institute, the Fraser Institute, and the Liberales Institut at the Friedrich Naumann Foundation for Freedom.

In this lab, you’ll be analyzing data from Human Freedom Index reports from 2008-2016. Your aim will be to summarize a few of the relationships within the data both graphically and numerically in order to find which variables can help tell a story about freedom.

Getting Started

Load packages

In this lab, you will explore and visualize the data using the **tidyverse** suite of packages. The data can be found in the companion package for OpenIntro resources, **openintro**.

Let’s load the packages.

```
library(tidyverse)
library(openintro)
data('hfi', package='openintro')
```

The data

The data we’re working with is in the openintro package and it’s called **hfi**, short for Human Freedom Index.

1. What are the dimensions of the dataset?

There are 1,458 observations with 123 variables.

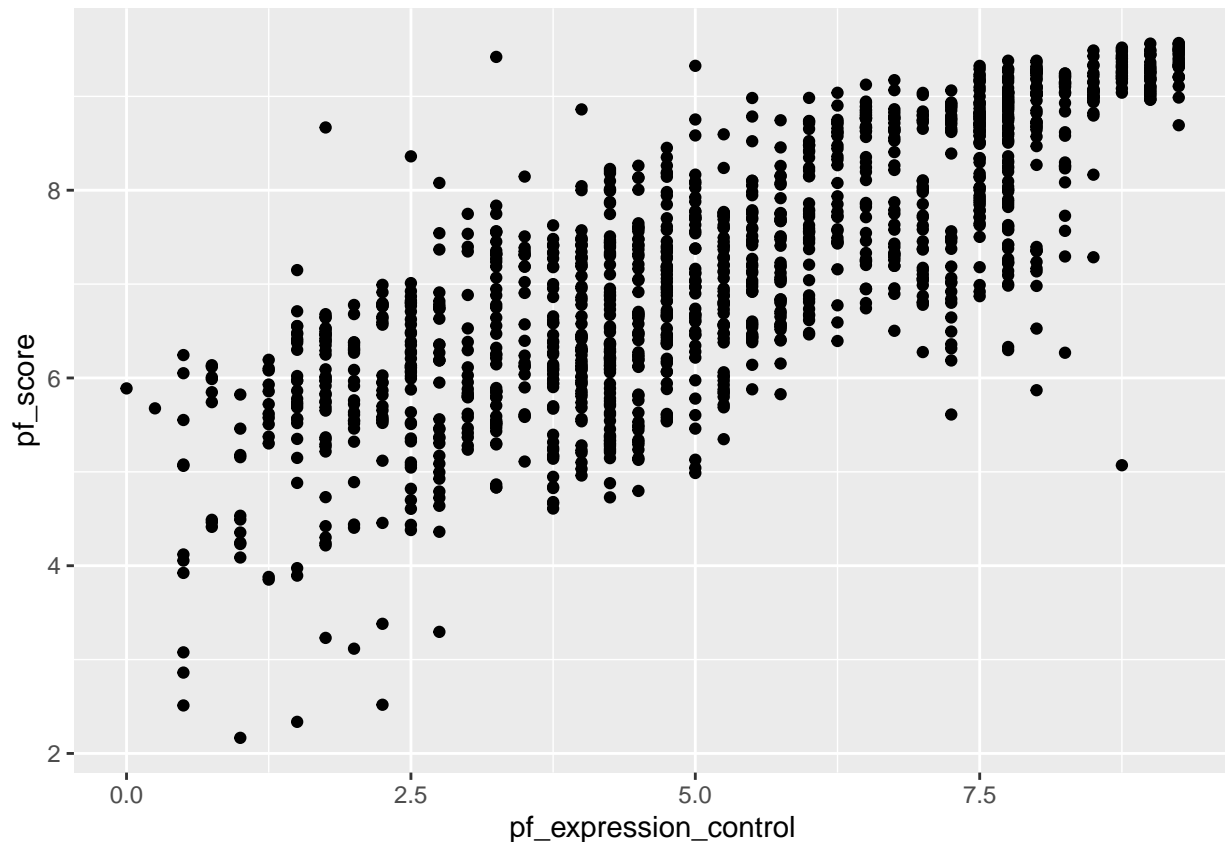
```
dim(hfi)
```

```
## [1] 1458 123
```

2. What type of plot would you use to display the relationship between the personal freedom score, **pf_score**, and one of the other numerical variables? Plot this relationship using the variable **pf_expression_control** as the predictor. Does the relationship look linear? If you knew a country’s **pf_expression_control**, or its score out of 10, with 0 being the most, of political pressures and controls on media content, would you be comfortable using a linear model to predict the personal freedom score?

Utilizing a scatterplot, we can see somewhat of a step pattern increase along the *pf_expression_control* variable. However, overall it does look to have a linear relationship and using a linear model to predict the personal freedom score may be viable.

```
hfi |>
  ggplot(aes(x = pf_expression_control, y = pf_score)) +
  geom_point()
```



If the relationship looks linear, we can quantify the strength of the relationship with the correlation coefficient.

```
hfi %>%
  summarise(cor(pf_expression_control, pf_score, use = "complete.obs"))

## # A tibble: 1 x 1
##   'cor(pf_expression_control, pf_score, use = "complete.obs")'
##                                     <dbl>
## 1                                     0.796
```

Here, we set the *use* argument to “complete.obs” since there are some observations of NA.

Sum of squared residuals

In this section, you will use an interactive function to investigate what we mean by “sum of squared residuals”. You will need to run this function in your console, not in your markdown document. Running the function also requires that the *hfi* dataset is loaded in your environment.

Think back to the way that we described the distribution of a single variable. Recall that we discussed characteristics such as center, spread, and shape. It's also useful to be able to describe the relationship of two numerical variables, such as `pf_expression_control` and `pf_score` above.

3. Looking at your plot from the previous exercise, describe the relationship between these two variables. Make sure to discuss the form, direction, and strength of the relationship as well as any unusual observations.

We can see from the scatterplot a moderate, positive increase where as `pf_expression_control` increases, the `pf_score` increases. However, the trend is not entirely linear as the tail ends when `pf_expression_control` is $\sim < 2.5$ and greater than 7.5, outliers begin to form a small concave down parabola.

Just as you've used the mean and standard deviation to summarize a single variable, you can summarize the relationship between these two variables by finding the line that best follows their association. Use the following interactive function to select the line that you think does the best job of going through the cloud of points.

```
# This will only work interactively (i.e. will not show in the knitted document)
hfi <- hfi %>% filter(complete.cases(pf_expression_control, pf_score))
DATA606::plot_ss(x = hfi$pf_expression_control, y = hfi$pf_score)
```

After running this command, you'll be prompted to click two points on the plot to define a line. Once you've done that, the line you specified will be shown in black and the residuals in blue. Note that there are 30 residuals, one for each of the 30 observations. Recall that the residuals are the difference between the observed values and the values predicted by the line:

$$e_i = y_i - \hat{y}_i$$

The most common way to do linear regression is to select the line that minimizes the sum of squared residuals. To visualize the squared residuals, you can rerun the plot command and add the argument `showSquares = TRUE`.

```
DATA606::plot_ss(x = hfi$pf_expression_control, y = hfi$pf_score, showSquares = TRUE)
```

Note that the output from the `plot_ss` function provides you with the slope and intercept of your line as well as the sum of squares.

4. Using `plot_ss`, choose a line that does a good job of minimizing the sum of squares. Run the function several times. What was the smallest sum of squares that you got? How does it compare to your neighbors?

The smallest sum of squares I got was 971.172. Another classmate got 1033.

The linear model

It is rather cumbersome to try to get the correct least squares line, i.e. the line that minimizes the sum of squared residuals, through trial and error. Instead, you can use the `lm` function in R to fit the linear model (a.k.a. regression line).

```
m1 <- lm(pf_score ~ pf_expression_control, data = hfi)
```

The first argument in the function `lm` is a formula that takes the form $y \sim x$. Here it can be read that we want to make a linear model of `pf_score` as a function of `pf_expression_control`. The second argument specifies that R should look in the `hfi` data frame to find the two variables.

The output of `lm` is an object that contains all of the information we need about the linear model that was just fit. We can access this information using the summary function.

```
summary(m1)
```

```
##
## Call:
## lm(formula = pf_score ~ pf_expression_control, data = hfi)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.8467 -0.5704  0.1452  0.6066  3.2060
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      4.61707    0.05745   80.36  <2e-16 ***
## pf_expression_control 0.49143    0.01006   48.85  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.8318 on 1376 degrees of freedom
## (80 observations deleted due to missingness)
## Multiple R-squared:  0.6342, Adjusted R-squared:  0.634
## F-statistic: 2386 on 1 and 1376 DF, p-value: < 2.2e-16
```

Let's consider this output piece by piece. First, the formula used to describe the model is shown at the top. After the formula you find the five-number summary of the residuals. The "Coefficients" table shown next is key; its first column displays the linear model's y-intercept and the coefficient of `pf_expression_control`. With this table, we can write down the least squares regression line for the linear model:

$$\hat{y} = 4.61707 + 0.49143 \times pf_expression_control$$

One last piece of information we will discuss from the summary output is the Multiple R-squared, or more simply, R^2 . The R^2 value represents the proportion of variability in the response variable that is explained by the explanatory variable. For this model, 63.42% of the variability in runs is explained by at-bats.

5. Fit a new model that uses `pf_expression_control` to predict `hf_score`, or the total human freedom score. Using the estimates from the R output, write the equation of the regression line. What does the slope tell us in the context of the relationship between human freedom and the amount of political pressure on media content?

$$\hat{y} = 5.153687 + 0.349862 \times pf_expression_control$$

As `pf_expression_control` increases by one unit, `human_freedom` increases by 0.349862.

```
lm_mod <- lm(hf_score ~ pf_expression_control, data = hfi)

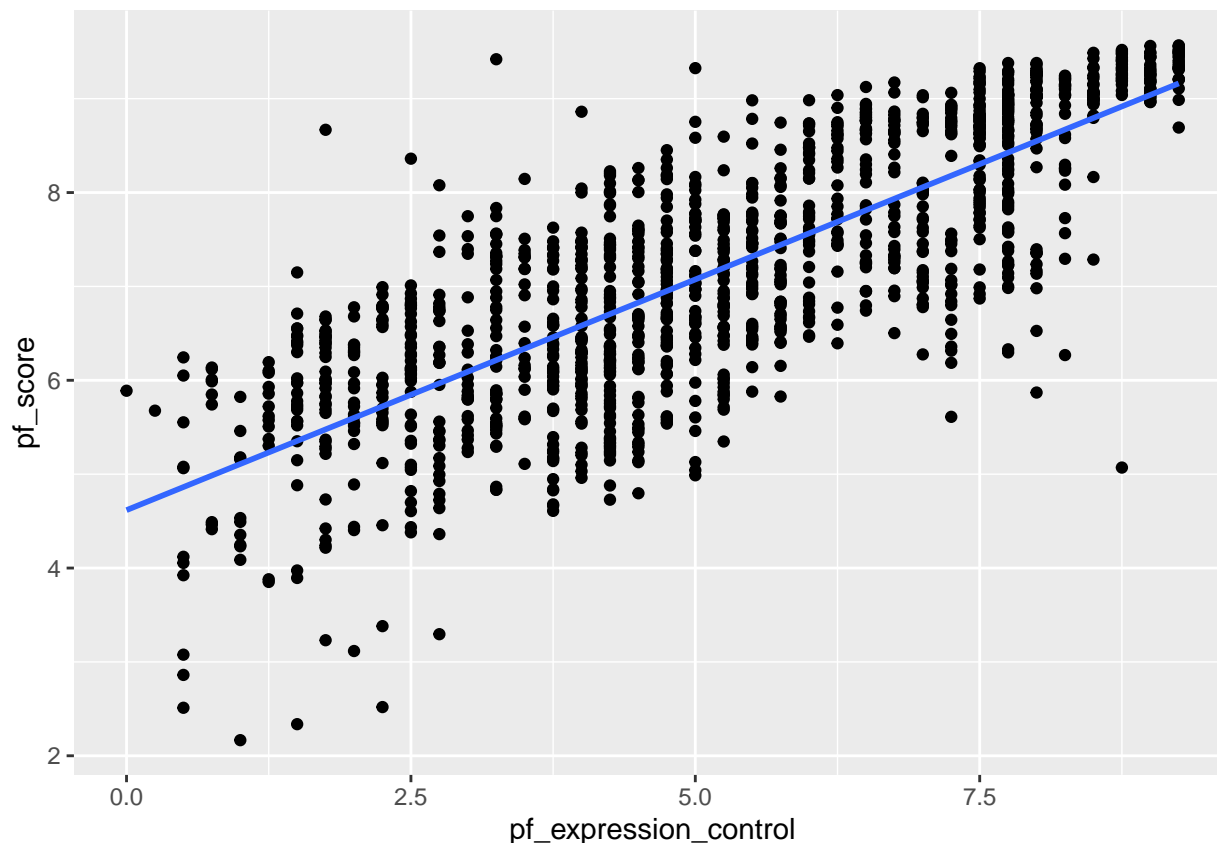
summary(lm_mod)

##
## Call:
## lm(formula = hf_score ~ pf_expression_control, data = hfi)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.6198 -0.4908  0.1031  0.4703  2.2933
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    5.153687   0.046070  111.87  <2e-16 ***
## pf_expression_control 0.349862   0.008067   43.37  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.667 on 1376 degrees of freedom
## (80 observations deleted due to missingness)
## Multiple R-squared:  0.5775, Adjusted R-squared:  0.5772
## F-statistic: 1881 on 1 and 1376 DF, p-value: < 2.2e-16
```

Prediction and prediction errors

Let's create a scatterplot with the least squares line for m1 laid on top.

```
ggplot(data = hfi, aes(x = pf_expression_control, y = pf_score)) +
  geom_point() +
  stat_smooth(method = "lm", se = FALSE)
```



Here, we are literally adding a layer on top of our plot. `geom_smooth` creates the line by fitting a linear model. It can also show us the standard error `se` associated with our line, but we'll suppress that for now.

This line can be used to predict y at any value of x . When predictions are made for values of x that are beyond the range of the observed data, it is referred to as *extrapolation* and is not usually recommended. However, predictions made within the range of the data are more reliable. They're also used to compute the residuals.

6. If someone saw the least squares regression line and not the actual data, how would they predict a country's personal freedom school for one with a 6.7 rating for `pf_expression_control`? Is this an overestimate or an underestimate, and by how much? In other words, what is the residual for this prediction?

We can find where 6.7 lands on the x-axis (`pf_expression_control`) and move up to where the regression line meets; roughly a 7.8 `pf_score`. However, let's be more accurate using the `predict` function and our regression equation which tells us it is 7.49776.

```
test_data <- data.frame(pf_expression_control = 6.7)

predict(lm_mod, test_data)
```

```
##          1
## 7.49776
```

Given that the dataset does not contain values for `pf_expression_control = 6.7`, as it has 0.25 increments, I will collect the next closest data points which is 6.75. Let's now determine the residuals for each one and calculate the average.

```
original_pts <-
  hfi |>
  select(pf_expression_control, pf_score) |>
  filter(pf_expression_control == 6.75)

resid_diff <-
  original_pts$pf_score - predict(lm_mod, test_data)

mean(resid_diff)

## [1] 0.5085551
```

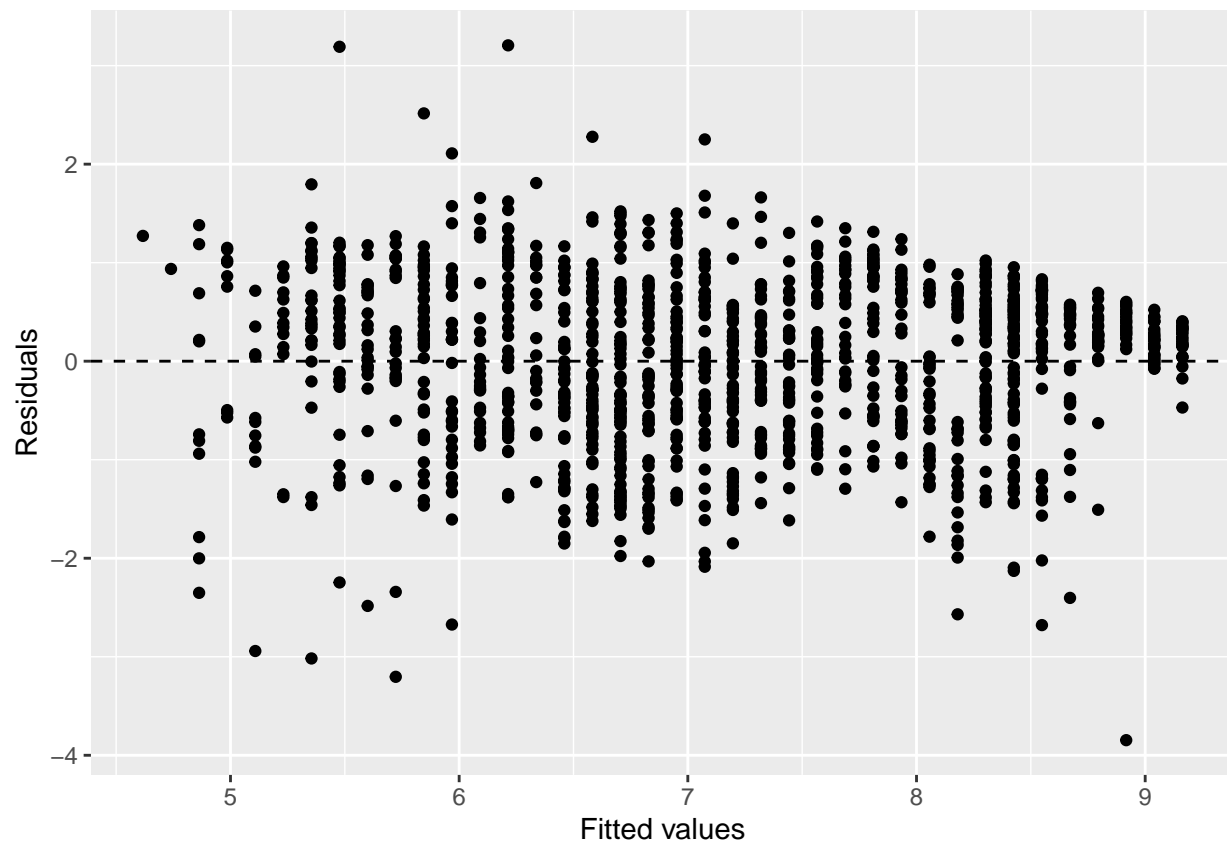
The residuals tell us that we will either overestimate or underestimate given that the linear regression model is trying to reduce the sum of squared errors. However, looking at the average of the residuals, we will be generally underestimating what 6.7 should be.

Model diagnostics

To assess whether the linear model is reliable, we need to check for (1) linearity, (2) nearly normal residuals, and (3) constant variability.

Linearity: You already checked if the relationship between `pf_score` and ‘`pf_expression_control`’ is linear using a scatterplot. We should also verify this condition with a plot of the residuals vs. fitted (predicted) values.

```
ggplot(data = m1, aes(x = .fitted, y = .resid)) +
  geom_point() +
  geom_hline(yintercept = 0, linetype = "dashed") +
  xlab("Fitted values") +
  ylab("Residuals")
```



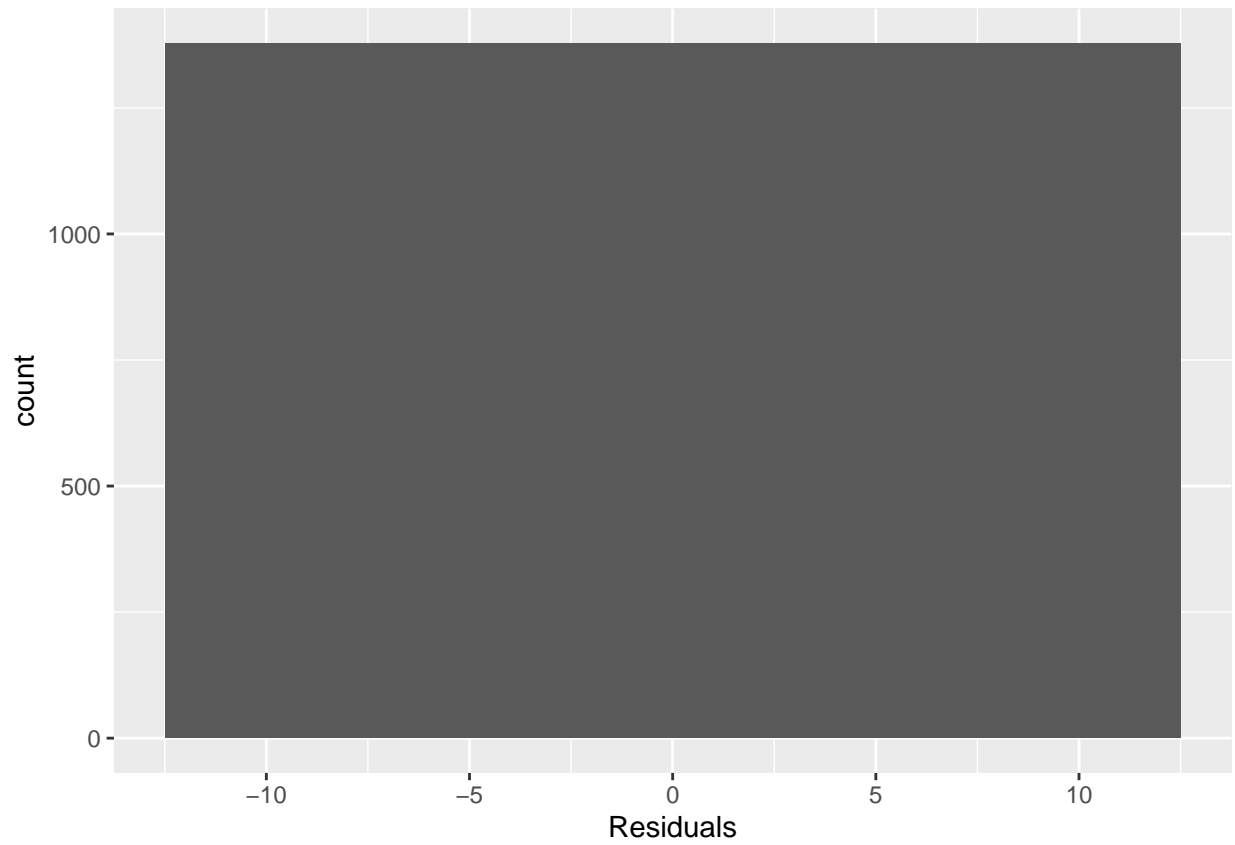
Notice here that `m1` can also serve as a data set because stored within it are the fitted values (\hat{y}) and the residuals. Also note that we're getting fancy with the code here. After creating the scatterplot on the first layer (first line of code), we overlay a horizontal dashed line at $y = 0$ (to help us check whether residuals are distributed around 0), and we also rename the axis labels to be more informative.

7. Is there any apparent pattern in the residuals plot? What does this indicate about the linearity of the relationship between the two variables?

Overall there does not seem to be any unusual pattern as the values are scattered around zero. This means that in general, there is a linear relationship between both variables.

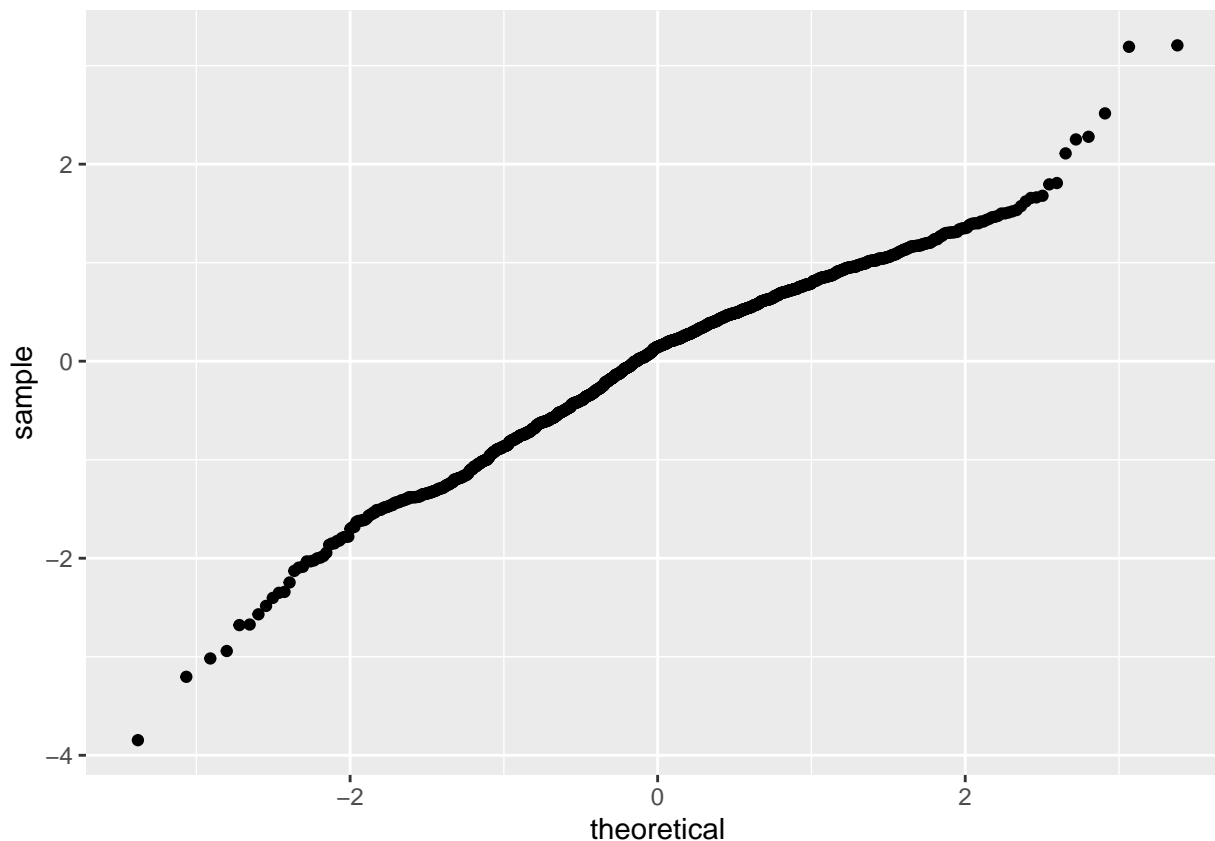
Nearly normal residuals: To check this condition, we can look at a histogram

```
ggplot(data = m1, aes(x = .resid)) +
  geom_histogram(binwidth = 25) +
  xlab("Residuals")
```

or a normal probability plot of the residuals.

```
ggplot(data = m1, aes(sample = .resid)) +  
  stat_qq()
```



Note that the syntax for making a normal probability plot is a bit different than what you're used to seeing: we set `sample` equal to the residuals instead of `x`, and we set a statistical method `qq`, which stands for “quantile-quantile”, another name commonly used for normal probability plots.

8. Based on the histogram and the normal probability plot, does the nearly normal residuals condition appear to be met?

Yes the normal condition appears to be met with some slight variation in the tails.

Constant variability:

9. Based on the residuals vs. fitted plot, does the constant variability condition appear to be met?

Yes, the constant variability condition is met as there is not a specific pattern being seen.

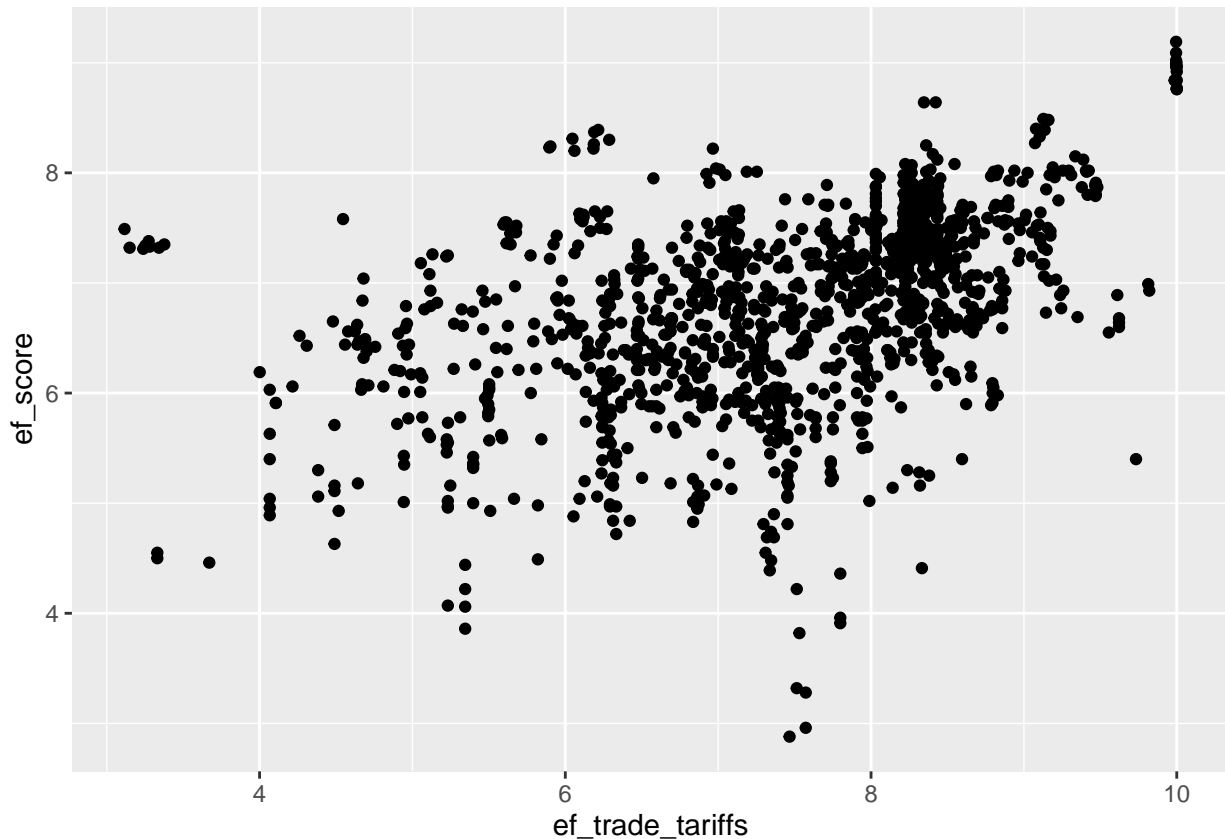
More Practice

- Choose another freedom variable and a variable you think would strongly correlate with it.. Produce a scatterplot of the two variables and fit a linear model. At a glance, does there seem to be a linear relationship?

`ef_trade_tariffs` vs. `ef_score`

There does seem to be a linear relationship, however quite weak.

```
hfi |>
  ggplot(aes(x = ef_trade_tariffs, y = ef_score)) +
  geom_point()
```



- How does this relationship compare to the relationship between `pf_expression_control` and `pf_score`? Use the R^2 values from the two model summaries to compare. Does your independent variable seem to predict your dependent one better? Why or why not?

```
lm_mod_2 <- lm(ef_score ~ ef_trade_tariffs, data = hfi)
```

```
summary(lm_mod_2)
```

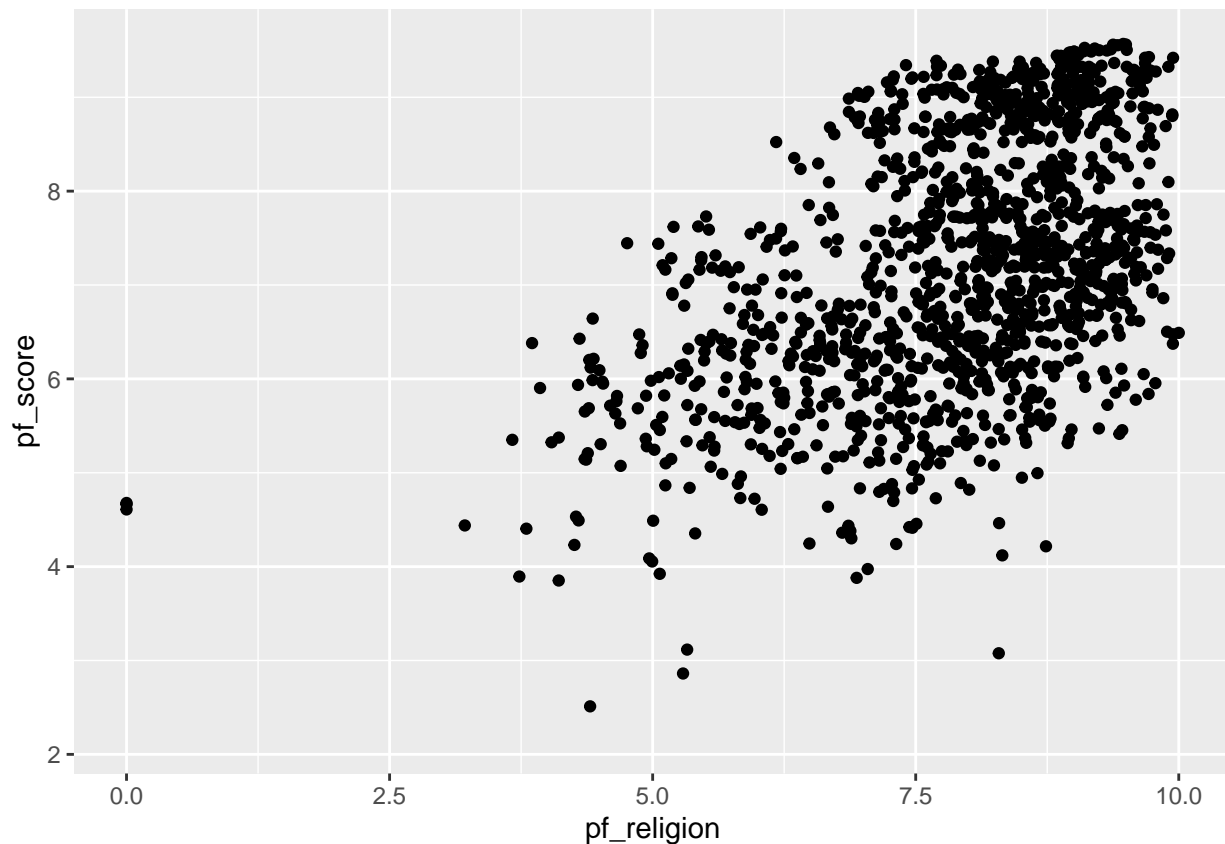
```
##
## Call:
## lm(formula = ef_score ~ ef_trade_tariffs, data = hfi)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.9200 -0.4577  0.1304  0.5004  2.1769
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    4.24787    0.12772   33.26  <2e-16 ***
## ef_trade_tariffs 0.34178    0.01693   20.19  <2e-16 ***
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.7715 on 1371 degrees of freedom
## (85 observations deleted due to missingness)
## Multiple R-squared:  0.2292, Adjusted R-squared:  0.2286
## F-statistic: 407.6 on 1 and 1371 DF,  p-value: < 2.2e-16
```

The linear model between `ef_trade_tariffs` vs. `ef_score` is significantly weaker with an R^2 of 0.2292 compared to `pf_expression_control` and `pf_score` and an R^2 of 0.6342

- What's one freedom relationship you were most surprised about and why? Display the model diagnostics for the regression model analyzing this relationship.

```
hfi |>
  ggplot(aes(x = pf_religion, y = pf_score)) +
  geom_point()
```



```
lm_mod_2 <- lm(pf_score ~ pf_religion, data = hfi)
summary(lm_mod_2)
```

```
##
## Call:
```

```
## lm(formula = pf_score ~ pf_religion, data = hfi)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.3515 -0.8728 -0.0744  1.0532  2.3520
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  3.30981    0.18744   17.66  <2e-16 ***
## pf_religion  0.49689    0.02346   21.18  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.17 on 1366 degrees of freedom
## (90 observations deleted due to missingness)
## Multiple R-squared:  0.2472, Adjusted R-squared:  0.2467
## F-statistic: 448.6 on 1 and 1366 DF,  p-value: < 2.2e-16
```

Looking into the relationship between pf_score vs. pf_religion I would have thought it would have a stronger association then it shows. This may be due to the skew viewpoints living in the United States and how religion is portrayed within media. Overall throughout the globe, religion does not seem to be a major point between personal freedoms of religion to their overall personal freedom score.
