

Probability Distributions

2018-01-17

Preparation

```
# Load libraries
library(tidyverse)
library(broom)
library(sm)
```

Normal Distribution

The probability distribution of a normal distribution is defined as

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left[-\frac{(x - \mu)^2}{2\sigma^2} \right]$$

for $-\infty \leq x \leq \infty$. Consider a normal distribution with a mean (μ) of 50, and a standard deviation (σ) of 10. We can compute the probability density ($p(x)$) for a particular x value by using this equation. For example, the probability density for $x = 65$ can be found using,

$$p(65) = \frac{1}{10\sqrt{2\pi}} \exp \left[-\frac{(65 - 50)^2}{2 \times 10^2} \right] = 0.01295176$$

Using R, we can carry out the computation,

```
(1/(10*sqrt(2*pi))) * exp(-(225)/200)
```

```
## [1] 0.01295176
```

There is also a more direct way to compute this using the `dnorm()` function. This function computes the density of x from a normal distribution with a specified mean and sd.

```
dnorm(x = 65, mean = 50, sd = 10)
```

```
## [1] 0.01295176
```

If we compute the density for several x values and plot them, we get the familiar normal shape; the graphical instantiation of the mathematical equation.

```
library(tidyverse)
data.frame(
  X = seq(from = 10, to = 90, by = 0.01)
) %>%
  rowwise() %>%
  mutate( Y = dnorm(x = X, mean = 50, sd = 10) ) %>%
  ggplot(data = ., aes(x = X, y = Y)) +
    geom_line() +
    theme_bw() +
    geom_point(x = 65, y = 0.01295176, size = 3)
```

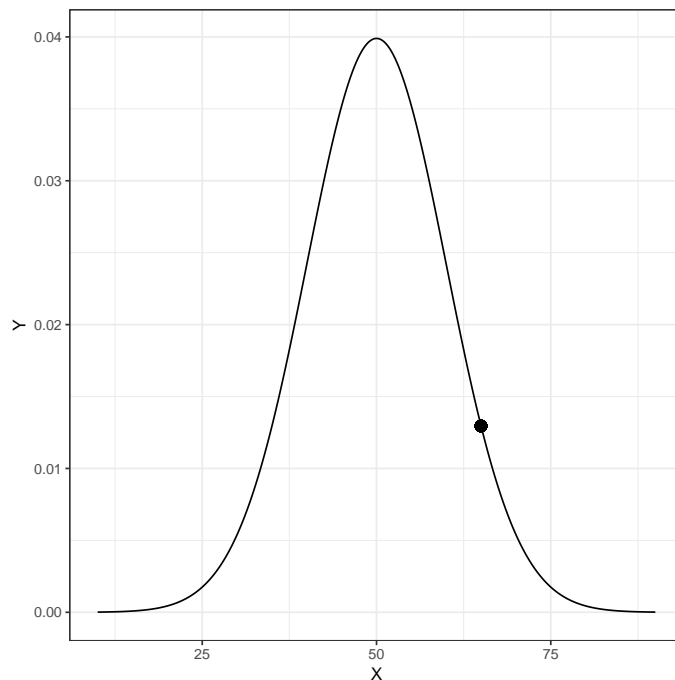


Figure 1. Plot of the probability density function (PDF) for a Normal distribution with mean of 50 and standard deviation of 10. The density value for $x = 65$, $p(65) = 0.01295176$, is also displayed on the PDF.

Other Useful R Functions for Working with Probability Distributions

There are four primary functions for working with the normal probability distribution:

- `dnorm()` : To compute the probability density (point on the curve)
- `pnorm()` : To compute the probability (area under the PDF)
- `qnorm()` : To compute the x value given a particular probability
- `rnorm()` : To draw a random observation from the distribution

Each of these requires the arguments `mean=` and `sd=`. Let's look at some of them in use.

Finding Cumulative Probability

The function `pnorm()` gives the probability x is less than or equal to some quantile value in the distribution; the cumulative probability. For example, to find the probability that $x \leq 65$ we would use,

```
pnorm(q = 65, mean = 50, sd = 10)
```

```
## [1] 0.9331928
```

This is akin to finding the proportion of the area under the normal PDF that is to the left of 65.

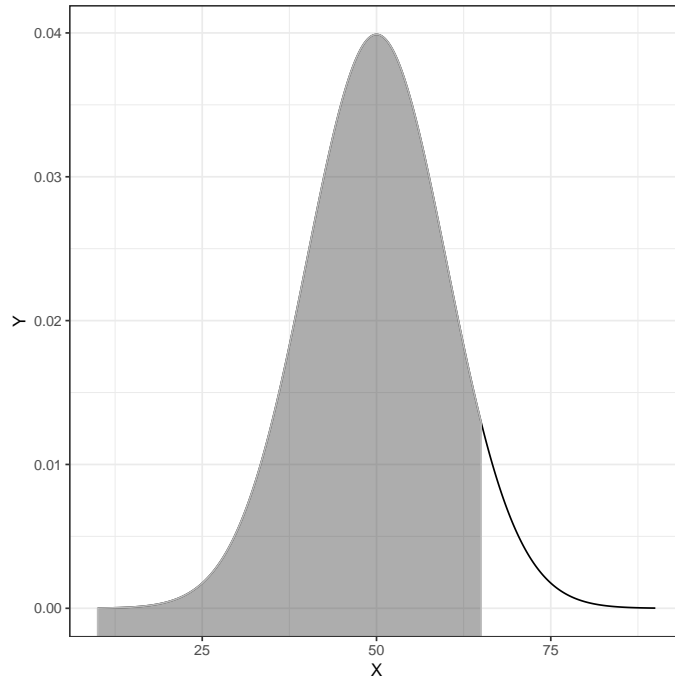


Figure 2. Plot of the PDF for a $\sim \mathcal{N}(50, 10)$ with the cumulative probability for $x \leq 65$ shaded.

For the mathematically inclined, the grey-shaded area is expressed as an integral

$$\int_{-\infty}^{65} p(x) dx$$

where $p(x)$ is the PDF for the normal distribution.

Cumulative Density and p -Value

This type of computation is used most commonly to find a p -value. The p -value is just the area under the distribution (curve) that is AT LEAST as extreme as some observed value. Consider a hypothesis test of whether a population parameter is equal to 0. Also consider that we observed a statistic (that has been standardized) of $z = 2.5$. Then, the p -value can be graphically displayed in the standard normal distribution as follows:

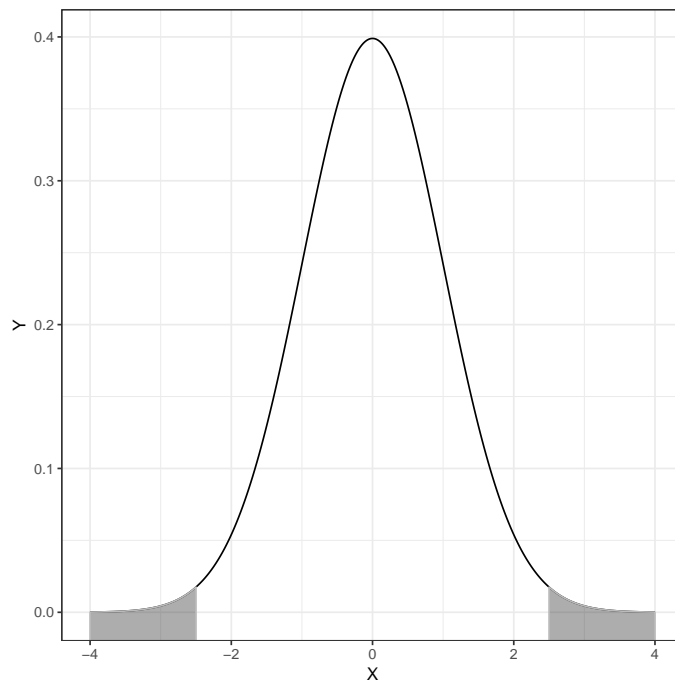


Figure 3. Plot of the probability density function (PDF) for the standard normal distribution ($M = 0$, $SD = 1$). The cumulative density representing the p -value for a two-tailed test evaluating whether $\mu = 0$ using an observed z of 2.5 is also displayed.

In most hypothesis tests, we test whether the parameter IS EQUAL to 0. Thus the values in the standard normal distribution more extreme than 2.5 encompass evidence against the hypothesis; those values greater than 2.5 and also those values less than -2.5 . (This is akin to testing a fair coin when both 8 heads OR 8 tails would provide evidence against fairness ...we have to consider evidence in both directions).

To compute this we use `pnorm()`. Remember, it computes the proportion of the area under the curve TO THE LEFT of a particular value. Here we will compute the area to the left of -2.5 and then double it to produce the actual p -value.

```
2 * pnorm(q = -2.5, mean = 0, sd = 1)
```

```
## [1] 0.01241933
```

Finding Quantiles

The `qnorm()` function is essentially the inverse of the `pnorm()` function. The `p` functions find the cumulative probability GIVEN a particular quantile. The `q` functions find the quantile GIVEN a cumulative probability. For example, in the normal distribution we defined earlier, half of the area is below the quantile value of 50 (the mean).

```
qnorm(p = 0.5, mean = 50, sd = 10)
```

```
## [1] 50
```

Student's t -Distribution

Student's t -distribution looks like a standard normal distribution. In the figure below, Student's t -distribution is depicted with a solid, black line and the standard normal distribution ($M = 0$, $SD = 1$) is depicted with a dotted, red line.

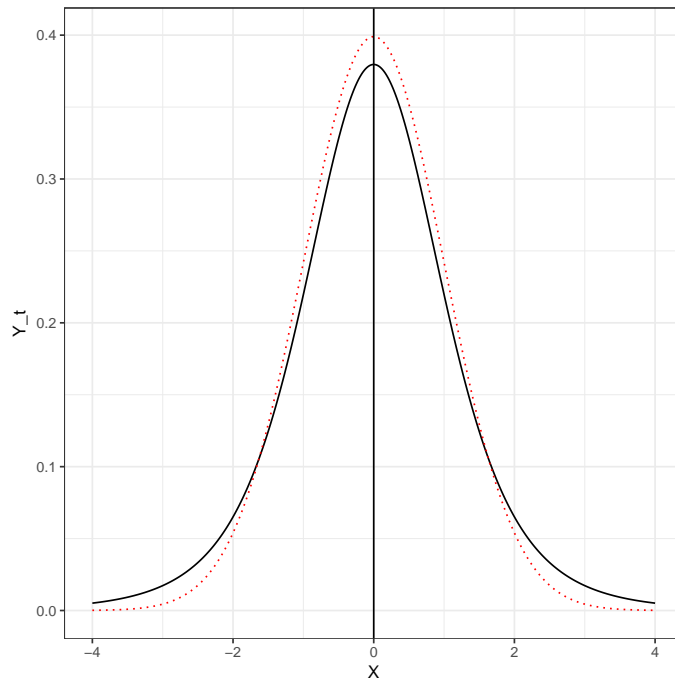


Figure 4. Plot of the probability density function (PDF) for the standard normal distribution (dotted, red line) and Student's t -distribution with 5 df (solid, black line).

Both the standard normal distribution and Student's t -distribution have a mean (expected value) of 0. The standard deviation for Student's t -distribution is larger than the standard deviation for the standard normal distribution ($SD > 1$). You can see this in the distribution because the tails in Student's t -distribution are fatter (more error) than the standard normal distribution.

In practice, we often use Student's t -distribution rather than the standard normal distribution when we are using sample data to estimate the population. This estimation increases the error and thus is typically modeled using Student's t -distribution.

Student's t -distribution constitutes a family of distributions—not just a single distribution. The specific shape (and thus probability density) is defined by the *degrees of freedom*; df . The plot below shows the standard normal distribution (purple) and four t -distributions with varying df -values.

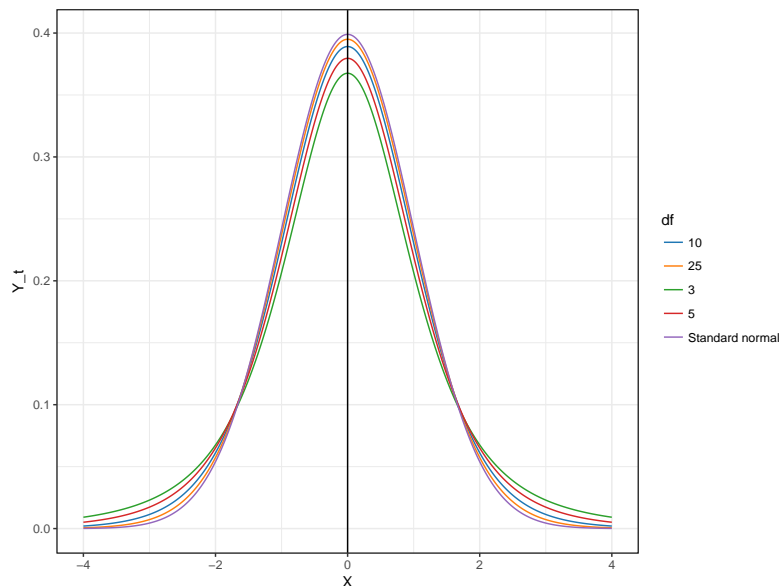


Figure 5. Plot of several t -distributions with differing df .

If we compare the means and SDs for these distributions,

df	M	SD
3	0	2.00
5	0	1.50
10	0	1.22
25	0	1.08
z	0	1.00

we find that the mean for all the t -distributions is 0, same as the standard normal distribution. All t -distributions are unimodal and symmetric around zero. The SD for every t -distribution is higher than the SD for the standard normal distribution. Student t -distributions with higher df values have less variation. It turns out that the standard normal distribution is a t -distribution with ∞ df . For the formula for the SD in a t -distribution, see Fox (2009).

There are four primary functions for working with Student's t -distribution:

- `dt()` : To compute the probability density (point on the curve)
- `pt()` : To compute the probability (area under the PDF)
- `qt()` : To compute the x value given a particular probability
- `rt()` : To draw a random observation from the distribution

Each of these requires the arguments `df=`. Let's look at some of them in use.

Comparing Probability Densities

How do the probability densities for a value of X compare across these distributions? Let's examine the X value of 2.

```
# Standard normal distribution
pnorm(q = 2, mean = 0, sd = 1)
```

```
## [1] 0.9772499
```

```
# t-distribution with 3 df
pt(q = 2, df = 3)
```

```
## [1] 0.930337
```

```
# t-distribution with 5 df
pt(q = 2, df = 5)
```

```
## [1] 0.9490303
```

```
# t-distribution with 10 df
pt(q = 2, df = 10)
```

```
## [1] 0.963306
```

```
# t-distribution with 25 df
pt(q = 2, df = 25)
```

```
## [1] 0.971762
```

We are essentially comparing the height of these distributions at $X = 2$.

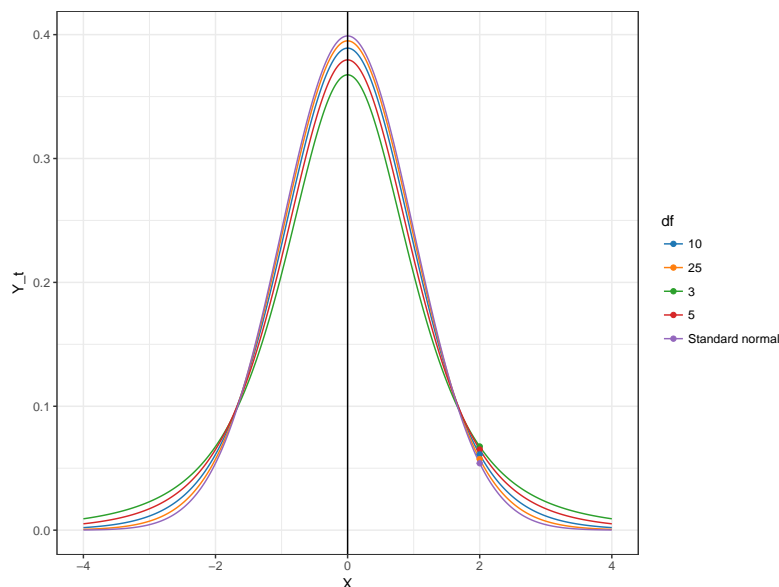


Figure 6. Plot of several t -distributions with differing df . The probability density for $t = 2$ is also displayed for each of the distributions.

Comparing Cumulative Densities

What if we wanted to look at cumulative density? Consider our hypothesis test of whether a population parameter is equal to 0. Also consider that we observed a statistic (that has been standardized) of 2.5 using a sample size of $n = 15$.

If we can assume that the SAMPLING DISTRIBUTION (or bootstrap distribution) is normally-distributed then we can use the cumulative density in a normal distribution to compute a p -value:

```
2 * pnorm(q = -2.5, mean = 0, sd = 1)
```

```
## [1] 0.01241933
```

If, however, the SAMPLING DISTRIBUTION (or bootstrap distribution) is t -distributed then we need to use the cumulative density for a t -distribution with the appropriate df to compute a p -value. For example if we use $df = n - 1$, the two-tailed p -value would be:

```
2 * pt(q = -2.5, df = 14)
```

```
## [1] 0.02546666
```

The p -value using the t -distribution is larger than the p -value computed based on the standard normal distribution. This is again because of the increased error (uncertainty) we are introducing when we estimate from sample. This added uncertainty makes it harder for us to reject a hypothesis.

Using the t -Distribution in Regression

```
# Read in data
city = read_csv(file = "~/Dropbox/epsy-8252/data/riverside.csv")
head(city)
```

```
## # A tibble: 6 x 6
##   education income seniority gender  male party
##   <int> <int> <int> <chr> <int> <chr>
## 1      8  37449      7 male      1 Democrat
## 2      8  26430      9 female    0 Independent
## 3     10  47034     14 male      1 Democrat
## 4     10  34182     16 female    0 Independent
## 5     10  25479      1 female    0 Republican
## 6     12  46488     11 female    0 Democrat
```

```
# Fit regression model
lm.1 = lm(income ~ 1 + education + seniority, data = city)

# Coefficient-level estimates
tidy(lm.1)
```

```
##           term estimate std.error statistic    p.value
## 1 (Intercept) 6769.1720 5372.8914  1.259875 2.177593e-01
## 2   education 2251.8456  334.6443  6.729073 2.202903e-07
## 3   seniority  738.7965  210.0954  3.516481 1.459777e-03
```

How do we obtain the p -value for each of the coefficients? Recall that the coefficients and SEs for the coefficients are computed directly from the raw data. Then we can compute a test-statistic by dividing the coefficient estimate by the SE. For example, to compute the education test statistic,

$$t = \frac{2252}{335} = 6.72$$

Since we are estimating the SE using sample data, our test statistic is likely t -distributed. Which value should we use for df ? Well, for that, statistical theory tells us that we should use the error df value. In our data,

$$\begin{aligned} n &= 32 \\ \text{Total } df &= 32 - 1 = 31 \\ \text{Model } df &= 2 \text{ (two predictors)} \\ \text{Error } df &= 31 - 2 = 29 \end{aligned}$$

Using the t -distribution with 29 df ,


```
2 * pt(q = -6.72, df = 29)
```

```
## [1] 2.257125e-07
```

For seniority (and the intercept), we would use the same t -distribution, but our test statistic would differ:

$$t_{\text{Intercept}} = \frac{6769}{5373} = 1.26$$
$$t_{\text{Seniority}} = \frac{739}{210} = 3.52$$

The associated p -values are:

```
# Intercept p-value
```

```
2 * pt(q = -1.26, df = 29)
```

```
## [1] 0.2177149
```

```
# Seniority p-value
```

```
2 * pt(q = -3.52, df = 29)
```

```
## [1] 0.001446316
```

Model-Level Inference: The F -Distribution

The model-level inference for regression was based on an F -statistic.

```
glance(lm.1)
```

```
##   r.squared adj.r.squared   sigma statistic    p.value df    logLik
## 1 0.7417857   0.7239778 7645.854   41.6549 2.976563e-09  3 -329.9724
##      AIC      BIC  deviance df.residual
## 1 667.9448 673.8077 1695313285         29
```

In this case, the test of

$$H_0 : \rho^2 = 0$$

is not statistically significant, $F(2, 29) = 41.7$, $p < .001$. In this test, our test statistic is R^2 . When we standardize this test statistic, we obtain an F -value. The F -distribution's shape is based on two df values—in this case 2 and 29. In the figure below, we show this F -distribution as a black line.

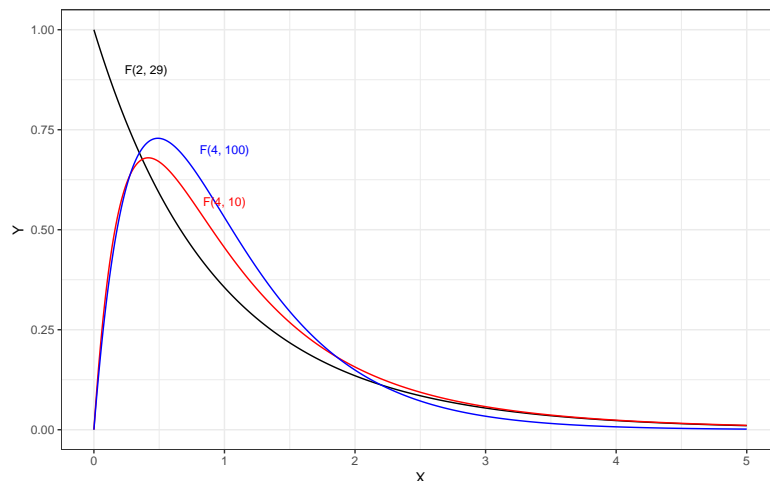


Figure 7. Plot of several F -distributions with differing df .

The F -distribution, like the t -distribution is a family of distributions. They are positively skewed and generally have a lower-limit of 0. Because of this, when we use the F -distribution to compute a p -value, we only compute the cumulative density GREATER THAN OR EQUAL TO the value of the standradized test statistic.

To compute the model-level F -statistic we use the sum of squares partitioning from the ANOVA table.

```
anova(lm.1)
```

```
## Analysis of Variance Table
##
## Response: income
##          Df      Sum Sq   Mean Sq F value    Pr(>F)
## education  1 4147330492 4147330492  70.944 2.781e-09 ***
## seniority  1  722883649  722883649  12.366  0.00146 **
## Residuals 29 1695313285   58459079
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

We first need to compute the model and error mean squares (MS). To compute a mean square we use the general computation

$$MS = \frac{SS}{df}$$

The model includes both the education and seniority predictor, so we combine the SS and df.

$$\begin{aligned} MS_{\text{Model}} &= \frac{SS_{\text{Model}}}{df_{\text{Model}}} \\ &= \frac{4147330492 + 722883649}{1 + 1} \\ &= \frac{4870214141}{2} \\ &= 2435107070 \end{aligned}$$

The error MS is

$$\begin{aligned} MS_{\text{Error}} &= \frac{SS_{\text{Error}}}{df_{\text{Error}}} \\ &= \frac{1695313285}{29} \\ &= 58459079 \end{aligned}$$

Then, we compute the F -statistic by computing the ratio of these two mean squares.

$$\begin{aligned} F &= \frac{MS_{\text{Model}}}{MS_{\text{Error}}} \\ &= \frac{2435107070}{58459079} \\ &= 41.7 \end{aligned}$$

This is the observed F -statistic for the model. To test whether $\rho^2 = 0$, we evaluate this in an F -distribution with the appropriate df for the model and error—namely, 2 and 29.

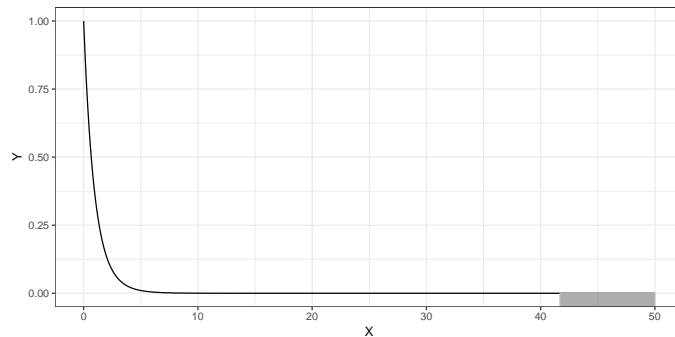


Figure 8. Plot of the probability density function (PDF) for the F -distribution with 2 and 29 df. The cumulative density representing the p -value for a two-tailed test evaluating whether $\rho^2 = 0$ using an observed F of 41.7 is also displayed.

The computation using the cumulative density `pf()` function is:

```
1 - pf(41.7, df1 = 2, df2 = 29)
```

```
## [1] 2.942114e-09
```

Mean Squares are Variance Estimates

Mean squares are estimates of the variance. Consider the computational formula for the sample variance,

$$\hat{\sigma}^2 = \frac{\sum(Y - \bar{Y})^2}{n - 1}$$

This is the total sum of squares divided by the total df . When we compute an F -statistic, we are finding the ratio of two different variance estimates—one based on the model (explained variance) and one based on the error (unexplained variance). Under the null hypothesis that $\rho^2 = 0$, we are assuming that all the variance is unexplained. In that case, our F -statistic would be close to zero. When the model explains a significant amount of variation, the numerator gets larger relative to the denominator and the F -value is larger.

Fox, J. (2009). *A mathematical primer for social statistics*. Thousand Oaks, CA: Sage.