
Linear Regression with One Regressor: Inference

Osman DOĞAN

Linear Regression with One Regressor: Inference

- Linear Regression with One Regressor: Inference
 - ① Simple hypothesis testing (or significance testing)
 - ② Confidence intervals for a regression coefficient
 - ③ Regression when the regressor is binary
 - ④ Heteroskedasticity vs Homoskedasticity
 - ⑤ The Gauss-Markov theorem
- Readings:
 - ① Stock and Watson (2020, Chapter 5),
 - ② Hanck et al. (2021, Chapter 5).

Simple hypothesis testing

- A hypothesis is simply a claim and a simple hypothesis in the regression analysis means a claim involving a single coefficient, often about the slope of the population regression function.
- The empirical analysis may require the researcher to test a hypothesis on β_1 .
- For example, in our class-size test-scores example, a taxpayer's claim could be that **cutting class size will not help boost test scores**.
- In other words, under this claim, the slope of the population regression function is zero, i.e., $\beta_1 = 0$.
- Can you reject the taxpayer's hypothesis that $\beta_1 = 0$, or should you accept it, at least tentatively pending further new evidence?

Simple hypothesis testing

- Our objective is to test a hypothesis, like $\beta_1 = 0$, using our sample to reach a tentative conclusion whether the hypothesis is correct or incorrect.
- To this end, we will start by stating the null and the alternative hypotheses.
- The null hypothesis will state the claim to be tested and we use H_0 to denote it. In our example, $H_0 : \beta_1 = 0$.
- The alternative hypothesis will state the scenario when the null hypothesis does not hold and we use H_1 to denote it:
 - 1 If $H_1 : \beta_1 \neq 0$ it is called a **two-sided** alternative.
 - 2 If $H_1 : \beta_1 < 0$ or $H_1 : \beta_1 > 0$, it is called a **one-sided** alternative.

Simple hypothesis testing

- In order to decide between H_0 and H_1 , we need a test statistics that will make use of the sample data and help us to reach a tentative conclusion.
- We will use the *t-statistic*, which is defined as

$$t = \frac{\text{estimator-hypothesized value}}{\text{standard error of the estimator}} \quad (1)$$

- Thus, the *t*-statistic for testing $H_0 : \beta_1 = c$, where c is the hypothesized value, takes the following form:

$$t = \frac{\hat{\beta}_1 - c}{SE(\hat{\beta}_1)}, \quad (2)$$

where $SE(\hat{\beta}_1)$ is the square root of the variance of $\hat{\beta}_1$.

- To compute *t*-statistic, we need to figure out $SE(\hat{\beta}_1)$.
- Also, since *t*-statistic will inherit the uncertainty due to sampling, we will have to figure out its sampling distribution to reach a tentative conclusion.

Simple hypothesis testing

- Recall from Chapter 4 that under the least squares assumptions, in large samples, we have

$$\hat{\beta}_1 \stackrel{A}{\sim} N\left(\beta_1, \sigma_{\hat{\beta}_1}^2\right), \text{ where } \sigma_{\hat{\beta}_1}^2 = \frac{1}{n} \frac{\text{var}((X_i - \mu_X)u_i)}{(\text{var}(X_i))^2}. \quad (3)$$

- The idea is to replace the unknown terms in $\sigma_{\hat{\beta}_1}^2$ with their sample counterparts. Thus, an estimator of $\sigma_{\hat{\beta}_1}^2$ is

$$\hat{\sigma}_{\hat{\beta}_1}^2 = \frac{1}{n} \frac{\frac{1}{n-2} \sum_{i=1}^n (X_i - \bar{X})^2 \hat{u}_i^2}{\left(\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2\right)^2} \quad (4)$$

- Then, $SE(\hat{\beta}_1)$ in (2) can be computed through $SE(\hat{\beta}_1) = \sqrt{\hat{\sigma}_{\hat{\beta}_1}^2}$.

Simple hypothesis testing

- Also, under the least squares assumptions, in large samples, we have

$$\hat{\beta}_1 \stackrel{A}{\sim} N\left(\beta_1, \sigma_{\hat{\beta}_1}^2\right) \Rightarrow \frac{\hat{\beta}_1 - \beta_1}{\sigma_{\hat{\beta}_1}} \equiv Z \stackrel{A}{\sim} N(0, 1) \quad (5)$$

- Also, (5) suggests that if H_0 is true then, the sampling distribution of the t -statistic can be approximated by $N(0, 1)$ in large samples, i.e., $t \stackrel{A}{\sim} N(0, 1)$.
- Recall the types of errors you can make in a testing problem (Chapter 3)
 - we reject a TRUE H_0 , i.e., **Type-I error**,
 - we fail to reject a FALSE H_0 , i.e., **Type-II error**.
- We choose a (significance) level for committing the Type-I error, and then minimize the likelihood of committing Type-II error.
- The conventional significance levels are 1%, 5% and 10%.
- We can now use the level and the sampling distribution of the test to define the rejection regions.

Simple hypothesis testing

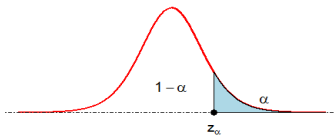
- Figure 1 illustrates the rejection regions according to the type of H_1 .
 - In each case, the shaded blue area is the level of test.
 - Depending on the type of H_1 , z_α , $-z_\alpha$, $-z_{\alpha/2}$ and $z_{\alpha/2}$ are critical values.
- Finally, we will reject H_0 if the t -statistic value falls in the rejection regions (the shaded blue areas).

Example 1

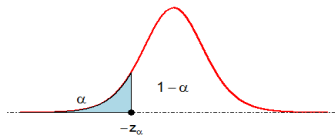
Consider $H_1 : \beta_1 \neq c$ in Figure 1c, and assume that $\alpha = 5\%$. Then, the critical value on the left tail is the 2.5th percentile of $N(0, 1)$, which is $-z_{\alpha/2} = -1.96$, and the critical value on the right tail is the 97.5th percentile of $N(0, 1)$, which is $z_{\alpha/2} = 1.96$. In R, we can use the following:

```
qnorm(0.025, mean = 0, sd=1, lower.tail = TRUE)=-1.96
```

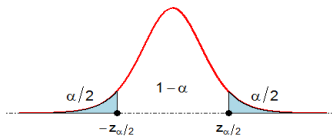
Assume that we obtained $t = \frac{\hat{\beta}_1 - c}{SE(\hat{\beta}_1)} = 2$. Then, we will reject H_0 because the test statistic value is in the rejection region. i.e., $t = 2 > z_{\alpha/2} = 1.96$.



(a) $H_0 : \beta_1 = c$ versus $H_1 : \beta_1 > c$



(b) $H_0 : \beta_1 = c$ versus $H_1 : \beta_1 < c$



(c) $H_0 : \beta_1 = c$ versus $H_1 : \beta_1 \neq c$

Figure 1: Rejection Regions

Hypothesis Testing

- Alternatively, we can also calculate a ***p*-value** to decide between H_0 and H_1 :

$$p\text{-value} = \begin{cases} P_{H_0} (t > |t_{\text{calc}}|) & \text{for } H_1 : \beta_1 > c, \\ P_{H_0} (t < -|t_{\text{calc}}|) & \text{for } H_1 : \beta_1 < c, \\ P_{H_0} (|t| > |t_{\text{calc}}|) & \text{for } H_1 : \beta_1 \neq c. \end{cases}$$

where t_{calc} is the value of the test statistic obtained from (2).

- P_{H_0} means this probability is calculated from the distribution of the test statistic under the null hypothesis.
- Since the asymptotic distribution of t -statistic is $N(0, 1)$, we have

$$p\text{-value} = \begin{cases} P_{H_0} (t > |t_{\text{calc}}|) = 1 - \Phi(|t_{\text{calc}}|) & \text{for } H_1 : \beta_1 > c, \\ P_{H_0} (t < -|t_{\text{calc}}|) = \Phi(-|t_{\text{calc}}|) & \text{for } H_1 : \beta_1 < c, \\ P_{H_0} (|t| > |t_{\text{calc}}|) = 2 \times \Phi(-|t_{\text{calc}}|) & \text{for } H_1 : \beta_1 \neq c. \end{cases}$$

Hypothesis Testing

- What does the p -value tell us?
- It gives the likelihood of obtaining a test statistic value that is more extreme than the actual one when the null is correct.
- Hence, the smaller p -value, the less likely that the null is correct. If the level is chosen as 5%, then reject H_0 if p -value is less than 5%.

Example 2

Consider $H_1 : \beta_1 \neq c$ in Figure 1c, and assume that $\alpha = 5\%$. Assume that we obtained $t = 2$ from (2). Then,

$$p\text{-value} = P_{H_0} (|t| > 2) = P_{H_0} (t > 2) + P_{H_0} (t < -2) = 2 \times \Phi(-2).$$

In R, we can compute the p -value in the following way:

```
2*pnorm(-2,mean = 0,sd=1,lower.tail = TRUE)= 0.045
```

Since $p\text{-value} = 0.045 < \alpha = 0.05$, we reject H_0 .

Simple hypothesis testing

- Consider the test score example:

$$\text{TestScore}_i = \beta_0 + \beta_1 \text{STR}_i + u_i. \quad (6)$$

- In R, the `lm` function computes the t -statistic for each regressor and provides the corresponding p -values for a two sided null hypothesis. The estimation results are given in Listing 1.
- Consider $H_0 : \beta_1 = 0$ versus $H_1 : \beta_1 \neq 0$ and assume $\alpha = 5\%$. Then, we can compute t -statistic as

$$t = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)} = \frac{-2.28}{0.48} \approx -4.75 \quad (7)$$

- The critical value is $z_{\alpha/2} = -1.96$. Since $t = -4.75 < z_{\alpha/2} = -1.96$. We reject H_0 and conclude that β_1 is statistically significant at 5% level.
- Alternatively, we can compute $p\text{-value} = 2 \times \Phi(-|t_{\text{calc}}|) = 2 \times \Phi(-4.75)$.

```
2*pnorm(-4.75, mean = 0, sd=1, lower.tail = TRUE)=2.034166e-06
```

Simple hypothesis testing

Listing 1: Estimation results

```
library(stargazer)
library(readxl)

CAschool = read_excel("caschool.xlsx", col_names = TRUE, skip = 0)
results = lm(testscr ~ str, data = CAschool)
summary(results)
```

```
Call:
lm(formula = testscr ~ str, data = mydata)
```

Residuals:

Min	1Q	Median	3Q	Max
-47.727	-14.251	0.483	12.822	48.540

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	698.9330	9.4675	73.825	< 2e-16 ***
str	-2.2798	0.4798	-4.751	2.78e-06 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 18.58 on 418 degrees of freedom

Multiple R-squared: 0.05124, Adjusted R-squared: 0.04897

F-statistic: 22.58 on 1 and 418 DF, p-value: 2.783e-06

Confidence Interval for β_1

- Recall that a 95% confidence is, equivalently:

- The set of points that cannot be rejected at the 5% significance level,
- A set-valued function of the data (an interval that is a function of the data) that contains the true parameter value 95% of the time in repeated samples.

- Suppose the level of the (two-sided) test is 5%. We know that the critical value (in absolute value) is 1.96.

- Then, the 95% CI for β_1 refers to the set of values for β_1 such that

$$P\left(\{\beta_1 : \left|\frac{\hat{\beta}_1 - \beta_1}{SE(\hat{\beta}_1)}\right| \leq 1.96\}\right) = 0.95. \text{ Hence,}$$

$$\begin{aligned} \{\beta_1 : \left|\frac{\hat{\beta}_1 - \beta_1}{SE(\hat{\beta}_1)}\right| \leq 1.96\} &= \{\beta_1 : -1.96 \leq \frac{\hat{\beta}_1 - \beta_1}{SE(\hat{\beta}_1)} \leq 1.96\} \\ &= \{\beta_1 : -\hat{\beta}_1 - 1.96 SE(\hat{\beta}_1) \leq -\beta_1 \leq -\hat{\beta}_1 + 1.96 SE(\hat{\beta}_1)\} \\ &= \{\beta_1 : \hat{\beta}_1 + 1.96 SE(\hat{\beta}_1) \geq \beta_1 \geq \hat{\beta}_1 - 1.96 SE(\hat{\beta}_1)\} \\ &= \{\beta_1 : \hat{\beta}_1 - 1.96 SE(\hat{\beta}_1) \leq \beta_1 \leq \hat{\beta}_1 + 1.96 SE(\hat{\beta}_1)\} \end{aligned}$$

- Thus, the 95% CI for β_1 is $\left[\hat{\beta}_1 - 1.96 SE(\hat{\beta}_1), \hat{\beta}_1 + 1.96 SE(\hat{\beta}_1)\right]$.

Confidence Interval for β_1

- Similarly, you can construct 90% and 99% CIs. You just need to change the critical value, 1.64 and 2.57, respectively.
- Notice that the CI is an estimator and is a random variable.
- In other words, from one sample to the next, it provides different values for the lower bound and upper bound of the interval.
- By construction, this interval estimator retains the true (unknown) value of β_1 in 95% of the intervals constructed under repeated sampling.
- It does not mean that for the sample at hand the likelihood of the true (unknown) value of β_1 being in the interval constructed is 95%.
- If the hypothesized value under H_0 is in the CI, we fail to reject the H_0 . Otherwise, we reject H_0 .

Confidence Interval for β_1

Example 3

Consider the test score example in (6). The estimation results are provided in Listing 1. Then, the 95% CI for β_1 is $\hat{\beta}_1 \pm 1.96 \times SE(\hat{\beta}_1)$, which gives

$$[-3.22, -1.34].$$

- The 95% CI for β_1 can be used to construct a 95% CI for the predicted effect of a general change in X .
- Recall that the expected change in Y associated with this change in X is $\beta_1 \Delta x$.
- The 95% CI for $\beta_1 \Delta x$ is

$$\left[\left(\hat{\beta}_1 - 1.96 SE(\hat{\beta}_1) \right) \Delta x, \left(\hat{\beta}_1 + 1.96 SE(\hat{\beta}_1) \right) \Delta x \right].$$

Regression when X is Binary

- Sometimes a regressor can be binary—that is, it takes on only two values, 0 and 1. For example, X can be
 - ① a worker's sex: $X = 1$ if female, $X = 0$ otherwise,
 - ② whether a school district is urban or rural: $X = 1$ if urban, $X = 0$ otherwise),
 - ③ whether the district's class size is small or large: $X = 1$ if small, $X = 0$ otherwise.
- A binary variable is also called an **dummy variable** or sometimes a **indicator variable**.
- The mechanics of regression with a binary regressor are the same as if it is continuous.

Regression when X is Binary

- Suppose you have a variable D_i that equals either 0 or 1, depending on whether the student-teacher ratio is less than 20:

$$D_i = \begin{cases} 1, & \text{if the student-teacher ratio in } i\text{'th district} < 20 \\ 0, & \text{if the student-teacher ratio in } i\text{'th district} \geq 20. \end{cases}$$

- Let the population regression model with D_i as the regressor is

$$Y_i = \beta_0 + \beta_1 D_i + u_i, \quad i = 1, 2, \dots, n. \quad (8)$$

- Because D_i is not continuous, it is not useful to think of β_1 as a slope; indeed, because D_i can take on only two values.
- The best way to interpret β_1 in a regression with a binary regressor is to consider, one at a time, the two possible cases, $D_i = 0$ and $D_i = 1$.
- Under the least squares assumptions,
 - ① if the student-teacher ratio is high, then $D_i = 0$, and we have $E(Y_i | D_i = 0) = \beta_0$,
 - ② if the student-teacher ratio is low, then $D_i = 1$, and we have $E(Y_i | D_i = 1) = \beta_0 + \beta_1$.

Regression when X is Binary

- Hence,

$$E(Y_i|D_i = 1) - E(Y_i|D_i = 0) = \beta_0 + \beta_1 - \beta_0 = \beta_1$$

- Thus, β_1 is the difference between two population means.
- In our example, it is the difference between the mean test score in districts with low student-teacher ratios and the mean test score in districts with high student-teacher ratios.
- Listing 2 provides the estimation result for (8). The results show that $\hat{\beta}_1 = 7.37$.
- Thus, the students in districts with the STR less than 20 on average have 7.37 more points than the students in districts with the STR larger or equal to 20.

Regression when X is Binary

Listing 2: Dummy variable regression results: R

```

CASchool$D = (CASchool$str < 20) # create the dummy variable
results = lm(testscr ~ D, data = CASchool)
summary(results)

Call:
lm(formula = testscr ~ D, data = mydata)

Residuals:
    Min       1Q   Median       3Q      Max
-50.601 -14.047  -0.451  12.841  49.399

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)   649.979      1.388  468.380 < 2e-16 ***
DTRUE          7.372       1.843   3.999 7.52e-05 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 18.72 on 418 degrees of freedom
Multiple R-squared:  0.03685,    Adjusted R-squared:  0.03455
F-statistic: 15.99 on 1 and 418 DF,  p-value: 7.515e-05

```

Regression when X is Binary

Listing 3: Dummy variable regression results: Python

```
import numpy as np
import pandas as pd
import statsmodels.api as sm
import statsmodels.formula.api as smf

# Import data
CAschool=pd.read_excel("data/caschool.xlsx")
# Create the dummy variable D
CAschool["D"] = CAschool["str"] < 20

# Specify the model
model2=smf.ols(formula='testscr~D',data=CAschool)
# Use the fit method to obtain parameter estimates
result2=model2.fit()
# Print the estimation results
print(result2.summary())
```

Heteroskedasticity vs. Homoskedasticity

Definition 1

The error term u_i is **homoskedastic** if the variance of the conditional distribution of u_i given X_i is constant for $i = 1, 2, \dots, n$, and in particular does not depend on X_i . Otherwise, the error term u_i is **heteroskedastic**.

- Recall the least squares assumptions did not say anything about the conditional distribution of u_i given X_i . In other words, we have been assuming heteroskedasticity implicitly.
- And we derived the variance of the sampling distribution of the OLS estimator under heteroskedasticity.

$$\sigma_{\hat{\beta}_1}^2 = \frac{1}{n} \frac{\text{var}((X_i - \mu_X)u_i)}{(\text{var}(X_i))^2}$$

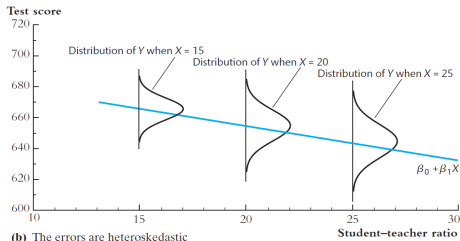
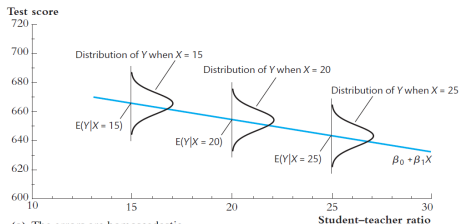
- Under homoskedasticity, this formula simplifies to

$$\sigma_{\hat{\beta}_1}^2 = \frac{1}{n} \frac{\text{var}((X_i - \mu_X)u_i)}{(\text{var}(X_i))^2} = \frac{1}{n} \frac{\text{var}((X_i - \mu_X)) \text{var}(u_i)}{(\text{var}(X_i))^2} = \frac{1}{n} \frac{\text{var}(u_i)}{\text{var}(X_i)}.$$

Heteroskedasticity vs. Homoskedasticity

FIGURE 5.2 Homoskedasticity and Heteroskedasticity

The figure plots the conditional distribution of test scores for three different class sizes (x). In figure (a), the spread of these distributions does not depend on x ; that is, $\text{var}(u|X = x)$ does not depend on x , so the errors are homoskedastic. In figure (b), these distributions become more spread out (have a larger variance) for larger class sizes, so $\text{var}(u|X = x)$ depends on x and the u is heteroskedastic.



Heteroskedasticity vs. Homoskedasticity

- Then, under homoskedasticity, $SE(\hat{\beta}_1)$ is given by

$$\hat{\sigma}_{\hat{\beta}_1} = \sqrt{\hat{\sigma}_{\hat{\beta}_1}^2} = \sqrt{\frac{1}{n} \frac{\frac{1}{n-2} \sum_{i=1}^n \hat{u}_i^2}{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2}} = \sqrt{\frac{\frac{1}{n-2} \sum_{i=1}^n \hat{u}_i^2}{\sum_{i=1}^n (X_i - \bar{X})^2}}.$$

- This formula is called the **homoskedasticity-only standard errors**. The `lm` function by default calculates standard errors according to this formula.
- Heteroskedasticity-robust standard errors (or Eicker-Huber-White standard errors) can be obtained using the **vcovHC** function from the **sandwich** package.

```
# install.packages("sandwich")
library(sandwich)
r1 <- lm(testscr ~ str, data = mydata)
vcov <- vcovHC(r1, type = "HC1")
robust_se <- sqrt(diag(vcov))

stargazer(r1, r1,
  se = list(NULL, robust_se),
  type = "text",
  column.labels = c("Homoskedastic",
    "Heteroskedastic"),
  align = TRUE)
```


Heteroskedasticity vs. Homoskedasticity

- According to the estimation results in Listing 4, the heteroskedastic standard errors are relatively larger.

Listing 4: Homoskedastic and heteroskedastic models

Dependent variable:		
	testscr	
	Homoskedastic	Heteroskedastic
	(1)	(2)
str	-2.280*** (0.480)	-2.280*** (0.519)
Constant	698.933*** (9.467)	698.933*** (10.364)
Observations	420	420
R ²	0.051	0.051
Adjusted R ²	0.049	0.049
Residual Std. Error (df = 418)	18.581	18.581
F Statistic (df = 1; 418)	22.575***	22.575***
Note:	*p<0.1; **p<0.05; ***p<0.01)	

Heteroskedasticity vs. Homoskedasticity

Listing 5: Homoskedastic and heteroskedastic models: Python

```
import numpy as np
import pandas as pd
import statsmodels.api as sm
import statsmodels.formula.api as smf
from statsmodels.iolib.summary2 import summary_col

# Specify the model
model=smf.ols(formula='testscr~str',data=CAAschool)
# Use the fit method to obtain parameter estimates
result=model.fit(cov_type="HC1")
# Print the estimation results
print(result.summary())
```

Heteroskedasticity vs. Homoskedasticity

- If the errors are either homoskedastic or heteroskedastic and you use heteroskedastic-robust standard errors, you are OK.
- If the errors are heteroskedastic and you use the homoskedasticity-only formula for standard errors, your standard errors will be wrong (the homoskedasticity-only estimator of the variance of $\hat{\beta}_1$ is inconsistent if there is heteroskedasticity).
- The two formulas coincide (when n is large) in the special case of homoskedasticity.
- Therefore, Stock and Watson (2020) suggest that we should always use heteroskedasticity-robust standard errors.

The Gauss-Markov theorem

- Recall the three least squares assumptions. Now, add the [homoskedasticity](#) assumption on the error term.
- Then, it can be shown analytically that the OLS estimator is the most precise estimator (conditional on X 's) within the class of linear unbiased estimators.
- Often stated as the OLS estimator is Best Linear (conditionally) Unbiased Estimator (BLUE).
- This results is known as the [Gauss-Markov theorem](#). See Appendix 5.2 for a proof.
- The Gauss-Markov theorem has two important limitations.
 - Homoskedasticity assumption is not likely to hold in practice. Then, the OLS estimator is no longer BLUE.
 - There are estimators that are not linear and conditionally unbiased, and under some conditions, these other estimators are more efficient than OLS.

Normality of the error term

- So far we have not specified the distribution of the error term, we just made assumptions on its moments.
- Addition to the four assumptions earlier, we assume also that the $u_i \sim N(0, \sigma_u^2)$.
- Then, the OLS estimator $\hat{\beta}_1$ is **exactly** distributed as normal with mean β_1 and $\sigma_{\hat{\beta}_1}^2$ under homoskedasticity.
- Also, the t -statistic now **exactly** distributed as Student t with $(n - 2)$ degrees of freedom.
- For large n , the Student t distribution with $(n - 2)$ degrees of freedom is almost the same as the standard normal distribution.
- For inference, this creates a difference if the sample size very small, say $n < 30$: the Student t critical values can be a fair bit larger than the standard normal critical values.

Practical implications

- If $n < 50$ and you really believe that, for your application, u_i is homoskedastic and normally distributed, then use the t_{n-2} instead of the $N(0, 1)$ critical values for hypothesis tests and confidence intervals.
- In most econometric applications, there is no reason to believe that u_i is homoskedastic and normal - usually, there are good reasons to believe that neither assumption holds.
- Fortunately, in modern applications, $n > 50$, so we can rely on the large- n results presented earlier, based on the CLT, to perform hypothesis tests and construct confidence intervals using the large- n normal approximation.

Bibliography I



Hanck, Christoph et al. (2021). *Introduction to Econometrics with R*. URL:
<https://www.econometrics-with-r.org/index.html>.



Stock, James H. and Mark W. Watson (2020). *Introduction to Econometrics*.
Fourth. Pearson.