

ØAMET4100 · Spring 2019

Lecture Note 2B

Instructor: Fenella Carpena

January 17, 2019

This lecture note provides a review of regression with a single regressor: hypothesis tests and confidence intervals (Stock & Watson, Chapter 5). This lecture note is not intended to be a comprehensive review of lecture or the textbook, since there is a lot more material than we have time to cover. However, I have tried to focus on the concepts which I believe are necessary to be successful in our class.

1 Hypothesis Testing for Regression Coefficients

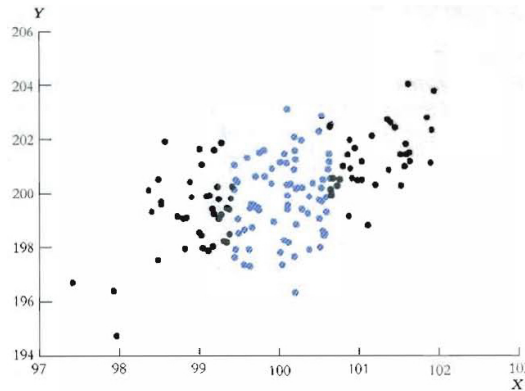
- Let us now consider how to make inferences, using our sample, about the relationship between X and Y in the population. For instance, some questions we will try to answer in this section are:
 - Is the observed relationship between X and Y in the sample strong enough to conclude that it also holds in the population?
 - How can we use the sample statistics $\hat{\beta}_0$ and $\hat{\beta}_1$ to determine a plausible range of values for β_0 and β_1 of the population regression line?
- To be able to use our sample estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ to say something about the population parameters β_0 and β_1 , we needed to know the standard errors of the sample estimates $SE(\hat{\beta}_0)$ and $SE(\hat{\beta}_1)$.
- **Why do $\hat{\beta}_0$ and $\hat{\beta}_1$ have standard errors? Aren't they just constants?** No, $\hat{\beta}_0$ and $\hat{\beta}_1$ are NOT constants. They are random variables. We already discussed this last week, but I repeat it again here since it is very important to understand this point. There is variability in $\hat{\beta}_0$ and $\hat{\beta}_1$ because if we draw a different sample, we would get a different value of $\hat{\beta}_0$ and $\hat{\beta}_1$. The standard errors describe this sample-to-sample variability of $\hat{\beta}_0$ and $\hat{\beta}_1$.
- In a regression with a single regressor, the formula for (homoskedastic) standard errors of $\hat{\beta}_1$ is

$$SE(\hat{\beta}_1) = \frac{s_{\hat{u}}}{\sqrt{n-1}} \cdot \frac{1}{s_X}$$

where $s_{\hat{u}}$ is the standard error of the regression (covered in the last section). n is the sample size, and s_X is the sample standard deviation of X . **Note that a smaller $SE(\hat{\beta}_1)$ means that our estimate of the population β_1 is more precise.** Here are several important features of $SE(\hat{\beta}_1)$:

- **As $s_{\hat{u}}$ falls, $SE(\hat{\beta}_1)$ falls.** This is because $s_{\hat{u}}$ is our estimate for the standard deviation of the population errors u , and if this standard deviation is small, then the data will have a tighter scatter around the population regression line. Hence, its slope will also be estimated more precisely.
- **As n increases, $SE(\hat{\beta}_1)$ falls.** That is, when our sample size is large, our estimate $\hat{\beta}_1$ will be close the true population coefficient β_1 with high probability. This is because the standard deviation of $\hat{\beta}_1$ decreases to zero as n increases, so the distribution of the $\hat{\beta}_1$ will be tightly centered around its mean β_1 when n is large. Another way to think about this is that the larger the sample, the more information we have about the population, so the more precise our estimate for β_1 is.
- **As s_X increases, $SE(\hat{\beta}_1)$ falls.** To get a better sense of why this is true, consider the figure below which presents a scatterplot of 150 artificial data points on X and Y . Suppose you were asked to draw a line as accurate as possible through *either* the blue dots or the black

dots—which would you choose? It would be easier to draw a precise line through the black dots, which have a larger variance than the blue dots. Similarly, the larger the variance of X , the more precise is $\hat{\beta}_1$.



- There is also a formula for $SE(\hat{\beta}_0)$ which you can find in the textbook. However I will not focus on it here since in most applications, we are likely to be more interested in $\hat{\beta}_1$ than $\hat{\beta}_0$ (e.g., we often want to know how Y changes when X changes).
- In practice, statistics software such as Stata calculate $SE(\hat{\beta}_0)$ and $SE(\hat{\beta}_1)$ automatically, so they are provided to you in the output of the regression table.
- The three approaches for hypothesis test we will construct are the following: (1) confidence interval, (2) test statistic (e.g., t -statistic), and (3) p -value.

1.1 Hypothesis Tests Using t -statistic

- **Step 1:** Construct a t -statistic. The general formula is:

$$t\text{-stat} = \frac{\text{sample statistic} - \text{hypothesized value}}{se(\text{sample statistic})}.$$

For example, if $\hat{\beta}_0 = 1.86534$, $SE(\hat{\beta}_0) = 0.40083$, and the null and alternative hypotheses are $H_0 : \beta_0 = 2$, $H_a : \beta_0 \neq 2$ (a two-sided hypothesis test), then

$$t\text{-stat} = \frac{1.86534 - 2}{0.40083} = -0.336.$$

- **Step 2:** Find the critical value. This critical value depends on: (1) whether we have a two-sided or one-sided test, and (2) the significance level α (given in the problem or chosen by you).
- **Step 3:** Compare the t -statistic with the critical value from the previous step, which will allow us to either **reject** or **fail to reject** H_0 .

1.2 Hypothesis Test: Using p -value

- **Step 1:** Construct the t -statistic as in the previous section.
- **Step 2:** Calculate the p -value, depending on whether we have a two-sided or one-sided test. For example, if you had a two-sided test, $p\text{-value} = P(|Z| > |t\text{-stat}|)$.
- **Step 3:** Compare the p -value obtained in the previous step with the level of α . We **reject** H_0 if $p\text{-value} < \alpha$, where α is the significance level. Otherwise, we **fail to reject** H_0 .

1.3 Confidence Intervals

- The general formula for the confidence interval that we learned previously for \bar{X} also applies to the case where we are constructing confidence intervals for β_0 and β_1 . That general formula is:

$$\text{sample statistic} \pm \text{two-sided critical value} * se(\text{sample statistic})$$

- The $(100 - \alpha)\%$ **confidence interval for the slope** β_1 is given by:

$$\widehat{\beta}_1 \pm z_{\alpha/2} * SE(\widehat{\beta}_1)$$

where α is the significance level (given in the problem or chosen by you), $z_{\alpha/2}$ is the two-sided critical value from the normal table, and n is the sample size. Note that the confidence interval for $\widehat{\beta}_0$ follows the same formula as above, but using $\widehat{\beta}_0$ instead of $\widehat{\beta}_1$.

- **How do we interpret the CI for β_1 ?** If $\alpha = 0.05$, then we have a 95% CI. This means that “we are 95% confident that the true population parameter β_1 lies between that interval.”

2 Dummy Variables

- So far we have been talking about X variables that are continuous (e.g., class size, age, etc.). However, regressions can also be used when X is binary, that is, when it can take on only two values (0 and 1). Such a variable is called **dummy variable** (also known as binary variable or indicator variable).
- To fix ideas, let us consider the following regression

$$\widehat{wage} = \widehat{\beta}_0 + \widehat{\beta}_1 \cdot master$$

where *wage* is measured annually in dollars, and *master* is a dummy variable equal to 1 if a worker has a master’s degree, and 0 otherwise.

- In the above regression, *master* is not continuous, so it is not useful to think of $\widehat{\beta}_1$ as a slope. In this case, **how should we interpret the coefficients?** We can easily do so by looking at the two cases when *master* = 1 or when *master* = 0. Specifically, we note the following.

- If *master* = 0, $\widehat{wage} = \widehat{\beta}_0$.
- If *master* = 1, $\widehat{wage} = \widehat{\beta}_0 + \widehat{\beta}_1$.
- Using the above two equations, we see that:
 - * $\widehat{\beta}_0$ is the average wage of workers without a master’s degree.
 - * $\widehat{\beta}_0 + \widehat{\beta}_1$ is the sample average wage of workers with a master’s degree.
 - * $\widehat{\beta}_1$ is the difference in the sample average wage between workers with master’s degree and workers without a master’s degree.

- Note that a hypothesis test for a difference in means between two groups can be carried out by regressing Y on a dummy variable X which defines groups. For example, let μ_{master} be the population average wage of workers with master’s degree, and $\mu_{nonmaster}$ be the population average wage of workers without a master’s degree. In the regression of *wages* on *master* shown above, the hypothesis test $H_0 : \mu_{master} - \mu_{Nonmaster} = 0$ vs. $H_1 : \mu_{master} - \mu_{NonMaster} \neq 0$ is equivalent to testing $H_0 : \beta_1 = 0$ vs. $H_1 : \beta_1 \neq 0$.

3 Homoskedasticity and Heteroskedasticity

- We say that the population error term u_i is **homoskedastic** if $var(u_i|X_i)$ is constant for all i . Otherwise, we say that u_i is **heteroskedastic**.
- To understand what heteroskedasticity means, let’s take a look at the following regression of *wages* on *schooling* (i.e., years of education)

$$wages_i = \beta_0 + \beta_1 schooling_i + u_i.$$

Note that $var(u_i|schooling_i) = var(wages_i - \beta_0 - \beta_1 schooling_i|schooling_i) = var(wages_i|schooling_i)$. Hence, this means that:

- If u_i is homoskedastic, both $var(u_i|schooling_i)$ and $var(wages_i|schooling_i)$ are constant.

- In particular, homoskedasticity would imply that $var(wages_i|schooling_i)$ is the same for all schooling levels. Another way of saying this is that the variability of wages around its mean is the same regardless of educational attainment.
 - Homoskedasticity is not realistic in this case because it is likely that people with more education have wider job opportunities, which could lead to more variability in wages. In contrast, people with low education levels have fewer opportunities and probably work minimum wage jobs, so there is less dispersion of wages among the uneducated.
 - In sum, we would expect that variability in wages is higher for the highly educated, and the variability in wages is low for those with low levels of schooling. Therefore, in this example, the errors u_i are likely heteroskedastic.
- **Why do we care about heteroskedasticity and homoskedasticity?** I can think of two reasons why we care.
 - If the 3 least squares assumptions hold **AND** u_i are homoskedastic, then the OLS estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ are BLUE (Best Linear Unbiased Estimators). This is the **Gauss-Markov Theorem**.
 - Heteroskedasticity and homoskedasticity have implications for calculating $SE(\hat{\beta}_1)$ and $SE(\hat{\beta}_0)$. For example, if the population errors are heteroskedastic but you use homoskedastic SEs, your hypothesis tests and confidence intervals will be invalid. However, since homoskedasticity is a special case of heteroskedasticity, the heteroskedastic-robust SEs will still be valid under homoskedasticity (see table below for a summary of the different cases). Hence, **a typical rule of thumb is to always use heteroskedasticity-robust SEs**. This is the “robust” option when running a regression in Stata.

| | | What you actually use | |
|-------------------|-----------------|--------------------------|------------------------|
| | | Homoskedastic | Heteroskedastic |
| Truth about u_i | Homoskedastic | VALID hyp. test and CI | VALID hyp. test and CI |
| | Heteroskedastic | INVALID hyp. test and CI | VALID hyp. test and CI |