

1 Population Regression

When we estimate things, our estimation is going to depend on whatever sample we happen to have obtained. That sample is usually not going to be a perfect representation of the population, and hence any given sample will differ from the population in random ways.

To illustrate, suppose you have a population of 100 people and you want to estimate their income. You take 20 random samples, someone else takes 20 random samples. Chances are you won't sample the exact same 20 people and hence your estimates will be a bit different. We need to account for that sampling variability.

In the context of regressions, we'd like a regression that best fits the population data. It will be given by the formula

$$y = \beta_1 + \beta_2 x + u,$$

which I will explain in detail momentarily. But think of this as being the line of best fit for the entire population where β_1 is the intercept and β_2 is the slope. We want to estimate β_1 and β_2 using a sample. The estimation method is called **ordinary least squares (OLS)**. To do so, we will rely on a number of assumptions in order to make sure our estimates for each β have nice properties.

2 Unbiased Estimators

It might help to refer to Figure 1 to get a visual feel for what's going on while reading these assumptions.

Assumption OLS1: Linear True Population Model. Again, a regression is just the line of *best* fit – it is not the line of *perfect* fit. When we talk about a specific data point i , we will assume that the true population model has form

$$y_i = \beta_1 + \beta_2 x_i + u_i. \tag{1}$$

What this says is we use the line $\beta_1 + \beta_2 x_i$ to best “predict” what y_i should be for a given value of x_i ; but since the regression line doesn't perfectly capture all data points, the prediction will be off by u_i . Accordingly, u_i is called the **error term**, sometimes called the **disturbance term**.

Assumption OLS2: Zero Conditional Mean. The zero conditional mean assumption states that

$$E[u_i|x_i] = 0 \quad \text{for all } i. \quad (2)$$

Remember, u_i represents how wrong the line of best fit is for data points (x_i, y_i) . The zero conditional mean assumption means that, given x_i , we expect the line of best fit to not be wrong *on average*.

This is important because we'd like to be able to answer the question, "what should I expect y to be, given any value of x ?" So consider a specific $x = x^*$, where x^* is just any old number for which we want to predict y .¹ This allows us to take the true population model and write

$$\begin{aligned} E[y|x = x^*] &= E[\beta_1|x = x^*] + E[\beta_2 x|x = x^*] + E[u|x = x^*] \\ &= \beta_1 + \beta_2 x^*. \end{aligned}$$

This is true because β_1 and β_2 are just numbers – there is nothing random about them – so we, uh, expect them to be themselves, regardless of what x is. And because of our zero conditional mean assumption, the error term drops out. Thus, the regression line is what we expect y to be for a given value of x .

For more intuition, suppose $E[u_i|x_i] \neq 0$. Then when we plug in some data point x_i , we expect the model to have some error, on average. It makes for a pretty lousy model when we expect it to be wrong, on average.

To summarize the population characteristics:

- The actual value y_i is given by $y_i = \beta_1 + \beta_2 x_i + u_i$.
- The regression line is what we expect y_i to be, given x_i . Expressed in the maths, $E[y_i|x = x_i] = \beta_1 + \beta_2 x_i$. This is a consequence of assumptions 1 and 2 combined.
- Hence the error term is given by $u_i = y_i - E[y_i|x = x_i]$.

OLS Assumptions 1 and 2 imply that OLS estimates (explained soon) b_1 and b_2 of β_1 and β_2 are unbiased, in other words, we expect the estimates to be their true values. In maths, $E[b_1] = \beta_1$ and $E[b_2] = \beta_2$.

¹For example, we might want to predict how many cavities a person has (y) if they eat 200 grams of sugar per day ($x = 200$).

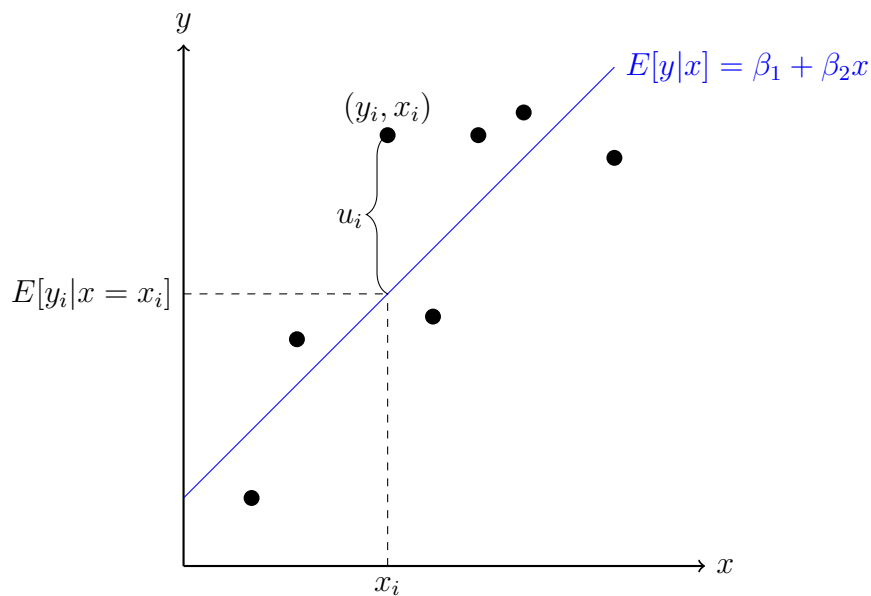


FIGURE 1: Suppose the dots here capture the entire population of data. Pick some arbitrary data point (x_i, y_i) . The regression line tells us $E[y_i|x = x_i]$, that is, what value we expect y_i to be for independent variable x_i . This is the **conditional mean** of y_i given x_i . But the regression line is a line of *best* fit, not a line *perfect* fit, so the actual value of y_i will in general be different than what we expect it to be based on the regression line. The difference between what y_i actually is and what we expect y_i to be based on the regression, $y_i - \beta_1 - \beta_2 x_i$, is the error term, u_i .

3 BLUE and BUE

3.1 BLUE

We can throw down two more assumptions to make analysis cleaner.

Assumption OLS3: Homoskedasticity. The variation of u_i given x_i is the same number σ_u^2 for any x_i . In math,

$$\text{Var}(u_i|x_i) = \sigma_u^2 \quad \text{for all } i. \quad (3)$$

This condition is illustrated in Figure 2.

Assumption OLS4: Uncorrelated Errors. Errors for different observations are uncorrelated, that is,

$$\text{Cor}(u_i, u_j) = 0 \quad \text{whenever } i \neq j. \quad (4)$$

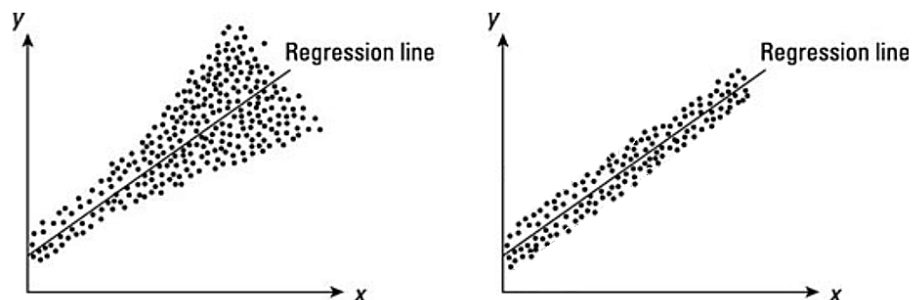


FIGURE 2: The figure on the left is an example of heteroskedasticity; the right an example of homoskedasticity. The left is heteroskedastic because the variation of errors around the regression line gets bigger as x increases.

Adding assumptions 3 and 4 allows us to say that the variation of y given x is also constant, and specifically, $\text{Var}(y|x) = \sigma_u^2$. OLS assumptions 1-4 imply also imply that estimates are **consistent**, provided the variances of the estimates go to zero as $n \rightarrow \infty$. Put somewhat crudely, this means that our estimates get arbitrarily close (in probability) to their true population values as the sample size increases. In math, we write $b \xrightarrow{p} \beta$.

We can go even further, however. Under OLS assumptions 1-4, the estimates are said to be **BLUE**, which stands for

- **B**est: estimates have the smallest standard errors...
- **L**inear: among linear models...
- **U**nbiased: that give unbiased...
- **E**stimator: um, estimates.

3.2 BLUE

We can make a fifth assumption for one more nice result.

Assumption OLS5: Normally Distributed Errors. Error terms have normal distribution with some variance σ^2 ,

$$u_i \sim \mathcal{N}(0, \sigma^2). \quad (5)$$

This allows us to say that OLS estimates are **BUE**, which means that they have the smallest standard errors among unbiased models, even when compared to nonlinear models. Also note that this is an essential condition if we want to do inference on small sample sizes.

3.3 Recap

We have five OLS assumptions that give the following implications:

- OLS1 (linear model) and OLS2 (zero conditional mean) imply unbiased OLS estimates.
- Adding OLS3 (homoskedasticity) and OLS4 (uncorrelated errors) imply that $\text{Var}(y|x) = \sigma_u^2$ is constant, and that the estimate b_i for each β_i is consistent.
- OLS assumptions 1-4 therefore imply that OLS estimates of β are BLUE.
- Adding OLS5 (normality of errors) implies that OLS estimates of β are BUE.

We will weaken many of these assumptions as we go further into the course, but it's almost always best to start with the easiest result and then break it down from there.

4 OLS Estimation of a Regression

Again, b_1 is the estimate of β_1 and b_2 the estimate of β_2 . Intuitively, we want a model that makes the fewest mistakes possible. We quantify “fewest” by considering the difference between the actual values y_i , and the **fitted values** as predicted by the model, given by $\hat{y}_i = b_1 + b_2x_i$; this is referred to as the **residual**, denoted e_i , defined as

$$e_i \equiv y_i - \hat{y}_i = y_i - [b_1 + b_2x_i]. \quad (6)$$

Think of the residual as capturing how wrong the estimated line of best fit is. If you think minimizing residuals sounds like a good idea, then you're on the right track.

We square each residual to ensure that it's positive, then we add the squared terms all up: this is the **residual sum of squares (RSS)**. We want the estimates that *minimize the residual sum of squares*. In math speak, we want to solve

$$(b_1, b_2) = \arg \min_{b_1, b_2} \sum_{i=1}^n (y_i - [b_1 + b_2x_i])^2. \quad (7)$$

The solution to this is the **ordinary least squares (OLS)** estimation for a linear regression. I omit the details, but explicitly solving OLS (using calculus to find critical points) gives

formulas

$$b_2 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} = \frac{s_{xy}}{s_x^2} = r_{xy} \times \frac{s_y}{s_x}, \quad (8)$$

$$b_1 = \bar{y} - b_2 \bar{x}, \quad (9)$$

where s_{xy} is the **sample covariance** defined by

$$s_{xy} \equiv \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}), \quad (10)$$

and r_{xy} is the **sample correlation coefficient** defined by

$$r_{xy} \equiv \frac{s_{xy}}{s_x s_y}. \quad (11)$$

I give several different expressions for b_2 because, depending on how a problem is worded, one expression might be more applicable than the others. Note that the first expression is the one given on the exam formula sheet.

Again, under OLS assumptions 1 and 2, the estimates will be unbiased: $E[b_1] = \beta_1$ and $E[b_2] = \beta_2$. That said, they will be different in generality than their population analogues because, well, they're estimates. Hence our estimated regression line will be more or less different than the population regression line, depending on how closely our sample reflects the population. This is illustrated in Figure 3.

Furthermore, OLS assumptions 3 and 4 imply that the variance of the slope estimate b_2 will be

$$\text{Var}(b_2) = \frac{\sigma_u^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \equiv \sigma_{b_2}^2. \quad (12)$$

OLS assumption 3 is most likely to break down in practice, in which case we will have **heteroskedasticity** – the variance of u_i will depend on x_i . In this case we need to use **heteroskedasticity-robust standard errors**, which mathematically are beyond the scope of this course, but are very easy to implement in Stata (so you should know that heteroskedasticity is a thing that we should be concerned about).

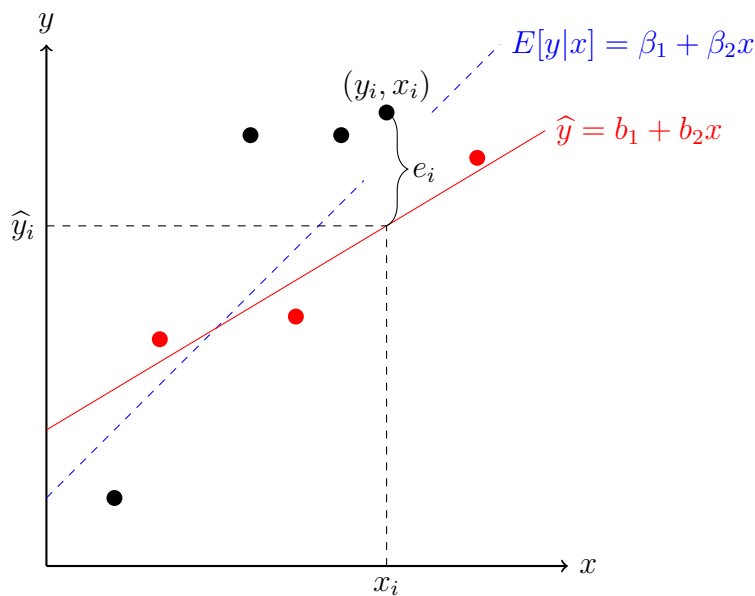


FIGURE 3: Suppose our sample consists of only the red dots. Thus the estimated regression line (in red) is different than the true population regression line (in blue). For x_i , it gives us a prediction for y_i , i.e. the fitted value \hat{y}_i . The fitted value will not in general be exactly the true value y_i , and the difference between the true value and the fitted value is the residual, $e_i = y_i - \hat{y}_i$. This example illustrates a positive residual, $e_i > 0$.

5 Explained and Unexplained Variation

To reiterate, we define the **residual sum of squares** to be

$$\text{RSS} \equiv \sum_{i=1}^n (y_i - \hat{y}_i)^2. \quad (13)$$

This captures the total error of the estimated regression line, squared so that the errors are positive. Dividing this by $n - 2$ (because we are estimating two parameters, one for each β) and taking the square root gives the **standard error of the regression**,

$$s_e \equiv \sqrt{\frac{\text{RSS}}{n-2}} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2}, \quad (14)$$

which is the sample analogue of σ_u as used in OLS3 and OLS4. You can think of this as being the variation of data around \bar{y} that cannot be explained by x . This is sometimes called the **standard error of the residuals** or **root mean squared error (RMSE)**.

On the other hand, the variation of data around \bar{y} that can be explained by x is the

explained sum of squares,

$$\text{ESS} \equiv \sum_{i=1}^n (\hat{y}_i - \bar{y})^2. \quad (15)$$

Finally, the total variation of data around \bar{y} is given by the **total sum of squares**,

$$\text{TSS} \equiv \sum_{i=1}^n (y_i - \bar{y})^2. \quad (16)$$

Based on the intuition it should not be surprising, and it is not difficult to show either, that

$$\text{TSS} = \text{ESS} + \text{RSS}. \quad (17)$$

Total variation is explained variation plus unexplained variation. Great.

The proportion of explained variation around \bar{y} is called the **R-squared** or **coefficient of determination**, defined as

$$R^2 \equiv \frac{\text{ESS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}}. \quad (18)$$

If R^2 is high, then x explains a lot about what's going on with y ; if R^2 is low, then it doesn't. There is no cutoff for what should be considered "high" or "low," however. Note that R^2 also equals the squared correlation between y and x , that is, $R^2 = r_{xy}^2$. Also note that R^2 is only valid if the regression includes the intercept.

6 Estimator Properties

We are primarily interested in β_2 because it captures the relationship between x and y . Under OLS assumptions 1-4, our slope estimator b_2 has expected value of β_2 because it is unbiased; and it also has variance $\sigma_{b_2}^2$. Thus we can write

$$b_2 \sim (\beta_2, \sigma_{b_2}^2). \quad (19)$$

For sufficiently large sample size (greater than 30), the z -score is approximately standard normal, that is,

$$z \equiv \frac{b_2 - \beta_2}{\sigma_{b_2}},$$

which is drawn from a $\mathcal{N}(0, 1)$ distribution, approximately.

But we don't actually know what σ_{b_2} because it is a function of σ_u , which is a unknown population parameter. So instead we must use the sample estimate of σ_u , given earlier as s_e . This then allows us to conclude that the sample standard error of b_2 is

$$\text{se}(b_2) = \frac{s_e}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} \quad (20)$$

So under OLS assumptions 1-4, for sufficiently large sample size (which does *not* have a clear cut rule-of-thumb in this case), we use the t -statistic

$$t \equiv \frac{b_2 - \beta_2}{\text{se}(b_2)}, \quad (21)$$

which is drawn from a $T(n-2)$ distribution, where the distribution is approximate. If we add an additional assumption that the error terms are normally distributed (OLS5), or if $n \rightarrow \infty$, then we can say that t is drawn from an exact $T(n-2)$ distribution.

7 Regression by Hand

Okay, so that's a lot to take in. At this point you should look at the formula sheet on one of the practice exams, because chances are you'll be relying on it when exam time comes around. To help you become familiar, I provide an example.

Consider the following data:

$$(y_1, x_1) = (2, 0),$$

$$(y_2, x_2) = (3, 3),$$

$$(y_3, x_3) = (4, 3).$$

Let us regress this entire thing manually.

Step 1: Estimate Regression Coefficients. We will need sample means,

$$\bar{y} = \frac{1}{3} [2 + 3 + 4] = 3,$$

$$\bar{x} = \frac{1}{3} [0 + 3 + 3] = 2,$$

sample variances,

$$s_y^2 = \frac{1}{2} [(2-3)^2 + (3-3)^2 + (4-3)^2] = 1,$$

$$s_x^2 = \frac{1}{2} [(0-2)^2 + (3-2)^2 + (3-2)^2] = 3,$$

and sample covariance and correlation,

$$s_{xy} = \frac{1}{2} [(2-3)(0-2) + (3-3)(3-2) + (4-3)(3-2)] = 1.5,$$

$$r_{xy} = \frac{1.5}{\sqrt{1}\sqrt{3}} \approx 0.866.$$

Now we can estimate the regression coefficients, given by

$$b_2 = \frac{s_{xy}}{s_x^2} = \frac{1.5}{3} = 0.5 \quad \text{or} \quad r_{xy} \times \frac{s_y}{s_x} = 0.866 \times \frac{1}{\sqrt{3}} = 0.5,$$

$$b_1 = \bar{y} - b_2\bar{x} = 3 - 0.5(2) = 2.$$

Therefore our estimated regression is $\hat{y}_i = 2 + 0.5x_i$.

Step 2: Calculate Residuals. We can find the standard error of the regression by finding the fitted values, i.e. by plugging each x_i into the regression formula and finding the corresponding \hat{y}_i . Doing so gives

$$\hat{y}_1 = 2 + 0.5(0) = 2,$$

$$\hat{y}_2 = 2 + 0.5(3) = 3.5,$$

$$\hat{y}_3 = 2 + 0.5(3) = 3.5.$$

The residuals are the difference between the actual y_i and what the regression line expects y_i to be based on x_i , which in our case are

$$e_1 = y_1 - \hat{y}_1 = 2 - 2 = 0,$$

$$e_2 = y_2 - \hat{y}_2 = 3 - 3.5 = -0.5,$$

$$e_3 = y_3 - \hat{y}_3 = -3.5 = 0.5.$$

Step 3: Calculate Standard Error of Residual. The residual sum of squares (RSS), uh, squares each residual and sums them up, so

$$\text{RSS} = (0)^2 + (-0.5)^2 + (0.5)^2 = 0.5.$$

Now we can find the standard error of the residuals using formula

$$s_e \equiv \sqrt{\frac{\text{RSS}}{n-2}} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2} = \sqrt{\frac{0.5}{3-2}} = 0.707.$$

Step 4: Calculate Standard Error of Slope Coefficient. The slope coefficient has standard error

$$\text{se}(b_2) = \frac{s_e}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}.$$

So we gotta figure out the denominator. Not a big deal, it's pretty much just the calculation for the standard deviation of x but without the division. The formula yields

$$\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} = \sqrt{(0-2)^2 + (3-2)^2 + (3-2)^2} = \sqrt{6}.$$

Therefore the standard error of b_2 is

$$\text{se}(b_2) = \frac{0.707}{\sqrt{6}} \approx 0.289,$$

Hence for a hypothesis test you would calculate the t -statistic

$$t \equiv \frac{b_2 - \beta_2^0}{\text{se}(b_2)},$$

where β_2^0 is the hypothesized value, and would perform inference about β_2 under the assumption that t was drawn from an approximate $T(3-2)$ distribution.