# The linear regression model, and model assumptions

Econometrics for minor Finance, Lecture 3

Tunga Kantarcı, Fall 2025

Can we improve student outcomes?

What are the policy tools?

Changing class size? Improving teacher quality? Investment in school?

Can smaller classes yield to higher test scores? What are the mechanisms?

Classroom management: Fewer students can mean fewer disruptions, leading to a better learning environment.

Teacher workload: With fewer students, teachers may be less overwhelmed and better able to prepare and assess.

Individual attention: Smaller classes may allow teachers to give more personalized instruction and feedback. Students might feel more engaged and less overlooked.

The California Standardized Testing and Reporting (STAR) dataset contains data on test performance and school characteristics. The data are from districts in California collected in 1999 by the California Department of Education. Test scores are the average of the reading and math scores on a standardized test administered to 5th grade students. The student-teacher ratio is the number of students divided by the number of teachers working full-time in a district. The data is analyzed in Kruger and Whitmore (Economic Journal, 2001).
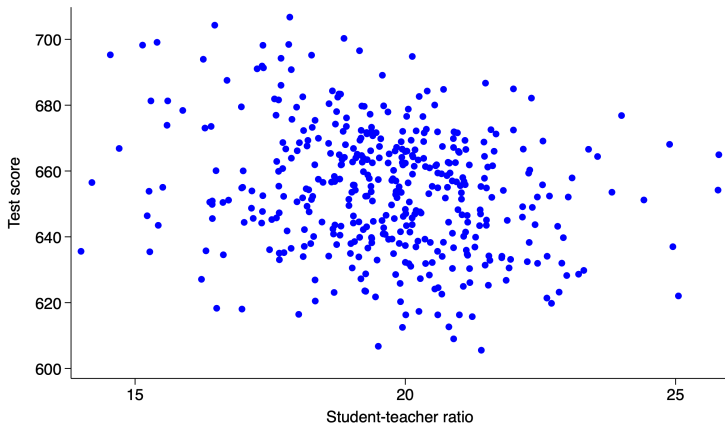
# An empirical challenge: Educational outcomes: Data: Summary statsitics

```
. summarize tstscr str

    Variable |        Obs        Mean    Std. dev.        Min        Max
-------------+-------------------------------------------------------------
      tstscr |        420    654.1565    19.05335     605.55     706.75
         str |        420    19.64043    1.891812         14       25.8
```

Our aim would be to formally analyze how test scores vary with the student-teacher ratio. For this, we need to assume a model, but not just any model. It must be one that fits the data well and captures the underlying relationship we care about.

So there are two things. First, is the model, which provides a structured lens for understanding the relationship. Second is fitting, which ensures that this lens reflects the actual patterns in the data.

An observational unit is the entity on which measurements are taken, such as a student, worker, firm, or country. An observational unit is random. It becomes non-random once we sample and observe it.

# The linear regression model: Model

Our model is

$$y_i = \beta_0 \cdot 1_i + \beta_1 x_i + u_i$$

for each of $n$ observational units.

- $y_i$ : outcome: e.g., test score
- $\beta_0$ : intercept term, unknown population parameter so needs to be estimated
- $1_i$ : a constant equal to 1, used to include the intercept term
- $\beta_1$ : slope coefficient, unknown population parameter so needs to be estimated
- $x_i$ : predictor: e.g., class size.
- $u_i$ : error term, to capture factors not observed in the data, unobserved but we can approximate it

# The linear regression model: Model

$y_i$ also called dependent variable, or response variable

$x_i$ also called independent variable, regressor, or explanatory variable

These names reflect convention, they do not imply causation. Regression describes conditional relationships, and sometimes causal ones.

The model

$$y_i = \beta_0 \cdot 1_i + \beta_1 x_i + u_i$$

is a population model. It provides the conceptual recipe for how the data we observe are generated in the population: the data generating process. That is, it assumes how the outcome $y$ is produced from the predictor

$$x$$

together with unobserved influences captured by the error term

$$u$$

The population parameters

$$\beta_0, \beta_1$$

summarize the systematic part of this relationship in the population. Randomness enters through the sampling process of $y_i$ and $x_i$.

The model for observational unit $i$ is

$$y_i = \beta_0 \cdot 1_i + \beta_1 x_i + u_i$$

As we typically have sample data on $n$ observational units, we can stack the observations to obtain a compact model for $n$ units

$$y = \beta_0 \cdot 1 + \beta_1 x + u$$

Moreover, although the intercept is technically multiplied by a vector of ones, we suppress it in notation and write:

$$y = \beta_0 + \beta_1 x + u$$

Like any other model, a regression model abstracts away complexity by reducing reality to a functional form. Models are approximations of the reality.

*"All models are wrong, but some are useful."*
— *George E. P. Box*

What all this implies is that we are making assumptions under which any regression model yields interpretable, credible, and sometimes causal insights about the world. A regression model is useful not because it is true, but because under the right assumptions it provides good approximations of relationships. If these assumptions fail, the model is not a good approximation of reality.

Linearity. The model is linear in the parameters.

The model

$$y = \beta_0 + \beta_1 x^2 + u$$

is

- linear in the parameters
- nonlinear in the regressor
- linear in the squared regressor

What makes a regression model linear is that it is linear in the parameters.

The model
$$y = x^\beta + u$$

is nonlinear in the parameter. This is not the linear regression model. Linearity means parameters enter additively and proportionally. If the true relationship looks more curved or complex, as it is here, this straight, proportional form is not a good assumption.

No perfect collinearity. If there are multiple explanatory variables, they must be linearly independent. In other words, no column vector corresponding to one variable can be expressed as a linear combination of the column vectors of the other variables.

If we include several explanatory variables in a regression model, each variable must contribute unique information. This means that no variable can be exactly predicted from the others. If one variable is just the sum or multiple of other variables, then it does not add anything new. When this happens, the regression model cannot separate the individual effects of the variables, and the estimator we want to use, to estimate the unknown parameters of the model, breaks down because it cannot distinguish which variable is responsible for changes in the outcome.

Consider the regression model

$$wage = \beta_0 \cdot 1 + \beta_1 d^{female} + \beta_2 d^{male} + u$$

where

$$d_i^{female} = \begin{cases} 1 & if & i = female \\ 0 & if & i = male \end{cases}$$

$$d_i^{male} = \begin{cases} 0 & if & i = female \\ 1 & if & i = male \end{cases}$$

Sum of the values in each row of $d^{female}$ and $d^{male}$ is equal to the value in that row of 1. Hence, one value can be perfectly predicted from other values. This is perfect collinearity.

$$\begin{bmatrix} 1 & d^{female} & d^{male} \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix}$$

If we include a constant and both male and female dummies, the female dummy is always equal to 1 – male dummy. This redundancy makes it impossible for the regression model to separate their effects. Each regressor must bring unique information to the model.

Exogeneity.

$$\mathrm{E}\,[u \mid x] = 0$$

This condition says that, once we account for $x$ in the model, the remaining error, $u$, is just random noise. It does not systematically lean up or down depending on $x$. In other words: given any value of $x$, the model is right on average. This assumption is also sometimes called the zero conditional mean assumption.

This assumption has two implications. First, the law of iterated expectations requires that

$$E[u] = E_X[E[u \mid x]]$$

Since

$$E[u \mid x] = 0$$

we have

$$E[u] = 0$$

Second, the law of iterated expectations requires that

$$E[ux] = E_X[E[ux \mid x]] = E_X[xE[u \mid x]]$$

Since

$$E[u \mid x] = 0$$

we have

$$E[ux] = 0$$

Now, from See Math refresher B in Wooldridge, we know that

$$\text{Cov}[u, x] = \text{E}[ux] - \text{E}[u]\text{E}[x]$$

Since

$$\text{E}[ux] = 0$$

and

$$\text{E}[u] = 0$$

we have

$$\text{Cov}[u, x] = 0$$

So the assumption implies that $u$ is not correlated with $x$.

Suppose we want to predict exam scores from study hours. Suppose students who study more also have hidden advantages, like natural ability. Since they are hidden, they show up in the error term. This makes the error correlated with study hours. The result is that the regression model mixes hidden ability with effort. So the slope is biased.

In Greek, homo means same, and skedasis means dispersion. Hence, homoskedasticity means same dispersion. Constant variance.

Homoskedasticity.

$$\text{Var}\left(u \mid x\right) = \sigma^2$$

This condition says that the error term $u$ has the same variance at all values of $x$. It means that, regardless of the level of $x$, the variability of $u$ is constant.

Suppose we want to predict wages from education. Individuals with little education face limited job options, so their wages vary less. The more educated can access a wider variety of jobs paying different wages, so their wages vary more. This means the variance of $u$ depends on $x$. The result is that the regression misjudges its own precision: when errors vary with education, it becomes harder to get a precise prediction of how wages vary with education.

The errors are normally distributed. $u$ has the mean and variance given by the exogeneity and homoskedasticity assumptions, and has a normal distribution:

$$u \mid x \sim N\left[0, \sigma^2\right]$$

We will use this assumption if $n$ is small but drop it if $n$ is large. This is a technical assumption. We will get back to this.

Random sampling. The data we collect

$$\{(x_i, y_i) : i = 1, 2, \ldots, n\}$$

is a random sample of the population that follows the population model. It says that each unit in the population has equal probability of selection, and selections are independent: observations are independently and identically distributed.

Suppose we want to know about the average education in the population. If only sample university students, the sample no longer represents the population. Then the model predictions will be systematically biased.

- Randomness does not eliminate sampling variation. Any single random sample may be unrepresentative, we can be unlucky. Representativeness is probabilistic, not absolute.
- Random sampling ensures the sample is expected to represent the population. Across repeated samples, the sampling distribution of statistics, e.g., means, proportions, etc. converges to the population values.
- Randomness protects against systematic bias, but not against chance variation in one particular sample.

A single random sample can be unrepresentative, but the distribution of a sample statistic across repeated samples centers on the true population value:

- In a 50/50 male-female population, one random sample of size $n = 10$ might yield 7 males and 3 females ($\hat{p} = 0.7$), which is not representative of $p = 0.5$.
- If we repeat this process hundreds of times, the sampling distribution of $\hat{p}$ clusters around the true proportion $p = 0.5$.
- Random sampling guarantees representativeness in expectation, not in every single draw.

Sampling Distribution of Sample Proportions (n=10, 500 samples)