# Standard linear regression model, model assumptions, OLS approximation

Empirical Methods, Lecture 2

Tunga Kantarcı, FEB, Groningen University, Spring 2025

Hereafter SLM stands for the standard linear regression model.
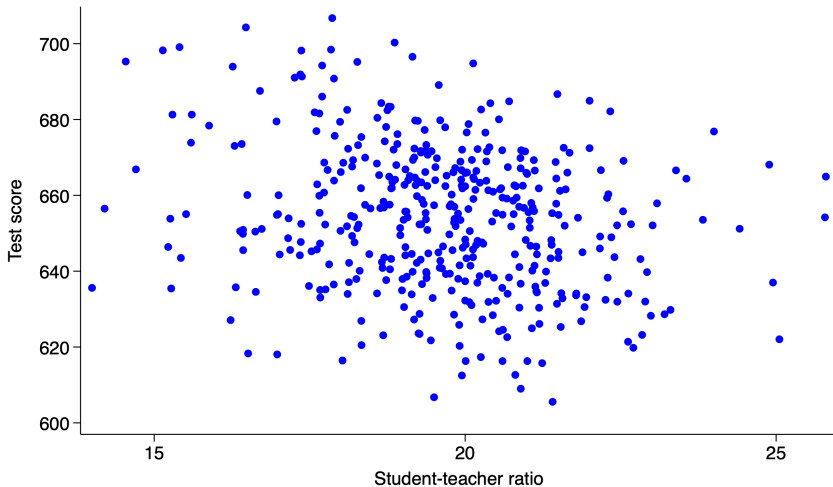
## SLM, an empirical question

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).

```
. 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
```

# SLM, an empirical question

Scatter plot of test scores vs student-teacher ratio:

# SLM, an empirical question

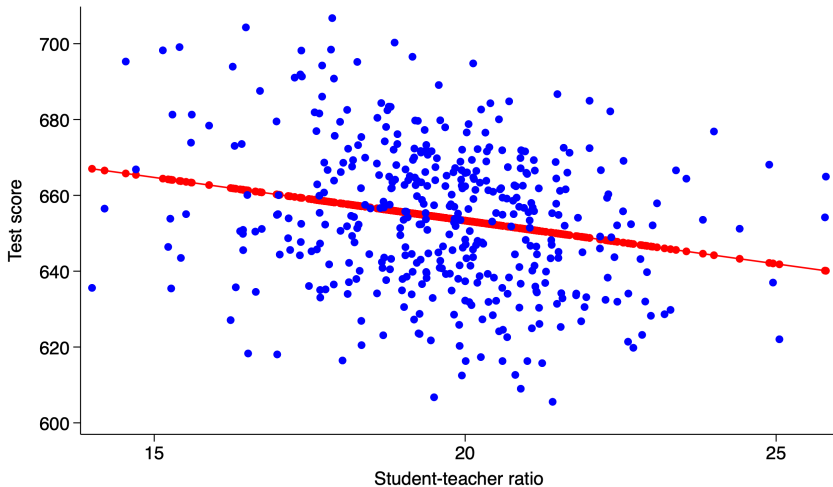We want to analyze how test scores vary with the student-teacher ratio.

To do this we will use the standard linear regression model.

We will first discuss the model and the assumptions we make while using the model.

We will then discuss the method to use to fit this model to the data.

# SLM, an empirical question

Scatter plot of test scores vs student-teacher ratio, fitted model:

But first we need to be clear about notation.

$$\boldsymbol{y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}.$$

$$
\underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}}_{\boldsymbol{y}}
=
\underbrace{\begin{bmatrix} 1 & x_{12} & \ldots & x_{1K} \\ 1 & x_{22} & \ldots & x_{2K} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n2} & \ldots & x_{nK} \end{bmatrix}}_{\boldsymbol{X}}
\underbrace{\begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_K \end{bmatrix}}_{\boldsymbol{\beta}}
+
\underbrace{\begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}}_{\boldsymbol{\varepsilon}}
$$

$$\boldsymbol{y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

$\boldsymbol{y}$: dependent variable. n×1. The bold font is for observations. A vector is always a column vector.

$y_i$: an observation in a row of $\boldsymbol{y}$.

i: unit of study.

$$\boldsymbol{y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

$\boldsymbol{X}$: matrix of variables. $n \times K$. The bold font indicates multiple observations. The big font indicates multiple variables.

$\boldsymbol{x}_k$: a column in $\boldsymbol{X}$. $n \times 1$. It contains $n$ observations for variable $k$. $k$, $l$, $m$ are used to indicate different columns. The bold font indicates multiple observations.

$\boldsymbol{x}'_i$: a row in $\boldsymbol{X}$. $1 \times K$. It contains observations for $K$ variables for unit $i$. $i$, $j$, $t$, $s$ are used to indicate different rows. The bold font indicates multiple variables.

$\boldsymbol{x}_i$: column vector formed by the transpose of a row in $\boldsymbol{X}$. K$\times$1.

$x_{ik}$: an observation in row $i$, column $k$ of $\boldsymbol{X}$.

$$\boldsymbol{y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

$\boldsymbol{\beta}$: true, or population, coefficient vector. K×1. Unobserved.

$\beta_k$: a coefficient in a row of $\boldsymbol{\beta}$. These are slope parameters.

If you let $\boldsymbol{x}_0$ be a column of 1s, $\beta_0$ is the constant term, or the intercept.

$$\boldsymbol{y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

$\boldsymbol{\varepsilon}$: error. n×1. Unobserved.

$\varepsilon_i$: an element in a row of $\boldsymbol{\varepsilon}$.

So for $i$, we have

$$y_i = \mathbf{x}'_i \boldsymbol{\beta} + \varepsilon_i$$

and if

$$\mathbf{x}' = \begin{bmatrix} 1 & x \end{bmatrix}$$

and

$$\boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix},$$

we have

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i.$$

What is a linear regression model?

Which assumptions make a liner regression model 'standard'?

A1. Linearity: the model is linear in the parameters.

The model

$$y_i = \beta_1 + \beta_2 x_{i2}^2 + \varepsilon_i$$

is linear in the parameters, nonlinear in the regressor, linear in the squared regressor

The model

$$y_i = x_{i2}^{\beta_2} + \varepsilon_i$$

is nonlinear in the parameter.

A2. Full column rank: $rank\,(\mathbf{X}) = K$. Remember that $\mathbf{X}$ is $n \times K$ matrix. It contains $K$ columns. Hence, A2 means $\mathbf{X}$ has full column rank.

This means that the columns of X are linearly independent, meaning no column can be written as a linear combination of other columns.

A2 is not satisfied in two cases.

First, if $n < K$. Note that $rank\,(\boldsymbol{X}) \leq min\,(n, K)$. Hence, $rank\,(\boldsymbol{X})$ cannot be K if $n < K$. In practice this is not likely.

Second is the case where there is an exact relationship among any of the columns of $\boldsymbol{X}$.

E.g., consider the regression

$$wage_i = x_{0i}\beta_0 + d_i^{female}\beta_1 + d_i^{male}\beta_2 + \varepsilon_i$$

where $x_{0i} = 1$, and

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

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

Sum of the values in each row of $\boldsymbol{d}^{female}$ and $\boldsymbol{d}^{male}$ is equal to the value in that row of $\boldsymbol{x}_0$. Hence, one value can be perfectly predicted from other values. $rank(\boldsymbol{X}) \neq K$. This is perfect multicollinearity.

$$\begin{bmatrix} \boldsymbol{x}_0 & \boldsymbol{d}^{female} & \boldsymbol{d}^{male} \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix}$$

# SLM, assumptions, full column rank

We will learn later in this lecture about the OLS method to estimate $\boldsymbol{\beta}$.

Perfect multicollinearity is a problem for estimating $\boldsymbol{\beta}$. The OLS estimator of $\boldsymbol{\beta}$ is given by

$$\hat{\boldsymbol{\beta}}_{OLS} = (\boldsymbol{X}'\boldsymbol{X})^{-1}\boldsymbol{X}'\boldsymbol{y}.$$

Rank of $\boldsymbol{X}$ is not $K$. Hence rank of $\boldsymbol{X}'\boldsymbol{X}$ is not $K$. Square matrices are invertible if they have full rank. $\boldsymbol{X}'\boldsymbol{X}$ does not have full rank and hence it is not invertible. This implies that $\hat{\boldsymbol{\beta}}$ has multiple solutions.

This means that when there is perfect multicollinearity between variables, we cannot use the OLS method. This is why Stata will drop a variable to avoid perfect multicollinearity with that variable.

A3. Strict exogeneity:

$$E\left[\varepsilon_i \mid \boldsymbol{x}_k\right] = 0.$$

What does this moment condition say? Recall that $\boldsymbol{x}_k$ contains $n$ observations for variable $k$. The stated condition says that the expected value of the error at observation $i$ in the sample is independent of the explanatory variable $k$ observed at any observation, including observation $i$. It says that the average of the error is the same across all observations of the independent variable, and that this average is 0. More on this later.

Why is it strict? Take a look at the definition of weak exogeneity:

$$E\left[\varepsilon_i \mid x_{ik}\right] = 0.$$

$x_{ik}$ is observation $i$ for variable $k$. That is, we do not consider all $n$ observations of variable $k$, but just observation $i$. That is why

$$E\left[\varepsilon_i \mid \mathbf{x}_k\right] = 0$$

is strict, since all $n$ observations of variable $k$ are considered.

Note that strict exogeneity,

$$E[\varepsilon_i \mid \mathbf{x}_k] = 0,$$

can be considered to apply to all $K$ variables as

$$E[\varepsilon_i \mid \mathbf{X}] = 0.$$

But this is beside the point. What makes it strict is about $n$ not $K$.

Why do we need A3? The model is

$$y = X\beta + \varepsilon.$$

Taking the expectation conditional on $X$,

$$\begin{aligned} E\left[y \mid X\right] &= E\left[X\beta \mid X\right] + E\left[\varepsilon \mid X\right] \\ &= X\beta. \end{aligned}$$

That is, A3 gives the conditional expectation function, or the population regression function.

Later we will study the other reasons for A3, and how we can relax it.

A3 has two implications. First, by the LIE,

$$E\left[\varepsilon_i\right] = E_{\boldsymbol{X}}\left[E\left[\varepsilon_i \mid \boldsymbol{X}\right]\right] = 0.$$

Second, note that

$$\text{Cov}\left[\varepsilon_i, \boldsymbol{X}\right] = \text{E}\left[\varepsilon_i \boldsymbol{X}\right] - \text{E}\left[\varepsilon_i\right]\text{E}\left[\boldsymbol{X}\right]$$

and

$$\text{E}\left[\varepsilon_i \boldsymbol{X}\right] = \text{E}_{\boldsymbol{X}}\left[\text{E}\left[\varepsilon_i \boldsymbol{X} \mid \boldsymbol{X}\right]\right] = \text{E}_{\boldsymbol{X}}\left[\boldsymbol{X}\text{E}\left[\varepsilon_i \mid \boldsymbol{X}\right]\right].$$

Hence, if

$$\text{E}\left[\varepsilon_i \mid \boldsymbol{X}\right] = \boldsymbol{0},$$

then

$$\text{Cov}\left[\varepsilon_i, \boldsymbol{X}\right] = \boldsymbol{0}.$$

It says that $\varepsilon_i$ is not correlated with $\boldsymbol{X}$, or any function of $\boldsymbol{X}$.

$$E[\varepsilon_i \mid \boldsymbol{X}] = \boldsymbol{0}$$

can be easily violated. E.g., suppose

$$\varepsilon_i^* = \varepsilon_i + \boldsymbol{x}_k \beta_k,$$

where $\beta_k \neq 0$, and $\boldsymbol{x}_k$ is correlated with $\boldsymbol{X}$. Then, $\varepsilon_i^*$ is correlated with $\boldsymbol{X}$ because

$$E[\varepsilon_i^* \mid \boldsymbol{X}] \neq 0.$$

This is restrictive in practice. We would want to include $\boldsymbol{x}_k$ in the model so that

$$E[\varepsilon_i^* \mid \boldsymbol{X}] = 0.$$

But what if $\boldsymbol{x}_k$ is unobserved? We cannot include it.

A4. Errors are homoskedastic and non-autocorrelated.

Homoskedasticity: each $\varepsilon_i$ has the same variance $\sigma^2$ conditional on $\boldsymbol{X}$:
$$\text{Var}\left[\varepsilon_i \mid \boldsymbol{X}\right] = \sigma^2, \ \forall \ i.$$

Nonautocorrelation: each $\varepsilon_i$ is uncorrelated with every other disturbance $\varepsilon_j$ conditional on $\boldsymbol{X}$:

$$\text{Cov}\left[\varepsilon_i, \varepsilon_j \mid \boldsymbol{X}\right] = 0, \ \forall \ i \neq j.$$

Later we will study how we can relax this assumption.

If $E[\varepsilon_i \mid \boldsymbol{X}] = 0$,

$$\text{Var}[\varepsilon_i \mid \boldsymbol{X}] = E[\varepsilon_i^2 \mid \boldsymbol{X}] - (E[\varepsilon_i \mid \boldsymbol{X}])^2 = E[\varepsilon_i \varepsilon_i \mid \boldsymbol{X}] = \sigma^2,$$

and

$$\text{Cov}[\varepsilon_i, \varepsilon_j \mid \boldsymbol{X}] = E[\varepsilon_i \varepsilon_j \mid \boldsymbol{X}] - E[\varepsilon_i \mid \boldsymbol{X}]E[\varepsilon_j \mid \boldsymbol{X}] = E[\varepsilon_i \varepsilon_j \mid \boldsymbol{X}] = 0.$$

The variance-covariance matrix for $n$ errors is

$$\text{Var}[\varepsilon \mid \boldsymbol{X}] = E[\varepsilon\varepsilon' \mid \boldsymbol{X}] - E[\varepsilon \mid \boldsymbol{X}]E[\varepsilon' \mid \boldsymbol{X}] = E[\varepsilon\varepsilon' \mid \boldsymbol{X}].$$

Note that $\varepsilon$ is $n \times 1$, and hence $\varepsilon\varepsilon'$ is $n \times n$. This implies that

$$\text{Var}[\varepsilon \mid \boldsymbol{X}] = E[\varepsilon\varepsilon' \mid \boldsymbol{X}] = \sigma^2 I_n = \sigma^2 \boldsymbol{I}$$

which is a $n \times n$ matrix.

$$
\mathsf{E}\left[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}' \mid \boldsymbol{X}\right] =
\begin{bmatrix}
\mathsf{E}\left[\varepsilon_1\varepsilon_1 \mid \boldsymbol{X}\right] & \mathsf{E}\left[\varepsilon_1\varepsilon_2 \mid \boldsymbol{X}\right] & \dots & \mathsf{E}\left[\varepsilon_1\varepsilon_n \mid \boldsymbol{X}\right] \\
\mathsf{E}\left[\varepsilon_2\varepsilon_1 \mid \boldsymbol{X}\right] & \mathsf{E}\left[\varepsilon_2\varepsilon_2 \mid \boldsymbol{X}\right] & \dots & \mathsf{E}\left[\varepsilon_2\varepsilon_n \mid \boldsymbol{X}\right] \\
\vdots & \vdots & \vdots & \vdots \\
\mathsf{E}\left[\varepsilon_n\varepsilon_1 \mid \boldsymbol{X}\right] & \mathsf{E}\left[\varepsilon_n\varepsilon_2 \mid \boldsymbol{X}\right] & \dots & \mathsf{E}\left[\varepsilon_n\varepsilon_n \mid \boldsymbol{X}\right]
\end{bmatrix}
$$

$$
=
\begin{bmatrix}
\sigma^2 & 0 & \dots & 0 \\
0 & \sigma^2 & \dots & 0 \\
& & \vdots & \\
0 & 0 & \dots & \sigma^2
\end{bmatrix}
$$

$$
=
\underbrace{\begin{bmatrix}
1 & 0 & \dots & 0 \\
0 & 1 & \dots & 0 \\
& & \vdots & \\
0 & 0 & \dots & 1
\end{bmatrix}}_{I_n} \sigma^2
$$

A5. Random sampling: the data $\{(\boldsymbol{x}_i, y_i) : i = 1, 2, \ldots, n\}$ is a random sample following the population model $\boldsymbol{y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$. It says that all elements of the data have the same probability of being selected from the population. That is, the observations are i.i.d. This implies that the data have been chosen to be representative of the population.

The sample selection model deals with situations where this assumption fails. This course does not study this model.
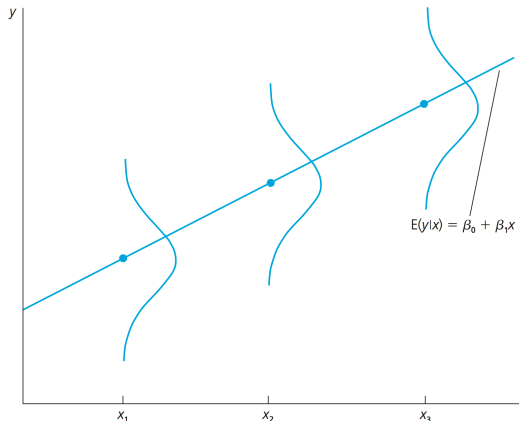
A6. $\varepsilon_i$ is normal. That is, $\varepsilon_i$ has the mean and variance given by A3 and A4, and has a normal distribution. That is,

$$\varepsilon \mid \boldsymbol{X} \sim N\left[\boldsymbol{0}, \sigma^2 \boldsymbol{I}\right].$$

We will use this assumption if $n$ is small. We will drop this assumption if $n$ is large.

# SLM, summary of assumptions

A1: regression line is linear in $\boldsymbol{\beta}$. A3: the conditional expectation function. A4: errors have a constant variance conditional on $\boldsymbol{X}$, and hence so do $\boldsymbol{y}$. The latter because $\mathrm{Var}\left[\varepsilon_i \mid \boldsymbol{X}\right] = \mathrm{Var}\left[y_i \mid \boldsymbol{X}\right]$. A6: errors are normal, and hence so do $\boldsymbol{y}$. The following figure demonstrates all of these assumptions:



$E(y|x) = \beta_0 + \beta_1 x$

## OLS approximation

Consider the SLM

$$\boldsymbol{y} = \boldsymbol{X}\boldsymbol{\beta} + \varepsilon.$$

$\boldsymbol{\beta}$ is unknown and we want to estimate it. The best estimate is the one that makes $\boldsymbol{y}$ as close to $\boldsymbol{X}\boldsymbol{\beta}$ as possible since our aim is to explain $\boldsymbol{y}$ with $\boldsymbol{X}\boldsymbol{\beta}$ as much as possible. Let $\hat{\boldsymbol{\beta}}$ be a candidate for $\boldsymbol{\beta}$ that intends to minimise the sum of squared residuals

$$S(\hat{\boldsymbol{\beta}}) = (\boldsymbol{y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})'(\boldsymbol{y} - \boldsymbol{X}\hat{\boldsymbol{\beta}}).$$

The necessary condition for a minimum is

$$\frac{\partial S(\hat{\boldsymbol{\beta}})}{\partial \hat{\boldsymbol{\beta}}} = -2\boldsymbol{X}'\boldsymbol{y} + 2\boldsymbol{X}'\boldsymbol{X}\hat{\boldsymbol{\beta}} = \boldsymbol{0}.$$

$\partial S(\hat{\boldsymbol{\beta}})/\partial \hat{\boldsymbol{\beta}}$ is calculated using matrix differentiation. If A2 holds, $S(\hat{\boldsymbol{\beta}})$ attains a minimum at $\hat{\boldsymbol{\beta}}_{OLS}$ which takes the form

$$\hat{\boldsymbol{\beta}}_{OLS} = (\boldsymbol{X}'\boldsymbol{X})^{-1}\boldsymbol{X}'\boldsymbol{y}.$$

# OLS approximation

The derivation above used matrix differentiation and notation. In this course we will hardly use matrix algebra and notation.

That is, the model of interest in standard form is

$$y_i = \beta_0 + \beta_1 x_i + u_i.$$

The sum of squared residuals in standard form is then

$$\sum_{i=1}^{N}(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2.$$

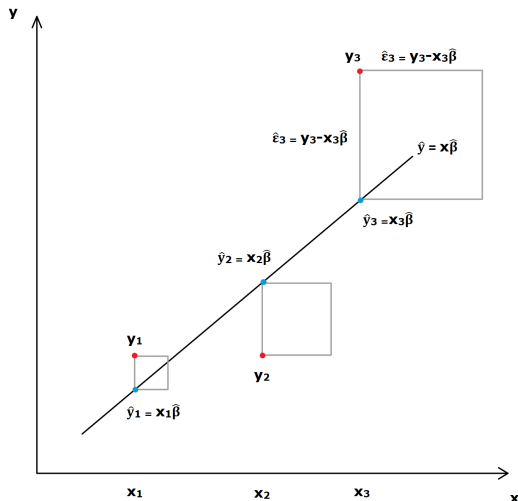Taking the derivatives with respect to $\hat{\beta}_0$ and $\hat{\beta}_1$ and setting them equal to 0, and using some algebra tricks, the OLS estimator of $\beta_0$ is

$$\hat{\beta}_{0,OLS} = \bar{y} - \hat{\beta}_{1,OLS}\bar{x}$$

and the OLS estimator of $\beta_1$ is

$$\hat{\beta}_{1,OLS} = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n}(x_i - \bar{x})}.$$

# OLS approximation

By minimizing the sum of "squared" residuals, OLS fits as good as possible a regression line to the data points.

# OLS approximation

OLS is an approximation, or estimation, method. It is not a model.

The model is the standard liner regression model.

We estimate the population, or slope, parameters of the model using the OLS method.

# OLS approximation

The solution to the least squares problem is

$$\boldsymbol{X}'\boldsymbol{X}\hat{\beta}_{OLS} - \boldsymbol{X}'\boldsymbol{y} = -\boldsymbol{X}'(\boldsymbol{y} - \underbrace{\boldsymbol{X}\hat{\beta}_{OLS}}_{\hat{\boldsymbol{y}}}) = -\boldsymbol{X}'\hat{\varepsilon} = 0.$$

$$\boldsymbol{X}'\boldsymbol{X}\hat{\beta}_{OLS} = \boldsymbol{X}'\boldsymbol{y}$$

are also called the normal equations.

Recall the population regression function given by

$$\mathrm{E}\left[\boldsymbol{y} \mid \boldsymbol{X}\right] = \boldsymbol{X}\boldsymbol{\beta}.$$

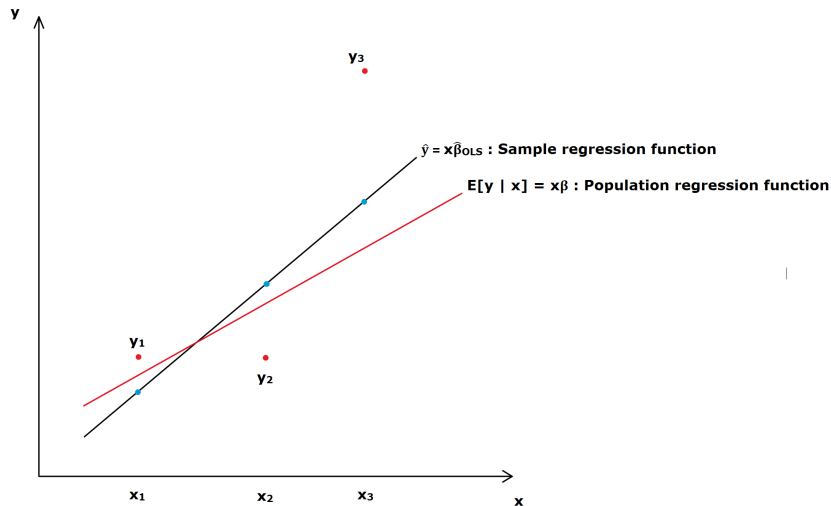The solution to the least squares problem gives

$$\hat{\boldsymbol{y}} = \boldsymbol{X}\hat{\boldsymbol{\beta}}_{OLS}$$

which represent the predictions of the regression model. This is the sample regression function. It is the estimate of the population regression function:

$$\widehat{\mathrm{E}\left[\boldsymbol{y} \mid \boldsymbol{X}\right]} = \boldsymbol{X}\hat{\boldsymbol{\beta}}_{OLS}.$$

# OLS approximation

Compare the unknown population regression function to the sample regression function.

Consider the sample regression function

$$\widehat{\mathsf{E}[\boldsymbol{y} \mid \boldsymbol{X}]} = \boldsymbol{X}\hat{\boldsymbol{\beta}}_{OLS}.$$

Considering discrete changes in the predicted dependent variable and the independent variable, we have:

$$\frac{\Delta \widehat{\mathsf{E}[\boldsymbol{y} \mid \boldsymbol{X}]}}{\Delta \boldsymbol{X}} = \hat{\boldsymbol{\beta}}_{OLS}.$$

This tells that when the independent variable changes by some unit, the dependent variable changes on average by the OLS estimate.
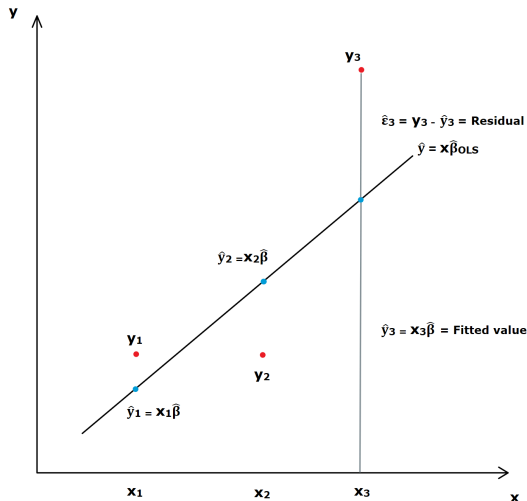
The solution to the least squares problem gives

$$\boldsymbol{y} = \hat{\boldsymbol{y}} + \hat{\varepsilon}$$

which shows that we explain the dependent variable by the prediction of our model and our error.

# OLS approximation

An observation of the dependent variable is explained by the prediction of our model and our error: $y_3 = \hat{y}_3 + \hat{\epsilon}_3$.

## OLS approximation, implications

The solution has three implications:

1. If the first column of $\boldsymbol{X}$, $\boldsymbol{x}_0$, is a column of 1s, i.e. the regression includes a constant, the residuals, or deviations from the regression line, sum to zero:
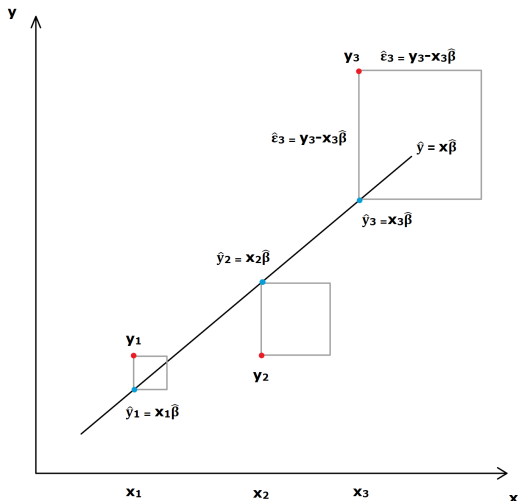
$$\boldsymbol{x}_0'\hat{\boldsymbol{\varepsilon}} = \sum_i^n \hat{\varepsilon}_i = 0.$$

2. $\bar{y} = \bar{\boldsymbol{x}}'\hat{\boldsymbol{\beta}}_{OLS} + \bar{\hat{\varepsilon}}$, and since $\bar{\hat{\varepsilon}} = 0$ by the first implication, $\bar{y} = \bar{\boldsymbol{x}}'\hat{\boldsymbol{\beta}}_{OLS}$. This says that the regression hyperplane passes through the point of means of the data.
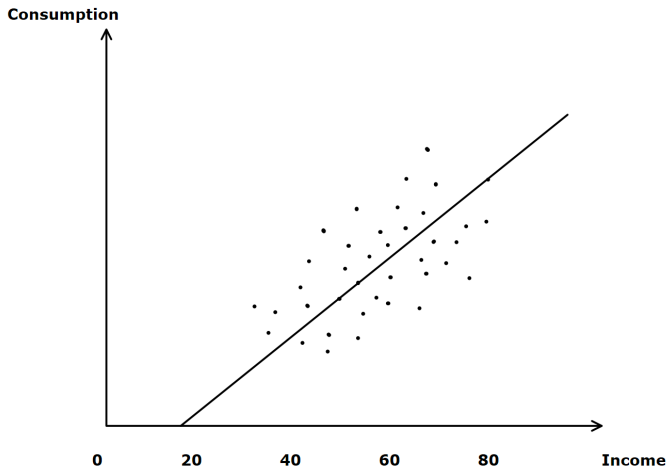
3. $\boldsymbol{y} = \hat{\boldsymbol{y}} + \hat{\boldsymbol{\varepsilon}}$ from the solution. Taking the means, we obtain $\bar{y} = \bar{\hat{y}} + \bar{\hat{\varepsilon}}$. Since $\bar{\hat{\varepsilon}} = 0$ by the first implication, we obtain $\bar{y} = \bar{\hat{y}}$.

# OLS approximation, insights

OLS is not robust to outliers. $y_3$ contributes too much to the minimization problem of OLS: it pulls the regression towards itself.

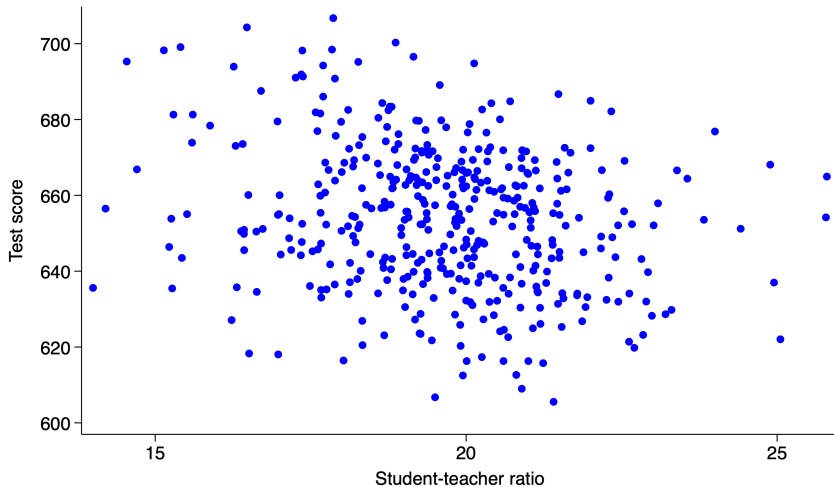For incomes between 40 and 80, the consumption function can be approximated by the line (model). Does the line describe the consumption-income relationship for all incomes, or only for the those in the center? Only in the center! What is the predicted consumption when income is 10? A negative value! Models are approximations. Approximations do not work well if we move too far away from the point of approximation. OLS is a good approximator around the average value of x.

# OLS approximation, example

Scatter plot of test scores vs student-teacher ratio:

Our population regression model is:

$$testscr = \beta_0 + \beta_1 str + u.$$

We want to estimate the population parameters of the model.

# OLS approximation, example

We use the OLS method

```
. regress testscr str

      Source |       SS           df       MS      Number of obs   =       420
-------------+----------------------------------   F(1, 418)       =     22.58
       Model |  7794.11012         1  7794.11012   Prob > F        =    0.0000
    Residual |  144315.484       418  345.252353   R-squared       =    0.0512
-------------+----------------------------------   Adj R-squared   =    0.0490
       Total |  152109.594       419  363.030056   Root MSE        =    18.581

------------------------------------------------------------------------------
     testscr | Coefficient  Std. err.      t    P>|t|     [95% conf. interval]
-------------+----------------------------------------------------------------
         str |  -2.279808   .4798256    -4.75   0.000    -3.22298   -1.336637
       _cons |   698.933    9.467491    73.82   0.000     680.3231   717.5428
------------------------------------------------------------------------------
```

which gives the sample regression model, or the fitted model:

$$\widehat{testscr} = 698.93 - 2.28 \; str$$

# OLS approximation, example

For each $i$, we have a prediction on the regression line:

Interpretation of the estimated coefficient of the independent variable:

$$\frac{\Delta E\,[\widehat{testscr \mid str}]}{\Delta str} = -2.28.$$

On average, a unit increase in student-teacher ratio is associated with a $-2.28$ points decrease in test scores.