

ScPoEconometrics

Simple Linear Regression

Florian Oswald, Gustave Kenedi and Pierre Villedieu
SciencesPo Paris
2021-01-19

Quick "Quiz" on Last Week's Material

1. From your *computer* ↗ connect to www.wooclap.com/SCPOTIDY2

OR

2. From your *phone* ↗ flash QR code below



Today - Real 'metrics finally



- Introduction to the *Simple Linear Regression Model* and *Ordinary Least Squares (OLS) estimation*.
- Empirical application: *class size* and *student performance*
- Keep in mind that we are interested in uncovering **causal** relationships



Class size and student performance

- What policies *lead* to improved student learning?
- Class size reduction has been at the heart of policy debates for *decades*.



Class size and student performance

- What policies *lead* to improved student learning?
- Class size reduction has been at the heart of policy debates for *decades*.
- We will be using data from a famous paper by **Joshua Angrist and Victor Lavy (1999)**, obtained from **Raj Chetty and Greg Bruich's course**.
- Consists of test scores and class/school characteristics for fifth graders (10-11 years old) in Jewish public elementary schools in Israel in 1991.
- National tests measured *mathematics* and (Hebrew) *reading* skills. The raw scores were scaled from 1-100.



Task 1: Getting to know the data

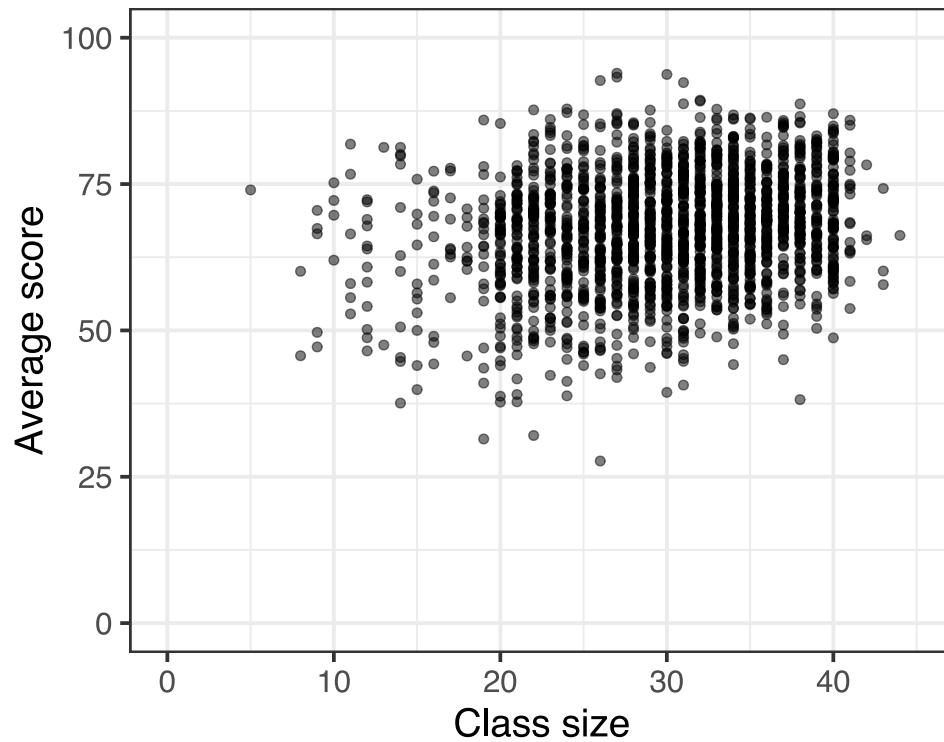
07 : 00

1. Load the data from [here](#) as `grades`. You need to find the function that enables importing `.dta` files. (FYI: `.dta` is the extension for data files used in *Stata*)
2. Describe the dataset:
 - What is the unit of observations, i.e. what does each row correspond to?
 - How many observations are there?
 - What variables do we have? View the dataset to see what the variables correspond to.
 - What do the variables `avgmath` and `avgverb` correspond to?
 - Use the `skim` function from the `skimr` package to obtain common summary statistics for the variables `classize`, `avgmath` and `avgverb`. (*Hint: use `dplyr` to select the variables and then simply pipe (%>%) `skim()`.*)
3. Do you have any priors about the actual (linear) relationship between class size and student achievement? What would you do to get a first insight?
4. Compute the correlation between class size and math and verbal scores. Is the relationship positive/negative, strong/weak?

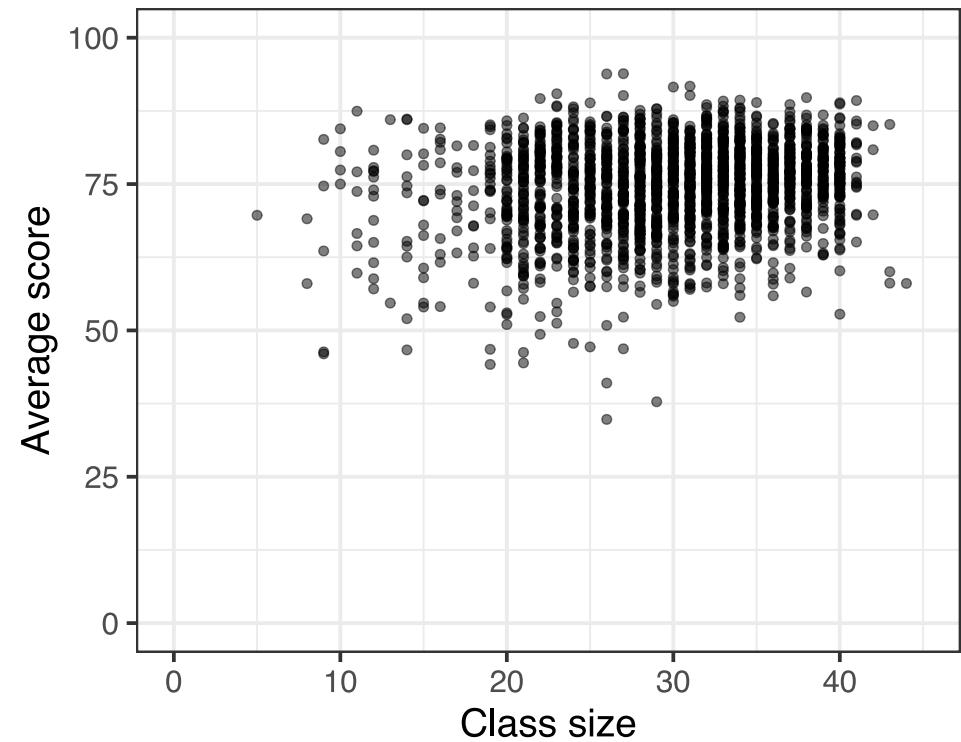


Class size and student performance: Scatter plot

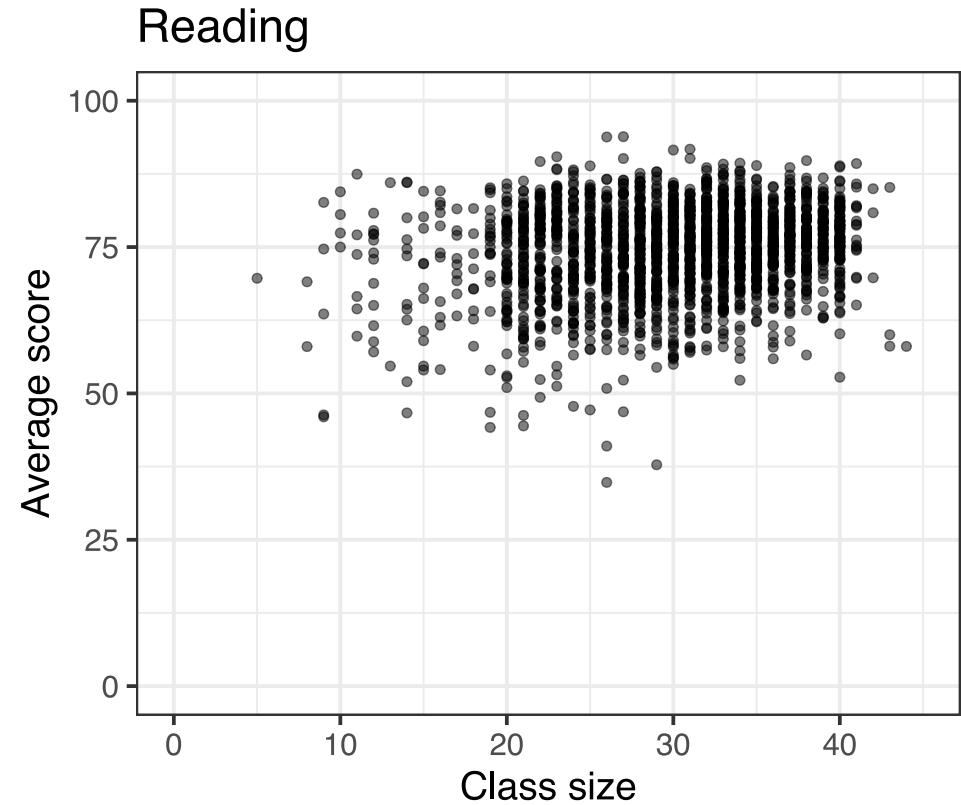
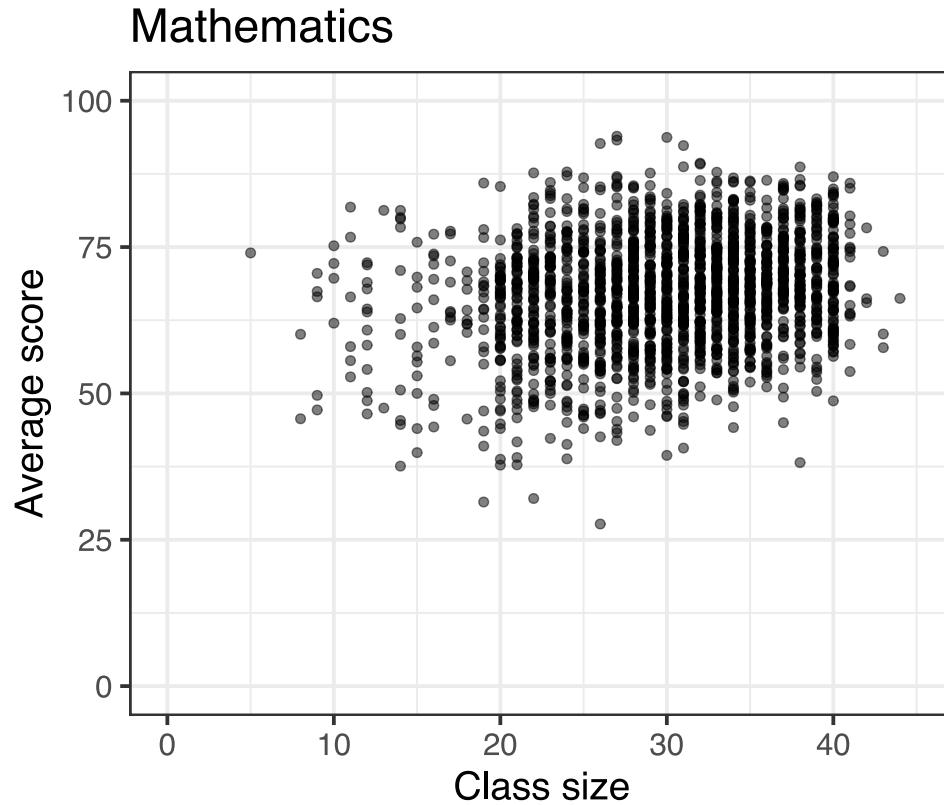
Mathematics



Reading



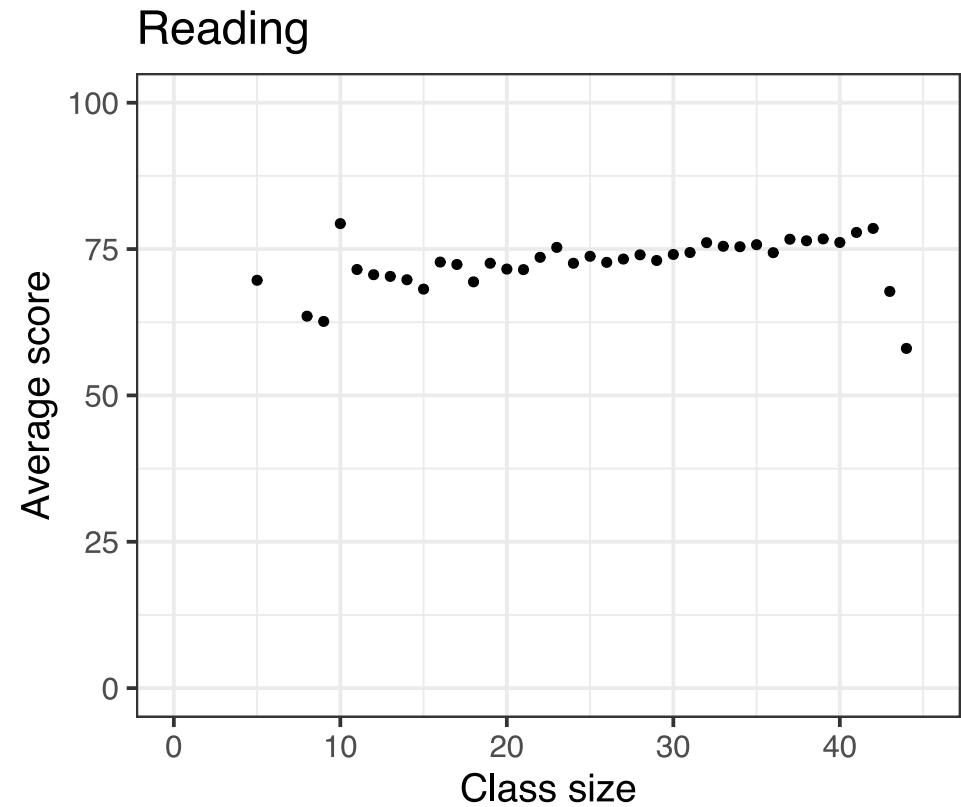
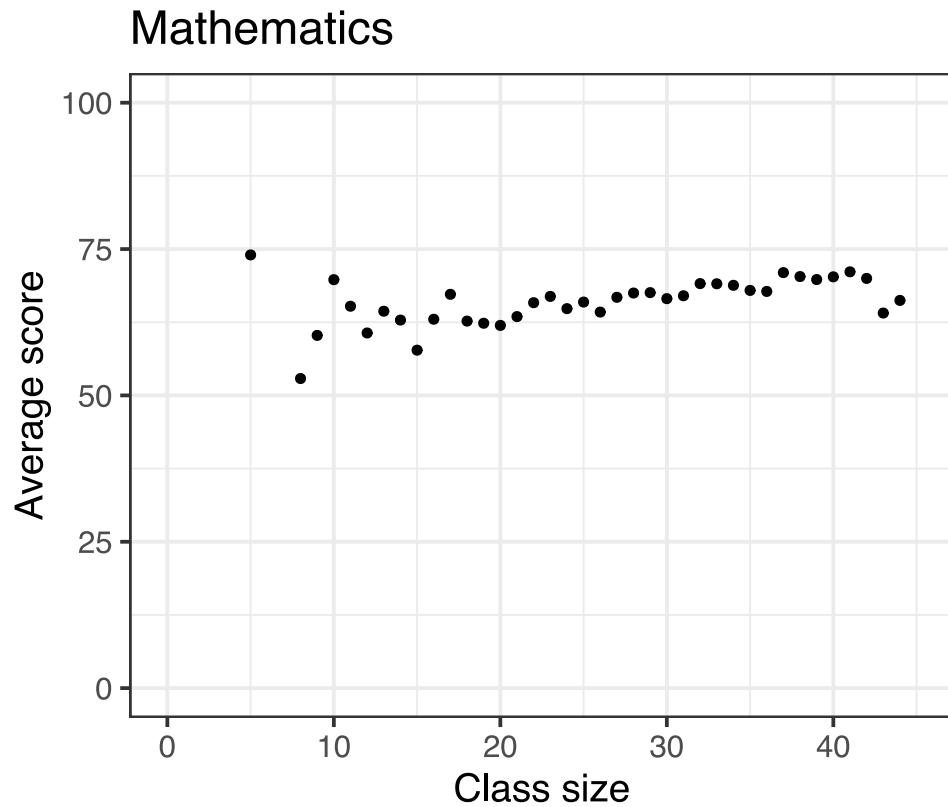
Class size and student performance: Scatter plot



- Somewhat positive association as suggested by the correlations. Let's compute the average score by class size to see things more clearly!

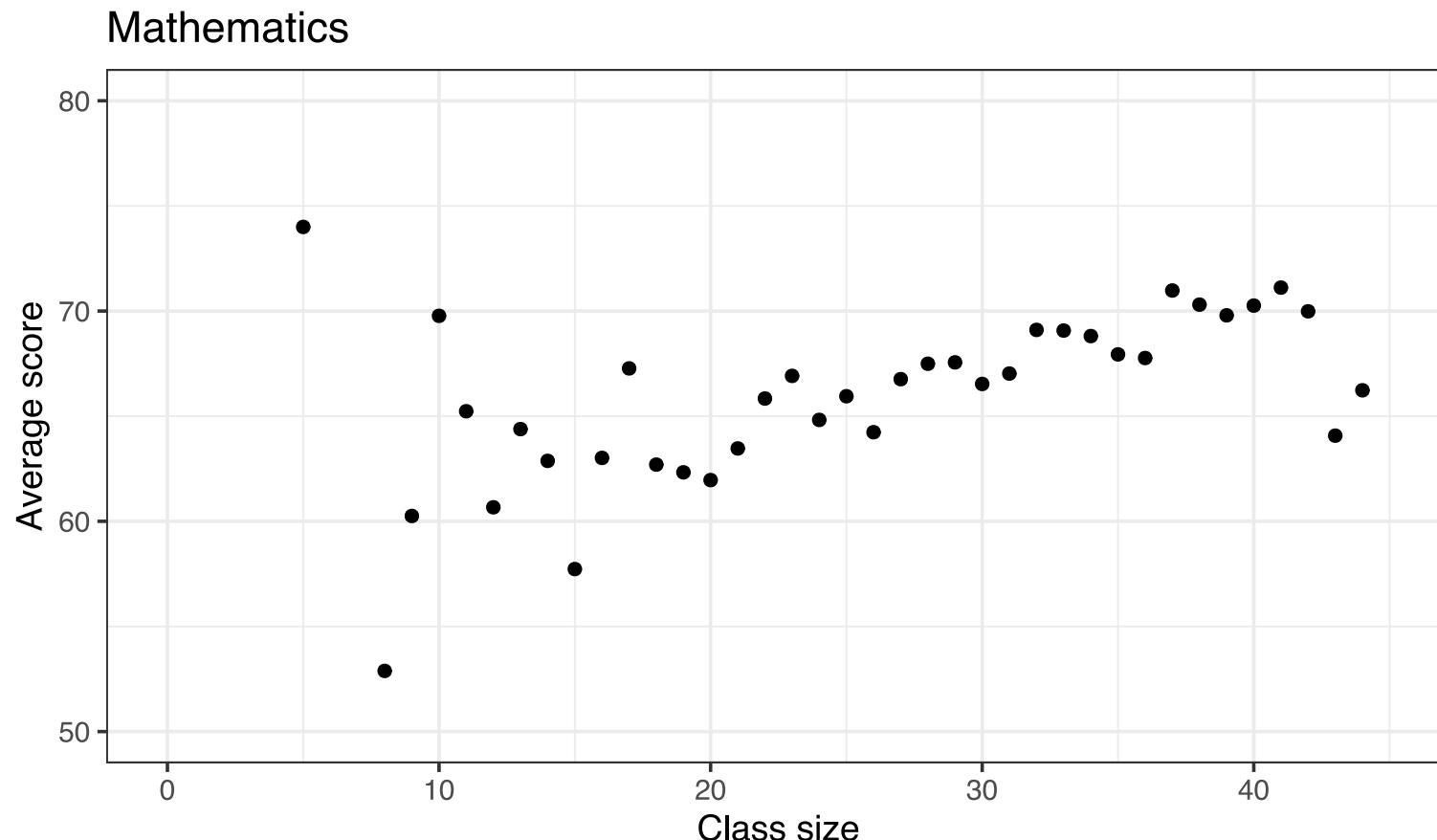


Class size and student performance: Binned scatter plot



Class size and student performance: Binned scatter plot

- We'll first focus on the mathematics scores and for visual simplicity we'll zoom in



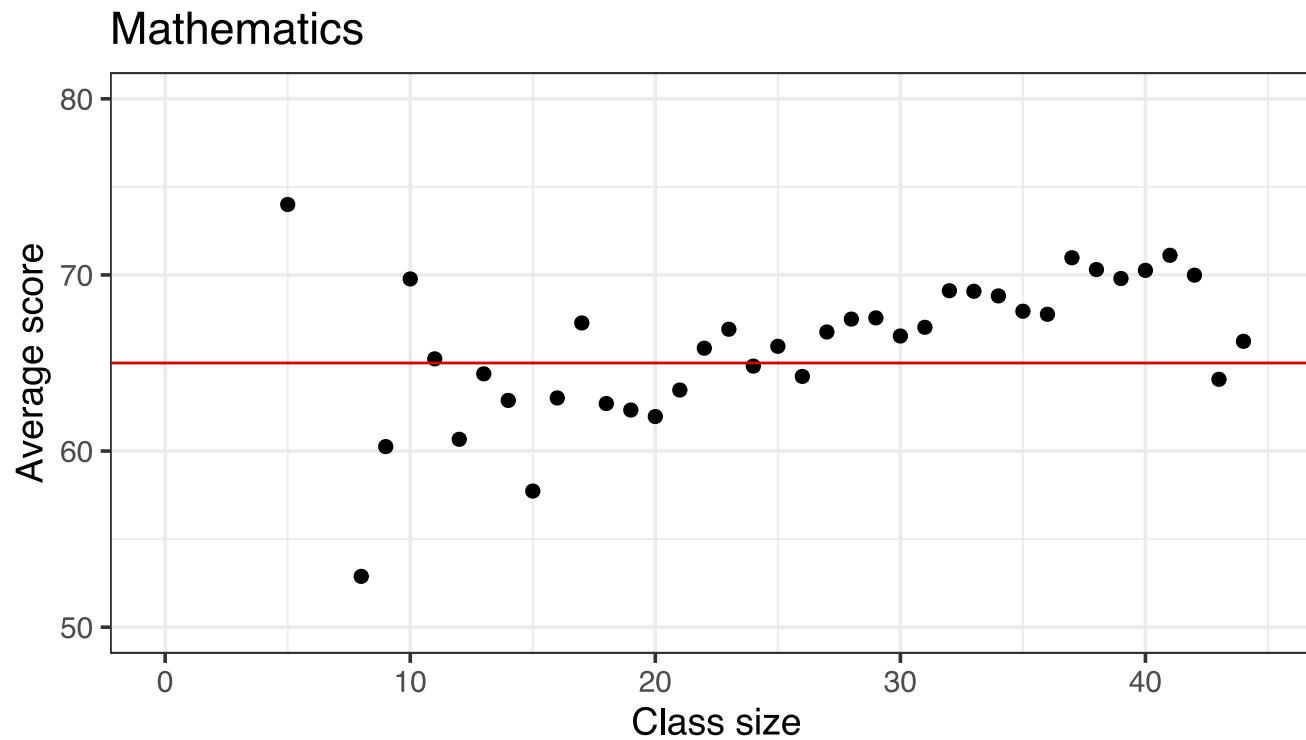
Class size and student performance: Regression line

How to visually summarize the relationship: **a line through the scatter plot**



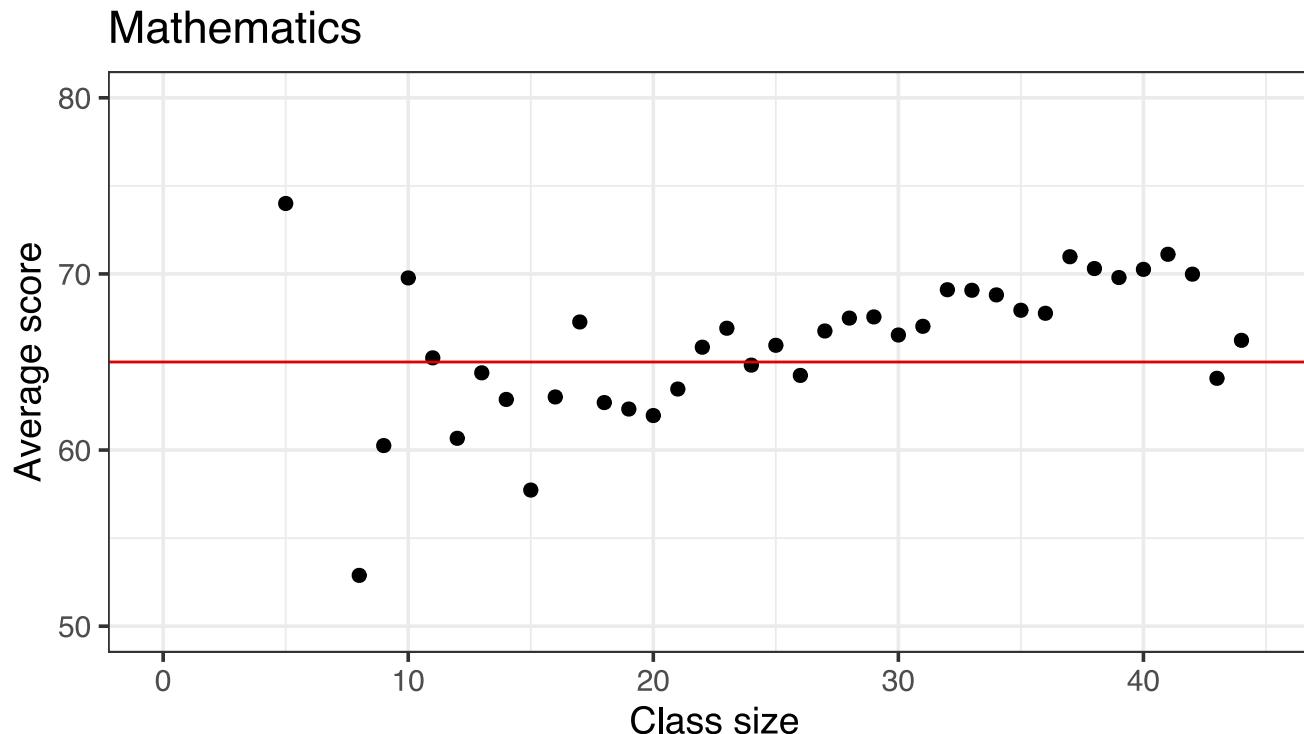
Class size and student performance: Regression line

How to visually summarize the relationship: **a line through the scatter plot**



Class size and student performance: Regression line

How to visually summarize the relationship: **a line through the scatter plot**

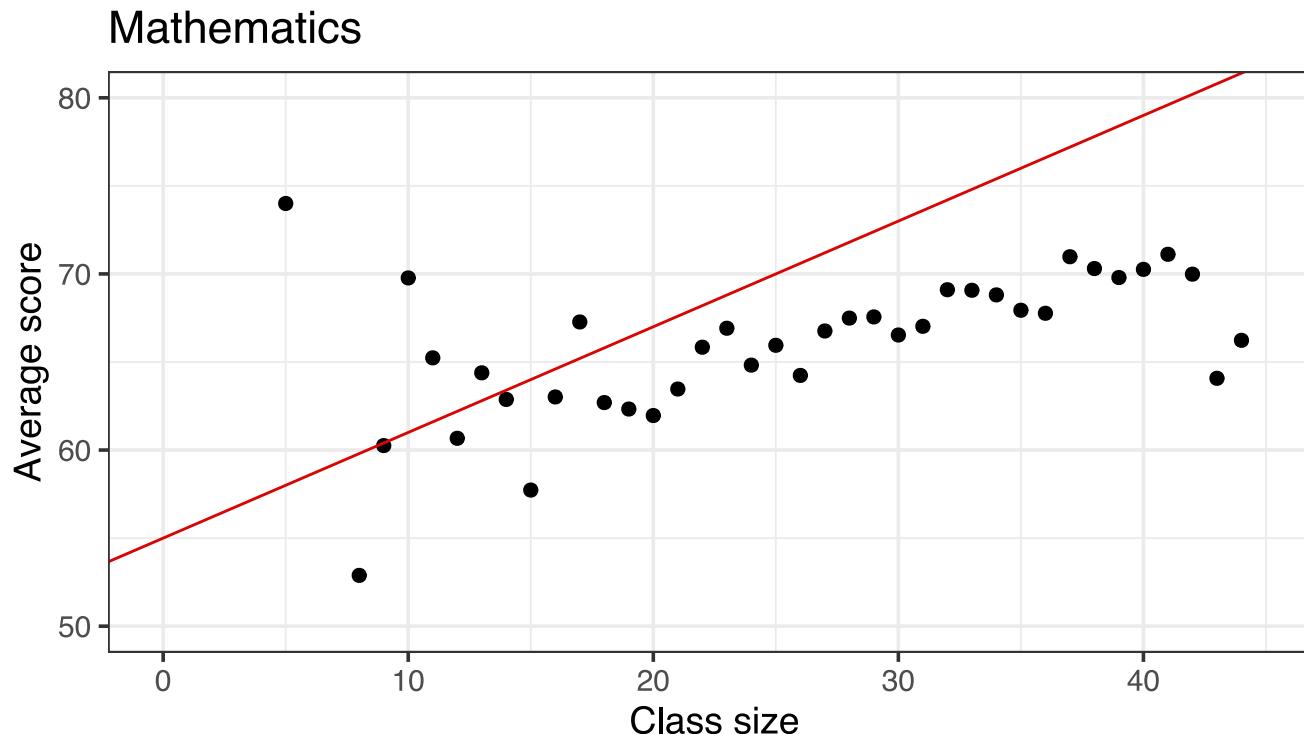


- A *line!* Great. But **which** line? This one?
- That's a *flat* line. But average mathematics score is somewhat *increasing* with class size 😞



Class size and student performance: Regression line

How to visually summarize the relationship: **a line through the scatter plot**



- **That one?**
- Slightly better! Has a **slope** and an **intercept** 😐
- We need a rule to decide!



Simple Linear Regression

Let's formalise a bit what we are doing so far.

- We are interested in the relationship between two variables:



Simple Linear Regression

Let's formalise a bit what we are doing so far.

- We are interested in the relationship between two variables:
 - an **outcome variable** (also called **dependent variable**):
average mathematics score (y)



Simple Linear Regression

Let's formalise a bit what we are doing so far.

- We are interested in the relationship between two variables:
 - an **outcome variable** (also called **dependent variable**):
average mathematics score (y)
 - an **explanatory variable** (also called **independent variable** or **regressor**):
class size (x)



Simple Linear Regression

Let's formalise a bit what we are doing so far.

- We are interested in the relationship between two variables:
 - an **outcome variable** (also called **dependent variable**):
average mathematics score (y)
 - an **explanatory variable** (also called **independent variable** or **regressor**):
class size (x)
- For each class i we observe both x_i and y_i , and therefore we can plot the *joint distribution* of class size and average mathematics score.



Simple Linear Regression

Let's formalise a bit what we are doing so far.

- We are interested in the relationship between two variables:
 - an **outcome variable** (also called **dependent variable**):
average mathematics score (y)
 - an **explanatory variable** (also called **independent variable** or **regressor**):
class size (x)
- For each class i we observe both x_i and y_i , and therefore we can plot the *joint distribution* of class size and average mathematics score.
- We summarise this relationship with a line (for now). The equation for such a line with an intercept b_0 and a slope b_1 is:

$$\hat{y}_i = b_0 + b_1 x_i$$



Simple Linear Regression

Let's formalise a bit what we are doing so far.

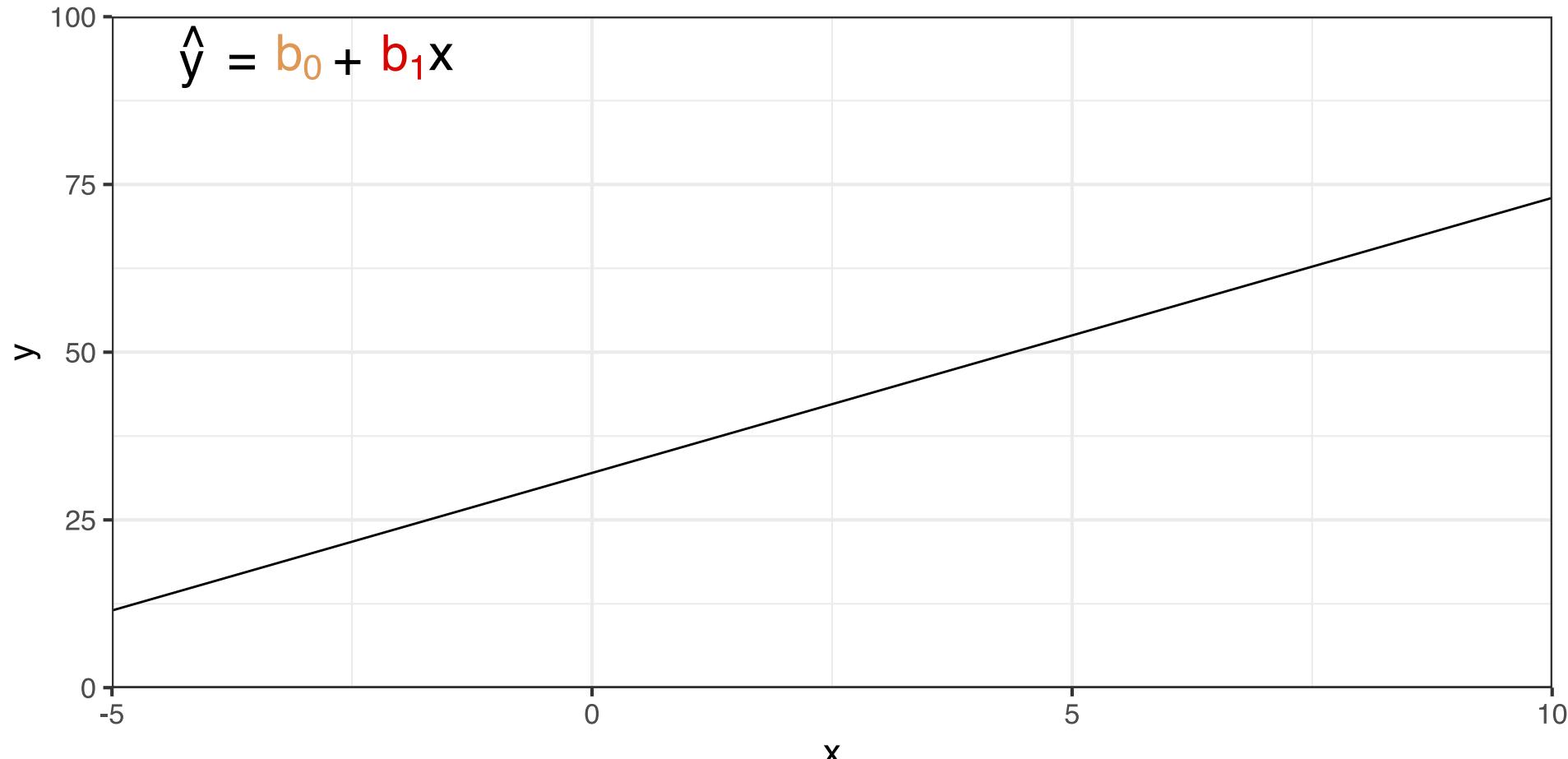
- We are interested in the relationship between two variables:
 - an **outcome variable** (also called **dependent variable**):
average mathematics score (y)
 - an **explanatory variable** (also called **independent variable** or **regressor**):
class size (x)
- For each class i we observe both x_i and y_i , and therefore we can plot the *joint distribution* of class size and average mathematics score.
- We summarise this relationship with a line (for now). The equation for such a line with an intercept b_0 and a slope b_1 is:

$$\hat{y}_i = b_0 + b_1 x_i$$

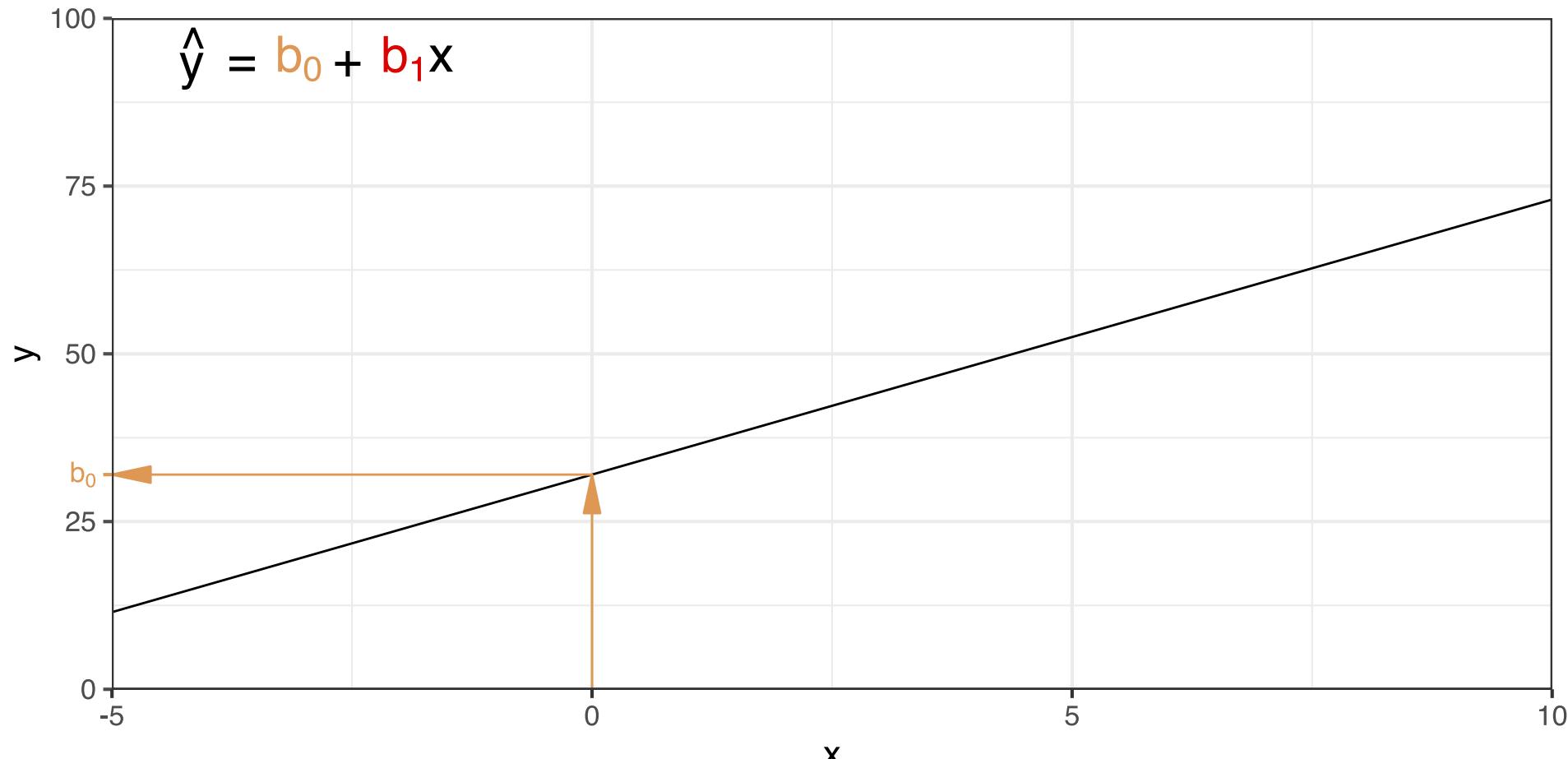
- \hat{y}_i is our *prediction* for y at observation i (y_i) given our model (i.e. the line).



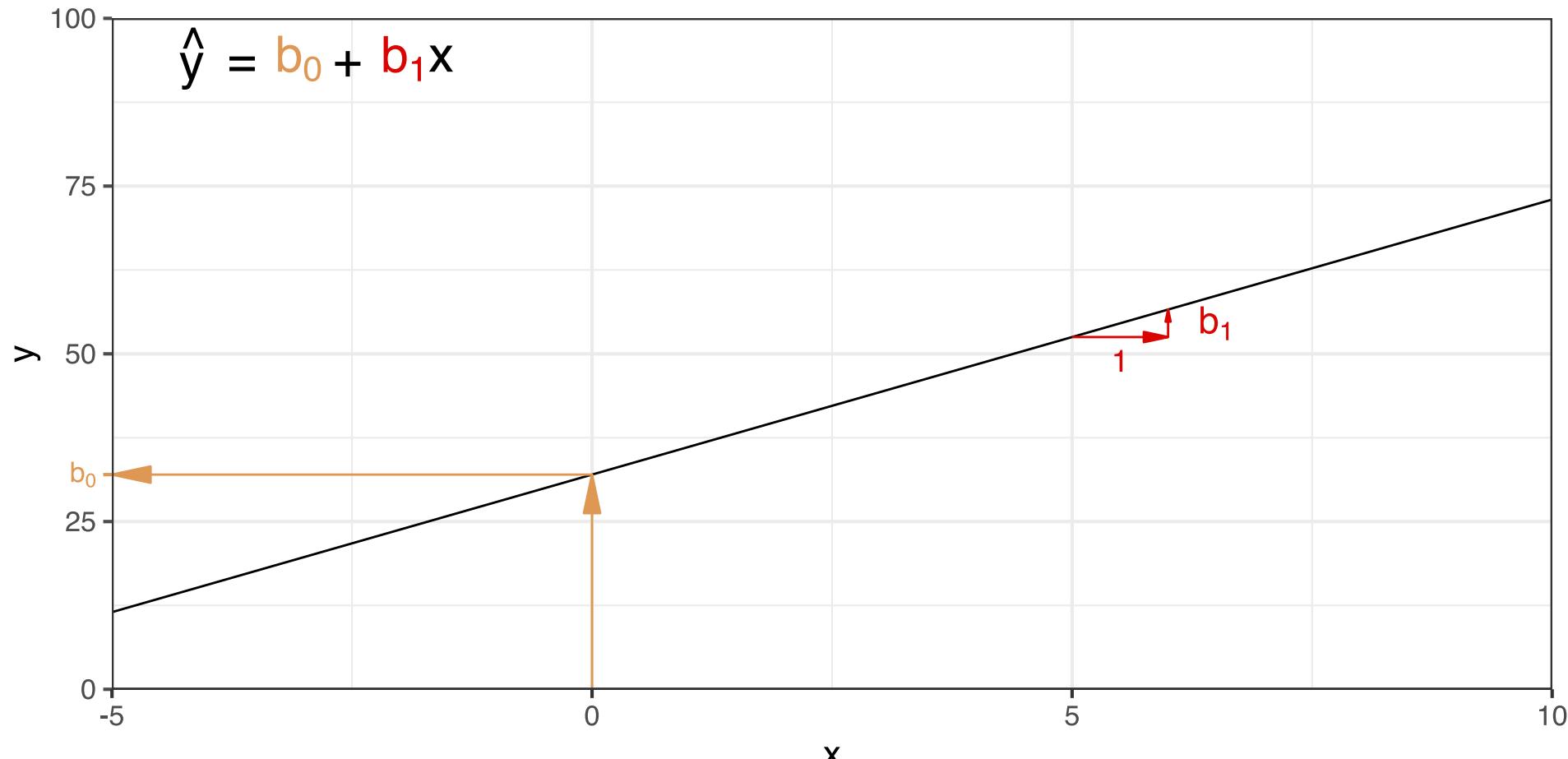
What's A Line: A Refresher



What's A Line: A Refresher



What's A Line: A Refresher



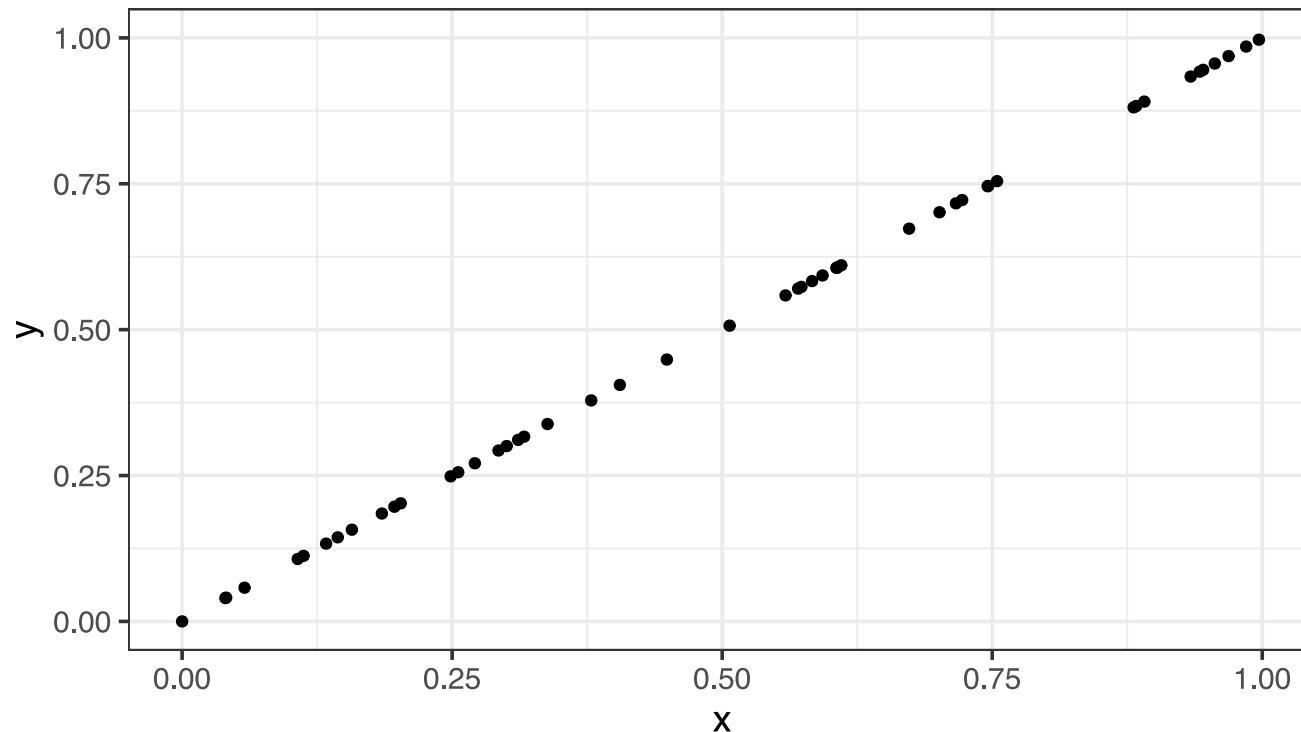
Simple Linear Regression: Residual

- If all the data points were **on** the line then $\hat{y}_i = y_i$.



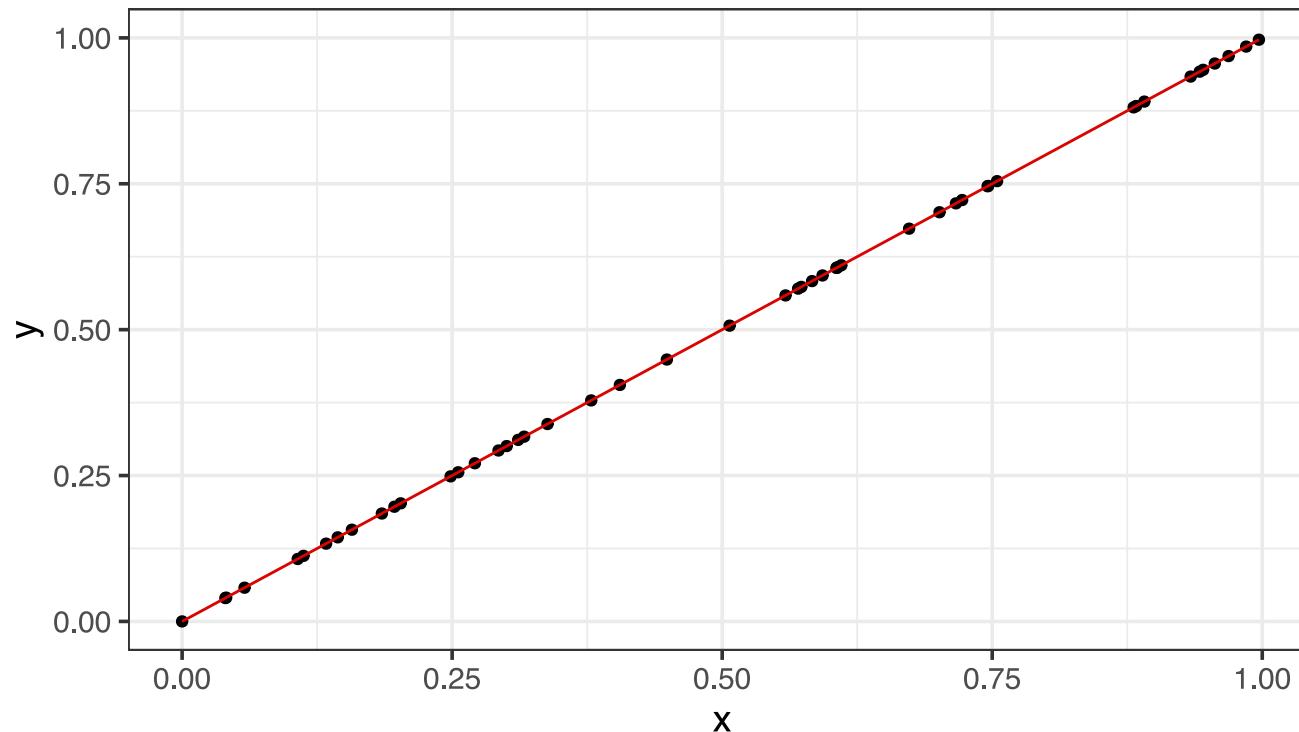
Simple Linear Regression: Residual

- If all the data points were **on** the line then $\hat{y}_i = y_i$.



Simple Linear Regression: Residual

- If all the data points were **on** the line then $\hat{y}_i = y_i$.



Simple Linear Regression: Residual

- If all the data points were **on** the line then $\hat{y}_i = y_i$.
- However, since in most cases the *dependent variable* (y) is not *only* explained by the chosen *independent variable* (x), $\hat{y}_i \neq y_i$, i.e. we make an **error**.
This **error** is called the **residual**.



Simple Linear Regression: Residual

- If all the data points were **on** the line then $\hat{y}_i = y_i$.
- However, since in most cases the *dependent variable* (y) is not *only* explained by the chosen *independent variable* (x), $\hat{y}_i \neq y_i$, i.e. we make an **error**.
This **error** is called the **residual**.
- At point (x_i, y_i) , we note this residual e_i .



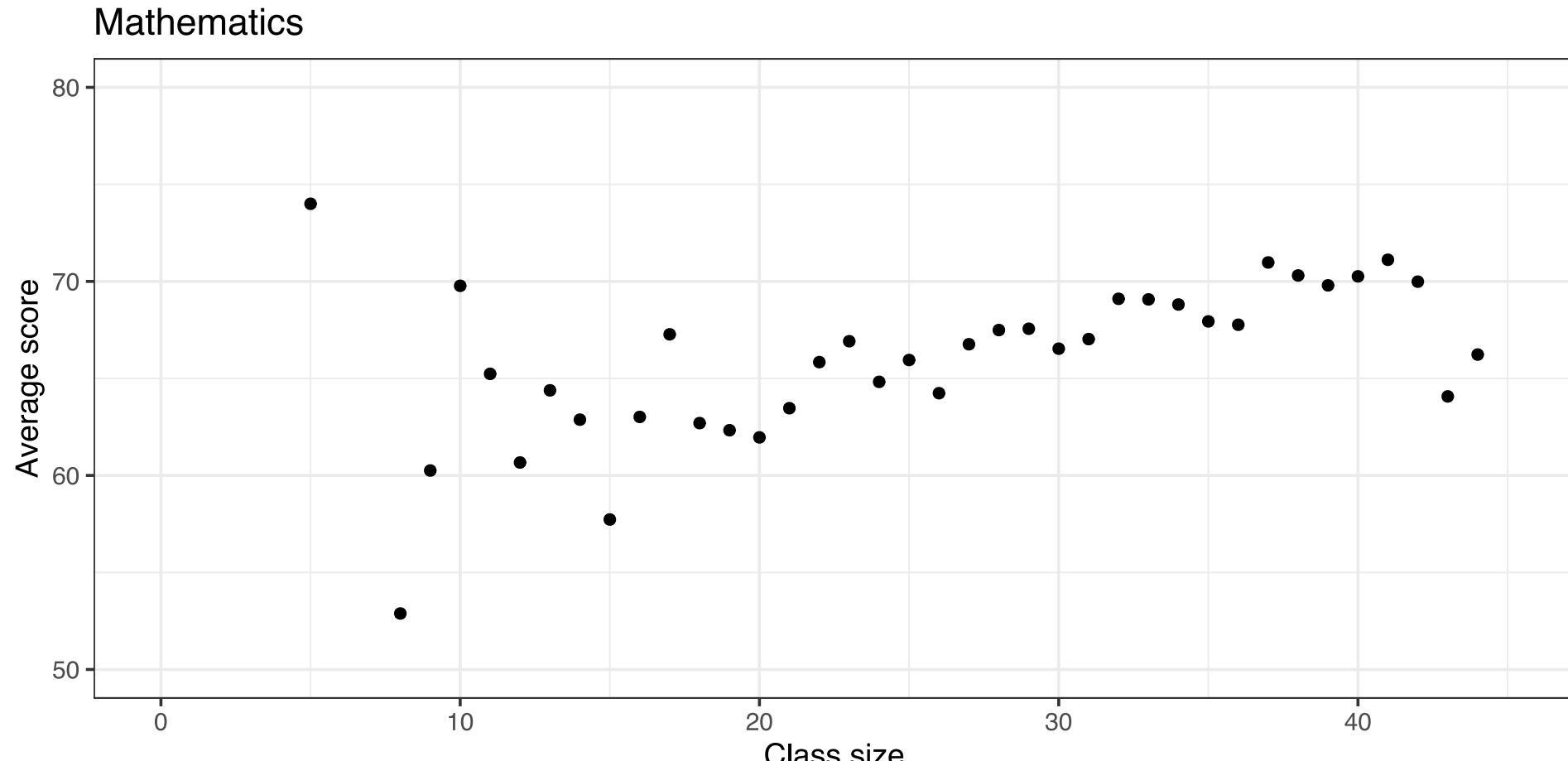
Simple Linear Regression: Residual

- If all the data points were **on** the line then $\hat{y}_i = y_i$.
- However, since in most cases the *dependent variable* (y) is not *only* explained by the chosen *independent variable* (x), $\hat{y}_i \neq y_i$, i.e. we make an **error**.
This **error** is called the **residual**.
- At point (x_i, y_i) , we note this residual e_i .
- The *actual data* (x_i, y_i) can thus be written as *prediction + residual*:

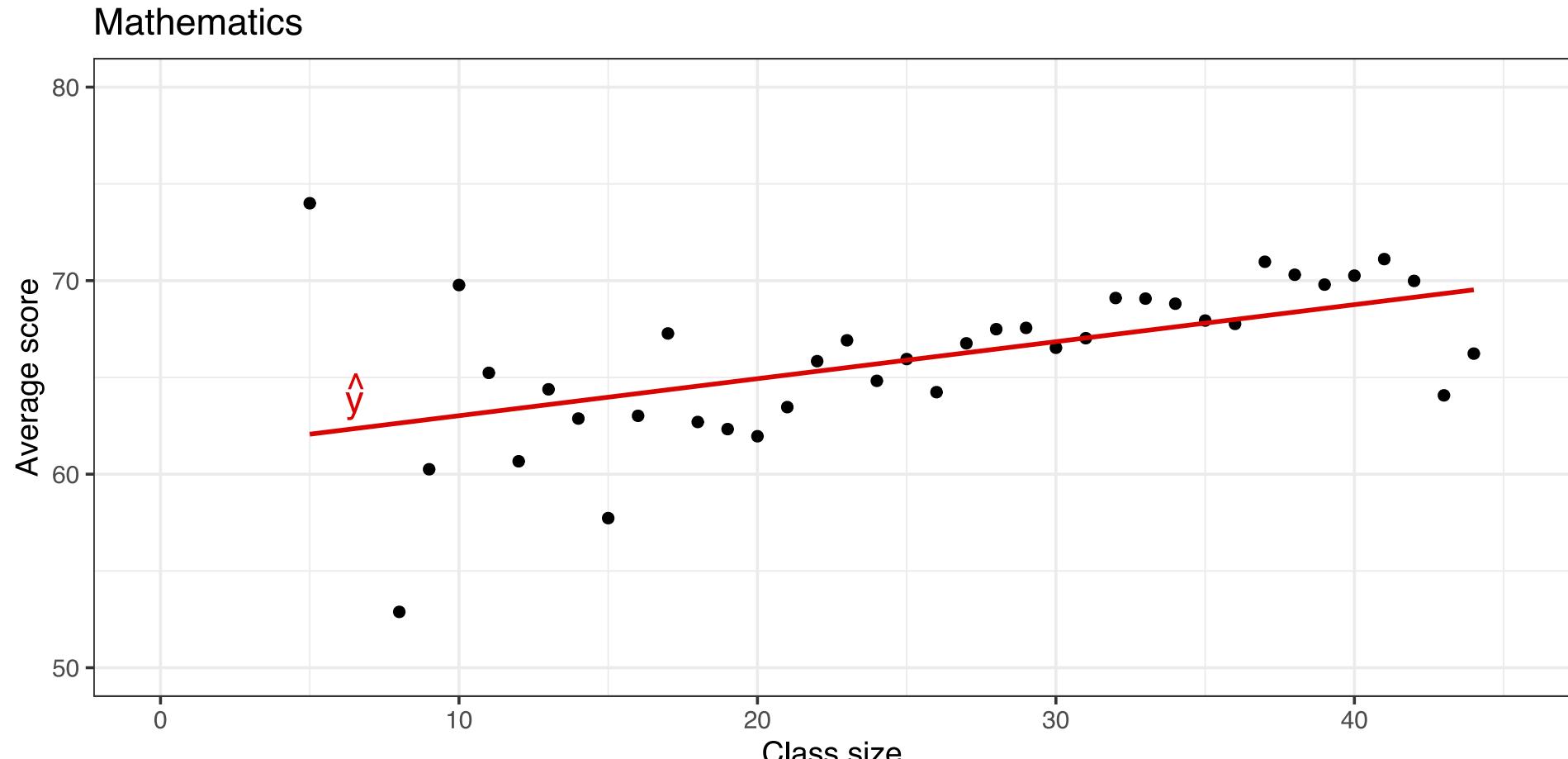
$$y_i = \hat{y}_i + e_i = b_0 + b_1 x_i + e_i$$



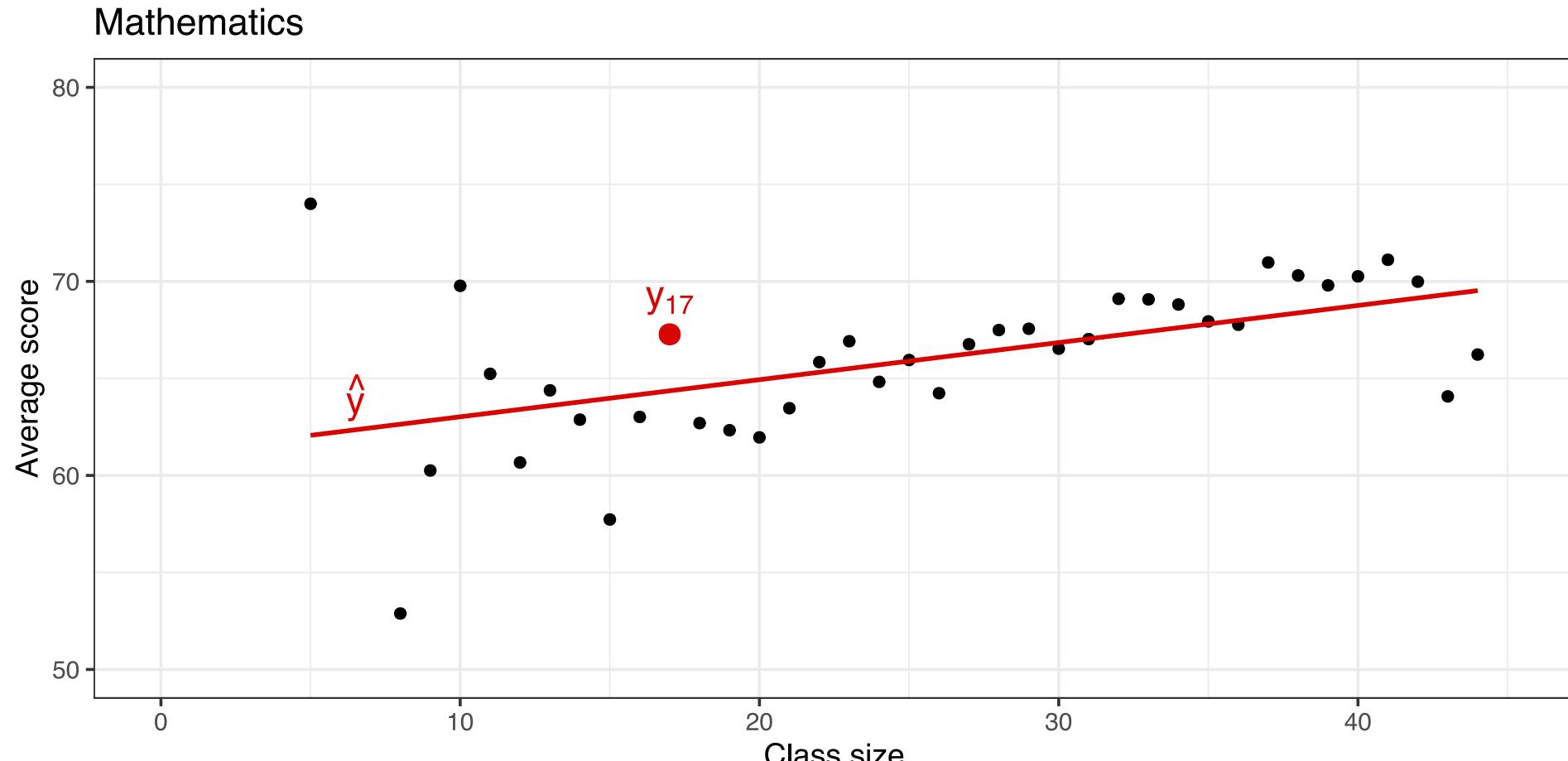
Simple Linear Regression: Graphically



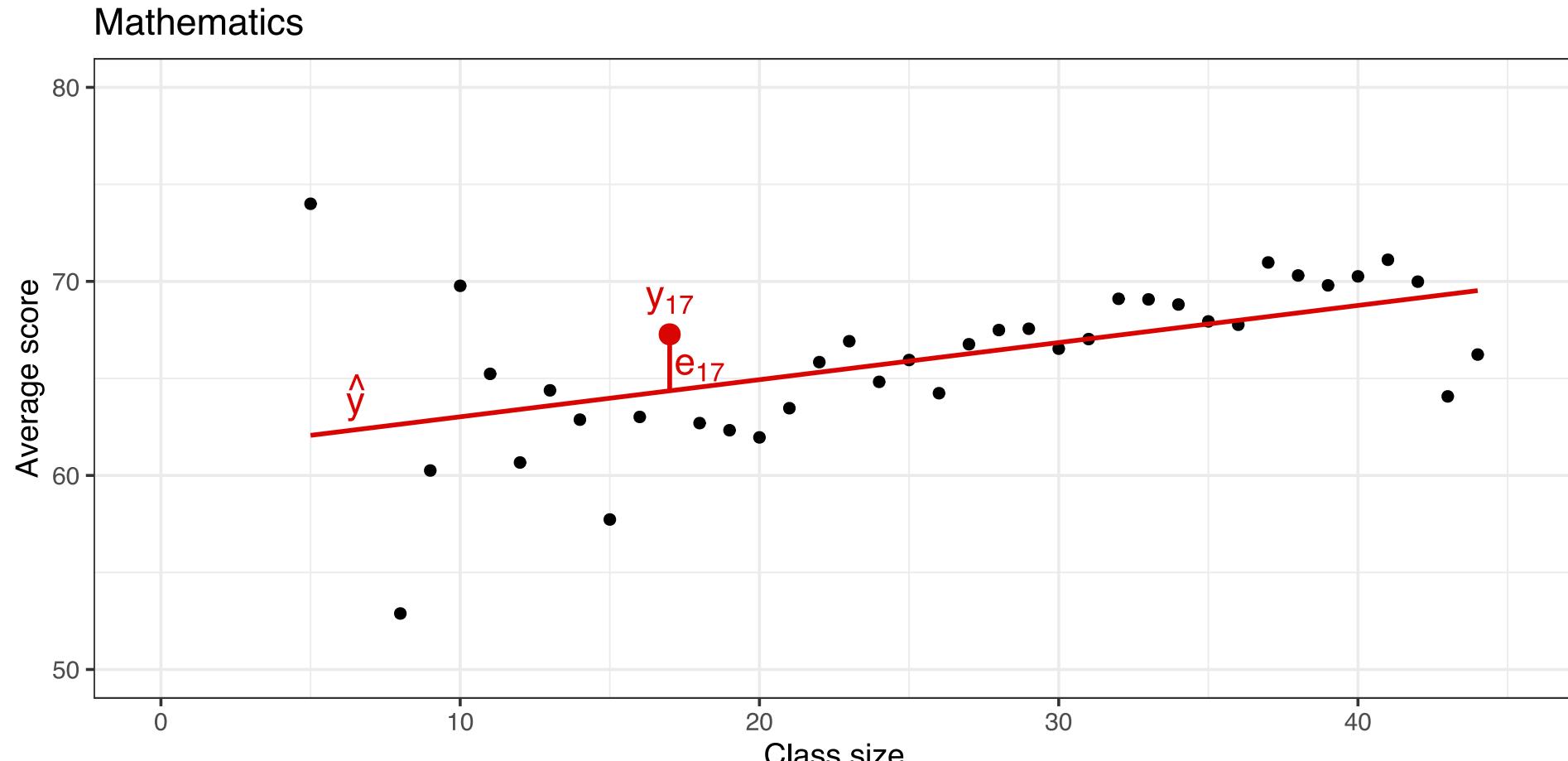
Simple Linear Regression: Graphically



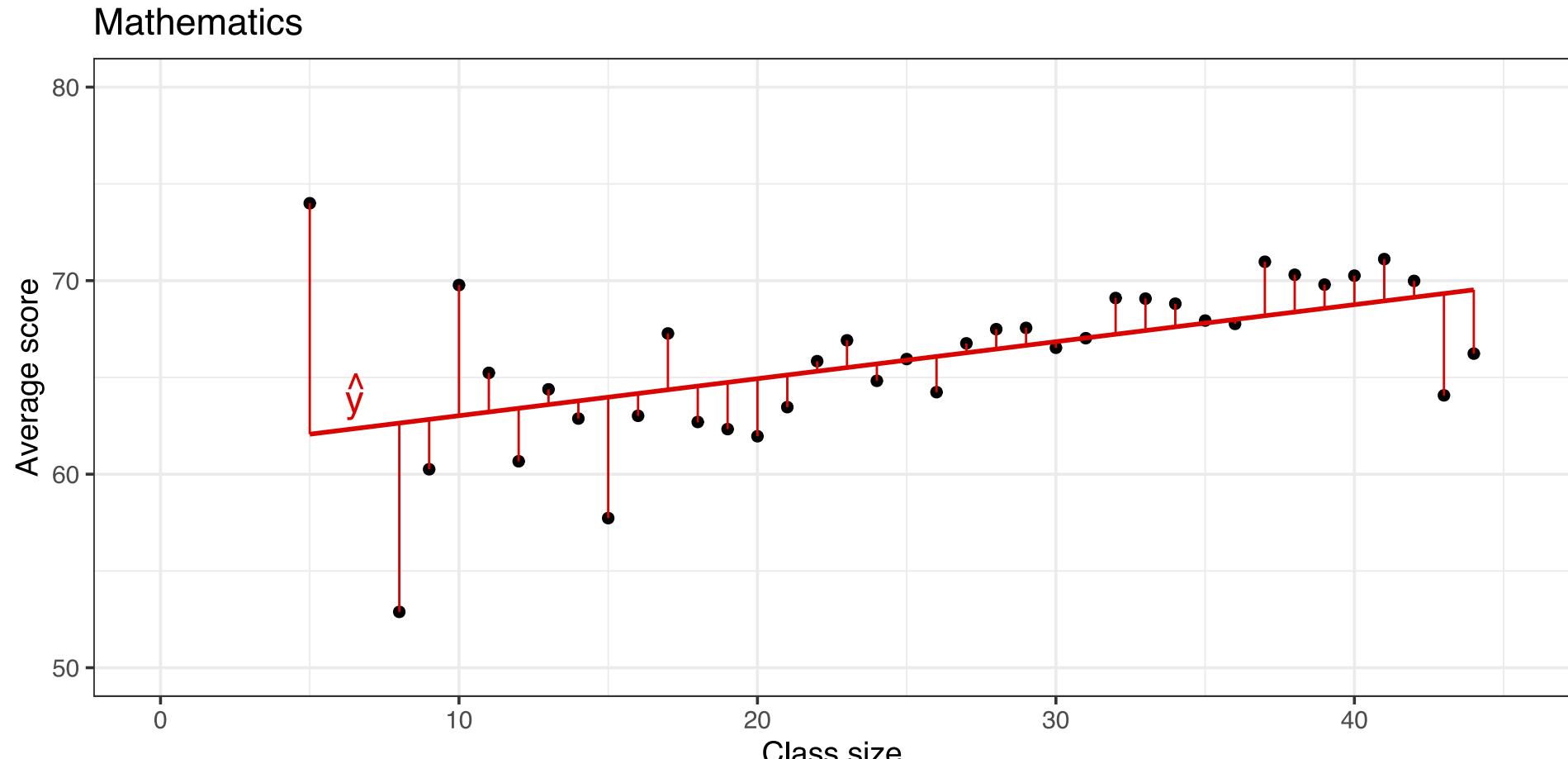
Simple Linear Regression: Graphically



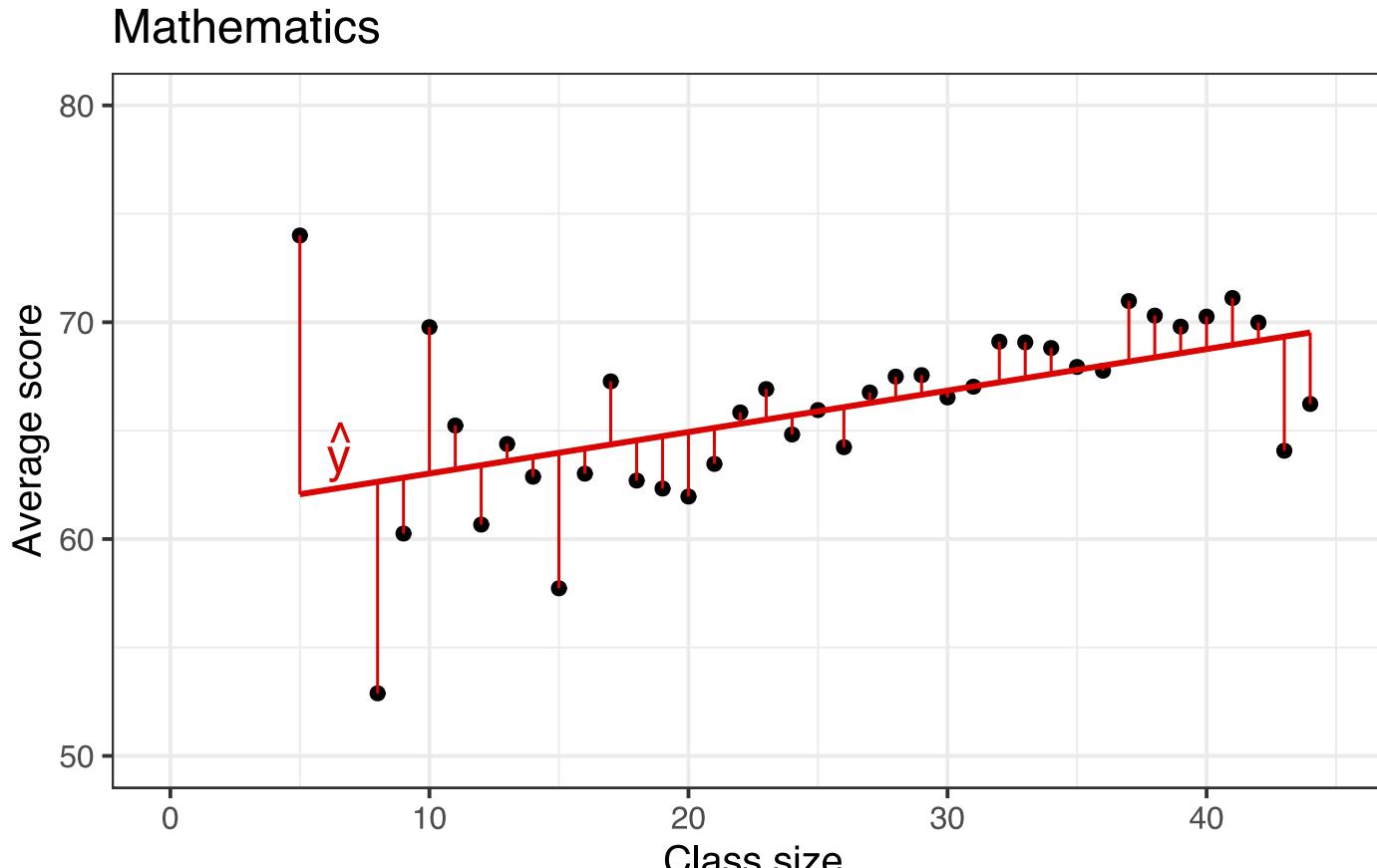
Simple Linear Regression: Graphically



Simple Linear Regression: Graphically



Simple Linear Regression: Graphically



Which
"minimisation"
criterion
should (can)
be used?



Ordinary Least Squares (OLS) Estimation



Ordinary Least Squares (OLS) Estimation

- Errors of different sign (+/-) cancel out, so we consider **squared residuals**

$$\forall i \in [1, N], e_i^2 = (y_i - \hat{y}_i)^2 = (y_i - b_0 - b_1 x_i)^2$$

- Choose (b_0, b_1) such that $\sum_{i=1}^N e_i^2 + \dots + e_N^2$ is **as small as possible**.

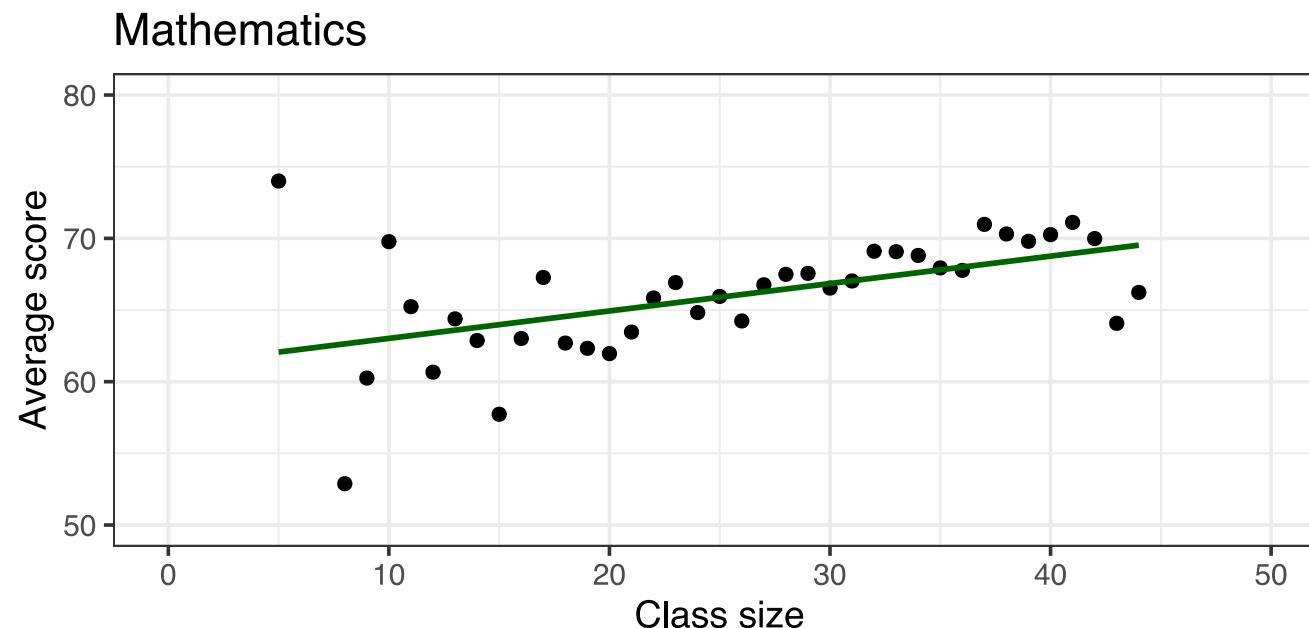


Ordinary Least Squares (OLS) Estimation

- Errors of different sign (+/-) cancel out, so we consider **squared residuals**

$$\forall i \in [1, N], e_i^2 = (y_i - \hat{y}_i)^2 = (y_i - b_0 - b_1 x_i)^2$$

- Choose (b_0, b_1) such that $\sum_{i=1}^N e_i^2 + \dots + e_N^2$ is **as small as possible**.

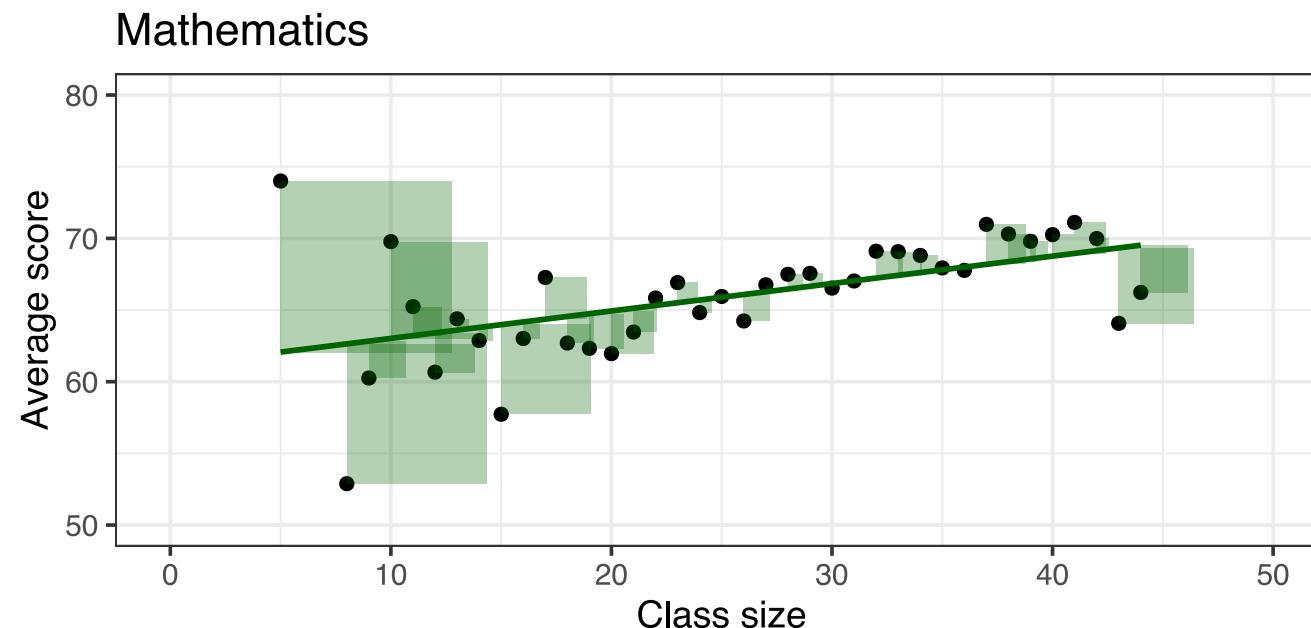


Ordinary Least Squares (OLS) Estimation

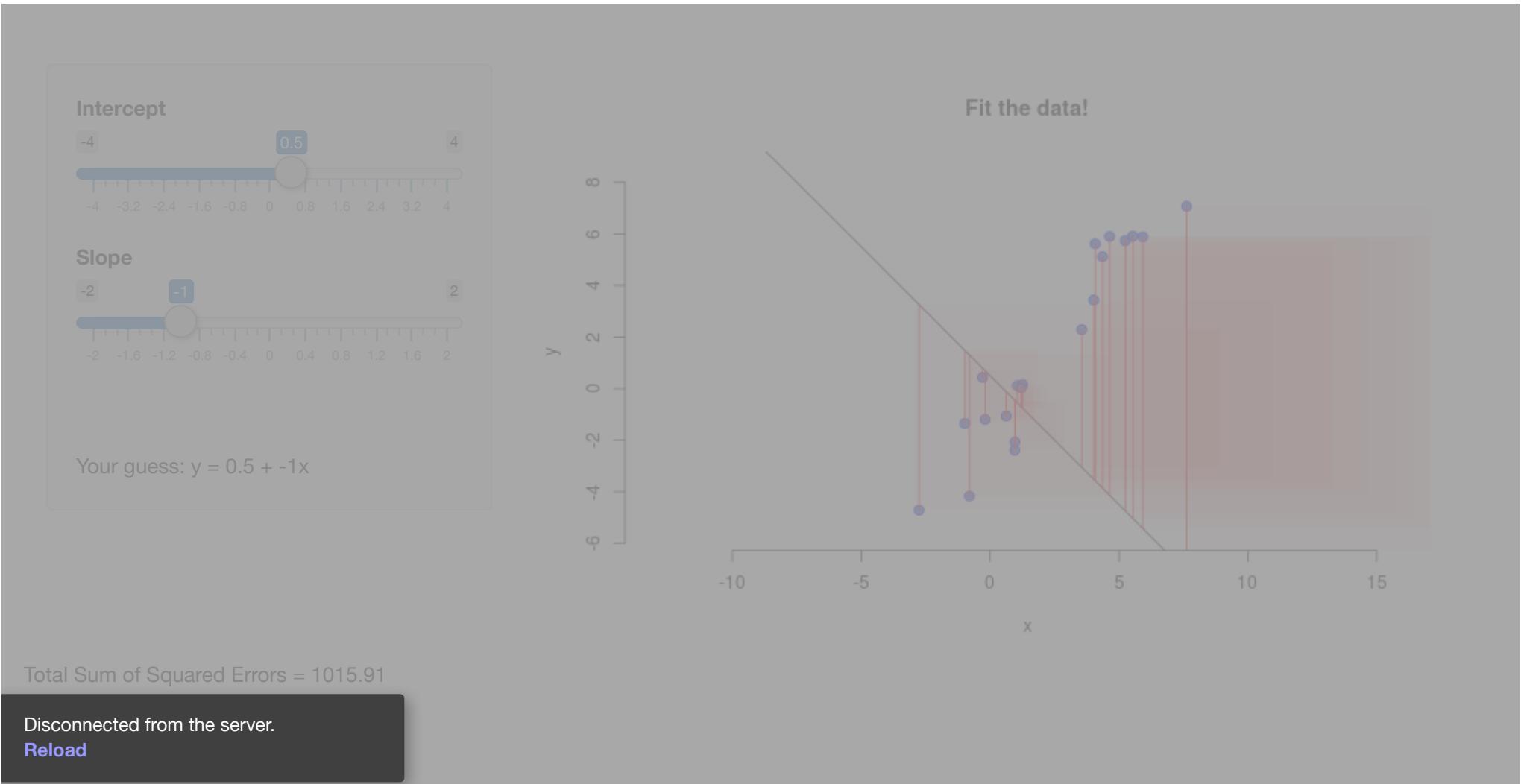
- Errors of different sign (+/-) cancel out, so we consider **squared residuals**

$$\forall i \in [1, N], e_i^2 = (y_i - \hat{y}_i)^2 = (y_i - b_0 - b_1 x_i)^2$$

- Choose (b_0, b_1) such that $\sum_{i=1}^N e_i^2 + \dots + e_N^2$ is **as small as possible**.



Ordinary Least Squares (OLS) Estimation



Ordinary Least Squares (OLS) Estimation



Ordinary Least Squares (OLS): Coefficient Formulas

- **OLS**: *estimation* method consisting in minimizing the sum of squared residuals.
- Yields **unique** solutions to this minimization problem.
- So what are the formulas for b_0 (intercept) and b_1 (slope)?



Ordinary Least Squares (OLS): Coefficient Formulas

- **OLS**: estimation method consisting in minimizing the sum of squared residuals.
- Yields **unique** solutions to this minimization problem.
- So what are the formulas for b_0 (intercept) and b_1 (slope)?
- In our single independent variable case:

$$\text{Slope: } b_1^{OLS} = \frac{\text{cov}(x,y)}{\text{var}(x)}$$

$$\text{Intercept: } b_0^{OLS} = \bar{y} - b_1 \bar{x}$$



Ordinary Least Squares (OLS): Coefficient Formulas

- **OLS**: estimation method consisting in minimizing the sum of squared residuals.
- Yields **unique** solutions to this minimization problem.
- So what are the formulas for b_0 (intercept) and b_1 (slope)?
- In our single independent variable case:

Slope: $b_1^{OLS} = \frac{cov(x,y)}{var(x)}$

Intercept: $b_0^{OLS} = \bar{y} - b_1 \bar{x}$

- ! You should know these formulas, especially the one for b_1^{OLS} !



Ordinary Least Squares (OLS): Coefficient Formulas

- **OLS**: estimation method consisting in minimizing the sum of squared residuals.
- Yields **unique** solutions to this minimization problem.
- So what are the formulas for b_0 (intercept) and b_1 (slope)?
- In our single independent variable case:

Slope: $b_1^{OLS} = \frac{cov(x,y)}{var(x)}$

Intercept: $b_0^{OLS} = \bar{y} - b_1 \bar{x}$

- ! You should know these formulas, especially the one for b_1^{OLS} !
- These formulas do not appear from magic. They can be found by solving the minimisation of squared errors. The maths can be found [here](#) for those who are interested.



Ordinary Least Squares (OLS): Interpretation

For now assume both the dependent variable (y) and the independent variable (x) are numeric.



Ordinary Least Squares (OLS): Interpretation

For now assume both the dependent variable (y) and the independent variable (x) are numeric.

| Intercept (b_0): **The predicted value of y (\hat{y}) if $x = 0$.**



Ordinary Least Squares (OLS): Interpretation

For now assume both the dependent variable (y) and the independent variable (x) are numeric.

Intercept (b_0): **The predicted value of y (\hat{y}) if $x = 0$.**

Slope (b_1): **The predicted change, on average, in the value of y associated to a one-unit increase in x .**



Ordinary Least Squares (OLS): Interpretation

For now assume both the dependent variable (y) and the independent variable (x) are numeric.

Intercept (b_0): **The predicted value of y (\hat{y}) if $x = 0$.**

Slope (b_1): **The predicted change, on average, in the value of y associated to a one-unit increase in x .**

- ⚠ Note that we use the term *associated*, **clearly avoiding interpreting b_1 as the causal impact of x on y** . To make such a claim, we need some specific conditions to be met.
(Next week!)



Ordinary Least Squares (OLS): Interpretation

For now assume both the dependent variable (y) and the independent variable (x) are numeric.

Intercept (b_0): **The predicted value of y (\hat{y}) if $x = 0$.**

Slope (b_1): **The predicted change, on average, in the value of y associated to a one-unit increase in x .**

- ⚠ Note that we use the term *associated*, **clearly avoiding interpreting b_1 as the causal impact of x on y** . To make such a claim, we need some specific conditions to be met.
(Next week!)
- Also notice that the units of x will matter for the interpretation (and magnitude!) of b_1 .



Ordinary Least Squares (OLS): Interpretation

For now assume both the dependent variable (y) and the independent variable (x) are numeric.

Intercept (b_0): **The predicted value of y (\hat{y}) if $x = 0$.**

Slope (b_1): **The predicted change, on average, in the value of y associated to a one-unit increase in x .**

- ⚠ Note that we use the term *associated*, **clearly avoiding interpreting b_1 as the causal impact of x on y** . To make such a claim, we need some specific conditions to be met.
(Next week!)
- Also notice that the units of x will matter for the interpretation (and magnitude!) of b_1 .
- You need to be explicit about what the unit of x is!**



OLS with R

- In R, OLS regressions are estimated using the lm function.
- This is how it works:

```
lm(formula = dependent variable ~ independent variable, data = data.frame containing the data)
```



OLS with R

- In R, OLS regressions are estimated using the lm function.
- This is how it works:

```
lm(formula = dependent variable ~ independent variable, data = data.frame containing the data)
```

Class size and student performance

Let's estimate the following model by OLS: average math score $_i = b_0 + b_1 \text{class size}_i + e_i$

```
# OLS regression of class size on average maths score
lm(avgmath_cs ~ classize, grades_avg_cs)
## 
## Call:
## lm(formula = avgmath_cs ~ classize, data = grades_avg_cs)
## 
## Coefficients:
## (Intercept)    classize
##       61.1092      0.1913
```



Ordinary Least Squares (OLS): Prediction

```
##  
## Call:  
## lm(formula = avgmath_cs ~ classize, data = grades_avg_cs)  
##  
## Coefficients:  
## (Intercept)    classize  
##       61.1092      0.1913
```



Ordinary Least Squares (OLS): Prediction

```
##  
## Call:  
## lm(formula = avgmath_cs ~ classize, data = grades_avg_cs)  
##  
## Coefficients:  
## (Intercept)    classize  
##       61.1092      0.1913
```

This implies (abstracting the i subscript for simplicity):

$$\hat{y} = b_0 + b_1 x$$

$$\text{average } \widehat{\text{math score}} = b_0 + b_1 \cdot \text{class size}$$

$$\text{average } \widehat{\text{math score}} = 61.11 + 0.19 \cdot \text{class size}$$



Ordinary Least Squares (OLS): Prediction

```
##  
## Call:  
## lm(formula = avgmath_cs ~ classize, data = grades_avg_cs)  
##  
## Coefficients:  
## (Intercept)    classize  
##       61.1092      0.1913
```

This implies (abstracting the i subscript for simplicity):

$$\hat{y} = b_0 + b_1 x$$

$$\text{average } \widehat{\text{math score}} = b_0 + b_1 \cdot \text{class size}$$

$$\text{average } \widehat{\text{math score}} = 61.11 + 0.19 \cdot \text{class size}$$

What's the predicted average score for a class of 26 students? (Using the *exact* coefficients.)

$$\text{average } \widehat{\text{math score}} = 61.11 + 0.19 \cdot 26$$

$$\text{average } \widehat{\text{math score}} = 66.08$$



Task 2: OLS Regression

10 : 00

Run the following code to aggregate the data at the class size level:

```
grades_avg_cs <- grades %>%
  group_by(classize) %>%
  summarise(avgmath_cs = mean(avgmath),
            avgverb_cs = mean(avgverb))
```

1. Compute the OLS coefficients b_0 and b_1 of the previous regression using the formulas on slide 25. (*Hint*: you need to use the cov, var, and mean functions.)
2. Regress average verbal score (dependent variable) on class size (independant variable). Interpret the coefficients.
3. Is the slope coefficient similar to the one found for average math score? Was it expected based on the graphical evidence?
4. What is the predicted average verbal score when class size is equal to 0? (Does that even make sense?!)
5. What is the predicted average verbal score when the class size is equal to 30 students?



OLS Variations / Restrictions

- All are described **in the book**. Optional 😊
- There is an app for each of them:

type	App
No Intercept, No regressors	<code>launchApp('reg_constrained')</code>
Centered Regression	<code>launchApp('demeaned_reg')</code>
Standardized Regression	<code>launchApp('reg_standardized')</code>



Predictions and Residuals: Properties

- The average of \hat{y}_i is equal to \bar{y} .

$$\begin{aligned}\frac{1}{N} \sum_{i=1}^N \hat{y}_i &= \frac{1}{N} \sum_{i=1}^N b_0 + b_1 x_i \\ &= b_0 + b_1 \bar{x} = \bar{y}\end{aligned}$$

- The average (or sum) of residuals is 0.

$$\begin{aligned}\frac{1}{N} \sum_{i=1}^N e_i &= \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i) \\ &= \bar{y} - \frac{1}{N} \sum_{i=1}^N \hat{y}_i \\ &= 0\end{aligned}$$



Predictions and Residuals: Properties

- The average of \hat{y}_i is equal to \bar{y} .

$$\begin{aligned}\frac{1}{N} \sum_{i=1}^N \hat{y}_i &= \frac{1}{N} \sum_{i=1}^N b_0 + b_1 x_i \\ &= b_0 + b_1 \bar{x} = \bar{y}\end{aligned}$$

- The average (or sum) of residuals is 0.

$$\begin{aligned}\frac{1}{N} \sum_{i=1}^N e_i &= \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i) \\ &= \bar{y} - \frac{1}{N} \sum_{i=1}^N \hat{y}_i \\ &= 0\end{aligned}$$

- Regressor and residuals are uncorrelated (by definition).

$$Cov(x_i, e_i) = 0$$

- Prediction and residuals are uncorrelated.

$$\begin{aligned}Cov(\hat{y}_i, e_i) &= Cov(b_0 + b_1 x_i, e_i) \\ &= b_1 Cov(x_i, e_i) \\ &= 0\end{aligned}$$

Since $Cov(a + bx, y) = bCov(x, y)$.



Linearity Assumption: Visualize your Data!

- It's important to keep in mind that covariance, correlation and simple OLS regression only measure **linear relationships** between two variables.
- Two datasets with *identical* correlations and regression lines could look *vastly* different.



Linearity Assumption: Visualize your Data!

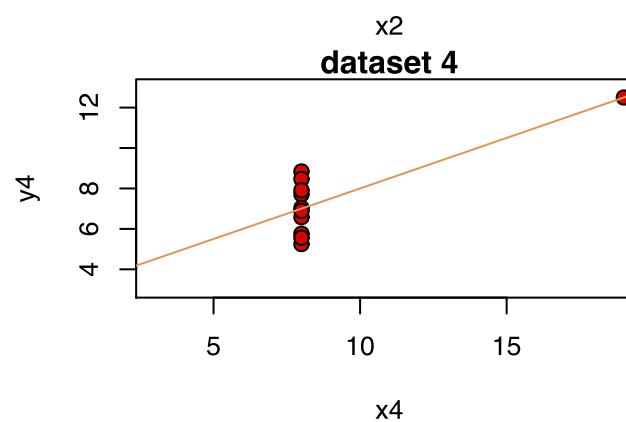
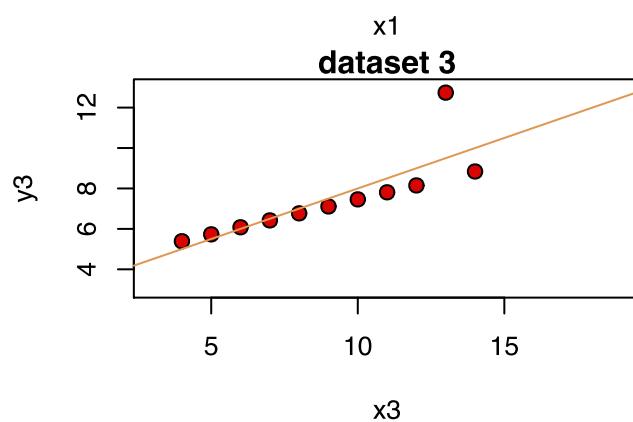
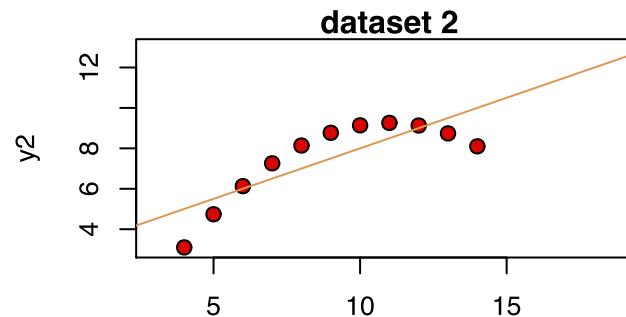
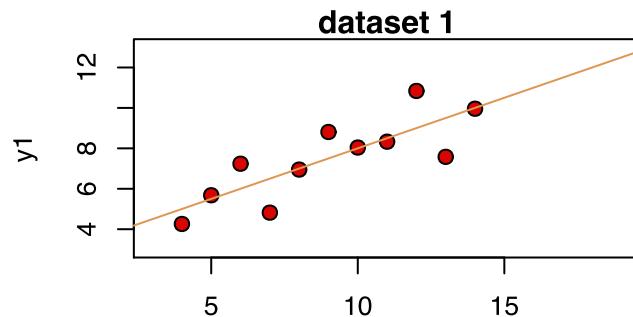
- It's important to keep in mind that covariance, correlation and simple OLS regression only measure **linear relationships** between two variables.
- Two datasets with *identical* correlations and regression lines could look *vastly* different.

- Is that even possible?



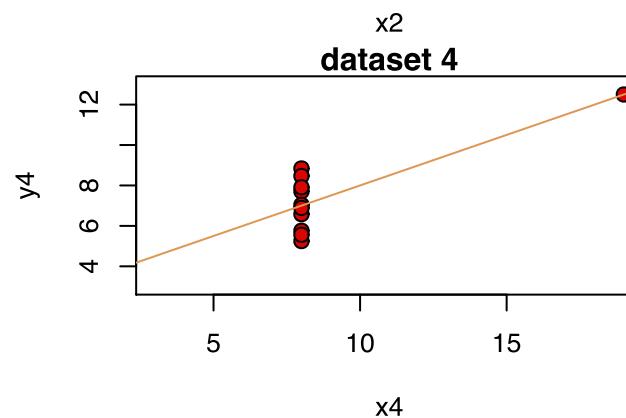
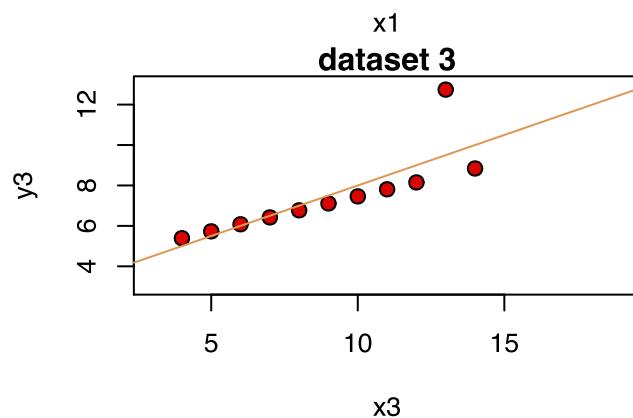
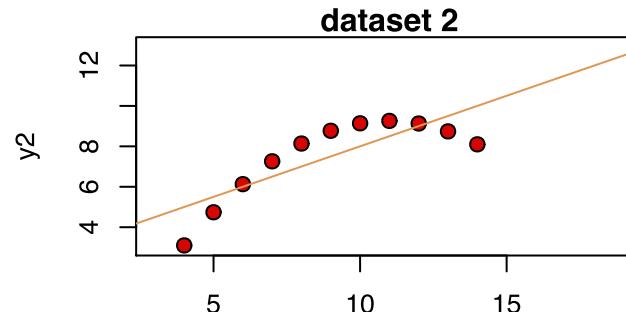
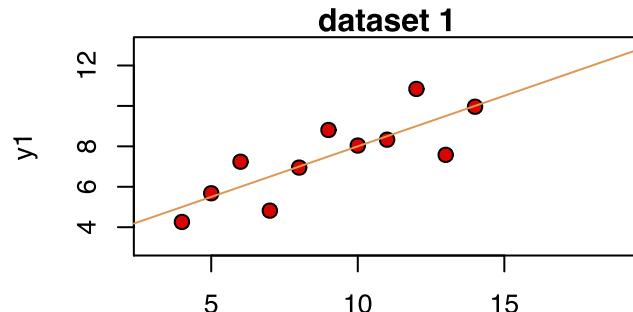
Linearity Assumption: Anscombe

- Francis Anscombe (1973) came up with 4 datasets with identical stats. But look!



Linearity Assumption: Anscombe

- Francis Anscombe (1973) came up with 4 datasets with identical stats. But look!



dataset	cov	var(y)	var(x)
1	5.501	4.127	11
2	5.500	4.128	11
3	5.497	4.123	11
4	5.499	4.123	11



Nonlinear Relationships in Data?

- We can accomodate non-linear relationships in regressions.
- Just add a *higher order* term like this:

$$y_i = b_0 + b_1 x_i + b_2 x_i^2 + e_i$$

- This is **multiple regression** (in 2 weeks!)



Nonlinear Relationships in Data?

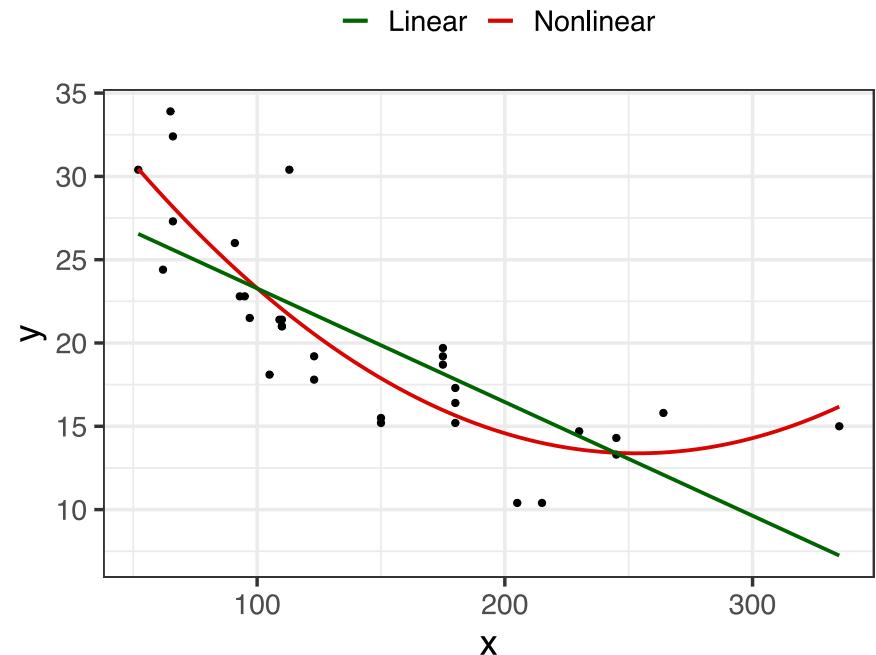
- We can accommodate non-linear relationships in regressions.
- Just add a *higher order* term like this:

$$y_i = b_0 + b_1 x_i + b_2 x_i^2 + e_i$$

- This is **multiple regression** (in 2 weeks!)

- For example, suppose we had this data and fit the previous regression model:

Nonlinear relationship between x and y



Analysis of Variance

- Remember that $y_i = \hat{y}_i + e_i$.
- We have the following decomposition:

$$\begin{aligned}Var(y) &= Var(\hat{y} + e) \\&= Var(\hat{y}) + Var(e) + 2Cov(\hat{y}, e) \\&= Var(\hat{y}) + Var(e)\end{aligned}$$

- Because:
 - $Var(x + y) = Var(x) + Var(y) + 2Cov(x, y)$
 - $Cov(\hat{y}, e) = 0$
- **Total variation (SST) = Model explained (SSE) + Unexplained (SSR)**



Goodness of Fit

- The R^2 measures how well the **model fits the data**.



Goodness of Fit

- The R^2 measures how well the **model fits the data**.

$$R^2 = \frac{\text{variance explained}}{\text{total variance}} = \frac{SSE}{SST} = 1 - \frac{SSR}{SST} \in [0, 1]$$



Goodness of Fit

- The R^2 measures how well the **model fits the data**.

$$R^2 = \frac{\text{variance explained}}{\text{total variance}} = \frac{SSE}{SST} = 1 - \frac{SSR}{SST} \in [0, 1]$$

- R^2 close to 1 indicates a **very high explanatory power** of the model.
- R^2 close to 0 indicates a **very low explanatory power** of the model.



Goodness of Fit

- The R^2 measures how well the **model fits the data**.

$$R^2 = \frac{\text{variance explained}}{\text{total variance}} = \frac{SSE}{SST} = 1 - \frac{SSR}{SST} \in [0, 1]$$

- R^2 close to 1 indicates a **very high explanatory power** of the model.
- R^2 close to 0 indicates a **very low explanatory power** of the model.
- Interpretation:* an R^2 of 0.5, for example, means that the variation in x "explains" 50% of the variation in y .



Goodness of Fit

- The R^2 measures how well the **model fits the data**.

$$R^2 = \frac{\text{variance explained}}{\text{total variance}} = \frac{SSE}{SST} = 1 - \frac{SSR}{SST} \in [0, 1]$$

- R^2 close to 1 indicates a **very high explanatory power** of the model.
- R^2 close to 0 indicates a **very low explanatory power** of the model.
- *Interpretation:* an R^2 of 0.5, for example, means that the variation in x "explains" 50% of the variation in y .
- ⚠ Low R^2 does **NOT** mean it's a useless model! Remember that econometrics is interested in causal mechanisms, not prediction!



Task 3: R^2 and goodness of fit

10 : 00

1. Regress `avgmath_cs` on `classize`. Assign to an object `math_reg`.
2. Pass `math_reg` in the `summary()` function. What is the (multiple) R^2 for this regression? How can you interpret it?
3. Compute the squared correlation between `classize` and `avgmath_cs`. What does this tell you about the relationship between R^2 and the correlation in a regression with only one regressor?
4. Install and load the `broom` package. Pass `math_reg` in the `augment()` function and assign it to a new object. Use the variance in `avgmath_cs` (SST) and the variance in `.fitted` (predicted values; SSE) to find the R^2 using the formula on the previous slide.
5. Repeat steps 1 and 2 for `avgverb_cs`. For which exam does the variance in class size explain more of the variance in students' scores?



SEE YOU IN TWO WEEKS!

✉ florian.oswald@sciencespo.fr

🔗 Slides

🔗 Book

🐦 @ScPoEcon

�� @ScPoEcon

