

Econometrics Discussion Section 2

John Green

Spring 2024

Linear regression

We now want to move on to examining how to understand the relationship between two variables.

- Distinguish *prediction* from *causation*
- Linear regression model
- R^2 , SER, F-test
- Necessary assumptions

Minimizing error with one variable

- Suppose we have data $X_1, X_2, X_3, \dots, X_n = \{X_i\}_{i \leq n}$
- What is the best guess for the value of an arbitrary X_j ?
 - The best fit will be the mean, $E[X]$
- But wait: which error are we minimizing?

Minimizing error with two variables

- Which error are we minimizing?
 - Usually the *squared* error — but not necessarily!
- Now what if we have data on two variables:

$$(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n) = \{(X_i, Y_i)\}_{i \leq n}$$

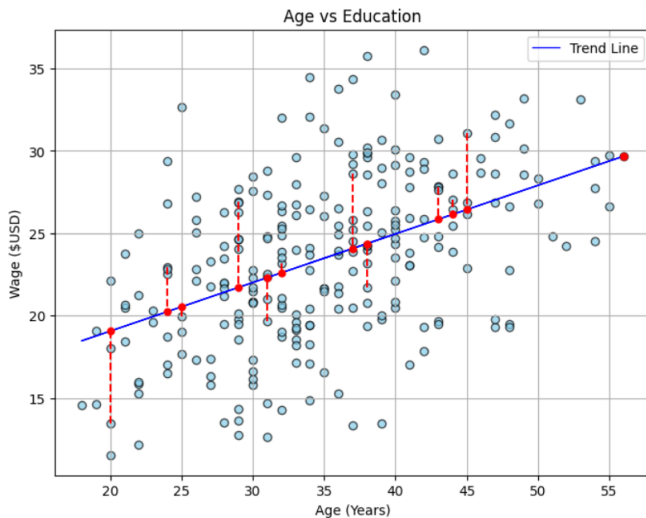
- We want to understand their relationship! Suppose we think that X causes Y : then we might want to know what is our best guess for an arbitrary Y_j , given its X_j ?
 - X is years of education and Y is income; if I tell you that someone has 12 years of education, what is your best guess for their income?
 - This will be the *conditional expectation*: $E[Y|X = 12]$

Minimizing error with two variables

- We want to understand their relationship! Suppose we think that X causes Y : then we might want to know what is our best guess for an arbitrary Y_j , given its X_j ?
 - X is years of education and Y is income; if I tell you that someone has 12 years of education, what is your best guess for their income?
 - This will be the *conditional expectation*: $E[Y|X = 12]$
- Let's assume a *linear* relationship:

$$Y_i = \beta_0 + \beta_1 X_i$$

- Then our problem is simply to draw a line through the 2D data which minimizes the errors
 - β_0 is the intercept, β_1 is the slope



Ordinary least squares

- The line which minimizes the errors is the *ordinary least squares* regression line
- Found by solving

$$\min_{\beta_0, \beta_1} \sum_{i=1}^n (Y_i - (\beta_0 + \beta_1 X_i))^2$$

- In practice, we only have a **sample**, and so we have to estimate these two parameters:
 - Add an error term: $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$
 - We will denote these estimates as $\hat{\beta}_0$ (intercept) and $\hat{\beta}_1$ (slope)

Regression output

- R^2 measures the variance in Y that is explained by X
 - $R^2 = \frac{ESS}{TSS}$
- Standard error of regression (SER): spread of the residual ϵ
- Root mean squared error (RMSE): similar, just calculated with n instead of $n - 2$ in the denominator

Assumptions

- All we've done so far is draw a line — no assumptions needed!
- What if we want to think about X as having a causal effect on Y , as in the case of more years of school mechanically leading to a higher income?
 - eg if we were to hold all other factors constant (as in a controlled and randomized experiment) what do we expect an extra year of education to do to income?

Assumptions

- $E[\epsilon|X = x] = 0$ so that $\hat{\beta}_1$ is unbiased estimator
- $\{(X_i, Y_i)\}_{i \leq n}$ are i.i.d. (will give us sampling distributions for coefficient estimates)
- Large outliers are rare