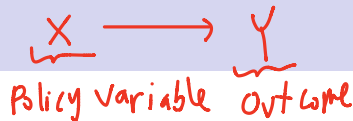


# Introduction to Simple Linear Regression

Hammad Shaikh

September 30, 2019

## Regression Overview



- Empirical analysis in economics is to provide precise quantitative answers to questions of economic interest
  - What is the effect of reducing class size on test scores?
- Economic model relates economic variables of interest to one another using a equation
  - Achievement =  $f(\text{effort, class size, parental investment})$
- Econometric model completes an economic model by specifying any additional uncertainty
  - Achievement =  $f(\text{effort, class size, parental investment, } \epsilon)$

stochastic error

## Linear regression model

$$X \rightarrow Y$$

- $Y$  = dependant / outcome / response variable
  - What are plausible  $Y$ 's in class size reduction policy?

↳  $Y \in \{\text{test score, Parent satisfaction}\}$

- $X$  = independent / explanatory / predictor variable
  - Contains treatment of interest and other factors that effect  $Y$
  - What are the  $X$ 's in class size reduction policy?

↳  $X \in \{\text{class size, student-teacher ratio}\}$

- Simple regression:  $Y = \beta_0 + \beta_1 X + \epsilon$

↳ i) data  $(Y, X)$ , ii) parameters  $\beta_0, \beta_1$ , iii) error

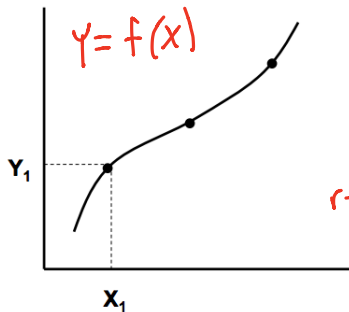
- Multiple regression:  $Y = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k + \epsilon$

↳  $\underbrace{\text{class size}}$   $\underbrace{\text{Hour study}}$

## Functional vs. Statistical Relationship

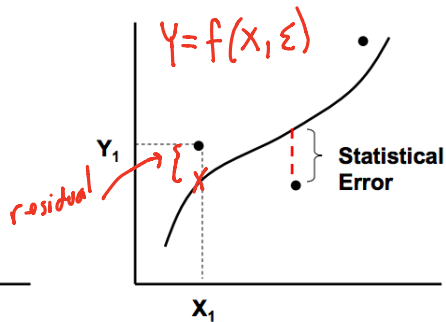
- Regression model describes the statistical relationship between outcome  $Y$  and response variable(s)  $X$

### Functional Relationship



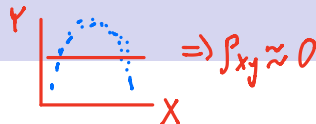
Economic model

### Statistical Relationship



Econometric

## Relationship Between X and Y



- The covariance is a measure of the linear association between X (class size) and Y (test score)

- $S_{xy} = \widehat{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$
- Units are Units of X  $\times$  Units of Y (No. of students  $\times$  Score)

$\hookrightarrow$  Hard to interpret magnitude

- $Cov(X, Y) > 0$  means a positive relation between X and Y

- Correlation is a unit less measure of the strength of linear relationship between X and Y

- $\rho_{xy} = \frac{S_{xy}}{S_x S_y}$  is a number between -1 and 1
- $\rho_{xy} = 1$  means perfect positive linear relationship

$$\text{units } \text{Corr}(X, Y) = \frac{\text{Units } X \cdot \text{Units } Y}{\text{Units } X \cdot \text{Units } Y}$$

## Simple Regression Example $X \rightarrow Y$

- Question: What is the relationship between  $\overset{\text{X}}{\text{class size}}$  and  $\underset{\text{Y}}{\text{test scores}}$  in California?

$\rightarrow$  unit of obs. is school districts ( $d$ )

- Data available from 420 California school districts
  - 5th grade district average math and reading score
  - Student to Teacher Ratio (STR): number of students divided by number of teachers (within district)

- What is the regression model of interest?

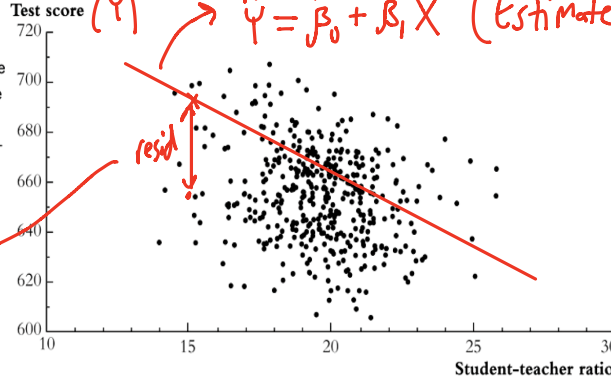
$$\text{Math}_{5d} = \beta_0 + \beta_1 \text{STR}_d + \varepsilon_d \quad \left( \begin{array}{l} \text{Population} \\ \text{reg. model} \end{array} \right)$$

$d \in \{1, 2, \dots, 420\}$

## Test Score and Student to Teacher Ratio

**FIGURE 4.2** Scatterplot of Test Score vs. Student-Teacher Ratio (California School District Data)

Data from 420 California school districts. There is a weak negative relationship between the student-teacher ratio and test scores: the sample correlation is  $-0.23$ .



420 data points

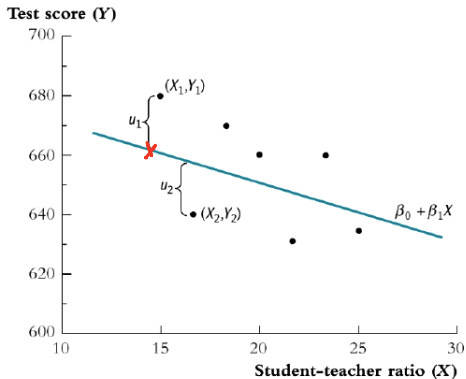
$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X \quad (\text{Estimated})$$

- We want to model above relationship with a simple linear regression

## Estimating Simple Regression

**FIGURE 4.1** Scatter Plot of Test Score vs. Student-Teacher Ratio (Hypothetical Data)

The scatterplot shows hypothetical observations for seven school districts. The population regression line is  $\beta_0 + \beta_1 X$ . The vertical distance from the  $i^{\text{th}}$  point to the population regression line is  $Y_i - (\beta_0 + \beta_1 X_i)$ , which is the population error term  $u_i$  for the  $i^{\text{th}}$  observation.



$$\min_{\beta_0, \beta_1} \sum_{i=1}^n u_i^2(\beta_0, \beta_1)$$

- Simple regression estimates:  $\hat{\beta}_1 = \frac{\widehat{Cov}(X, Y)}{\widehat{Var}(X)}$ ,  $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$ 
  - Known as Ordinary Least Squares (OLS) estimator



## Effect of STR on Achievement

- $\overset{Y}{\text{TestScore}_d} = \beta_0 + \beta_1 \overset{X}{\text{STR}_d} + \epsilon_d$  (Population)
  - We want to estimate  $\beta_1 = \frac{\Delta \text{TestScore}}{\Delta \text{STR}}$ . Interpret  $\beta_1$ ?  
↳ when STR ↑ by 1 the impact on test scores is  $\beta_1$  points on avg.
- Line of best fit:  $\widehat{\text{TestScore}}_d = \hat{b}_0 + \hat{b}_1 \text{STR}_d$  (Estimated)
  - $(\hat{b}_0, \hat{b}_1)$  found by minimizing  $\sum_{i=1}^n (\text{TestScore}_d - \widehat{\text{TestScore}}_d)^2$   
↳ OLS  
 $e(\hat{b}_0, \hat{b}_1)$
- $\hat{b}_1 = \frac{\widehat{\text{Cov}}(\text{TestScore}_d, \text{STR}_d)}{\widehat{\text{Var}}(\text{STR}_d)}$  and  $\hat{b}_0 = \overline{\text{TestScore}} - \hat{b}_1 \overline{\text{STR}}$

## Effect of STR on Achievement Cont.

- Estimated model:  $\widehat{TestScore}_d = \underbrace{698.9}_{\hat{b}_0} - 2.28 \underbrace{STR_d}_{\hat{b}_1}$
- Primary estimate of interest is  $\hat{b}_1 = -2.28$ 
  - Districts with one more student per teacher on average are associated with 2.28 points lower test scores
- How to interpret intercept of  $\hat{b}_0 = 698.9$ ?

↳  $STR=0$  is not in data and hence no interpretation.