

Note: Annotations start at slide 12

Introduction to Econometrics: Linear Regression

Hammad Shaikh

January 18, 2018

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 an equation
 - Achievement = $f(\text{effort, class size, parental investment})$
- Econometric model completes an economic model by specifying any additional uncertainty
 - Achievement = $f(\text{effort, class size, parental investment, } \epsilon)$

Linear regression model

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

- 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?

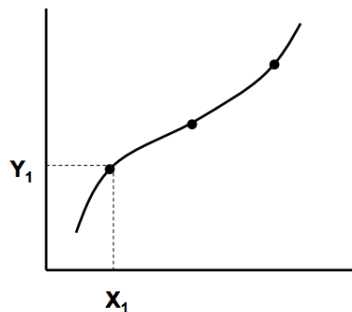
- Simple regression: $Y = \beta_0 + \beta_1 X + \epsilon$

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

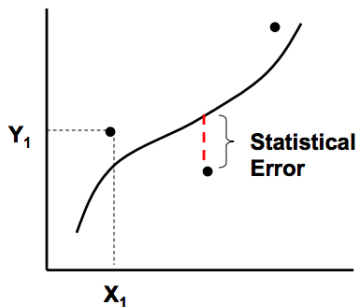
Functional vs. Statistical Relationship

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

Functional Relationship



Statistical Relationship



Relationship Between X and Y

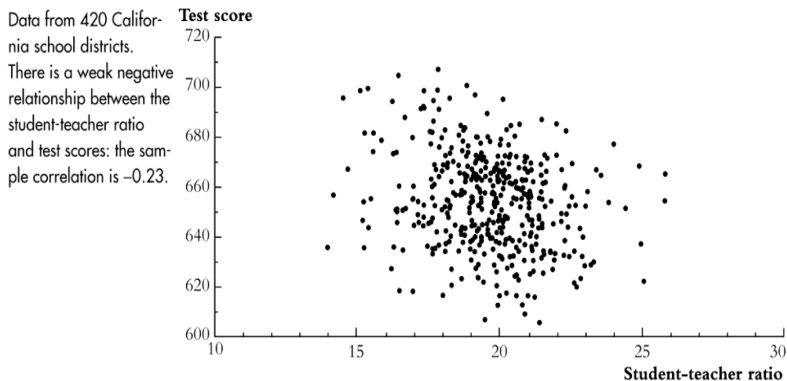
- The covariance is a measure of the linear association between X (class size) and Y (test score)
 - $S_{xy} = \widehat{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$
 - Units are Units of X \times Units of Y (No. of students \times Score)
- $Cov(X, Y) > 0$ means a positive relation 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

Simple Regression Example

- Question: What is the relationship between class size and test scores in California?
- 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?

Test Score and Student to Teacher Ratio

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

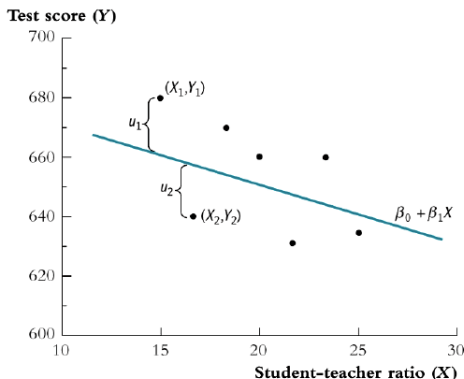


- We want to model above relationship with a simple linear regression

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^{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^{th} observation.



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

Effect of STR on Achievement

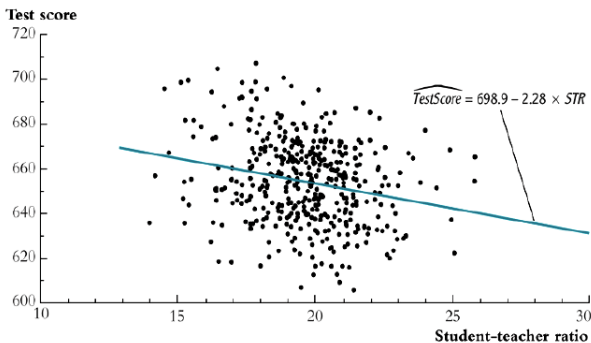
- $TestScore_d = \beta_0 + \beta_1 STR_d + \epsilon_d$
 - We want to estimate $\beta_1 = \frac{\Delta TestScore}{\Delta STR}$. Interpret β_1 ?
- Line of best fit: $\widehat{TestScore}_d = \widehat{b}_0 + \widehat{b}_1 STR_d$
 - $(\widehat{b}_0, \widehat{b}_1)$ found by minimizing $\sum_{i=1}^n (TestScore_d - \widehat{TestScore}_d)^2$
- $\widehat{b}_1 = \frac{\widehat{Cov}(TestScore_d, STR_d)}{\widehat{Var}(STR_d)}$ and $\widehat{b}_0 = \overline{TestScore} - \widehat{b}_1 \overline{STR}$

Effect of STR on Achievement Cont.

- Districts with larger class sizes (higher STR) are associated with lower test scores

FIGURE 4.3 The Estimated Regression Line for the California Data

The estimated regression line shows a negative relationship between test scores and the student-teacher ratio. If class sizes fall by 1 student, the estimated regression predicts that test scores will increase by 2.28 points.



Effect of STR on Achievement Cont.

- Estimated model: $\widehat{TestScore}_d = 698.9 - 2.28STR_d$
- 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 $\hat{b}_0 = 698.9$?

Properties of Slope Estimator

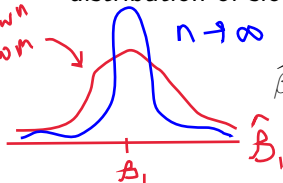
Recall $\begin{cases} \text{Rand. var } \hat{\beta}_1 \\ \text{Estimate } \hat{\beta}_1 \end{cases}$

- We generally want estimators to be unbiased and consistent
 - Slope estimator $\hat{\beta}_1$ unbiased if $E(\hat{\beta}_1) = \beta_1$
 - Slope estimator $\hat{\beta}_1$ consistent if $\hat{\beta}_1 \xrightarrow{P} \beta_1$ as n grows large

$\Sigma \perp X$ is assumption discussed later

- It can be shown (using CLT) that if ϵ independent of X then the distribution of slope estimator $\hat{\beta}_1 = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$ is:

$\hat{\beta}_1$ drawn from



$$\hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma_\epsilon^2}{\sum_{i=1}^n (X_i - \bar{X})^2}\right)$$

Notice:
 $n \cdot \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$
 $n \sigma_x^2$

- Show that $\hat{\beta}_1$ is unbiased and consistent

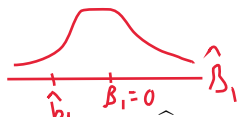
unbiased: $E(\hat{\beta}_1) = \beta_1$
 \hookrightarrow mean of normal

consistent: $\text{Var}(\hat{\beta}_1) \rightarrow 0$ as $n \rightarrow \infty$
 so $\hat{\beta}_1 \rightarrow \beta_1$

Simple Linear Regression and Hypothesis Testing

- Simple linear regression: $TestScore_d = \beta_0 + \beta_1 STR_d + \epsilon_d$
 - β_0 (intercept) and β_1 (slope) are unknown parameters
 - Use sample $(STR_d, TestScore_d)_{d=1}^n$ to make inference about the simple linear regression parameters

Note: It is possible $\beta_1 = 0$ but $\hat{b}_1 < 0$
no effect



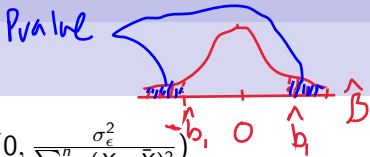
- Question: How much can we trust the primary estimate \hat{b}_1 ?
 - $\hat{b}_1 < 0$ is more trustable if $Var(\hat{\beta}_1)$ is small
 - Want to check $\hat{b}_1 < 0$ not just due to sampling error
- Null Hypothesis: Class size has no effects on achievement

$$H_0: \beta_1 = 0$$

- Alternative Hypothesis: Class size effects achievement

$$H_1: \beta_1 \neq 0 (\beta_1 < 0 \text{ or } \beta_1 > 0)$$

SLR and Hypothesis Testing Cont.



- Under $H_0 : \beta_1 = 0$ we have $\hat{\beta}_1 \sim N(0, \frac{\sigma_\epsilon^2}{\sum_{i=1}^n (X_i - \bar{X})^2})$
 - Since ϵ unknown, σ_ϵ^2 is also unknown. The solution is to replace it with s_e^2 , the sample variance of the residuals

$$s_e^2 = \frac{1}{n} \sum_i e_i^2, \quad e_i \text{ is residual, } e_i = y_i - \hat{y}_i$$

- If $H_1 : \beta_1 \neq 0$ we compute p-value = $2 * Pr(\hat{\beta}_1 > \hat{b}_1)$
 - \hat{b}_1 is very significant if p-value < 0.01 , significant if p-value < 0.05 , and marginally significant if p-value < 0.1

- Computing p-value involves $SE(\hat{b}_1) = \sqrt{Var(\hat{\beta}_1)}$

- Typically (not always) if $|\frac{\hat{b}_1}{SE(\hat{b}_1)}| > 2$ then \hat{b}_1 is significant

$$t_{stat} = \frac{\hat{b}_1 - 0}{SE(\hat{b}_1)} > 2 \implies \underbrace{\hspace{10em}}_{\text{p-value} < 0.05}$$

- P-value ≈ 0 for class size application

\rightarrow class size related to test scores

Fitness of Regression Model

- R^2 measures the proportion of variation in the outcome (Y) explained by the independent variable(s) (X)
 - R^2 is a number between 0 and 1 (R^2 is unitless)
 - $R^2 = 1$ means regression model perfectly fits the data

- $R^2 = \frac{SSR}{SST}$; SST = Sum of Square Total, SSR = Sum of Square Regression

- $SST = \sum_{i=1}^n (y_i - \bar{y})^2$ and $SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$

variation in y

Var. in y explained by reg. \hat{y}

- R^2 applies to both simple and multiple linear regression

Ex. $R^2 = 0.3$ for $\text{Mârk}_i = 40 + 5 \text{Lecture Attend}_i$
→ 30% of variation in grades is explained by lecture attendance

Simple Linear Regression Summary

- The population linear regression model

- $Y = \beta_0 + \beta_1 X + \epsilon$

β_1 is param. of interest

- Line of best fit and OLS estimator

- $\hat{\beta}_1 = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$ and $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$

Slope Estimator

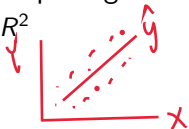
- Hypothesis testing

- $H_0 : \beta_1 = 0$ and $H_1 : \beta_1 \neq 0$

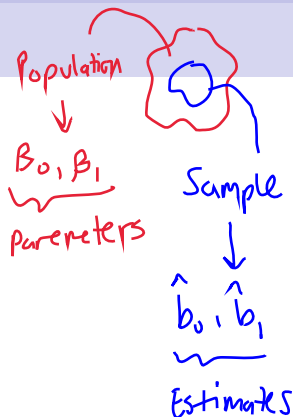
X and Y not related

- Measures of fit for simple regression: $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

- Correlation and R^2



Line of best fit



Extending to Multiple Regression

$$Y = \beta_0 + \beta_1 X + \epsilon$$

(Handwritten red and blue annotations: a red arrow points from $\beta_1 X$ to ϵ , and a blue arrow points from ϵ to $\beta_1 X$)

- Results from simple linear regression are usually not causal
 - Many other factors that affect both X and Y are not accounted for in the model
 - Can bias slope estimates (omitted variable bias)

ϵ related to X is a problem

- Returns to education: $AdultIncome_i = \beta_0 + \beta_1 YrsEduc_i + \epsilon_i$
 - What are some variables in ϵ_i that may bias $\hat{\beta}_1$?

\uparrow Parent Inc $\left\{ \begin{array}{l} \uparrow Yrs Educ \\ \uparrow Income \end{array} \right.$, \uparrow Ability $\left\{ \begin{array}{l} \uparrow Yrs Educ \\ \uparrow Income \end{array} \right.$
Hard to measure

- Two solutions to help obtain causal result:
 - 1) Randomized control trial, or 2) Multiple regression

Best method to get causality

Randomized Control Trial (RCT)

$\uparrow X_i$ not related to ϵ_i

RCT is gold standard

- Simple regression model: $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$

We want $\epsilon \perp X$

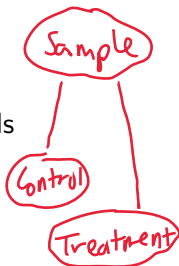
- In a RCT the X s are randomly assigned to individuals
 - No omitted variable bias since X_i independent to ϵ_i
 - Now \hat{b}_1 has a causal interpretation

→ Only diff. b/w control & Treatment is X

- Correlation does not imply causation?
 - Generally true for observational data, but false for experimental data where treatment variable is randomly assigned

- Returns of education: $Y_i = \beta_0 + \beta_1 \text{YrsEduc}_i + \epsilon_i$ to Yrs Educ
- Can we randomly assign years of education to individuals?

→ Not ethical, so no



Problem ϵ_i related to Yrs Educ

Multiple Regression

- Slope estimate in simple regression can be biased from omitted variables related to X and Y

- Solution is to include the omitted variables into the model

→ Compare people with diff. X_1 , but same X_2, X_3, \dots, X_k

- Multiple regression: $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + \epsilon$

- β_1 = effect of changing X_1 on Y holding X_2, \dots, X_k constant
- \hat{b}_1 can be causal if all relevant variables are included
 - Conditional independence: ϵ indep. to X_1 given X_2, \dots, X_k

→ Cond. indep: Given X_2, \dots, X_k , ϵ unrelated to X_1

- Returns to education:

$$Y_i = \beta_0 + \beta_1 \text{YrsEduc}_i + \beta_2 \text{Exp}_i + \beta_3 \text{ParentIncome}_i + \epsilon_i$$

→ Cond. indep: Yrs Educ indep. of ϵ given Exp & Parent Income

↳ Must likely false as ability & motivation are still omitted

Regression Table Example

→ 4 regressions below
↳ Two Outcomes

→ ① $wage_i = 11 + 2 \text{ Yrs Educ}_i$

Table: Income and Health Returns to Education (Fake Data)

	①	②	③	④
	Hourly Wage	Hourly Wage	Years Lived	Years Lived
Constant	11*** (2.5)	10*** (0.1)	65*** (10)	66*** (10)
Years of Educ	2*** (0.5)	1*** (0.1)	2*** (0.25)	3*** (0.3)
Experience		3*** (0.8)		0.5** (0.245)
Parent income (\$1000)		0.1** (0.048)		0.15* (0.075)
R-square	0.15	0.30	0.10	0.20
No. of individuals	15000	15000	15000	15000

Stars denote level of significance *10%, ** 5%, and ***1%.

$N=15000$ indiv. in data

- Regression table generally contain coefficient estimates, standard errors, no. of observations, and R^2

\hat{b}_1 in ④: People with an extra year of educ. are associated with living 3 yrs longer on avg. after controlling for work exp. and their parent income

Summary of Linear Regression

- Goal: examine causal relationship between outcome Y and explanatory variable X
- Simple linear regression is a good starting point
 - Slope estimate is likely biased due to omitted variables that effect both X and Y
- Experiments (RCTs) are ideal for determining causal relationship between X and Y
 - Costly and sometimes unfeasible
- Multiple regression can control for several relevant variables
 - Obtain causal relationship under conditional independence