Tutorial: Multiple Linear Regression

Hammad Shaikh

October 6, 2019

Fitness of Regression Model

R² relavant for prediction but not param informa

*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

Lo R² ~ model fits data well

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

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

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

Simple Linear Regression Summary

- The population linear regression model
 - $Y = \beta_0 + \beta_1 X + \epsilon$ • porch: β_0, β_1
- Line of best fit and OLS estimator • $\hat{\beta}_1 = \frac{Cov(X,Y)}{Var(X)}$ and $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$ ($\hat{Y} = \hat{\beta}_1 + \hat{\beta}_1 X$)
- Hypothesis testing
- $H_0: \beta_1 = 0$ and $H_1: \beta_1 \neq 0$ \implies prode < X X does not impact Y
- Measures of fit for simple regression: $\widehat{y} = \widehat{b}_0 + \widehat{b}_1 x$
 - Correlation and R^2

 $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + ... + \beta_R X_R + \Sigma$ shark Extending to Multiple Regression Results from simple linear regression are usually not causal • Other factors that affect both X and Y are not considered • Can bias slope estimates (omitted variable bias) Y= Shork attacks, X= Ice Gran Sold =) B, >> 0 • Returns to education: AdultIncome_i = $\beta_0 + \beta_1$ YrsEduc_i + ϵ_i • What are some variables in ϵ_i that may bias \hat{b}_1 ? setfort, ability, motivation) $\omega_{v}(X, \varepsilon) \neq 0$ R ias • Two solutions to help obtain causal result: • 1) Randomized control trial, or 2) Multiple regression 4 X LE Ly Hold Observed

Randomized Control Trial $(X \text{ random assigh } \rightarrow) X \perp E)$ • Simple regression model: $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$ • In a RCT the Xs are randomly assigned to individuals • No omitted variable bias since X_i independent to ϵ_i • Now \hat{b}_1 has a causal interpretation

Ly Unbias $E(\hat{\beta}_1) = \beta_1$

Correlation does not imply causation?

- Generally true for observational data, but false for experimental data where treatment variable is randomly assigned
- Returns of education: $Y_i = \beta_0 + \beta_1 YrsEduc_i + \epsilon_i$
 - Can we randomly assign years of education to individuals?

Lo No, experiments not always feasible

hold - X2, X3 Unst X - Y 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 • Multiple regression: $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \ldots + \beta_k X_k + \epsilon$ • β_1 = effect of changing X_1 on Y holding X_2, \ldots, X_k constant • b_1 can be causal if all relevant variables are included • Conditional independence: ϵ indep. to X_1 given X_2, \ldots, X_k 4 ELX, X2,...,XR constant • Returns to education: $Y_{i} = \beta_{0} + \beta_{1} YrsEduc_{i} + \beta_{2} Exp_{i} + \beta_{3} ParentIncome_{i} + \epsilon_{i}$ 5 \$1 is change in Y when PyrsEdue holding constant Exp, Paul Inc.

	Table: Income and I	Health	Returns to E	Education ((Fake Dat	a) 🕖 🔪
{			Hourly Wage	Hourly Wage	Years Lived	Years Lived
	Constant		11***	10***	65***	66***
	Years of Educ	1	(2.5) 2***	(0.1) 1***	(10) 2***	(10) 3***
	Experience	SE	(0.5)	(0.1) 3***	(0.25)	(0.3) 0.5**
	Parent income (\$1000)		t=4	(0.8) 0.1** (0.048)		(0.245) 0.15* (0.075)
	R-square No. of indivisuals		0.15 15000	0.30	0.10 15000	0.20 15000
	Stars denote level of significance $^*10\%,^{**}5\%,$ and $^{***}1\%.$		N=15000			

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



• Suppose we want to make scatter plot of earnings on education but adjust for parental education

 $Eorn = \beta_0 + \beta_1 Educ + \beta_2 PE + \Sigma$

• Step 1: We need to obtain variation in education that is independant of the parental education Ly $E_{dvc} = d_0 + d_1 PE + \overline{\xi} = \overline{\xi}(nsidvals)$ Vor. in Educ. not explained by PE • Step 2: Related earnings on the variation obtained in step 1. 4 Earn = Yo+Y, Z + M Earn $\beta \hat{\gamma} = \hat{k}$

Summary of Linear Regression

- Goal: examine causal relationship between outcome Y and explanatory variable X
- Simple linear regression is a good starting point
 - $\bullet\,$ 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 unfeasable
- Multiple regression can control for several relevant variables
 - Obtain causal relationship under conditional independance

 $L_{1} \geq L X_{1} | X_{2}, \dots, X_{R}$