

What is Regression?

Jeffrey Arnold

April 12, 2016

What is a relationship and why do we care?

- ▶ Most of what we want to do in the social science is learn about how two or more variables are related
- ▶ Examples:
 - ▶ Does turnout vary by types of mailers received?
 - ▶ Is the quality of political institutions related to average incomes?
 - ▶ Does conflict mediation help reduce civil conflict?

Notation and conventions

- ▶ Y - the dependent variable or outcome or regressand or left-hand-side variable or response
 - ▶ Voter turnout
 - ▶ Log GDP per capita
 - ▶ Number of battle deaths
- ▶ X - the independent variable or explanatory variable or regressor or right-hand-side variable or treatment or predictor
 - ▶ Social pressure mailer versus Civic Duty Mailer
 - ▶ Average Expropriation Risk
 - ▶ Presence of conflict mediation
- ▶ Generally our goal is to understand how Y varies as a function of X :

$$Y = f(X) + \text{error}$$

Conditional expectation review

- ▶ How to describe relationship between X and Y ?
- ▶ **Definition** The **conditional expectation function** (CEF) or the **regression function** of Y given X , denoted $E[Y|X = x]$ is the function that gives the mean of Y at various values of x .
- ▶ Note that this is a function of the *population* distributions.
- ▶ Regression at its most fundamental is about how the mean of Y changes as a function of X

Difference in means as difference in conditional expectation

- ▶ We've been writing μ_y and μ_x for the means in different groups.
- ▶ Each of these are conditional expectations. Define Y to be the loan amount, $X = 1$ to indicate a man, and $X = 0$ to indicate a woman and then we have:

$$\mu_m = E[Y|X = 1]$$

$$\mu_w = E[Y|X = 0]$$

- ▶ Notice here that since X can only take on two values, 0 and 1, then these two conditional means completely summarize the CEF.
- ▶ How do we calculate this? We've already done this: it's just the usual sample mean among the men and then the usual sample mean among the women:

$$\hat{E}[Y_i|X_i = 1] = \frac{1}{n_1} \sum_{i: X_i=1} Y_i$$

$$\hat{E}[Y_i|X_i = 0] = \frac{1}{n_0} \sum_{i: X_i=0} Y_i$$

Discrete covariate: sample conditional expectations

- ▶ In the last section we had a binary covariate. What if X is discrete?
- ▶ The same logic applies, we can still estimate $E[Y|X = x]$ with the sample mean among those who have $X_i = x$:

$$\hat{E}[Y_i|X_i = x] = \frac{1}{n_x} \sum_{i: X_i = x} Y_i$$

Continuous covariate (I): each unique value gets a mean

- ▶ What if X is continuous? Can we calculate a mean for every value of X ?
- ▶ Not really, because remember the probability that two values will be the same in a continuous variable is 0.
- ▶ Thus, we'll end up with a very “jumpy” function, $\hat{E}[Y_i|X_i = x]$, since n_x will be at most 1 for any value of x .
- ▶ You can imagine that this will jump around a lot from sample to sample. The estimates, $\hat{E}[Y_i|X_i = x]$, will have high sampling variance.
- ▶ For some values of x we never observe anything

Continuous covariate (II): stratify and take means

- ▶ So, that seems like each value of X won't work, but maybe we can take the continuous variable and turn it into a discrete variable. We call this **stratification**.
- ▶ Once it's discrete, we can just calculate the means within each **strata**.

Continuous covariate (III): model relationship as a line

- ▶ The stratification approach was fairly crude: it assumed that means were constant within strata, but that seems wrong.
- ▶ Can we get a more global model for the regression function? Well, maybe we could assume that it is linear:

$$E[Y_i|X_i = x] = \beta_0 + \beta_1 x$$

- ▶ Why might we do this? Parsimony, first and foremost: 2 numbers to predict any value.
- ▶ Some other nice properties we'll talk about in the coming weeks.
- ▶ Here is the linear regression function for the weight-active minutes relationships:
- ▶ We'll see soon how we estimate this line. It's a bit more complicated than the stratify and calculate means.

Parametric vs. nonparametric models

- ▶ The conditional mean approach for discrete independent variables are **nonparametric** because they make no assumptions about the functional form of $E[Y_i|X_i = x]$.
- ▶ We just estimate the mean among each value of x .
- ▶ With continuous independent variables, this approach breaks down because of the number of values.
- ▶ Need to make **parametric** assumptions about the functional form of $E[Y_i|X_i = x]$ in order to make progress
- ▶ These are parametric because they involve writing the functional form in terms of parameters, like the linear model.

Bias-variance tradeoff

- ▶ How we model the regression function, $E[Y_i|X_i = x]$, affects our the behavior of our estimates:
- ▶ Low bias (function “nails” every point)
- ▶ High variance (drastic changes from sample to sample)

Bias-variance tradeoff

- ▶ How we model the regression function, $E[Y_i|X_i = x]$, affects our the behavior of our estimates:
- ▶ Higher bias (misses “local” variation)
- ▶ Low variance (slope and intercept will only change slightly from sample to sample)