



The World's Largest Open Access Agricultural & Applied Economics Digital Library

This document is discoverable and free to researchers across the globe due to the work of AgEcon Search.

Help ensure our sustainability.

Give to AgEcon Search

AgEcon Search

<http://ageconsearch.umn.edu>

aesearch@umn.edu

*Papers downloaded from **AgEcon Search** may be used for non-commercial purposes and personal study only. No other use, including posting to another Internet site, is permitted without permission from the copyright owner (not AgEcon Search), or as allowed under the provisions of Fair Use, U.S. Copyright Act, Title 17 U.S.C.*

No endorsement of AgEcon Search or its fundraising activities by the author(s) of the following work or their employer(s) is intended or implied.

The Stata Journal (2016)
16, Number 1, pp. 23–24

Regressions are commonly misinterpreted: Comments on the article

James W. Hardin
University of South Carolina
Department of Epidemiology and Biostatistics
Columbia, SC
jhardin@sc.edu

How much should we really care whether a description says “comparison” or “change”? Will it lead to a mistake if we say “held constant” instead of “all other things being equal” or “clamping the other variables”? Is there a single best phrase that should be prescribed to all textbooks that discuss the interpretation of coefficients in a regression model? The article by Dr. Hoaglin presents advice on interpretation of regression models and criticism of some commonly applied interpretations.

The author states that 1) the correct interpretation of regression coefficients is evident in added-variable plots; 2) the correct interpretation should be based on an examination of coefficients in multivariate normal distributions and the geometry of least squares; and 3) the proper application of multivariable models requires caution in calculating predictions that average over other variables.

To motivate the discussion, the author uses a notation that places two pieces of information in the subscript of regression parameters. The first part of the subscript identifies the outcome variable and the subscript of the associated covariate. The second part of the subscript identifies the concomitant covariates in the model. In most textbooks (and nearly all articles), this notation is simplified to identify the associated covariate only because the rest of the information is available in context.

Dr. Hoaglin advocates the phrase “adjusting for” instead of “controlling for” when identifying concomitant covariates in discussion of the interpretation of a particular covariate of interest. I agree with the author’s assertion that “controlling for” could imply that randomization rules were applied over those covariates in the collection of data.

To illustrate added-variable plots, the author uses Stata’s ubiquitous automobile dataset. In the example, two correlated covariates are specified in a linear regression model of gallons per 100 miles: the total weight of the car in pounds (weight) and the cubic inch displacement of the car’s engine (displacement). The example illustrates that even though each covariate has a strong relationship as evidenced in a scatterplot, neither appears to be significantly associated with the outcome variable in a multivariable regression. Indeed, in the univariate models, each covariate is found to have a significant association with the outcome variable.

In sections 3 and 4, the author cautions that in most cases, independent variables take on a limited set of values. As such, obtaining predicted values for which one

covariate is allowed to change while other covariates are fixed at their mean may not be meaningful. Not only may mean values of certain covariates lack meaning, any specific covariate pattern used in a prediction could be unrepresented in the data. In the vast majority of cases in public health, predictions are obtained for groups, and the model assumptions adequately address this given that it is the comparison of the predictions that is germane.

In section 4, Dr. Hoaglin cautions against the “holding other covariates constant” phrase when interpreting a coefficient associated with a variable that enters the model in multiple ways either as part of an interaction or in an additional function form as is the case in polynomial models. All authors (of texts to which I had access in my bookshelf) go to great lengths to cover interpretation in such models, and none of them advocate using this phrase in those instances.

The added-variable plot and the initial sections lead Dr. Hoaglin to favoring the phrase “the change in the outcome per unit increase in the covariate of interest after adjusting for simultaneous linear change in the data at hand” over “the change in the outcome per unit increase in the covariate of interest holding the other variables constant”. There is nothing I can say against this preference. Although accurate, the author’s phrase leaves me unsatisfied. But I would never use that phrase. It is dull and lifeless; it fails to illuminate or highlight, and it forces an awkward overly wordy presentation of something that, frankly, I do not mind leaving a little vague.

At worst, “held constant” is a placeholder with which we have become a little too comfortable. That phrase is a nod toward the fitted model and the model’s underlying assumptions. Once the model is fit, and under the assumptions of that model, a researcher can make calculations despite any lack of covariate pattern representation of the particular sample. If we cannot interpret those calculations in terms of the model, then what good is it?

Rather than what specific phrasing was used to describe coefficients and marginal means, I found the examples far more compelling for a different prescription: the inclusion of detailed tabulations, summaries, and regression plots. I contend that the “held constant” phrase is not as relevant as the inclusion of important contextual information.