



Electronic Research Archive of Blekinge Institute of Technology
<http://www.bth.se/fou/>

This is an author produced version of a journal paper. The paper has been peer-reviewed but may not include the final publisher proof-corrections or journal pagination.

Citation for the published Journal paper:

Title:

Author:

Journal:

Year:

Vol.

Issue:

Pagination:

URL/DOI to the paper:

Access to the published version may require subscription.

Published with permission from:

Assessing variable importance through mean value decomposition

H. E. T. Holgersson^{a,b,*}, T. Norman^a and S. Tavassoli^c

Jönköping International Business School, Box 1026, 55 111 Sweden

^b*Linnæus University, 351 95, Växjö, Sweden*

^c*Blekinge Institute of Technology, Department of Industrial Economics, 371 79 Karlskrona, Sweden*

Abstract. Regression analysis is usually concerned with inference of the mean value of a dependent variable conditioned on a set of independent variables. In this paper we propose using a mean value decomposition to assess economic significance as a supplementary tool to statistical hypothesis testing. Unlike many previously suggested methods the proposed decomposition is simple to conduct and requires no specific software. The technique is demonstrated and contrasted with hypothesis tests by an empirical example involving the income of Mexican children, which shows that the two inference approaches provide different and supplementary pieces of information.

JEL classifications: C54, C51, I32

Keywords: Conditioning, Inference, Regression analysis, Mean Value Decomposition, Goodness-of-Fit

* Corresponding author. E-mail: thomas.holgersson@jibs.hj.se

I. Introduction

The linear multiple regression model is one of the most commonly applied dependence techniques used in econometrics. The common text-book approach to empirical regression analysis (e.g., Gujarati, 2003; Neter *et al*, 1996; Theil, 1971; Greene, 2008) involves point estimation and hypothesis testing of individual regression parameters, usually with the purpose of making statements about elasticities. The strong focus on hypothesis testing, however, seems to have its origin in controlled experiments and quality control, and may not always correspond to the implicit questions connected to economics research. McClosky and Ziliak (1996) raise a concern regarding the use of statistical significance through hypothesis testing for science and policy. Their hypothesis is that in conventional economics research, statistical significance is taken to be the same as ‘economic significance’. In a later paper, McClosky and Ziliak (2004) re-investigate this hypothesis and come to the conclusion that most published works in economics persistently ignore the magnitude of the impact of variables, hence causing economic damage. A rare exception to this is the papers by Solon (1992) and Zimmerman (1992) in which the results are analysed in terms of economic significance in favour of statistical significance. Further, both Carver (1978) and Johnson (1999) argue against the use of statistical significance testing since it is too often erroneously applied and should therefore be abandoned. Carver (1978) advises researchers to accompany the minimum p -value with a statistic reflecting the size of the effect or the strength of the association between X and Y , i.e. not only considering ‘statistical significance’, but also ‘economic (practical) significance’. In other words, regardless of whether or not a variable is statistically significant, the magnitude of its value should be interpreted with respect to the research hypothesis in question.

In this paper we propose a simple approach to deal with the relative importance of explanatory variables through an unconditional inference approach. While traditional regression analysis is

concerned with inference of the mean value of a dependent variable Y conditioned on a set of explanatory variables, a decomposition of the unconditional mean value can better help to understand the phenomenon in hand, especially concerning the level of economic relevance of the variables. This is particularly important when the magnitude of the estimated parameter of a significant variable is very small while the mean of the explanatory variable (say X_j) is large. In this situation, relying merely on conventional conditional inference procedures may underestimate the economical, or practical, effect of X_j on Y , since it does not take into account the size of X_j . This paper proposes a mean value decomposition of the dependent variable defined through the law of iterated expectations, to be used as a supplement to traditional significance tests. This decomposition is simple to conduct and explores the economic (practical) significance of the economic phenomenon of interest in a sense which cannot be fully explored with traditional hypothesis testing.

II. Decomposition of the regression response function

A theoretical economic model is usually written in very general terms. For example, a common production function may be defined by

$$Y(i) = F(K(i), L(i))$$

where Y is output, K is capital, L is labour, i is an index over, for example, different states, and $F(\square)$ is an unspecified function (e.g. Romer, 2006, Mankiw, 2000). In contrast to such a general model, empirical research usually makes very precise statements. An econometric analysis may first specify a model in terms of variables only, e.g. *GDP/Capita = Population + Education + Unemployment* or similar, then present a table of estimated regression coefficients, and finally determine the relative importance of the explanatory variables by parameter significances or t -statistics. In other words, the research process starts off with a very general problem and ends up with very precise conclusions. The point made in this paper is not that hypothesis testing is unimportant but that the standard econometric regression analysis should be supplemented with a broader type of statement concerning the importance of a specific input variable. Of particular interest is the unconditioned mean value of the response variable, since it does not restrict the analysis to the specific values of the observed explanatory variables. This is further explored below.

The expected value of a dependent variable Y conditioned on an explanatory variable \mathbf{X} is denoted as $E[Y|\mathbf{X}]$. The unconditional expected value $E[Y]$ is different from the conditional expected value since the specific values of \mathbf{X} are disregarded. These two expectations are conveniently linked to each other through the law of iterated expectations (Stuart, Ord and Arnold, 1994; Greene, 2008), defined by

$$E[Y] = E_x \left[E[Y|\mathbf{X}] \right], \quad (1)$$

where E_x denotes the expectation taken with respect to \mathbf{X} . In terms of a linear regression model we may write the conditioned expected values as

$$E[Y|\mathbf{X}] = E_y \left[\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}|\mathbf{X} \right] = E_y \left[\mathbf{X}\boldsymbol{\beta}|\mathbf{X} \right] + E_y \left[\boldsymbol{\varepsilon}|\mathbf{X} \right] = \mathbf{X}\boldsymbol{\beta},$$

where \mathbf{X} is the observed vector of explanatory variables and $\boldsymbol{\beta}$ is the parameter vector to be estimated. Hence by (1) we may decompose the unconditioned mean as

$$E[Y] = E_x \left[E_y \left[Y|\mathbf{X} \right] \right] = E_x \left[\mathbf{X} \right] \boldsymbol{\beta} = \boldsymbol{\mu}_x \boldsymbol{\beta}, \quad (2)$$

where $\boldsymbol{\mu}_x$ is the average value vector of the explanatory variables and $\boldsymbol{\beta}$ is the same as above. Inference of the regression model is usually drawn by adapting significance tests or confidence intervals with respect to the elements of $\boldsymbol{\beta}$. The conventional inference of the conditional expectation, however, does not take the level of $\boldsymbol{\mu}_x$ into account. For example, a statistically significant yet very small $\boldsymbol{\beta}$ in terms of magnitude might be concluded as having an important economic impact on $E[Y]$. However, according to equation (2), the level or magnitude of $\boldsymbol{\mu}_x$ also has an important meaning. In traditional econometric inference this part is completely disregarded and the focus is set on beta parameters only. On the other hand, the impact on $E[Y]$ in terms of $\boldsymbol{\mu}_x \boldsymbol{\beta}$ is explicitly described by the decomposition in (2). An empirical counterpart of this is readily obtained by

$$\bar{Y} = \bar{\mathbf{X}}' \hat{\boldsymbol{\beta}} = \hat{\beta}_0 + \hat{\beta}_1 \bar{X}_1 + \dots + \hat{\beta}_p \bar{X}_p, \quad (3)$$

where $\bar{\mathbf{X}}$ is the sample mean vector of the explanatory variables in \mathbf{X} and $\hat{\boldsymbol{\beta}}$ is an appropriate estimate of $\boldsymbol{\beta}$, such as the ordinary least square (OLS)- or maximum likelihood estimate. The

mean value decomposition in (3) allows the analyst to assess each term's effect on the mean value of the dependent variable. In order to express the right hand side components of (3) in terms of percentage contribution to \bar{Y} , we may divide equation (3) with \bar{Y} on the left- and right hand side as follows:

$$1 = \hat{\beta}_0/\bar{Y} + \hat{\beta}_1\bar{X}_1/\bar{Y} + \dots + \hat{\beta}_p\bar{X}_p/\bar{Y} \quad (4)$$

For example, the total contribution of X_1 to the mean value of Y is determined by $\hat{\beta}_1\bar{X}_1/\bar{Y}$. This measure, in contrast to other measures of variable importance obtained through orthogonal rotation or other eigenvalue transformations which have been proposed in the literature, is linear and retains the one-to-one relationship between original variables and the individual components of the decomposition, which in turn makes it easy to interpret and apply for the average user. Moreover, the decomposition in (4) also provides a goodness-of-fit measure with respect to each explanatory variable, in the sense that in a model with good fit, the absolute value of each $\hat{\beta}_j\bar{X}_j/\bar{Y}$ should be large relative to $\hat{\beta}_0/\bar{Y}$. In the case that the intercept term dominates the expected value of \bar{Y} , i.e. if $\hat{\beta}_0/\bar{Y}$ is much larger than a corresponding term $\hat{\beta}_j\bar{X}_j/\bar{Y}$, $j = 1, 2, \dots, p$, this indicates that the explanatory variable X_j only makes a small contribution to Y . It may also be shown that $\hat{\beta}_j\bar{X}_j$ unlike $\hat{\beta}_j$ is invariant to scale changes in X_j , e.g. from euros to dollars.

The empirical mean-value decomposition in (3) and its scaled version in (4) may in many cases be more intuitive and easier to interpret than, for example, significance tests, standardised beta coefficients or quadratic measures such as the extra sum of squares (Neter *et al*, 1996) or measures involving rotation of original coordinates (Fabbris, 1980). The next section presents an empirical economic application where the proposed mean-value decomposition is demonstrated along with some alternative measures for the purpose of comparison.

III. Empirical Analysis: Mexican Children Income

In order to demonstrate the use of the above discussed mean value decomposition, we use a data set concerning the income of poor children in Mexican urban areas, their age, the number of children aged under six years in the family, sex, education, working hours, and family income. The data originate from a socioeconomic experiment known as PROGRESA conducted by the Mexican government and the World Bank (source: ENCELURB database). Survey data have been assembled based on interviews of households consisting of a total of five million people over a number of years. The data set used in this section represents a subpopulation of these people interviewed in 2003. Descriptive statistics of the variables are supplied in Table 1 below.

Table 1: Variables and descriptive statistics

Variables	Obs	Mean	Std. Dev.	Min	Max
Income of child (log)	3037	5.298	1.106	1.609	9.212
Child education: in school(1), otherwise (0)	3037	0.294	0.456	0	1
No. of children aged under six in family	3037	0.630	0.897	0	7
Age	3037	15.775	2.144	6	18
Sex: male (1), female (0)	3037	0.658	0.475	0	1
Working hours per week	3037	42.751	21.932	1	98
Family income	3037	1256.624	1178.655	12	30000

The model to be used for inference is the conventional log-linear multiple regression model defined as follows:

$$\ln(Y_i) = \beta_0 + \beta_1 ChEduc_i + \beta_2 Und6_i + \beta_3 Age_i + \beta_4 Sex_i + \beta_5 CWHours_i + \beta_6 FInc_i + \varepsilon_i \quad (5)$$

where Y_i is the income of child i in urban areas of Mexico, β_0 is the constant intercept term, and the explanatory variables corresponding to each β_j , $j=1, \dots, 6$ are defined in Table 1 in order of appearance. Table 2 contrasts the conditional versus unconditional mean value of the Mexican child income model. It reports the OLS estimates of the beta parameters of model (5) with standard errors and, using the empirical decomposition in (4), the table also reports the percentage effect of the term $\hat{\beta}_j \bar{X}_j$ on mean value of Y .

Table 2: Mean value decomposition in Mexican children's income (Model 5)

Variables \bar{X}_j	$\hat{\beta}_j$	\bar{X}_j	$\hat{\beta}_j \bar{X}_j$	$\frac{\hat{\beta}_j \bar{X}_j}{\bar{Y}}$
Child education	-0.393*** (0.0391)	0.294	-0.116	-2.2%
No. children aged under six in family	0.00945 (0.0169)	0.630471	0.006	0.1%
Age	0.154*** (0.00802)	15.77493	2.429	45.9%
Sex	0.00498 (0.0319)	0.657782	0.003	0.1%
Working hours per week	0.0103*** (0.000753)	42.75123	0.440	8.3%
Family income	0.000342*** (1.29e-05)	1256.624	0.430	8.1%
Constant	2.099*** (0.136)	1	2.099	39.6%
Observations	3037			
R-squared	0.434			
Sum				100%

Standard errors in parentheses
 *** p<0.01, ** p<0.05, * p<0.1

Four regression parameters turn out to be highly significant in terms of their t -statistics, i.e. those corresponding to *Child education*, *Age*, *Working hours*, and *Family income*. According to Table 2, if a child works one more extra hour, then his (log) income is expected to increase by about 0.01 units ($\hat{\beta}_5 = 0.01$) holding other variables constant, while one unit increase in the *Family income* of the child would increase his income by only 0.0003 units ($\hat{\beta}_6 = 0.00034$). Hence $\hat{\beta}_5$ is about 34 times larger than $\hat{\beta}_6$. On the other hand, by mean value decomposition, which makes a statement about the effect on the typical or average Child income, a different conclusion is made. It is seen that the average value of *Working hours* has an 8.3% effect on the average (log) income of the Mexican children, which is almost identical to the *Family income* effect (8.1%). This means that, in terms of the mean value decomposition, *Working hours* is equally important as *Family income* even though the corresponding beta parameters differ by a factor of 34. Hence, the point estimates of the regression parameters together with the traditional variable importance measures give contrary results compared with the variable importance measures of the mean value decomposition, for instance in the above comparison of *Working hours* and *Family income*. It should, however, be noted that the two views of the regression model address different questions. The individual regression parameters express the marginal effect of changing the corresponding explanatory variable by one unit, holding all other variables constant, while the individual components of the mean value decomposition explicitly take the levels of the explanatory variables into account and thereby express the corresponding variables' contribution to the average value of the dependent variable.

Additional measures for variable importance have been proposed by, for example, Fabbris (1980), Budescu (1993) and Johnson (2000). These methods no doubt have their own merits but at the same time are technically involved and difficult to interpret. Because simplicity is an important concern in this paper they will not be further considered here.

For example, if policy makers attempted to increase the income of the “average” Mexican school child, then increasing the *Family income* would be just as important as increasing the *Working hours*, even though the latter has a much higher beta value (sampling errors disregarded). In other words, traditional regression methods should, of course, not be replaced but supplemented with the proposed decomposition.

IV. Summary

This paper argues that standard econometric regression analysis does not fully explore the economic dynamics analysed. By linking the conditioned and unconditioned mean values of the response variable we propose a decomposition of the regression model that is argued to provide additional information not given by individual regression parameters or t -tests. Moreover, the components within the mean value decomposition also provide a linear goodness-of-fit measure for the model, which should be compared with previously proposed measures such as changes in R -square or extra sums of squares, which are quadratic and therefore less straightforward to interpret. The proposed method, unlike most other measures of variable importance, does not require any special software facilities. A data set involving income data for poor Mexican children is explored using the proposed mean value decomposition along with other techniques. It is shown that the mean value decomposition gives additional and non-overlapping information when compared with traditional analysis from, for example, standardised beta coefficients or t -statistics, and should therefore be a simple and useful complement to standard methods.

References

- Budescu, D. V. (1993) Dominance analysis: A new approach to the problem of relative importance of predictors in multiple regression, *Psychological Bulletin*, **114**, 542-551.
- Carver, R. P. (1978) The case against statistical significance testing, *Harvard Educational Review* **48**, 378-399.
- Fabbris, L. (1980) Measures of predictor variable importance in multiple regression: An additional suggestion, *Quality & Quantity*, **4**, 787-792.
- Greene, W. H. (2008) *Econometric Analysis*, 6th edn. Pearson Education, New Jersey.
- Gujarati, D. (2003) *Basic Econometrics*, McGraw Hill, New York.
- Johnson, D. H. (1999) The Insignificance of Statistical Significance Testing, *The Journal of Wildlife Management*, **63**, 763-772.
- Johnson, J. W. (2000) A heuristic method for estimating the relative weight of predictor variables in multiple regression, *Multivariate Behavioral Research*, **35**, 1-19.
- Mankiw, N. G. (2000) *Macro Economics*, 4th edn, Worth Publishers, New York.
- McClosky, D. and Ziliak, S. (1996) The Standard Error of Regressions, *Journal of Economic Literature*, **34**, 97-114.
- McClosky, D. and Ziliak, S. (2004) Size matters: the standard error of regressions in the American Economic Review, *The Journal of Socio-Economics*, **33**, 527-546.
- Neter, J, Kutner, M. J., Nachsheim, C. J. and Easserman, W. (1996) *Applied Linear Regression Models*, 3rd edn, Irwin Inc., Illinois.
- Romer, D. (2006) *Advanced Macroeconomics*, 3rd edn, McGraw-Hill Irwin, Boston.
- Solon, G. (1992) Intergenerational income mobility in the United States. *American Economic Review*, **82**, 393-408.
- Stuart, A., Ord, K. and Arnold, S. (1994) *Kendall's Advanced Theory of Statistics Vol 2A*, 6th edn, Wiley, New York.
- Theil, H. (1971) *Principles of Econometrics*, Wiley, New York.
- Zimmerman, D. J. (1992) Regression toward mediocrity in economic stature, *American Economic Review*, **82**, 409-429.