



WORKING PAPER SERIES

Aggregated vs. Disaggregated Data in Regression Analysis: Implications for Inference

Thomas A. Garrett

Working Paper 2002-024B
<http://research.stlouisfed.org/wp/2002/2002-024.pdf>

November 2002

FEDERAL RESERVE BANK OF ST. LOUIS
Research Division
411 Locust Street
St. Louis, MO 63102

The views expressed are those of the individual authors and do not necessarily reflect official positions of the Federal Reserve Bank of St. Louis, the Federal Reserve System, or the Board of Governors.

Federal Reserve Bank of St. Louis Working Papers are preliminary materials circulated to stimulate discussion and critical comment. References in publications to Federal Reserve Bank of St. Louis Working Papers (other than an acknowledgment that the writer has had access to unpublished material) should be cleared with the author or authors.

Photo courtesy of The Gateway Arch, St. Louis, MO. www.gatewayarch.com

Aggregated vs. Disaggregated Data in Regression Analysis: Implications for Inference

Thomas A. Garrett
Federal Reserve Bank of St. Louis
Research Division
411 Locust Street
St. Louis, Missouri 63102

phone: 314-444-8601
email: tom.a.garrett@stls.frb.org

Abstract

This note demonstrates why regression coefficients and their statistical significance differ across degrees of data aggregation. Given the frequent use of aggregated data to explain individual behavior, data aggregation can result in misleading conclusions regarding the economic behavior of individuals.

Aggregated vs. Disaggregated Data in Regression Analysis: Implications for Inference

I. Introduction

Every field of economics uses aggregated data to test hypotheses about the behavior of individuals. Examples in macroeconomics include the use of aggregate consumption and income to test the permanent income hypothesis (Hall, 1978), and forecasting national personal consumption expenditures using consumer sentiment indices (Carroll, et al. 1994; Bram and Ludvigson, 1998). The use of aggregated data to explain individual behavior makes the assumption that the hypothesized relationship between the economic variables in question is homogenous across all individuals. When the behavior of economic agents is not the same, a regression analysis using aggregated data can provide conclusions regarding economic relationships that are different than if less aggregated data were used. Correcting for this ‘aggregation bias’ has received careful attention in the literature.¹

This note develops a simple framework to show how coefficient estimates and their statistical significance can differ using aggregated versus less aggregated data. The analysis explores the effects of dependent variable data aggregation on coefficient estimates and standard errors. A classic example of the framework presented here is the empirical work of Carroll, et al. (1994) and Bram and Ludvigson (1998) that explores the ability of consumer sentiment (a measure of an individual’s perception of economic conditions) to forecast national personal consumption expenditures (a highly aggregated measure of consumption). The analysis presented here is useful to both economics graduate students and applied economists as it provides general insights into the impact of data aggregation on regression estimates and conclusions made from statistical inference.

II. Analysis

Statistical inference on regression coefficients is traditionally done using a t -statistic. The t -statistic value for testing $H_0: \beta=0$ is the estimated coefficient, $\hat{\beta}_k$, divided by its standard error. Thus, for a chosen critical value, the significance of any coefficient depends upon the size of the coefficient and its variance. The difference in coefficient size and variance (and thus statistical significance) from regressions using various levels of data aggregation in the dependent variable is the focus of this analysis.

Consider m regression equations each having T observations and data matrix \mathbf{X} that is assumed identical across equations:

$$\begin{aligned} Y_1 &= \mathbf{X}\beta_1 + U_1 \\ Y_2 &= \mathbf{X}\beta_2 + U_2 \\ &\dots \\ &\dots \\ Y_m &= \mathbf{X}\beta_m + U_m \end{aligned} \tag{1}$$

where \mathbf{Y}_m is a $T \times 1$ vector, \mathbf{X} is $T \times K$ with $K-1$ explanatory variables (assume a constant term), β_m is the $K \times 1$ vector of estimated coefficients, and \mathbf{U}_m is the $T \times 1$ residual vector. As an example of the above framework that follows the consumer sentiment literature, each \mathbf{Y}_m could be personal consumption expenditures for each state ($m = 50$) and \mathbf{X} consists of lagged consumer sentiment values.

Summing equations provides the aggregated regression equation

$$Y_M = \mathbf{X}\beta_M + U_M \tag{2}$$

where the aggregated regression coefficient vector β_M is²

$$\beta_M = \sum_{m=1}^M \beta_m \quad (3)$$

It then follows that

$$Y_M = \sum_{m=1}^M Y_m ; \quad U_M = \sum_{m=1}^M U_m \quad (4)$$

It is clear from (3) that any estimated coefficient from the aggregated regression is simply equal to the sum of the corresponding coefficients from the less aggregated regressions. Depending on the signs and magnitudes of each β_m , the estimated impact of changes in each explanatory variable from the aggregated and less aggregated regressions can be quite different. For example, if all β_m are positive, then β_M will be larger than any β_m . The estimated impact of changes in explanatory variable on the dependent variable is thus much greater in the aggregated regression than in the less aggregated regression.

Now consider the variance of a coefficient vector. The variance of a coefficient vector from any OLS regression is [see Greene (1990), page 184]

$$\text{var}[\beta_m] = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1} \quad (5)$$

where σ^2 is the sum of squared residuals (RSS). Because \mathbf{X} is the same for each regression equation, the variance of a specific coefficient in the aggregated regression relative to the

variance of each corresponding coefficient in the less aggregated regressions thus depends solely on differences in the RSS between the less aggregated regressions and the aggregated regression.

Using the expression for U_M in (4), the aggregated regression RSS is

$$U_M^2 = (\sum_m U_m)^2 \quad (6)$$

The questioned relationship between the RSS from the aggregated regression and the summed RSS from the less aggregated regressions ($\sum_m U_m^2$) can be written as:

$$(\sum_m U_m)^2 \cong \sum_m U_m^2 \quad (7)$$

or

$$(\sum_m U_m) \cdot (\sum_n U_n) \cong \sum_m U_m \cdot U_m \quad (8)$$

Giving the final result

$$\sum_m U_m \cdot U_m + 2 \cdot \sum_{m=1}^{N-1} \sum_{n=m+1}^N U_m \cdot U_n \cong \sum_m U_m \cdot U_m \quad (9)$$

Let $Z = 2 \cdot \sum_{m=1}^{N-1} \sum_{n=m+1}^N U_m \cdot U_n$. Z is the sum of the product of cross equation residuals. If $Z < 0$,

then the RSS from the aggregated regression will be less than the sum of RSS from the less aggregated regressions. The sign of Z depends upon the covariance (and thus correlation)

between the cross-equation residuals.³ If $Z < 0$, then the cross-equation residuals have a

negative covariance and are negatively correlated -- the residuals from the individual regression

equations have opposing signs (a result of differing signs on corresponding slope coefficients), and summing positive and negative residuals reduces the absolute magnitude of the aggregated regression residuals. Thus, the spread of the residuals in the aggregate regression is less than the sum of the residual spreads from the less aggregated regressions. Conversely, a $Z > 0$ implies that the summed residuals from the less aggregated regressions result in aggregated residuals that are more than the sum of the individual residual spreads -- residuals from the individual regression equations have the same signs (a result of same signs on corresponding slope coefficients), and summing same sign residuals increases the absolute magnitude of the aggregated regression residuals.

In closing, the analysis finds the following connection between data aggregation and statistical inference: the size of the RSS from the aggregated regression relative to the sum of the RSS from less aggregated regressions depends on the correlation of residuals across the less aggregated regressions. The difference in the RSS then translates into differences in coefficient standard errors as seen in (5). This, combined with the additive relationship between the less aggregated and aggregated coefficients shown in (3), gives different conclusions from statistical inference across levels of data aggregation.

III. Summary

This analysis has shown why the sign and significance of coefficient estimates from regressions using aggregated data can differ from regressions using less aggregated data. The size of a coefficient estimate from aggregated data is shown to be the sum of each coefficient from the less aggregated regressions. More importantly, it is shown that the RSS from the

aggregated regression can be larger or smaller than the sum of RSS from less aggregated regressions depending upon the covariance and correlation between cross-equation residuals. Because, given a chosen critical value, statistical significance of a coefficient is a positively related to coefficient size and negatively related to its standard error, it is likely that coefficients from an aggregated regression are statistically significant whereas identical coefficients from less aggregated regressions are statistically insignificant, and vice versa. This results in different conclusions regarding economic behavior depending upon the level of data aggregation.

Endnotes

1. See Goodfriend (1992), Thomas and Tauer (1994), Mittlehammer, et. al (1996), Davis (1997), and Cherry and List (2002).

2. This is derived from the fact that

$$\begin{aligned}\beta_m &= (X'X)^{-1}X'Y_m \\ \sum_m \beta_m &= \sum_m (X'X)^{-1}X'Y_m \\ \sum_m \beta_m &= (X'X)^{-1}X'\sum_m Y_m\end{aligned}$$

3. This is seen by taking the expectation of (9). $\text{Cov}(U_m, U_n) = \rho_{U^m, U^n} \cdot \sigma_{U^m} \cdot \sigma_{U^n}$ where ρ is the correlation between residual vectors and σ is the standard error of each residual vector.

References

- Bram, Jason and Sydney Ludvigson, "Does Consumer Confidence Forecast Household Expenditures? A Sentiment Index Horse Race." *Federal Reserve Bank of New York Economic Policy Review*, vol. 4 (June 1998): 59-78.
- Carroll, Christopher, Jeffrey Fuhrer, and David Wilcox, "Does Consumer Sentiment Forecast Household Spending? If So, Why?" *American Economic Review*, vol. 84, no. 5 (December 1994): 1397-1408.
- Cherry, Todd and John List, "Aggregation Bias in the Economic Model of Crime." *Economics Letters*, vol. 75, no. 1 (March 2002): 81-86.
- Davis, George, "Product Aggregation Bias as a Specification Error in Demand Systems," *American Journal of Agricultural Economics*, vol. 79, no. 1 (February 1997): 100-109.
- Goodfriend, Marvin, "Information-Aggregation Bias." *American Economic Review*, vol. 82, no. 3 (June 1992): 508-519.
- Greene, William H. Econometric Analysis. MacMillan Publishing, New York, 1990.
- Hall, Robert E. "Intertemporal Substitution and Consumption." *Journal of Political Economy*, vol. 96 (April 1998): 339-357.
- Mittlehammer, Ron, Hongqi Shi and Thomas Wahl, "Accounting for Aggregation Bias in Almost Ideal Demand Systems." *Journal of Agricultural and Resource Economics*, vol. 21, no. 2 (December 1996): 247-262.
- Thomas, Arthur and Loren Tauer, "Linear Input Aggregation Bias in Nonparametric Technical Efficiency Measurement." *Canadian Journal of Agricultural Economics*, vol. 42, no. 1 (March 1994): 77-76.