

# Exploration of the Use of Empirical Bayes Procedures for Estimating Changes in Occupancy Rate and Persons per Household

Lynn Weidman, Robert Creecy, Donald Malec, Julie Tsay  
Statistical Research Division, U. S. Census Bureau, Washington, DC 20233<sup>1</sup>

## Abstract

The Census Bureau is carrying out research with the intent of improving the housing unit-based method for estimating population totals. One approach is to use the Decennial Census as the baseline for measures and update them annually using American Community Survey (ACS) data to estimate their change since the census year. This study looks at the possibility of using an Empirical Bayes approach to produce county estimates of change with smaller variances than direct ACS county estimates. Results from national and state models are compared. Data from the 1990 and 2000 long form samples are used to represent ACS and data from the 1990 and 2000 short form are used as independent variables in the models. In the actual application, sources for the latter would be the Master Address file and Administrative Records.

## 1. Introduction

Staff of the Population Division and others in the U. S. Census Bureau have been carrying out research into methodologies for improving housing unit-based county population estimates. The housing unit-based method estimates the number of housing units (HUs), percent of occupied housing units (%OCC), and persons per occupied housing unit (PPH), and uses their product (HUs\*%OCC\*PPH) to estimate the number of persons living in HUs. The American Community Survey (ACS) is a very large new household survey that provides updated information annually on %OCC, PPH, and related measures, as well as on a large number of demographic, economic, and social characteristics. All of these characteristics were formerly available only once a decade following the decennial census. The approach of the project we report on is to use the ACS as the source of direct estimates of change in %OCC and PPH between the current year and the previous census year, then apply Empirical Bayes (EB) methods to find estimates of change in %OCC and PPH that have smaller variances than these direct estimates. These are then combined with the previous census values to get the current year estimates.

The sample size and continuous nature of the ACS make it an obvious source to use as the basis for producing current year estimates with relatively small sampling variability. The variables correlated with %OCC and PPH that are required for the EB approach would in practice be derived from administrative records and the Census Bureau's Master Address File (MAF). This paper reports on the initial results of this research project using decennial census data and introduces some of the practical and statistical issues that must be addressed before this approach can be applied.

We give a short presentation of the general approach and some estimation issues in the following section, summarize the EB method in section 3, describe the data used in section 4, give selected regression models in section 5, summarize results of the EB estimation in section 6, and present a brief conclusion in section 7.

## 2. General Approach

Consider the situation where the Census Bureau is producing housing unit-based estimates of county populations using %OCC and PPH in the years following the 2010 Census (C2010). This project looks at a method of estimating the changes in these variables between the 2010 ACS and a future year ACS. For each of the variables, the change estimate would be combined with the Census 2010 value to obtain an estimate of the current year value. Why don't we just use the current year value of ACS as our estimate? There are the following differences between the data collection methodologies used in the census and the ACS that result in their not estimating exactly the same %OCC and PPH parameters.

---

<sup>1</sup> This report is released to inform interested parties of research and to encourage discussion. The views expressed are those of the authors and not necessarily those of the U.S. Census Bureau.

- (1) The reference date for each ACS interview is the day on which it is begun, which can be any day throughout the year. The census determines its reference date in a similar manner, but it is always close to April 1.
- (2) Because the ACS represents an average of characteristics that can change over time, the persons included as residents of a HU are generally only those who are staying there for longer than two months. We refer to this as the current residence rule. This differs from the census which includes all persons who are considered to stay at the HU ‘most of the time.’ (All sample persons currently staying at a non-HU residence, referred to as group quarters (GQ), are interviewed as in the census.)
- (3) An ACS interview of the persons in a HU can take place during a three-month period. A HU that is initially vacant during this time can become occupied and an interview completed later, so the HU is measured as occupied although it is only occupied during part of the period. A HU can also be initially occupied and become vacant later, but if the interview was completed during the initial period then it is also measured as occupied although it was not occupied during the entire period. There are additional scenarios, but overall they result in the number of occupied houses being over-counted compared to the census.

As a result of these differences in data collection methods, changes between estimates from ACS in two different years will need to be scaled to correctly represent changes based on the Census 2010 residence rule.

There is an additional issue with using ACS county estimates that needs to be addressed. Single year estimates and their variances are currently only produced for counties (and other geographical entities) with a population larger than 65,000, due to larger variances for smaller populations. But it would be easy to produce them for all smaller counties because the necessary weighting has been done. And the larger sampling variances for %OCC and PPH in these smaller counties is not a concern for the EB approach, as they are just used in determination of the relative contributions of the direct and modeled estimates to the EB estimate.

### 3. Empirical Bayes Method

The EB method requires a direct estimate of the characteristic of interest and variables correlated with that characteristic to be available. The direct estimate of the characteristic is regressed on these correlated variables and for each observation a regression estimate of the characteristic is calculated in addition to the direct estimate. The EB estimator is then a weighted average of the direct estimate and the linear regression estimate of the characteristic. The weights on the two estimates depend on the sampling variance of the direct estimates and how well the model fits the data.

For purposes of illustration we will use change in PPH as the variable being estimated. Let

$$\theta_i = (\text{current year PPH under current residence rule for county } i) - (\text{2010 PPH under current residence rule for county } i), \quad i=1,2,\dots,k, \tag{1}$$

and

$$Y_i = \theta_i + \eta_i \tag{2}$$

be a direct estimate of  $\theta_i$  from a survey, where the  $\eta_i$ 's are independent sampling errors with  $E(\eta_i | \theta_i) = 0$  and  $V(\eta_i | \theta_i) = V_i$ , so the  $Y_i$ 's are design unbiased for  $\theta_i$  (their expected values over all possible samples are equal to the quantities they are attempting to estimate).

Now suppose that we have a set of  $r$  model variables available for each county  $i$ , represented by the vector  $\mathbf{z}_i' = (z_{i1}, z_{i2}, \dots, z_{ir})$ , such that

$$\theta_i = \mathbf{z}_i' \underline{\beta} + \varepsilon_i, \tag{3}$$

is a model that represents PPH change, where  $\underline{\beta}$  is a vector of unknown coefficients,  $E(\varepsilon_i) = 0$ , and  $V(\varepsilon_i) = A_i$  with the  $\varepsilon_i$ 's independent.

Then, substituting (3) into (2), we can write a model for our direct estimate as

$$Y_i = \theta_i + \eta_i = \mathbf{z}_i' \underline{\beta} + \varepsilon_i + \eta_i. \tag{4}$$

This is called a mixed model for  $Y_i$  because it includes the fixed parameters  $\underline{\beta}$  and the random term  $\varepsilon_i$  with  $V(\varepsilon_i) = A_i$ . (This use of a single variable  $\varepsilon_i$  is the simplest form for the random part of a mixed model. The random part could be a vector, a regression, or

another function of random variables.) If all the  $V_i + A_i$  are known, then we can compute a regression estimate  $\underline{b}$  of the coefficient vector  $\underline{\beta}$  via generalized least squares as

$$\underline{b} = (Z'DZ)^{-1}Z'D\underline{Y},$$

with  $Z$  the  $k \times r$  matrix with rows  $z_i'$ ,  $D^{-1} = \text{Diagonal}(V_1 + A_1, \dots, V_k + A_k)$ , and  $\underline{Y} = (Y_1, \dots, Y_k)$ .

If we make the assumption that both  $\varepsilon_i$  and  $\eta_i$  have normal distributions, then we have

$$Y_i | \theta_i \sim N(\theta_i, V_i) \text{ and } \theta_i | \underline{\beta}, \eta_i \sim N(z_i' \underline{\beta}, A_i). \quad (5)$$

(We can make other distributional assumptions but normality makes the problem simpler, so we will work with it for now.)

We then look at estimators of the  $\theta_i$ , which are weighted averages of the  $Y_i$  and the  $z_i' \underline{\beta}$ .

$$\hat{\theta}_i = C_i Y_i + (1 - C_i) z_i' \underline{\beta}. \quad (6)$$

Then the current year PPH is estimated as  $g(\text{C2010 PPH}, \hat{\theta}_i)$ , where  $g(\cdot)$  is a function that adjusts for definitional and operational differences between the ACS and the decennial census. (If there were no such differences, then  $g(\cdot) = (\text{C2010 PPH}) + \hat{\theta}_i$ .) The best linear unbiased predictor (BLUP) in (6) when the  $A_i$ 's and  $V_i$ 's are known requires  $C_i = A_i / (V_i + A_i)$ . Since the  $A_i$ 's and  $V_i$ 's are not known, we substitute estimates  $\hat{A}_i$  and  $\hat{V}_i$  for them to get  $\hat{C}_i$  and the empirical BLUP (EBLUP). Under our assumption of normality of the variance components, the EBLUP estimators are the same as the EB estimators, where the posterior distributions of the  $\hat{V}_i$  are normal with

$$\text{mean } \hat{C}_i Y_i + (1 - \hat{C}_i) z_i' \underline{\beta}, \text{ and variance } \hat{C}_i V_i. \quad (7)$$

$\hat{C}_i V_i$  is actually an underestimate of the variance because it does not take account of the variability in  $\hat{A}_i$  and  $\hat{V}_i$ . It may be desirable to find a better estimator of the true posterior variance than  $\hat{C}_i V_i$  and there are approaches in the literature for doing so, e.g. in Rao, 2003. (See Morris, 1983 for additional details about the basic EB set-up.)

## 4. Data Used in the Study

Our task was to carry out an exploration of the feasibility of using the EB procedure with data similar to what would be available in actual implementation. The 2010 ACS data are represented by the C90 longform and the 2010 shortform data are represented by the C90 shortform (100% data). Future year ACS data are represented by the C2K longform, while the MAF and administrative record variables for the future year are represented by the C2K shortform.

### 4.1 Model Variables

In the actual application of this methodology, the sources of the EB model variables should have very small sampling variances. If we tried to use variables from the ACS, they would have sampling variances of a magnitude, especially for the less populous geographies, that would need to be incorporated into the estimation procedure. We would prefer to select variables from the shortform of C2010, the MAF, and administrative records. Here we use %single units, %10+ units, %urban units, %Hispanic, and %non-Hispanic white from the C90 shortform, the change in these variables between the C90 and C2K shortforms, and appropriate interactions between these two sets as model variables.

### 4.2 Estimation Variables

It is not appropriate to use the census longform estimates of %OCC and PPH as the estimation variables when shortform estimates are used as the model variables. This is because the final stage of longform weighting controls combinations of several characteristics to be equal to their shortform values, which results in (some of) the  $z_{ij}$  in the mixed model having higher correlation with the estimation variables than in the actual application. To avoid this, we create estimates of %OCC and PPH by using only the HU sampling rates and a single 'nonresponse' adjustment as follows. Both the Census 1990 and 2000 samples are selected from multiple strata. On the

long form files there are sample occupied HUs that have data but no final weights. This is because they were determined not to meet the definition of a HU. We treat them as nonresponding units and apply a nonresponse ratio adjustment by stratum for each county.

The variance for the number of occupied HUs in a county for each census year is estimated using the appropriate formula for a stratified sample. It is divided by the square of the number of HUs to get variance(%OCC). (Our assumption is that for each county the total number of HUs in the sampling frame is equal or very close to the sum over strata of the number of units in sample times the inverse of the stratum sampling rate. This sum is used for the number of HUs when calculating %OCC and variance estimates.) The variance formula for Census 1990 is given by

$$\hat{V}(\% \hat{OCC}_{90}) = \left( \frac{1}{N_{90}} \right) \sum_{h,90} N_{h,90}^2 \left( \frac{p_{h,90}(1-p_{h,90})(1-f_{h,90})}{(n_{h,90}-1)} \right), \quad (8)$$

where

h denotes stratum,

$N_{90}$  is the number of HUs in the county,

$N_{h,90}$  is the number of HUs in stratum h,

$n_{h,90}$  is the number of sample HUs in stratum h,

$p_{h,90}$  is the fraction of sample HUs in stratum h that are occupied, and

$f_{h,90}$  is the proportion of HUs sampled in stratum h

In full notation there would be additional subscripts for county and state but for simplicity we consider a given county and state. For Census 2000 the subscripts 90 are replaced by 00. The estimated variance of the difference in %OCC and PPH between Census 1990 and Census 2000 is the sum of their individual variance estimates.

The variance of the estimated number of persons in a county is also estimated based on the appropriate stratified sampling formula, with adjustment for the nonresponding occupied units and other assumptions required. Derivation of the following approximate variance formula is given in Appendix 2 of Weidman, et al. (2008).

$$\hat{V}(\text{pers\`ons}_{90}) \approx \sum_{h,90} N_{h,90}^2 \left\{ \frac{p_{h,90}(1-p_{h,90})(1-f_{h,90})}{(n_{h,90}-1)} \left[ \frac{s_{hor,90}^2(1-\rho_{ho,90}f_{h,90})}{\rho_{ho,90}n_{ho,90}} + \bar{x}_{hor,90}^2 \right] + p_{ho,90}^2 \frac{s_{hor,90}^2(1-\rho_{ho,90}f_{h,90})}{\rho_{ho,90}n_{ho,90}} \right\} \quad (9)$$

where

$n_{ho,90}$  is the number of occupied sample HUs in stratum h,

$\rho_{ho,90}$  is the response rate for occupied HUs in stratum h,

$\bar{x}_{hor,90}$  is the mean number of persons in occupied respondent HUs in stratum h, and

$s_{hor,90}^2$  is the estimated variance of the number of persons per occupied HU in stratum h based on the responding HUs.

An estimated variance for PPH is obtained using the standard formula for the approximate variance of the ratio of two random variables,

$$\hat{V}(\hat{PPH}_{90}) \approx \frac{\hat{V}(\text{pers\`ons}_{90})}{N_{90}^2 (\% \hat{OCC}_{90})^2} + \frac{(\text{pers\`ons}_{90})^2}{N_{90}^4 (\% \hat{OCC}_{90})^4} N_{90}^2 \hat{V}(\% \hat{OCC}_{90}), \quad (10)$$

and substituting  $\hat{V}(\% \hat{OCC}_{90})$  from (8) and  $\hat{V}(\text{pers\`ons}_{90})$  from (9).

The estimated variance of a difference between estimates from separate census years is the sum of the individual year variances, due to the independence of their samples.

## 5. Regression Models

The initial step in development of mixed models was to identify variables correlated with change in %OCC and PPH via forward selection stepwise linear regression. (DC, DE, HI, and RI are excluded from the state models because they have five or fewer counties and no degrees of freedom for estimating the random county effect in the mixed model. In future work we would pursue how to handle these states, perhaps by grouping them together or combining them individually with neighboring states.) In most cases only a few variables are included in the regression model and  $R^2$  is of a reasonable size. These results suggest that modeling change in %OCC and PPH with independent variables similar to those used in this task should provide some improvement in estimates via the EB procedure.

## 6. Empirical Bayes Modeling Results

We attempted to fit mixed models with a single random county term for %OCC and PPH for the nation and most states, using the model variables selected in the stepwise regressions. The county sampling variances  $V_i$  of the estimation variables were treated as fixed at their estimated values. The variances  $A_i$  for the random county effects in a given model were assumed to have the same value  $A$ . We used both the SAS procedure MIXED and a custom-written R program using an E-M algorithm (Creedy, 2008) to find maximum likelihood estimates of the parameters in the mixed models. For the national models SAS gave the message that there was not enough memory to estimate them. Even consultation with SAS staff did not lead to a solution for this problem. We were able to fit these models using R. For %OCC both programs obtained solutions for all states. For PPH complete convergence was not obtained for 17 states. For some of these states MIXED was not able to start the solution procedure and for others  $\hat{A} = 0$  but the final Hessian was not positive definite. For these states, plus states 22 and 25, R was not able to obtain convergence but  $\hat{A}$  was close to zero when it stopped. As a double check on these results, we searched the likelihood surface for a few of the 17 states using a Fortran program and  $\hat{A}$  was equal to zero, so we use  $\hat{A} = 0$  as the estimate for these 17 states. (Because of the variability in the distributions of the estimates  $\hat{A}$ , in a given state the true value that is being estimated may not be close to zero. Seven of these states have more than 50 counties, so the variability is probably fairly small and the true values are likely to be near 0. The remaining 10 states have fewer than 40 counties and their true values are more likely to not be near 0.) None of the states for which both programs converged to an estimate without any warning messages found  $\hat{A} = 0$ .

Figure 1 shows some percent reductions in variance for the national and state EB estimates of change in county PPH compared to the direct sample estimates, using the estimated coefficient  $\hat{C}_i = \hat{A} / (\hat{V}_i + \hat{A})$ . To give an idea of the range of reductions, they are calculated for the counties with the minimum, median, and maximum estimated sampling variances in each state.

*%OCC variance reduction summary.* There is very little reduction for the minimum variance counties using either model, except for one state. This is not surprising since these variances are so small. For the median variance counties most of the national model reductions are less than 20% and about half the states have additional improvement of at least 10% for the state model. The national model gives reductions of at least 30% for most states and the state models usually show substantial additional reduction – more than 40% for some states.

*PPH variance reduction summary.* The variance reductions are in general much larger for PPH than for %OCC. Twenty-two states show more than 10% reduction for the minimum variance counties with the national model, and many show an additional 10% or more reduction with the state model. For two states this additional improvement is more than 70%, so the national model is not appropriate for them. Most states show at least 40% reduction for their median variance counties with noticeable additional reductions for the state models. There are large reductions with the national model in most states, so the state models do not usually offer substantial additional reduction.

Note two things about these comparisons. First, care should be used when interpreting the amount of improvement of a state model over the national model. Estimates for the state models are based on many fewer degrees of freedom than are those for the national model, so the variances of the variance component estimates are larger. But 30 of the 47 states have more than 50 counties and these variances are probably suitably small to allow valid comparison of the two values. Secondly, there are a few states for which the national model gives more variance reduction than the state model. We might expect that this would happen in states with smaller numbers of counties, where the national model would supply more degrees of freedom for estimation. Upon examination of the number of counties, we see that this is not the case and the state with the most counties, Texas, is in this group.

Figure 2 shows, as an example from a single state, MS, how the estimates of change in PPH differ across the sample, the state regression model, the EB procedure, and the full Bayes procedure. The counties are ordered from smallest to largest sampling variance. For each county the EB estimate lies between the sample and model estimate since it is a linear combination of them. The most important thing to note is that the Full Bayes and EB estimates are close together for most counties, so that in these cases the EB estimator obtains most of the benefits of the full Bayes estimator without the additional assumptions the latter requires about prior distributions. However, there are some counties where the full Bayes estimator does not lie between the sample and model estimates, so the EB estimate is not close to it.

## 7. Conclusion

The purpose of this paper is to present research into the feasibility of using the EB approach to reduce the variability of direct estimates of change in %OCC and PPH. Overall we see that the EB approach can give noticeable reductions in variance from the direct sample estimates, especially as sampling errors get larger, even with the simple models used. For most states, using state-specific information in the models gives additional improvement over using just national information. The size of the variance reductions and the closeness of EB estimates to the full Bayes estimates for most counties suggest that further research into the application of the EB methodology to estimating %OCC and PPH from the ACS with auxiliary data from the MAF and administrative records would be worthwhile.

Before this methodology can be applied to the situation introduced at the beginning of this paper, there are multiple avenues of investigation that would need to be pursued. Several issues concerning the relationship of estimates between the decennial census and the ACS were introduced but not pursued, as well as the issue of how to handle states with few counties. In addition, we have not attempted to look at more complex mixed models to find additional reductions in variances of EB estimates.

Of course, the EB approach can be applied with any sample estimator, not just the ones used here. So it may be possible to use it with estimators investigated in other HU-based method research projects after determining a set of correlated auxiliary variables.

## References

- Cochran, William (1977). *Sampling Techniques*, John Wiley and Sons.
- Creedy, Robert (2008). "Computation of Empirical Bayes Estimates Using Single Level Mixed Models", *Research Report: Statistics #2008-2*, Statistical Research Division, U.S. Census Bureau.
- Morris, Carl (1983). "Parametric Empirical Bayes Inference: Theory and Applications", *Journal of the American Statistical Association*, 73, pp. 47-55.
- Rao, J.N.K. (2003), *Small Area Estimation*, John Wiley and Sons.
- Weidman, L., Creedy, R., Malec, D., and Tsay, J. (2008). "Exploration of the Use of Empirical Bayes Procedures for Estimating Changes in Occupancy Rate and Persons per Household", *Research Report: Statistics #2008-7*, Statistical Research Division, U.S. Census Bureau.

**Figure 1. Percent Variance Reduction in PPH by State: National and State Models**

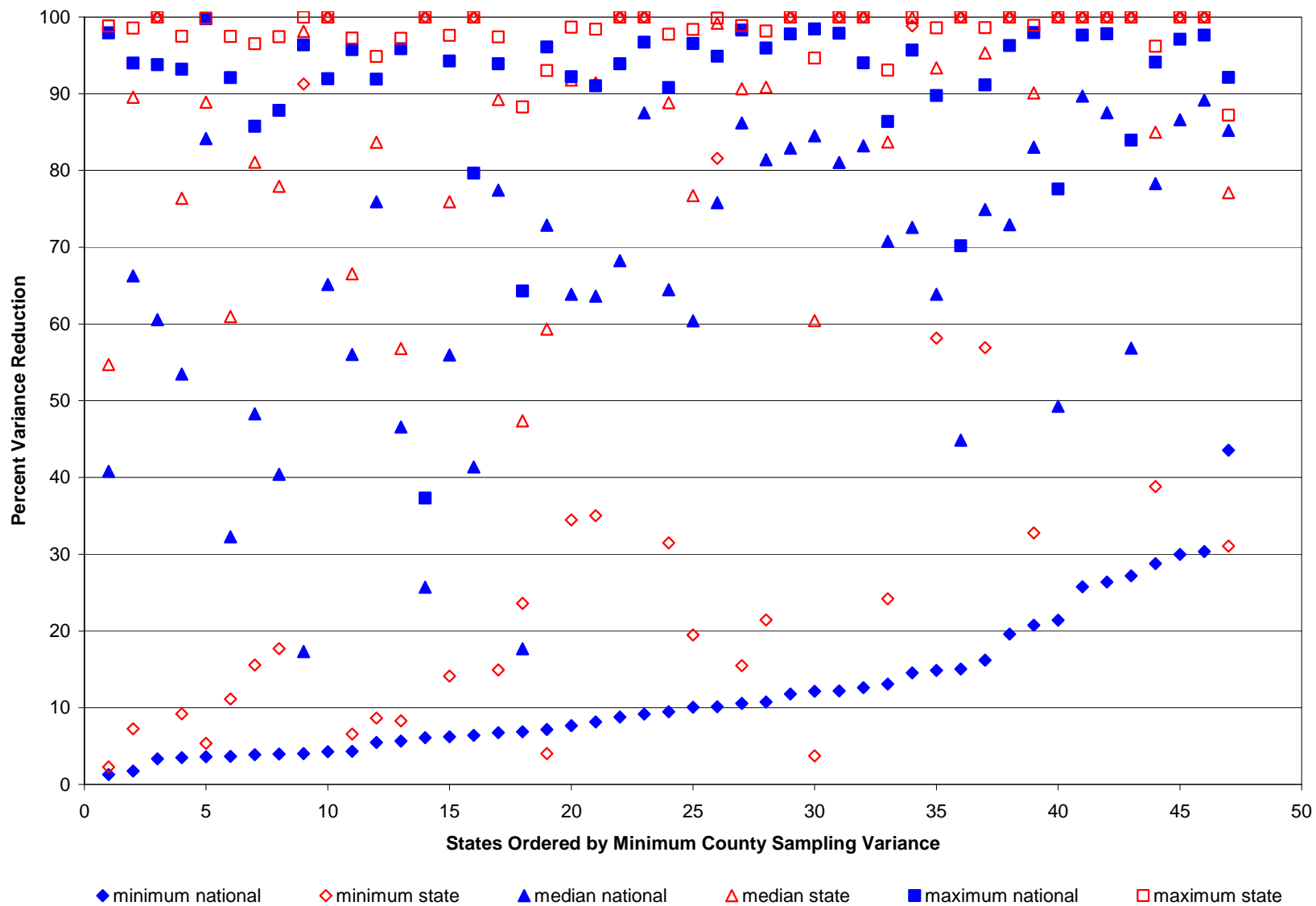


Figure 2. Comparison of County Estimates of Change in PPH for Mississippi

