

## CREATING IMPROVED SURVEY DATA PRODUCTS USING LINKED ADMINISTRATIVE-SURVEY DATA

MICHAEL E. DAVERN\*

BRUCE D. MEYER

NIKOLAS K. MITTAG

Recent research linking administrative to survey data has laid the groundwork for improvements in survey data products. However, the opportunities have not been fully realized yet. In this article, our main objective is to use administrative-survey linked microdata to demonstrate the potential of data linkage to reduce survey error through model-based blended imputation methods. We use parametric models based on the linked data to create imputed values of Medicaid enrollment and food stamp (SNAP) receipt. This approach to blending data from surveys and administrative data through models is less likely to compromise confidentiality or violate the terms of the data sharing agreements among the agencies than releasing the linked microdata, and we demonstrate that it can yield substantial improvements of estimate accuracy. Using the blended imputation approach reduces root mean squared error (RMSE) of estimates by 81 percent for state-level Medicaid enrollment and by 93 percent for substate area SNAP receipt compared with estimates based

MICHAEL DAVERN is with the NORC at the University of Chicago, Health Care Research. BRUCE MEYER is with Harris School of Public Policy, University of Chicago. NIKOLAS MITTAG is with CERGE-EI. Any opinions and conclusions expressed here are those of the authors and do not necessarily represent the views of the New York Office of Temporary and Disability Assistance (OTDA), the US Census Bureau, or NORC at the University of Chicago. The ACS-OTDA data analysis was conducted at the Chicago Census Research Data Center by researchers with Special Sworn Status, and the results were reviewed to prevent the disclosure of confidential information. We would like to thank Katharine Abraham for comments on an earlier draft. Meyer and Mittag would like to thank the Alfred P. Sloan Foundation, the Charles Koch Foundation and the Russell Sage Foundation for support. Mittag is thankful for financial support from the Czech Science Foundation (through grant no. 16-07603Y) the UNCE project (UNCE/HUM/035) and the Czech Academy of Sciences (through institutional support RVO 67985998).

\*Address Correspondence to Michael Davern, Executive Vice President of Research, NORC at the University of Chicago, Health Care Research, 55 E Monroe Suite 3000, Chicago IL 60603 USA. Email: Davern-Michael@norc.org. Bruce Meyer, Professor, Harris School of Public Policy, University of Chicago, 1155 E. 60th Street, Chicago IL 60637. Email: Meyer1@uchicago.edu. Nikolas Mittag, Assistant Professor, CERGE-EI, Politických vězňů 7, 110 00 Prague, Czech Republic.

on the survey data alone. Given the high level of measurement error associated with these important programs in the United States, data producers should consider blended imputation methods like the ones we describe in this article to create improved estimates for policy research.

**KEYWORDS:** Administrative data; Data linkage; Measurement error; Survey methods.

## 1. INTRODUCTION

The goal of survey design is to make good decisions about resources spent to reduce different types of survey error given the level of funding and the specific research question the survey sponsor would like the data to answer. Survey researchers catalog the potential sources of survey errors that can influence survey-based estimates into five basic sources of error: (1) sampling error, (2) sample coverage error, (3) nonresponse error, including both unit and item nonresponse, (4) measurement error, and (5) processing error ([Federal Committee on Statistical Methodology \(FCSM\) 2001](#)). In this article, our main objective is to improve overall estimate accuracy by reducing measurement error in survey estimates of program participation by means of using models developed on linked administrative and survey data to impute program receipt status. We discuss methods to blend information from the linked data with the public use data that neither require access to the restricted linked data nor compromise confidentiality in the public use data.<sup>1</sup> We show that estimates blending administrative and survey data substantially reduce error that has been observed in critical policy relevant survey estimates of Medicaid enrollment and Supplemental Nutrition Assistance Program (SNAP) receipt. We argue that such utilization of data linkage to improve estimates is a more cost-effective approach to increase survey accuracy than many other current practices used to reduce survey error.

Past research has demonstrated substantial measurement error and bias in survey estimates of Medicaid enrollment and SNAP receipt. Work using the Current Population Survey (CPS) found that 43 percent of those linked to administrative data with Medicaid coverage did not report having coverage (false negatives). On the other hand, only 1 percent of respondents in the CPS reported having Medicaid coverage that could not be confirmed through the linkage. This pattern results in a large net undercount ([Davern, Klerman, Ziegenfuss, Lynch, and Greenberg 2009a](#)).<sup>2</sup> Research on survey misreporting

1. The additional risk to disclosure from synthetic data or model parameters is minimal. See [Reiter \(2003\)](#) for a discussion.

2. Note that the net undercount is smaller than the difference between the false negative and false positive rate, as the denominator of the false positive rate (true nonrecipients) is much larger than the denominator of the false negative rate (true participants).

of SNAP has also found that a substantial share of recipients do not report receipt. For New York state, [Celhay, Meyer and Mittag \(2017\)](#) found false negative rates of 42 and 26 percent in the CPS and American Community Survey (ACS), respectively. [Meyer, Goerge, and Mittag \(2014\)](#) found even higher rates in the same surveys for Illinois (48 and 32 percent) and Maryland (53 and 37 percent). On the other hand, the false positive rates (i.e., reported SNAP receipt that cannot be verified) are low at around 1 percent (e.g., 1 percent for the New York ACS), resulting in the substantial net underreporting of food assistance that is documented in [Meyer, Mok, and Sullivan \(2015a, 2015b\)](#) and [Meyer and Mittag \(forthcoming\)](#).

Survey error and the bias it causes for Medicaid and SNAP is a serious problem for the policy research community and the Federal Statistical system as these survey estimates are used for critical purposes. Medicaid and SNAP are two important noncash benefits provided by states and funded through a federal-state partnership. It is critical for surveys to measure them accurately as benefit recipients are better off than similar nonrecipients, because they have more resources. For example, when calculating the Supplemental Poverty Measure (which includes SNAP benefits but not Medicaid) or related measures, having accurate benefit receipt information is critical to capture the full resources available to a person or family ([US Census 2015](#)). The impact of these benefits on understanding poverty in the United States is large overall and even larger for specific demographic groups ([US Census 2015](#)). The adjustments for noncash benefits often rely on survey responses for SNAP that are known to have significant measurement error and that undercount the participation in these programs (and as a result may overcount poverty). The errors also result in understatement of the poverty reduction of the various programs and the mistakes in the relative importance of the individual programs in poverty reduction.

In addition to measuring poverty, these data are critical for (1) providing general knowledge and statistics on the programs, (2) evaluating these programs to see whether specific policy objectives are met, and (3) aiding official Congressional Budget Office legislative “scoring” to provide cost estimates for critical legislative initiatives such as the Affordable Care Act ([Congressional Budget Office 2007](#)) and simulation models used by federal agencies such as the Urban TRIM model ([Urban Institute 2015](#)). They are also used for official purposes by agencies to develop important health expenditure estimates for the country and states ([Cuckler and Sisko 2013](#)). Given these important uses of the survey data and the evidence that these data have considerable measurement error and bias, improvements to the data could substantially aid policy-making.

In this article, we estimate the magnitude of data quality gains that would be possible if agencies or policy researchers began to routinely use imputations from models that rely on linked survey and administrative data. We model the relationship between “true” and reported receipt status in the linked data (using

only public use file variables as predictors) and use this model to impute a more accurate receipt indicator in the public use data for all respondents. The models are estimated using the confidential linked data so that the parameters from the models can be released to the public. We demonstrate that such methods can substantially increase estimate quality in studies of program receipt, focusing on Medicaid and SNAP.

The methods could also be applied to other government programs important to understanding poverty on which administrative records exist, including Temporary Assistance for Needy Families (TANF), Supplemental Security Income and the Earned Income Tax Credit, as well as other programs such as social security or unemployment insurance. We focus on binary variables here, but the approach also works for continuous variables known to be measured with error (e.g., self-reported height and weight data, housing prices, and employer characteristics). More generally, the methods can be used to amend or improve survey data whenever additional data can be linked to the survey, but making the linked data publicly available is infeasible due to confidentiality concerns or data license restrictions.

## 2. LINKING DATA AS A WAY TO REDUCE MEASUREMENT ERROR IN ESTIMATES

One way to try to ameliorate the potential limitations of any data system is to combine it with other sources of data through linkage ([National Academies of Sciences, Engineering and Medicine 2017a, 2017b](#)). In this study, we do this by combining survey reported data with program administrative data and use the linked data to create models we can use to partially correct for measurement error in survey reported program participation. Previous studies have used linked data to (1) examine sample coverage and accuracy ([Bee, Gathright, and Meyer 2015; Meyer, and Mittag 2017](#)), (2) impute variables ([Davern et al. 2009a](#)), (3) substitute administrative values for reported values ([Nicholas and Wiseman 2010; Hokayem, Bollinger, and Ziliak 2015; Meyer and Mittag forthcoming](#)), (4) supplement survey reported data ([Abowd, Stinson, and Benedetto 2006](#)), and (5) correct estimates ([Davern et al. 2009a; Schenker, Raghunathan, and Bondarenko 2010; Mittag forthcoming](#)). These studies demonstrate that data linkage can improve survey accuracy by illuminating problems in survey design and methodology. Data linkage can also be used to improve survey accuracy more directly by combining information from the linked variables and the survey responses.

In this article, we explore one of the potential benefits of combining administrative data with survey data. Specifically, our central objective is to show the reduction in measurement error that results when a model-based blended imputation model is used to impute values to replace the fallible survey responses. We use reported survey data that have been linked to

external administrative data to estimate the relationship between “true” and reported values of the variable of interest (e.g., being enrolled in Medicaid or receiving SNAP). Then, we use this model to impute improved receipt indicators. Using these imputed values instead of the misreported survey measures yields partially corrected estimates of the statistics of interest. An example of the method is [Schenker et al. \(2010\)](#) who start with a set of data from the National Health and Nutrition Examination Survey (NHANES) that has both the reported survey items and clinically measured items to diagnose hypertension, diabetes, and obesity. They then model the clinically diagnosed values using the survey reported values along with potential covariates of measurement error as predictors. Once the model is developed on the NHANES with both survey-reported and clinically measured variables, they use the model to multiply impute clinical outcomes in data for which they do not have the actual clinical measures in the National Health Interview Survey. For a Bayesian approach, see [He and Zaslavsky \(2009\)](#). [Blackwell, Honaker and King \(2017\)](#) review such (multiple) imputation as a method to correct for measurement error. They propose a similar imputation strategy that relies on independence assumptions for situations in which validation data are not available.

Our analyses build on applications to Medicaid in [Davern et al. \(2009a\)](#) and SNAP in [Mittag \(forthcoming\)](#). [Davern et al. \(2009a\)](#) link Medicaid administrative data to CPS data and then use the model developed on the linked data to impute an indicator of Medicaid enrollment in the CPS for subsequent years. [Mittag \(forthcoming\)](#) employs a similar approach using SNAP administrative data linked to ACS data to correct estimates of food stamp receipt out of sample in the ACS.

We use a mean squared error (MSE) metric to measure the gains in estimate quality that can be realized by using linked data to create blended imputed estimates. The MSE is defined by:

$$\text{MSE} = \text{Bias Squared} + \text{Variance}.$$

The MSE of an estimate is the expected value of the square of its deviation from the true parameter of interest (i.e., it combines estimator bias and variance). Thus, when evaluating the quality of different survey estimates, preference is given to the one with the smaller MSE. In our tables, we take the square root of the MSE or the root mean squared error (RMSE) in order to put the measure on the same scale as the original estimate. It is important to note that our RMSE is itself likely a biased estimate due to the fact that the measure of “truth” we use to compute it contains error as well. We discuss these points in the next section, when we review limitations of our approach.

Table 1 presents results from our application to Medicaid receipt. The first four columns of numbers in table 1 are drawn from Davern et al. (2009a).<sup>3</sup> They used the 2001–2002 CPS linked to Medicaid Statistical Information System (MSIS) data from 2000–2002 to create a person-level, logistic, regression model of Medicaid receipt (for detailed tables on the data linkage and the evaluation, see US Census Bureau 2007, 2008a, and 2008b). Of the CPS respondents linked to MSIS who show Medicaid enrollment in MSIS at some point during the reference period, roughly 43 percent do not report having Medicaid, resulting in a Medicaid undercount (Davern, Klerman, Baugh, Call, and Greenberg 2009b). However, because 47 percent of those linked do correctly report, survey-reported Medicaid enrollment is a critical predictor of enrollment in the imputation model.

Stratifying on survey reported Medicaid status, Davern et al. (2009a) estimated two models to partially correct for survey measurement error. The first model used a logistic regression to predict whether a person received Medicaid in MSIS given that they did not report having Medicaid in the survey (i.e., a false-negative model). The second model predicted whether a person received Medicaid given that they had reported Medicaid coverage in the survey (i.e., a true positive model). Both models condition on key covariates such as age, sex, and income and include state fixed effects. The models also include reported Medicaid enrollment, which is a key predictor of actual enrollment. Fitting a well-specified model is crucial for the accuracy of the correction, but it is a standard specification question that is beyond the scope of this article. Davern et al. (2009a) discuss estimation of the Medicaid models we use here. See Appendix A of the supplementary data online of this article for the estimated model parameters.

We use the coefficients from these two logistic regression models from 2000 and 2001 to predict each person's probability of being enrolled in Medicaid in the 2007 and 2008 CPS (covering calendar years 2006 and 2007) given their self-reported coverage and other key covariates such as age, sex, income, and state of residence. We then use these person-level predicted probabilities to estimate the number of people having Medicaid by state, which are reported in table 1.

The point of this article is to add the evaluation of estimate accuracy in the last four columns of table 1. The bias is estimated as the difference between the state estimate of enrollment in 2006–2007 and the Medicaid enrollment numbers found on Kaiser State Health Facts (Ellis et al. 2008). The Kaiser number is likely imperfect as well, and bias can vary from state to state given that each state has different ways they compile the data for Kaiser. In addition,

3. The standard errors for the imputed Medicaid enrollment estimates in Davern et al. (2009a) were incorrect and did not appropriately adjust for the design effect of the CPS complex sample design. The standard errors in table 1 of this article for imputed Medicaid by state have been adjusted for the design effect of the CPS survey.

Table 1. Percent Change in Root Mean Squared Error of Medicaid Enrollment Estimates from the Blended Imputation Model to the Regular CPS Estimates, by State: Average of Calendar Year 2006 and 2007

State	Medicaid Enrollment Estimate - CPS		Medicaid Enrollment Estimate - Imputed		Kaiser Medicaid Enrollment Estimate <sup>a</sup>	Root Mean Squared Errors (RMSE)		Percent Change from RMSE-CPS to RMSE-Imputed
	Rate	SE	Rate	SE		RMSE-CPS	RMSE-Imputed	
Alabama	11.2%	0.85%	13.9%	0.93%	14.7%	3.58%	1.24%	65.39%
Alaska	7.9%	0.68%	10.3%	0.77%	11.9%	4.06%	1.77%	56.40%
Arizona	15.0%	0.98%	17.5%	1.05%	15.8%	1.26%	2.00%	-58.20%
Arkansas	15.3%	0.93%	17.4%	0.98%	17.9%	2.75%	1.10%	59.91%
California	13.8%	0.35%	16.5%	0.38%	17.7%	3.92%	1.27%	67.69%
Colorado	7.6%	0.50%	8.7%	0.54%	8.0%	0.64%	0.94%	-48.03%
Connecticut	7.9%	0.56%	9.0%	0.59%	11.4%	3.59%	2.48%	30.98%
Delaware	10.0%	0.74%	13.7%	0.85%	16.9%	6.98%	3.37%	51.71%
District of Columbia	18.5%	1.09%	20.5%	1.14%	21.8%	3.53%	1.79%	49.38%
Florida	8.3%	0.40%	11.7%	0.47%	11.6%	3.35%	0.47%	85.91%
Georgia	9.8%	0.59%	12.9%	0.67%	13.4%	3.70%	0.86%	76.67%
Hawaii	9.6%	0.65%	12.6%	0.73%	14.5%	4.96%	2.09%	57.82%
Idaho	9.9%	0.78%	10.9%	0.82%	11.3%	1.56%	0.90%	42.19%
Illinois	10.3%	0.54%	13.4%	0.60%	15.3%	5.00%	1.99%	60.20%
Indiana	10.3%	0.71%	12.5%	0.78%	12.5%	2.33%	0.78%	66.57%
Iowa	11.0%	0.67%	12.2%	0.70%	10.7%	0.71%	1.59%	-125.22%
Kansas	8.5%	0.69%	10.8%	0.77%	9.1%	0.94%	1.87%	-98.70%
Kentucky	13.6%	0.84%	14.7%	0.86%	16.7%	3.17%	2.21%	30.35%
Louisiana	12.8%	1.00%	15.6%	1.08%	20.4%	7.63%	4.85%	36.38%
Maine	18.2%	0.84%	21.6%	0.89%	19.7%	1.75%	2.10%	-20.05%

Maryland	7.0%	0.51%	8.3%	0.55%	9.5%	2.53%	1.35%	46.57%
Massachusetts	14.7%	0.86%	13.9%	0.84%	16.1%	1.60%	2.32%	-44.66%
Michigan	11.9%	0.64%	12.7%	0.65%	15.1%	3.24%	2.47%	23.95%
Minnesota	10.3%	0.60%	12.2%	0.65%	11.3%	1.20%	1.06%	12.09%
Mississippi	16.7%	1.13%	16.5%	1.12%	18.0%	1.68%	1.84%	-9.40%
Missouri	11.5%	0.72%	15.8%	0.83%	12.5%	1.23%	3.43%	-178.24%
Montana	10.7%	0.90%	6.6%	0.72%	9.2%	1.78%	2.69%	-51.01%
Nebraska	7.8%	0.66%	11.6%	0.79%	10.0%	2.23%	1.81%	18.94%
Nevada	5.2%	0.54%	7.1%	0.62%	6.8%	1.71%	0.68%	60.39%
New Hampshire	5.6%	0.42%	7.3%	0.48%	8.3%	2.77%	1.16%	58.11%
New Jersey	7.4%	0.52%	8.6%	0.56%	8.8%	1.52%	0.59%	61.33%
New Mexico	14.7%	1.01%	18.1%	1.09%	20.3%	5.66%	2.46%	56.52%
New York	15.6%	0.53%	16.2%	0.54%	21.6%	6.02%	5.36%	10.92%
North Carolina	11.9%	0.66%	16.8%	0.76%	13.3%	1.47%	3.61%	-144.80%
North Dakota	8.0%	0.68%	10.3%	0.76%	8.4%	0.76%	2.11%	-178.57%
Ohio	12.0%	0.63%	13.5%	0.66%	14.1%	2.20%	0.93%	57.91%
Oklahoma	12.3%	0.84%	15.5%	0.92%	14.7%	2.59%	1.20%	53.56%
Oregon	10.0%	0.74%	11.8%	0.80%	9.1%	1.13%	2.79%	-146.14%
Pennsylvania	9.3%	0.51%	13.3%	0.60%	15.3%	6.03%	2.10%	65.12%
Rhode Island	17.1%	0.90%	16.7%	0.89%	15.7%	1.73%	1.35%	22.11%
South Carolina	13.2%	0.87%	16.8%	0.96%	14.6%	1.60%	2.40%	-49.70%
South Dakota	8.8%	0.70%	9.8%	0.73%	11.5%	2.84%	1.88%	33.65%
Tennessee	14.1%	0.95%	22.0%	1.13%	20.6%	6.54%	1.79%	72.60%
Texas	10.9%	0.38%	13.2%	0.42%	12.0%	1.23%	1.20%	3.09%
Utah	8.0%	0.77%	9.8%	0.84%	7.2%	1.08%	2.71%	-151.51%
Vermont	17.2%	0.93%	20.9%	1.00%	19.2%	2.21%	1.95%	11.74%

Continued



Table 1. Continued

State	Medicaid Enrollment Estimate - CPS		Medicaid Enrollment Estimate - Imputed		Kaiser Medicaid Enrollment Estimate <sup>a</sup>		Root Mean Squared Errors (RMSE)		Percent Change from RMSE-CPS to RMSE-Imputed
	Rate	SE	Rate	SE	Rate		RMSE-CPS	RMSE-Imputed	
Virginia	7.1%	0.52%	8.1%	0.55%	8.5%		1.46%	0.66%	54.60%
Washington	11.1%	0.69%	15.1%	0.79%	13.4%		2.36%	1.89%	20.06%
West Virginia	14.0%	0.91%	16.3%	0.97%	16.7%		2.88%	1.06%	63.38%
Wisconsin	11.5%	0.79%	12.0%	0.80%	12.2%		1.04%	0.83%	20.68%
Wyoming	7.5%	0.70%	9.0%	0.76%	10.9%		3.40%	1.99%	41.61%
<b>Total: United States</b>	11.4%	0.11%	13.8%	0.12%	14.3%		2.89%	0.54%	81.29%

SOURCE: 2007 and 2008 CPS ASEC data files.

<sup>a</sup>Independent Medicaid Enrollment Estimate derived from a two year average of the December 2006 and 2007 total Medicaid enrollment for each state (see Ellis et al. 2008) to match the 2007 and 2008 CPS estimate which corresponds to Calendar year 2006 and 2007.

concept alignment between the CPS measure and the Kaiser measure is not perfect: the Kaiser measure is an average monthly enrollment, and the CPS is a measure of Medicaid enrollment at any point in the last year. Thus, the CPS number should be higher and include more enrollees who churn on and off the program throughout the year. In general, this would mean the administrative data counts of Medicaid enrollment should be even higher than the Kaiser counts for the any time in the past year enrollment number. Finally, universes between CPS and Kaiser are not the same. Kaiser includes people in group quarters who may have died during the year who would not be counted in CPS. The impacts of these universe adjustments are important but will not significantly impact the findings of the article (see [US Census Bureau 2008a,b](#) to better understand the magnitude). Nevertheless, the Kaiser numbers are an independent estimate of enrollment in those years for comparison purposes. The first column of RMSEs are for the unadjusted CPS (i.e., what one would get if one simply tabulated the CPS public use file for those two years and created a two-year average). The second column of RMSEs compares the Kaiser rate to the CPS blended imputations on the individual-level predicted probabilities. The final column is the percent reduction (negative numbers are the percent increase) between the two RMSEs for any given state. The Medicaid example does not account for variance added by imputation so that the MSEs for the imputed model are too small.<sup>4</sup> However, the additional variance due to imputation modeling will likely be small relative to the reduction in bias (as we find in the SNAP example later on which adjusted for model variance).

For the United States as a whole, the RMSE for the model-based blended imputed estimate is 81 percent lower than the RMSE for the direct CPS estimate. This is a substantial reduction in RMSE, which is mainly due to bias being reduced. The direct CPS estimate of the Medicaid coverage rate for the United States is 11.4 percent, and the imputed estimate is 13.8 percent, which is much closer to the 14.3 percent in the Kaiser State Health Facts. In most states, the RMSE decreased between the CPS direct survey estimate and imputed estimate. There are, however, fourteen states that saw an increase in bias. One would expect estimates from some states to be closer to truth by chance due to the sampling error in the CPS. Due to the extrapolation, we cannot hold sampling error fixed here. The fact that RMSE increases for a smaller fraction of areas when we hold sampling error fixed in the second application suggests that this explains some but not all of these increases. Increases in RMSE may also stem from unobserved differences in reporting rates, which may lead us to overstate receipt in states with the most accurate reporting. The largest increases were in Utah, Arizona, North Dakota, North Carolina, Oregon, Iowa, Montana, Kansas, and Missouri. These are states in which the CPS fares particularly well: in these nine states, the difference between the CPS direct survey estimate and the Kaiser rate was only 0.8 percent on

4. We no longer have access to the Medicaid microdata to correct the variance estimates.

average, whereas in the thirty-seven states with positive reductions in RMSE, the average difference was 3.2 percent. When known, differences in reporting rates can be incorporated in the imputation model, making it desirable for future research to look for potential reasons and to attempt to improve on the fit of the model for these states.<sup>5</sup>

Our second illustration of how blending based on models from linking data can improve survey estimates examines SNAP receipt for small geographic areas<sup>6</sup> in New York state. The results in [table 2](#) are similar to those for Medicaid in [table 1](#). They are based on the model and results in [Mittag \(forthcoming\)](#), which uses administrative SNAP records linked to the ACS to develop a method of correcting survey estimates for measurement error. The validation data were created by linking administrative records on monthly SNAP payments for all recipients in New York state from the New York State Office of Temporary and Disability Assistance (OTDA) to the 2010 ACS survey data. The administrative records are based on actual, validated receipt, and the two data sources are linked with a high match rate at the household level. Thus, even though they are not free of error, the linked data appear accurate enough that we consider them to be the assumed best or unbiased measure of receipt. For further descriptions of data linkage and accuracy, see [Celhay, Meyer and Mittag \(2017\)](#), [Cerf Harris \(2014\)](#), [Mittag \(forthcoming\)](#), and [Scherpf, Newman, and Prell \(2014\)](#). As [Celhay, Meyer and Mittag \(2017\)](#) show, the linked data reveal substantial error in reported SNAP receipt and amounts. For example, 26 percent of administrative data recipient households do not report SNAP receipt in the ACS (false negatives). On the other hand, the false-positive rate (true nonrecipients reporting SNAP receipt) is low at 1.2 percent, resulting in the substantial net underreporting of government transfers that is documented in [Meyer et al. \(2015a, 2015b\)](#) and [Meyer and Mittag \(forthcoming\)](#).

The fifth column of [table 2](#) provides estimates of receipt rates that we consider to be unbiased from the linked data. We estimate these rates for the 39 county groups that can be identified in the ACS public use data. Comparing these receipt rates to the survey-based estimates in the first two columns underlines that there is net underreporting in all but one area and that reporting rates vary between these areas. [Cerf Harris \(2014\)](#) examines reporting rates at the county level in detail.

The main objective of this article is to assess how the survey estimates compare with the results in columns three and four, which contain estimates of the receipt rate using an imputation model to create blended imputed estimates that

5. For the state of Montana, the increase in the bias in the modeled results derives from the fact that over half of those on Medicaid were missing the linking information. Thus in Montana's case, too few people are imputed to have Medicaid as over half the enrollees were not linkable to the CPS ([US Census Bureau 2008a](#)).

6. We use the counties that can be identified in the public use ACS data and pool counties that cannot be separated in the public use data.

Table 2. Percent Change in Root Mean Squared Error of SNAP Receipt Estimates from the Blended Imputation Model to the Regular ACS Estimates, by New York State County or County Group: Average of Calendar Year 2010

Counties	SNAP receipt eate		SNAP Receipt Rate		Linked data	Root mean squared errors (RMSE)		Percent change from RMSE-ACS to RMSE-imputed
	Estimate - ACS		estimate - Imputed			RMSE-ACS RMSE-Imputed		
	Rate	SE	Rate	SE		Rate	SE	
Albany	11.6%	1.49%	15.6%	1.66%	14.6%	3.3%	1.9%	40.85%
Allegany, Cattaraugus	14.6%	1.83%	18.4%	1.97%	19.3%	5.0%	2.2%	56.95%
Bronx	41.1%	1.02%	47.7%	1.04%	52.1%	11.0%	4.5%	58.94%
Broome, Tioga	15.3%	1.59%	19.3%	1.80%	18.8%	3.9%	1.9%	51.69%
Cayuga, Madison, Onondaga	12.6%	0.92%	16.6%	1.03%	16.4%	4.0%	1.0%	73.98%
Chautauqua	15.9%	1.94%	19.2%	2.18%	21.5%	5.9%	3.2%	46.48%
Chemung, Schuyler	17.1%	2.35%	19.4%	2.47%	21.7%	5.2%	3.4%	35.55%
Chenango, Cortland	18.2%	2.28%	20.6%	2.47%	19.6%	2.7%	2.7%	1.77%
Clinton, Essex, Franklin, Hamilton	16.9%	2.05%	19.7%	2.17%	19.2%	3.1%	2.2%	27.05%
Columbia, Greene	9.4%	2.26%	14.5%	2.72%	10.2%	2.4%	5.0%	-107.30%
Delaware, Otsego, Schoharie	10.8%	1.61%	16.0%	2.06%	15.1%	4.6%	2.3%	50.68%
Dutchess	8.5%	1.34%	12.4%	1.62%	11.3%	3.1%	2.0%	36.72%
Erie	17.4%	0.95%	20.5%	1.00%	20.9%	3.6%	1.1%	70.42%
Fulton, Montgomery	14.6%	2.07%	19.5%	2.63%	28.3%	13.8%	9.2%	33.37%
Genesee, Orleans	11.3%	2.26%	15.6%	2.51%	16.4%	5.6%	2.7%	52.74%
Herkimer, Oneida	16.6%	1.55%	19.6%	1.65%	21.8%	5.4%	2.8%	48.86%
Jefferson, Lewis	19.1%	2.39%	22.0%	2.45%	19.5%	2.4%	3.4%	-41.32%
Kings (Brooklyn)	26.1%	0.68%	32.8%	0.73%	35.9%	9.9%	3.2%	67.74%
Livingston, Wyoming	12.2%	2.42%	16.1%	2.64%	13.8%	2.9%	3.5%	-19.51%
Monroe, Wayne	14.2%	0.82%	18.4%	0.92%	18.1%	4.0%	1.0%	75.95%

Continued

Table 2. Continued

Counties	SNAP receipt eate Estimate - ACS		SNAP Receipt Rate estimate - Imputed		Linked data	Root mean squared errors (RMSE)		Percent change from RMSE-ACS to RMSE-imputed
	Rate	SE	Rate	SE		RMSE-ACS	RMSE-Imputed	
Nassau	4.2%	0.42%	9.5%	0.67%	7.0%	2.8%	2.6%	9.78%
New York (Manhattan)	16.5%	0.69%	20.3%	0.73%	21.2%	4.7%	1.2%	75.07%
Niagara	15.9%	1.79%	18.7%	1.90%	19.4%	3.9%	2.0%	48.05%
Ontario	8.6%	2.78%	13.0%	2.98%	10.6%	3.4%	3.8%	-12.09%
Orange	11.1%	1.44%	17.2%	1.85%	11.0%	1.4%	6.4%	-342.91%
Oswego	18.1%	2.63%	20.9%	2.71%	22.1%	4.8%	3.0%	37.36%
Putnam, Westchester	6.0%	0.62%	10.6%	0.81%	9.4%	3.4%	1.5%	56.78%
Queens	17.2%	0.65%	24.4%	0.74%	23.9%	6.7%	0.9%	86.91%
Rensselaer	14.7%	2.32%	17.5%	2.42%	17.1%	3.3%	2.4%	26.25%
Richmond (Staten Island)	11.0%	1.21%	15.5%	1.41%	16.9%	6.0%	1.9%	67.51%
Rockland	13.8%	1.59%	18.1%	1.78%	15.1%	2.0%	3.5%	-68.77%
Saratoga	8.0%	1.65%	10.6%	1.76%	9.3%	2.1%	2.1%	-2.12%
Schenectady	10.2%	1.93%	13.9%	2.16%	17.3%	7.3%	4.0%	45.77%
Seneca, Tompkins	11.0%	2.70%	15.5%	2.91%	13.4%	3.6%	3.6%	-0.68%
St. Lawrence	16.4%	3.41%	20.3%	3.44%	21.5%	6.2%	3.6%	41.12%
Steuben, Yates	10.7%	1.72%	14.1%	1.90%	18.6%	8.1%	4.9%	39.79%
Suffolk	5.6%	0.53%	10.7%	0.73%	9.1%	3.6%	1.7%	51.93%
Sullivan, Ulster	13.1%	1.70%	17.4%	1.93%	18.1%	5.2%	2.0%	60.84%
Warren, Washington	11.1%	1.78%	13.8%	1.95%	15.9%	5.1%	2.8%	44.51%
Total: New York State	16.1%	0.20%	21.1%	0.22%	21.4%	5.3%	0.4%	92.86%

NOTE.—Source is the 2010 ACS. The measure in the first two columns and the parameters of the imputation model are from NY OTDA administrative data linked to the 2010 ACS.

partially correct the survey reports. The imputations are based on the method in [Mittag \(forthcoming\)](#), who uses the linked 2010 ACS data to estimate the conditional distribution of administrative SNAP receipt and amounts received given reported receipt and other covariates. The conditional distribution of SNAP amounts can be seen as a continuous distribution with a mass point at zero. However, we are only concerned with receipt and not with amounts received here. Therefore, we only use the estimate of the binary part of the distribution, which is a standard probit model. In addition to reported SNAP receipt, the model conditions on a large set of demographic and economic variables, including household composition, age, education, and income. The model does not condition on any geographic information, so that the variation between counties we examine here is only captured by the covariates. This makes the reduction in RMSE particularly noteworthy because accuracy could still be improved by incorporating geographic information. As stated previously, specification of the imputation model is crucial for the accuracy of the correction but beyond the scope of this article. It is discussed further in [Mittag \(forthcoming\)](#). Appendix Table A2 of the [supplementary data](#) online contains the estimated parameters of the conditional distribution.

We use the parameters of this model to predict a probability of SNAP receipt for each household as done with Medicaid. We then generate a receipt variable by taking twenty random draws from a Bernoulli distribution with the predicted probability for every household in the 2010 New York ACS sample. Taking multiple draws makes simulation error negligible and thus reduces SEs and avoids having to correct the SEs for simulation error.<sup>7</sup>

The last three columns of [table 2](#) contain RMSE defined the same way as for Medicaid above. We compute the bias in the survey and imputation-based estimates as the difference in the numbers from the linked data in the fifth column. The population totals to which we compare our estimates are affected by errors in the administrative data and linkage errors. However, these errors should be small and outweighed by the benefit that using the linked data ensures that the numbers are for the same population as our improved survey estimates (which exclude group quarters and the homeless) and that subject definitions are comparable. Contrary to the Medicaid application, the imputation model is estimated using the same sample. [Mittag \(forthcoming\)](#) further discusses extrapolation across time and geography. We are mainly interested in the percentage reduction in RMSE when replacing the survey reports by the imputations in the last column (i.e., by how much the imputations reduce error

7. A key difference between multiple imputation and the approach we take here is that we estimate the statistic of interest from the multiple stacked imputations rather than averaging estimates from repeated single imputations. For the subgroup means we estimate here, the two approaches are equivalent, but estimates and SEs differ in general. As discussed in [Mittag \(forthcoming\)](#), correlations and model parameters as in [Schenker et al. \(2010\)](#) may be inconsistent under single and standard multiple imputation, but the methods discussed here yields consistent estimates.

compared with uncorrected survey based estimates). The numbers for the entire state of New York in the last row show that the blended imputed estimates reduce RMSE by an impressive 93 percent. This is similar in magnitude to the reduction in RMSE for Medicaid and again driven by the reduction in bias. The standard errors are slightly larger than in the survey, but they are small in both cases due to the large sample. Thus, bias is the main determinant of RMSE. The reduction in bias more than makes up for the increase in standard errors. The survey understates receipt by 25 percent, while the imputations fall short of the actual share of recipients by 1 percent only.

This pattern also drives the results at the local level. The survey numbers underestimate receipt rates in all but one county, while the imputation-based numbers do not seem to be systematically biased. They are larger than the assumed “true” numbers in twenty-one out of thirty-nine areas and smaller in eighteen areas. The imputation-based rates are more accurate than the survey in terms of estimated RMSE in thirty-one out of thirty-nine areas. The reductions in RMSE are substantial: in twenty-nine of these thirty-one areas, RMSE is reduced by 25 percent or more, and in fifteen areas, the imputation-based measure cuts the error by more than half. However, RMSE of the imputed receipt rate is larger than the survey RMSE in eight of the thirty-nine areas. As with Medicaid, this result is primarily due to the fact that the survey closely replicates the numbers from the linked data for these eight areas (i.e., it is mainly driven by the good performance of the survey in these counties).

### 3. DISCUSSION

Recent federal data initiatives emphasize linking and combining data as a promising way to improve data for policy purposes. For example, key recommendations in the report of the [Commission on Evidence-Based Policymaking \(2017\)](#) call for producing higher-quality data by linking and combining data. After the commission report was released, the US Office of Management and Budget (OMB) put out a request for information to help improve federal statistics stating that “a priority has been placed on using new techniques and methodologies based on combining data from multiple sources” ([Federal Register 2018](#)). We believe our article demonstrates an operationally efficient way of accomplishing the goal of combining data from multiple sources to improve data quality and ensure data can be widely disseminated for evidence-based policy-making. This section first briefly discusses why data linkage is a cost-effective way to reduce MSE in surveys compared with other common approaches. We then illustrate how model-based blended imputations compare with two key alternatives: the status quo and directly replacing survey reports with the linked administrative values. We compare the strengths and weaknesses of these three approaches using the data quality criteria of the [FCSM \(2001\)](#).

### 3.1. Comparisons to Other MSE Reduction Approaches

Model-based blended imputation may not always yield the largest feasible error reduction, but we argue that it offers a more cost-effective way to lower the MSE of survey statistics than other commonly employed approaches on which large amounts of money are spent. Survey researchers often use tools such as larger sample sizes and/or reducing nonresponse as ways to reduce MSE in surveys. Larger sample sizes reduce MSE, but this option is both expensive and grows less effective at reducing variance (and MSE) with each additional case that is added to the sample (and it also adds respondent burden as more sample is added). Another common strategy to reduce bias is to reduce survey nonresponse. Reducing survey nonresponse through additional effort (more telephone calls, more in-person attempts to recruit a household, more mailings, etc.) and the use of incentives are costly, and there is little evidence they improve data quality. Research has shown that spending considerable funds on strategies aimed at increasing response rates can indeed increase response rates. However, survey research is concerned with response bias and not response rates per se. But nonresponse has been shown to have little impact on the bias survey estimates (Groves 2006; Groves and Peytcheva 2008). Also, linkages to administrative data have demonstrated that nonresponse bias is small for key policy-relevant variables such as income (Bee et al. 2015; Meyer and Mittag (2017)). However, as we show in our article, these same linkage studies often show significant amounts of measurement error in survey responses that often lead to sizeable bias in survey estimates. Thus, we believe that expensive attempts to reduce MSE (such as increasing sample sizes or increasing effort to convert nonrespondents) should be evaluated to make sure that they are cost-effective ways of reducing MSE relative to other alternatives such as data linkage.<sup>8</sup>

### 3.2. Data Quality of Model-Based Blended Imputations Estimates Versus Alternatives

To compare model-based blended imputations with key alternative approaches on criteria that are relevant to statistical agencies, we use the data quality framework developed in FCSM (2001). The FCSM identifies the four key elements of data quality as (1) accuracy, (2) relevance, (3) timeliness, and (4) accessibility. The two alternatives we explore are, first, not making any changes (i.e., having the agencies maintain current practice), and second, direct substitution (i.e., having the agencies link the data and directly replace the survey

8. Data linkage to administrative data can also facilitate other survey improvements besides reducing measurement error. For example, there is strong evidence that linking the sample frame to other sources of data can help surveys more efficiently allocate resources used in household listing (Montaquila et al. 2011).



report by the actual value from the administrative data—rather than imputing a value as was done in this study).

If the statistical agencies do not make any changes from current practice, we believe that data quality will continue to be a problem. A mounting literature demonstrates that the current approach is not *accurate*, by showing that estimates of key policy importance are biased by the substantial amount of measurement error (e.g., Davern et al. 2009a; Meyer and Mittag forthcoming). The results in this article show that the blended imputation approach would substantially improve estimate *accuracy* over the current practice. In addition to *accuracy*, the imputed estimates could also bring gains in *relevancy* in that the survey could be enhanced with additional information from the administrative data. For example, one could use the blended imputations to develop monthly enrollment flags instead of indicators of ever being enrolled during a 12-month reference period. Such an additional programmatic detail has the potential to improve the policy *relevance* of the blended imputations for policy research purposes.

On the other hand, the blended imputation approach cannot improve over current practice in terms of *timeliness* and *accessibility*. Our approach may slightly reduce timeliness if creating the imputations delays data release or if the imputations are produced after data release. In order to mitigate the impact *timeliness* and *accessibility*, the model coefficients could be created by the statistical agency themselves (similar to what appears in Appendix A of the [supplementary data](#) online) and distributed separately so as not to interrupt current data processing. Another approach could be to have the statistical agency grant access to the linked data to a third party using the Research Data Center (RDC) network. Interested third parties could include those working on micro simulation models that rely on the survey data such as the Urban Institute's TRIM (Urban Institute 2015), Congressional Budget Office simulation models (Congressional Budget Office 2007), or RAND's COMPARE (Eibner, Girosi, Price, Cordova, Hussey et al. 2010) simulation models, in addition to groups that disseminate the survey microdata such as Integrated Public Use Microdata Series (IPUMS), Inter-university Consortium for Political and Social Research (ICPSR), or National Bureau of Economic Research (NBER). If granted access, these third parties could estimate imputation model coefficients and/or the imputations themselves and distribute them through the current dissemination channels such as IPUMS, ICPSR, or NBER. This would impinge on *accessibility* of the imputed data because the imputations are not provided with the core data product, but this downside could be mitigated if data dissemination channels such as IPUMS, ICPSR, or NBER would include the imputed values in the version of the data they distribute.

The second alternative approach to model-based blended imputation is direct substitution of administrative data for survey data. Direct substitution is more *accurate*. In terms of *relevance*, direct substitution is better on some dimensions but worse on others. And finally, the direct substitution approach

is likely to be less *accessible* and *timely* than blended imputation. To elaborate on these points, the blended imputation approach is less *accurate* than the direct substitution approach for two main reasons. First, the model-based blending approach adds variance from estimated parameters of the model and imputation. The added variance from estimated parameters decreases the effective sample size of the linked data.<sup>9</sup> In addition to variance, the direct substitution method is more *accurate* because it does not rely on a potentially misspecified model. Most specification questions can be assessed with standard tests (see [Davern et al. 2009a](#) and [Mittag forthcoming](#) for discussions of the models we use here). For imputation models, particular attention should be paid to the choice of conditioning variables. As discussed in [Hirsch and Schumacher \(2004\)](#) and [Bollinger and Hirsch \(2006\)](#), the imputation model should condition on all covariates in the outcome model. The goal of the blended imputations is to reproduce the distribution of the accurately measured variable or its joint distribution with the relevant covariates. Researchers developing the imputation models have access to the linked data; so how closely a given model reproduces these distributions can be measured and tested using Kolmogorov- or Cramer-von Mises-type statistics. Most household surveys are used for a wide range of purposes, so the ideal imputation model may depend on the statistic of interest. This can be addressed by producing different models for different purposes, and this again presents a slight downside in terms of convenience compared with direct substitution. Also, as programs and reporting errors may change over time, the imputation models should constantly be evaluated and improved. And it is likely that one model may not be appropriate for all use cases, and the development of additional models for specific use cases is recommended.

When comparing the policy *relevance* of the direct substitution method to the model-based imputation method, there are pluses and minuses for each. For direct substitution (as with the model-based imputation method), relevant details from the administrative data can be included other than just enrollment or receipt. These details include the months the person was enrolled, the basis of eligibility, the exact program of enrollment (e.g., State Children's Health Insurance Program or limited benefits Medicaid program versus a full benefits Medicaid program), the services received, and the amount of benefits (among many others). The advantage of using direct substitution for these extra

9. Standard errors need to be adjusted for this additional variance, which makes the model-based imputation approach less convenient. Methods to do so are well developed but depend on details of the implementation. If one uses the imputation model to create multiple imputations or synthetic data, SEs can be estimated as discussed in [Rubin 1996](#), [Raghunathan, Reiter and Rubin \(2003\)](#) and [Reiter \(2003\)](#). When using the imputation model to integrate out the error ridden measure as in [Mittag \(forthcoming\)](#), SEs can be corrected for simulation error as discussed in [McFadden \(1989\)](#) and for estimated first stage parameters as described in [Newey and McFadden \(1994\)](#). If one is willing to specify prior distributions, Bayesian survey inference provides a compelling way to measure uncertainty, see [Little \(2012\)](#) for a discussion.

policy-relevant details is that this information is measured more *accurately* than in the blended imputation approach.

The policy *relevance* advantage for the blended imputation approach is that high-quality, administrative data are often not available for some geographic areas or time periods or some households cannot be linked (e.g., survey respondents opt out of linkage, or survey/administrative data are missing identifiers used for linkage). In such cases, model-based blended imputation can use geographic areas and time periods with linked data to develop models and then extrapolate. This approach, although susceptible to model variance and misspecification, can still lead to significant reductions in MSE. Mittag (forthcoming) discusses the required conditions and finds substantial improvements in accuracy even though the assumptions are at best approximations in his application. In the Medicaid empirical example, the model was created using MSIS data linked to 2000–2002 CPS data and was applied to microdata from the 2007–2008 CPS. There is likely to be some extrapolation error in this case since several states experienced changes in their Medicaid program over this time span. However, as we showed in the analysis presented in this article, the reductions in MSE are substantial nonetheless. A final benefit of the blended imputations is that the imputation model can be extended to impute true receipt for households that the agency is unable to link (e.g., the respondent opts out of linkage). Such extensions could also address the consequences or linkage errors (incorrect or incomplete linking identifiers on the survey or administrative data). Linking data on a regular basis will improve our understanding of the conditions under which extrapolation works and thereby help to validate and improve the imputation models and the policy *relevance* of the data that result. This can make the blended imputations more policy-relevant than direct substitution in cases where linkage is not possible or imperfect.

The possibility to extrapolate also gives the blended approach an advantage over direct substitution in terms of *timeliness*. The administrative data needed for direct substitution may sometimes not be ready for linking in a timely manner to allow for direct substitution. Using previous years of linked data for modeling, while potentially less accurate, will result in the production of more *timely* estimates for use in policy research.

The final and likely most important advantage of the blended imputation over direct substitution is *accessibility*. As pointed out by Bound, Brown, and Mathiowetz (2001), linked data or validation data are usually only accessible to a small group of researchers. The main reason for this situation is that making the linked data publicly available makes it easier to identify individuals and carries a confidentiality risk for both the survey and the administrative data. In the past, this risk has been deemed to be too large to allow for the release of directly substituted data. This situation may change in the future, possibly by adding noise, coarsening the linked variables, or creating an infrastructure for secure data access. But neither the model coefficients nor the model-based imputed values present as much de-identification risk as direct

substitution, as long as there is an imperfect model fit, although models also carry disclosure risk (Reiter and Mitra 2009). As Appendix A of the [supplementary data](#) online shows, statistical agencies are willing to release the required parameters, so the blended imputations approach is already feasible. Thus, the blended imputation approach has an advantage for public *accessibility* in that the model error may be large enough to pass a data-producing organization's confidentiality review.

## 4. CONCLUSION

All data (including survey and administrative data) have errors. However, it is critical that we move beyond acknowledging data limitations and create new data products that blend the strengths of each data system to reduce known errors. Such innovative methods to mitigate the flaws in any one data system have the potential to improve public policy decisions. From our two analyses of Medicaid and Food Stamps, we argue that in the realm of survey errors that we (1) can address and (2) have a measurable impact on data quality, reducing measurement error through linkage of administrative data to survey data is a way to achieve substantial MSE reductions. Current practice does not incorporate the results from linkage studies into the most widely used and circulated data products from data producers such as the US Census Bureau.<sup>10</sup> We believe that they can and should do more to correct for known survey measurement error. At a minimum, data producers (potentially in collaboration with the broader research community) should create alternatives to their standard data products that are known to have pronounced measurement error in policy-relevant variables.

We have demonstrated one approach for creating products that allows analysts to partially correct known measurement errors. The examples of Medicaid and SNAP receipt underline that the resulting improvements can be substantial as they reduced RMSE by 81 and 93 percent compared with the survey estimates for the geographic areas we examined. The model-based blended imputation approach has been found to work well for a wide range of use cases. It extends to multivariate analyses and more complex estimators. [Schenker et al. \(2010\)](#) and [Mittag \(forthcoming\)](#) impute both binary and continuous variables and find the blended imputation to work well for multivariate and nonlinear models.

We use the [FCSM \(2001\)](#) elements of data quality (accuracy, relevance, timeliness, and access) to evaluate the blended imputations. We provide evidence that a key advantage of model-based imputation is its improvement of *accuracy* compared with current practice. We argue that the improvement in *accuracy* is large enough to outweigh the disadvantages in *timeliness* and

10. Although we note a good recent example is in [Motro and Roth \(2017\)](#).

*access*. Direct substitution would be more accurate than the blended imputations and have similar *relevance*. If linked data can be made *accessible* to a wider audience in a *timely* fashion, then advantages of direct substitution may easily make this approach preferable. However, *access* to directly substituted data may not be able to be made public, making this route more difficult. And given the current state of affairs the blended imputation appears preferable on grounds of *accessibility* and *timeliness* despite the loss of *accuracy*. The method we propose does not pose as great a risk to data confidentiality and the privacy of respondents, nor does sharing it publicly violate the terms of some of data sharing agreements between agencies.

Additional data products that improve accuracy through model-based imputations could be created based on existing data linkage projects. These additions to current data products could consist of (1) a set of models complete with coefficients like the ones we generated in this article for Medicaid and SNAP so that users of the data could use them to create imputations themselves or (2) a separate imputed variable for all survey persons/households using models like the ones we have used that is included in future data products. Little (2012) discusses the advantages of these two options further.

The reasons why it is now imperative to use linked data in the creation of official statistics, reports, and data products are that (1) the foundational research for use of linked administrative data and survey data has been conducted for several potential sources, (2) there is clear evidence from these research projects that the amount of bias due to measurement error in the survey data could be significantly reduced, and (3) the necessary infrastructure for sharing data among federal agencies is in place, and directives have been supplied by the Office of Management and Budget (Burwell 2014; O'Hara 2016). Now is the time to start building the data products that use blended survey and administrative data in production as it will improve official statistics, reports, and data products. While not all linked administrative data and survey data are ready for production, we believe that there are substantive areas of policy research (such as Medicaid enrollment, Medicare enrollment, SNAP and other program receipt, and uninsurance calculations) that have the needed agreements in place and ongoing linkage projects. These projects can be leveraged to improve our ability to make policy-relevant estimates to evaluate and cost out policy proposals for use by organizations such as the Congressional Budget Office, and the Office of the Actuary at the Centers for Medicare & Medicaid Services.

## Supplementary Materials

Supplementary materials are available online at [academic.oup.com/jssam](https://academic.oup.com/jssam).

## REFERENCES

- Abowd, J. M., M. Stinson, and G. Benedetto (2006). "Final Report to the Social Security Administration on the SIPP/SSA/IRS Public Use File Project," US Census Bureau unpublished paper.
- Bee, C. A., G. Gathright, and B. D. Meyer (2015). "Bias from Unit Nonresponse in the Measurement of Income in Household Surveys," University of Chicago unpublished paper.
- Blackwell, M., J. Honaker, and G. King (2017). "A Unified Approach to Measurement Error and Missing Data: Overview and Applications," *Sociological Methods & Research*, 46(3), 303–341.
- Bollinger, C. R., and B. T. Hirsch (2006). "Match Bias from Earnings Imputation in the Current Population Survey: The Case of Imperfect Matching," *Journal of Labor Economics*, 24, 483–519.
- Bound, J., C. Brown, and N. Mathiowetz (2001). "Measurement error in survey data." In *Handbook of Econometrics*. Vol. 5, eds. James J. Heckman and Edward Leamer, Chapter 59, 3705–3843. Amsterdam: Elsevier.
- Burwell, S. (2014). "M-14-06: Memorandum for the Heads of Executive Departments and Agencies: Guidance for Providing and Using Administrative Data for Statistical Purposes," Office of Management and Budget, available at <https://obamawhitehouse.archives.gov/sites/default/files/omb/memoranda/2014/m-14-06.pdf> (last accessed August 15, 2018).
- Congressional Budget Office (2007). "Background Paper: CBO's Health Insurance Simulation Model a Technical Description," Congressional Budget Office, October 2007. Washington DC, available at <https://www.cbo.gov/publication/19224?index=8712> (last accessed August 15, 2018).
- Celhay, P., B. D. Meyer, and N. Mittag (2017). "Errors in Reporting and Imputation of Government Benefits and Their Implications," unpublished paper.
- Cerf Harris, B. (2014). "Within and Across County Variation in SNAP Misreporting: Evidence from Linked ACS and Administrative Records," CARRA unpublished paper #2014-05. US Census Bureau.
- Commission on Evidence-Based Policymaking (2017). "The Promise of Evidence-Based Policy," available at <https://www.cep.gov/content/dam/cep/report/cep-final-report.pdf> (last accessed August 15, 2018).
- Cuckler, G., and A. Sisko (2013). "Modeling per Capita State Health Expenditure Variation: State-Level Characteristics Matter," *Medicare and Medicaid Research and Review*, 3, E1–E24.
- Davern, M., J. A. Klerman, J. Ziegenfuss, V. Lynch, and G. Greenberg (2009a). "A Partially Corrected Estimate of Medicaid Enrollment and Uninsurance: Results from an Imputational Model Developed off Linked Survey and Administrative Data," *Journal of Economic and Social Measurement*, 34, 219–240.
- Davern, M., J. A. Klerman, D. Baugh, K. Call, and G. Greenberg (2009b). "An Examination of the Medicaid Undercount in the Current Population Survey (CPS): Preliminary Results from Record Linking," *Health Services Research*, 44, 965–987.
- Eibner, C., F. Girosi, C. C. Price, A. Cordova, P. S. Hussey, A. Beckman, and E. A. McGlynn (2010). *Establishing State Health Insurance Exchanges Implications for Health Insurance Enrollment, Spending, and Small Businesses*. Santa Monica, CA: RAND Corporation.
- Ellis, R. E., D. Roberts, D. M. Rousseau, and T. Schwartz (2008). "Medicaid Enrollment in 50 States: June 2008 Update." The Kaiser Commission on Medicaid and the Uninsured. Kaiser Family Foundation. Washington DC. <https://kaiserfamilyfoundation.files.wordpress.com/2013/01/7606-04.pdf> (Last Accessed August 9, 2018).
- Federal Committee on Statistical Methodology (FCSM) (2001). *Measuring and Reporting Sources of Error in Surveys*. Washington DC: Statistical Policy Office, Office of the Management and Budget, available at <https://nces.ed.gov/FCSM/pdf/spwp31.pdf> (last accessed August 15, 2018).
- Federal Register (2018). Office of Management and Budget. 83 FR 1634, P. 1634–35, available at <https://www.federalregister.gov/documents/2018/01/12/2018-00400/request-for-information> (last accessed August 15, 2018).
- Groves, R. M. (2006). "Nonresponse Rates and Nonresponse Bias in Household Surveys," *Public Opinion Quarterly*, 70, 646–675.

- Groves, R. M., and E. Peytcheva (2008). "The Impact of Nonresponse Rates on Nonresponse Bias: A Meta-analysis," *Public Opinion Quarterly*, 72, 167–189.
- He, Y., and A. M. Zaslavsky (2009). "Combining Information from Cancer Registry and Medical Records Data to Improve Analyses of Adjuvant Cancer Therapies," *Biometrics*, 65, 946–952.
- Hirsch, B. T., and E. J. Schumacher (2004). "Match Bias in Wage Gap Estimates Due to Earnings Imputation," *Journal of Labor Economics*, 22, 689–722.
- Hokayem, C., C. R. Bollinger, and J. P. Ziliak (2015). "The Role of CPS Nonresponse in the Measurement of Poverty," *Journal of the American Statistical Association*, 110, 935–945.
- Little, R. J. A. (2012). "Calibrated Bayes, an Alternative Inferential Paradigm for Official Statistics," *Journal of Official Statistics*, 28, 309–334.
- Meyer, B. D., R. Goerge, and N. Mittag (2014). "Errors in Survey Reporting and Imputation and Their Effects on Estimates of Food Stamp Program Participation," unpublished paper.
- Meyer, B. D., and N. Mittag (forthcoming). "Using Linked Survey and Administrative Data to Better Measure Income: Implications for Poverty, Program Effectiveness and Holes in the Safety Net," *American Economic Journal: Applied Economics*.
- Meyer, B. D., and N. Mittag (2017). "An Empirical Total Survey Error Decomposition Using Data Combination." Unpublished paper.
- Meyer, B. D., W. K. C. Mok, and J. X. Sullivan (2015a). "The Under-Reporting of Transfers in Household Surveys: Its Nature and Consequences," Harris School of Public Policy Studies, University of Chicago unpublished paper.
- . (2015b). "Household Surveys in Crisis," *Journal of Economic Perspectives*, 29, 199–226.
- Mittag, N. (forthcoming). "Correcting for Misreporting of Government Benefits," *American Economic Journal: Economic Policy*.
- Montaquila, J. M., V. Hsu, and J. M. Brick (2011). "Using a Match Rate Model to Predict areas where usps-Based Address Lists May Be Used in Place of Traditional Listing," *Public Opinion Quarterly*, 75, 317–335.
- Motro, J., and V. Roth (2017). "Using Administrative Records and Parametric Models in 2014 SIPP Imputations," unpublished paper, US Census Bureau.
- National Academies of Sciences, Engineering and Medicine (2017a). *Innovations in Federal Statistics: Combining Data Sources While Protecting Privacy*, Washington, DC: The National Academies Press.
- National Academies of Sciences, Engineering and Medicine (2017b). *Federal Statistics, Multiple Data Sources, and Privacy Protection: Next Steps*, Washington, DC: The National Academies Press.
- Newey, W. K., and D. L. McFadden (1994). "Large Sample Estimation and Hypothesis Testing." In *Handbook of Econometrics*. Vol. 4, ed. Robert F. Engle and Daniel L. McFadden, Chapter 36, 2111–2245. Amsterdam: Elsevier.
- Nicholas, J., and M. Wiseman (2010). "Elderly Poverty and Supplemental Security Income, 2002–2005," *Social Security Bulletin*, 70(2), 1–29.
- O'Hara, A. (2016). "Use of Administrative Records to Reduce Burden and Improve Quality," Committee on National Statistics Workshop on Respondent Burden March 8, 2016, Washington DC.
- Raghunathan, T. E., J. P. Reiter, and D. B. Rubin (2003). "Multiple Imputation for Statistical Disclosure Limitation," *Journal of Official Statistics*, 19, 1–16.
- Reiter, J. P. (2003). "Inference for Partially Synthetic, Public Use Microdata Sets," *Survey Methodology*, 29, 181–188.
- Reiter, J. P., and R. Mitra (2009). "Estimating Risks of Identification Disclosure in Partially Synthetic Data," *Journal of Privacy and Confidentiality*, 1(6), 99–110.
- Rubin, D. B. (1996). "Multiple Imputation after 18+ Years," *Journal of the American Statistical Association*, 91, 473–489.
- Schenker, N., T. E. Raghunathan, and I. Bondarenko (2010). "Improving on Analyses of Self-Reported Data in a Large-Scale Health Survey by Using Information from an Examination-Based Survey," *Statistics in Medicine*, 29, 533–545.



- Scherpf, E., C. Newman, and M. Prell (2014). "Targeting of Supplemental Nutrition Assistance Program Benefits: Evidence from the ACS and NY SNAP Administrative Records," unpublished paper.
- Urban Institute (2015). TRIM3 project website, [trim3.urban.org](http://trim3.urban.org), downloaded on November 13, 2015.
- US Census Bureau (2007). *Phase I Research Results: Overview of National Medicare and Medicaid Files. Report of the Research Project to Understand the Medicaid Undercount: The University of Minnesota's State Health Access Data Assistance Center, the Centers for Medicare and Medicaid Services, the Department of Health and Human Services Office of the Assistant Secretary for Planning and Evaluation, and the US Census Bureau*, Washington DC: US Census Bureau, available at [https://www.census.gov/did/www/snacc/docs/SNACC\\_Phase\\_I\\_Full\\_Report.pdf](https://www.census.gov/did/www/snacc/docs/SNACC_Phase_I_Full_Report.pdf) (last accessed September 25, 2018).
- US Census Bureau (2008a). *Phase II Research Results: Examining Discrepancies between the National Medicaid Statistical Information System (MSIS) and the Current Population Survey (CPS) Annual Social and Economic Supplement (ASEC)*, US Census Bureau: Washington DC, available at [https://www.census.gov/did/www/snacc/docs/SNACC\\_Phase\\_II\\_Full\\_Report.pdf](https://www.census.gov/did/www/snacc/docs/SNACC_Phase_II_Full_Report.pdf) (last accessed August 15, 2018).
- US Census Bureau (2008b). *Phase III Research Results: Refinement in the Analysis of Examining Discrepancies between the National Medicaid Statistical Information System (MSIS) and the Current Population Survey (CPS) Annual Social and Economic Supplement (ASEC)*, US Census Bureau: Washington DC, available at [https://www.census.gov/did/www/snacc/docs/SNACC\\_Phase\\_III\\_Executive\\_Summary.pdf](https://www.census.gov/did/www/snacc/docs/SNACC_Phase_III_Executive_Summary.pdf) (last accessed August 15, 2018).
- US Census Bureau (2015). The Supplemental Poverty Measure: 2014. P60-254, September 2015, available at <https://www.census.gov/content/dam/Census/library/publications/2015/demo/P60-254.pdf> (last accessed September 25, 2018).