



UNITED STATES DEPARTMENT OF COMMERCE
Economics and Statistics Administration
U.S. Census Bureau
Washington, DC 20233-0001

November 27, 2018

2018 AMERICAN COMMUNITY SURVEY RESEARCH AND EVALUATION REPORT MEMORANDUM SERIES #
ACS18-RER-07

MEMORANDUM FOR ACS Research and Evaluation Workgroup

From: Victoria A. Velkoff
Chief, American Community Survey Office

Prepared by: Sandra L. Clark and R. Chase Sawyer
American Community Survey Office

Amanda Klimek, Christopher Mazur,
William Chapin, and Ellen Wilson
Social Economic Housing Statistic Division

Subject: Housing Administrative Records Simulation

Attached is the final American Community Survey (ACS) Research and Evaluation report, "Housing Administrative Record Simulation." To improve survey response and reduce respondent burden, many surveys are turning to administrative records to replace or supplement survey data. This research simulated the use of administrative data to replace responses to the 2015 ACS for four survey items that showed promise in previous research. We created a simulated set of estimates that include administrative data and compared them to the 2015 ACS 1-year estimates. The study also calculated burden reduction estimates that showed how many ACS sample addresses could use data from administrative records in lieu of being asked the survey question. Additionally, the research measures improvements in item missing data rates by calculating the number of ACS households that left the ACS item blank but have data from administrative records. Finally, the research identified challenges that require further consideration and decision-making when we consider our next steps in the research on integrating administrative data into ACS production.

If you have any questions about this report, please contact Sandra Clark at 301-763-5884.

Attachment

cc:

Edward Castro (DSSD)
Dawn Nelson (CSM)
Mary Schwartz (SEHSD)

Nathan Walters (SEHSD)
Michaela Dillon (CARRA)
Shawn Bucholtz (HUD)

Intentionally Blank

Housing Administrative Record Simulation

FINAL REPORT



Sandra L. Clark and R. Chase Sawyer, American Community Survey
Amanda Klimek, Christopher Mazur, William Chapin, and Ellen
Wilson, Social Economic Housing Statistic Division

Intentionally Blank

TABLE OF CONTENTS

EXECUTIVE SUMMARY	iii
1. Introduction	1
2. Background	1
3. Literature Review	3
4. Research Questions and Methodology	4
4.1. Research Questions	4
4.2. Methodology.....	4
5. Assumptions and Limitations	11
6. Results	13
6.1. How do the simulated estimates supplemented with administrative records differ from the published estimates (published 2015 ACS data)?	13
6.1.1. Summary Metrics	13
6.1.2. Key Measures.....	15
6.2. How much does the use of administrative records reduce ACS respondent burden? Specifically, what is the number and percentage of housing units for which administrative records could replace ACS?.....	33
6.3. How much does the use of administrative data reduce item allocation rates? What effect does the edit process have on ACS response values and administrative record values? What effect does the selected value (ACS or Administrative record) have on other edits?.....	38
7. Conclusions.....	40
8. References.....	42
A. Distribution of Simulated vs. Published Estimate Differences by Geographic Area and Survey Item.....	A1

List of Tables

Table 1. Simulated vs. Published Key Measures for the United States.....	16
Table 2. Match Rate by Item	34
Table 3. Respondent Burden Reduction by Item.....	36
Table 4. Simulated vs. Published Allocation Rates by Item	39

List of Figures

Figure 1. Adaptive Design	5
Figure 2. Distribution of Simulated vs. Published Estimate Differences by Geographic Area	15
Figure 3. Simulated vs. Published - Total housing units built 2014 or later	18
Figure 4. Simulated vs. Published - Total housing units built 1939 or earlier	19
Figure 5. Simulated vs. Published - Total single-family homes and mobile homes on less than 1 acre	21
Figure 6. Simulated vs. Published - Total single-family homes and mobile homes on 1 to 9.9 acres	22
Figure 7. Simulated vs Published - Total single-family homes and mobile homes on 10 or more acres, state level geographies	23
Figure 8. Simulated vs Published - Total single-family homes and mobile homes on 10 or more acres, county level geographies	24
Figure 9. Simulated vs. Published - Median home value, state level geography.....	26
Figure 10. Simulated vs. Published - Median home value, county level geography.....	26
Figure 11. Simulated vs. Published - Median home value, place level geography.....	27
Figure 12. Simulated vs. Published - Home value less than \$10,000.....	28
Figure 13. Simulated vs. Published - Home value \$2,000,000 or more	29
Figure 14. Simulated vs. Published - Median real estate taxes paid for owner-occupied housing units.....	30
Figure 15. Simulated vs. Published - Real estate taxes less than \$800 for owner-occupied household with a mortgage.....	31
Figure 16. Simulated vs. Published - Real estate taxes of \$3,000 or more for owner-occupied households with a mortgage	32
Figure 17. Simulated vs. Published - No real estate taxes paid for owner-occupied households with a mortgage	33
Figure 18. Match Rate by Item, County Geographic Level, Contiguous U.S.....	35
Figure 19. Respondent Burden Reduction by Item by State.....	38

EXECUTIVE SUMMARY

The U.S. Census Bureau has been researching ways to use administrative data to replace or supplement survey data to improve data quality and reduce respondent burden. This research used a direct substitution method to simulate the use of administrative records in the 2015 American Community Survey (ACS). Four housing items were included in the test: year structure built, acreage, property value, and real estate tax.

We created a simulated dataset, which included data from administrative records, and used it to produce estimates that were compared to the 2015 ACS estimates (published). This allowed us to evaluate the impact of the direct replacement method on the survey estimates released to the public. We combined estimates related to the four studied housing items and created summary metrics to study overall differences. Additionally, we chose key measures for each topic to compare (e.g., median property value). Major findings from the comparison of summary metrics and key measures evaluated include:

- There were differences for a large proportion of summary metrics for all four items across geographic levels (nation, states, counties, and places).
- At the national level, all but one of the 15 key measures studied were statistically different.
- The simulated median property value estimates, which included administrative data, were lower than the 2015 ACS estimates at the state level and for the majority of the lower level geographies. For median property tax, we found similar results.
- The simulated estimates for the acreage item were generally lower than the published estimates.
- Compared to the published estimates, the simulated estimates have a larger number of housing units in the older year built categories and a smaller number of housing units in the more recently year built categories.

We also calculated burden reduction rates (i.e., not asking survey question) and item allocation rates (i.e., allocation of values due to missing data) to see how using administrative data could impact these measures. We found:

- For the housing characteristics tested, we were able to successfully replace responses to survey questions with administrative data for between 37.5 and 53.5 percent of ACS responding housing units.
- Using administrative data resulted in significantly lower item allocation rates (ranging from 2.3 to 12.4 percentage points).

Using administrative data can help improve item allocation rates and reduce respondent burden; however, many of the simulated estimates with the administrative data were significantly different from the published estimates. In addition, there were geographic disparities in coverage of the administrative data. Not having complete coverage of administrative data for all geographic areas and housing types means that data for some areas would contain mostly ACS response data, others mostly administrative data, and others with varying combinations of the two sources. Differential coverage and differences between estimates derived from administrative data versus self-response data are of particular concern, as these issues may compromise the ability of the ACS to represent all areas and housing units as equally as possible. We have concerns that we may not be able to achieve consistency strictly with a direct-replacement method.

Testing a direct-replacement model is an important step in our research of the use of administrative records in the ACS. While we do not believe that a direct replacement method would work for the ACS, possible research in the future can explore other methods such as hybrid, modeled, and data harmonization approaches. Administrative records provide a vast amount of data that we would like to use in the ACS. This research supports the great potential of administrative records, but also revealed some of the challenges we must address to move forward with administrative records in ACS production in a responsive manner that meets the Census Bureau's high quality standards.

1. INTRODUCTION

To improve survey response and reduce respondent burden, many surveys are turning to administrative records to replace or supplement survey data. This research simulates the use of administrative records to replace responses to the 2015 American Community Survey (ACS) for four survey items about housing: year structure built, acreage, property value, and real estate tax.

Using a direct-replacement methodology, we created a simulated set of estimates that include administrative data and compared them to the 2015 ACS 1-year estimates. While the primary goal was to see the impact that the administrative data could have on the published ACS data products, the research also measured improvements in missing data rates by calculating the number of ACS households that left the ACS item blank but have administrative data. Additionally, the study calculated burden reduction estimates to show how many ACS sample addresses could potentially use data from administrative records in lieu of being asked the survey question.

2. BACKGROUND

The ACS is a nationwide survey that collects information on demographic, social, economic, and housing characteristics about the nation's population every year. Data collected through the ACS provide important statistics used by communities, businesses, government entities, and researchers.

The U.S. Census Bureau attempts to contact over 3.5 million U.S. housing units every year to participate in the ACS. Initially, we ask for response through the internet. Next, we send a mail questionnaire to addresses that do not respond via internet. Finally, we use computer-assisted telephone interviews (CATI) and computer-assisted personal interviews (CAPI) to follow-up with addresses that do not self-respond through the internet or mail modes.^{1,2} The distribution of total 2015 ACS responses by mode, weighted, was 32 percent internet, 21 percent mail, 5 percent CATI, and 43 percent CAPI.

The Census Bureau has been actively looking for ways to reduce the burden placed on respondents who participate in the American Community Survey (ACS). In 2014, the ACS conducted a comprehensive review of all of the questions on the ACS to determine their federal

¹ Computer assisted telephone follow-up interviews were discontinued by the American Community Survey in October 2017. They were used to gather data in 2015 though, and are part of the simulated and control datasets.

² For households that do not self-respond via internet or mail, we would first attempt a CATI interview if we have a phone number for the sampled address. CAPI follow-up is our last attempt to reach nonresponders and this operation is conducted for a subsample of nonresponding addresses.

needs and measure their burden on the public. While the Census Bureau found almost all of the questions had significant value to the Federal government, it also found that sometimes some of the most important questions took longer to answer than others, and therefore were more burdensome to the public.

The Census Bureau set forth a research agenda to improve the survey and reduce burden (Chappell and Obenski, 2014). The most recent vintage of this report is, “Agility in Action 2.0: A Snapshot of Enhancements to the American Community Survey” (U.S. Census Bureau, 2017). One proposal is to use administrative records to replace and/or supplement response data. Administrative data are data collected from sources other than survey respondents. The source could be federal, state, or local governments, or commercial entities (who collect their own data (e.g., electric companies) or serve as vendors who compile administrative data from other sources to sell). If administrative data are able to meet the need for data on a particular topic, it may no longer be necessary to ask some questions on the ACS.

The Census Bureau has already conducted several research projects to assess the feasibility of using administrative records in the ACS. In contract with the Census Bureau’s American Community Survey Office, NORC at the University of Chicago researched the availability of data sources and their potential for use in the ACS (Ruggles, 2015). The Census Bureau also identified several topics for further study, and for each topic they reviewed available administrative data values and compared them to ACS self-reported and imputed responses (Moore, 2016; Dillon, forthcoming).

Previous research was able to match (or link) administrative property tax records to 64 percent of 2014 ACS households. Data values for survey items are not available for all matched households, though the availability of data values from the administrative records was high for several housing items. However, the match rates and availability of data values were scarce for some geographies and housing characteristics. For example, property tax data for some jurisdictions are not available and therefore are not included in administrative records. Additionally, due to the way tax records are recorded it is difficult to match administrative records and ACS data for some structures, such as multi-unit rental buildings that do not have individual tax records. For these reasons, previous research concluded that the tax data from administrative records cannot sufficiently replace the questions on the ACS in their entirety. However, the findings suggested that administrative data may be suitable as a direct substitution of a survey question for a subset of ACS households or to assist with the edit and imputation of missing values. Of the ACS housing topics researched, the following were found to be the best candidates for further study: year structure built, acreage, property value, and property taxes.

This research continued our efforts to assess administrative records for use in ACS production by testing direct replacement of four housing items. The primary goals of this study were to:

- Determine the impact that using the administrative records might have on the ACS estimates.
- Understand the potential benefits and risks of using administrative records for the ACS program.

The Census Bureau plans to conduct future research to study specific details pertaining to the costs and other resources required to implement potential changes. Additionally, it would be useful to conduct future research to explore how administrative records could be used to inform our editing and imputation procedures.

3. LITERATURE REVIEW

Researchers have acknowledged that survey participation is continually declining (Brick and Williams, 2013) and we are constantly looking for ways to improve. Using data from administrative records may reduce respondent burden and increase overall survey participation, resulting in more complete response data. In 2014, the ACS program conducted the ACS Content Review (Chappell and Obenski, 2014), which researched ways the ACS could strengthen the survey, improve its content, and reduce respondent burden. Ruggles (2015) followed-up with a review of administrative data sources that could be used to replace or improve questions on the ACS. Recently, the Census Bureau has taken this research to the next level by matching administrative data to ACS addresses and studying match rates and comparing the presence of comparable administrative data and its agreement with ACS response data (Moore, 2015; Dillon, forthcoming).

Moore (2015) and Dillon (forthcoming) found that some items may be candidates for using administrative data to supplement ACS response data due to high match rates and available data values. This research is intended to build on the previous research, see how the 2015 ACS data products would have differed if we used administrative data in place of some of the ACS response data, and serve as a systems test to assess the feasibility of implementing administrative data in our production process.

Kingkade (2013) evaluated the differences in self-reported ACS home values with administrative data from a commercial vendor and found that the differences can be associated in clear cut ways to characteristics of the household, householder, and the location of the property. This finding could support the use of administrative records if their data were considered more reliable than a respondent's self-reported value.

This research proposes using an adaptive design in our automated modes of data collection. Increases in the use of administrative data and combining multiple data sources has amplified the complexities surrounding survey implementation. To accommodate these challenges surveys are turning towards adaptive and responsive design strategies (Chun, et al., 2017).

4. RESEARCH QUESTIONS AND METHODOLOGY

4.1. Research Questions

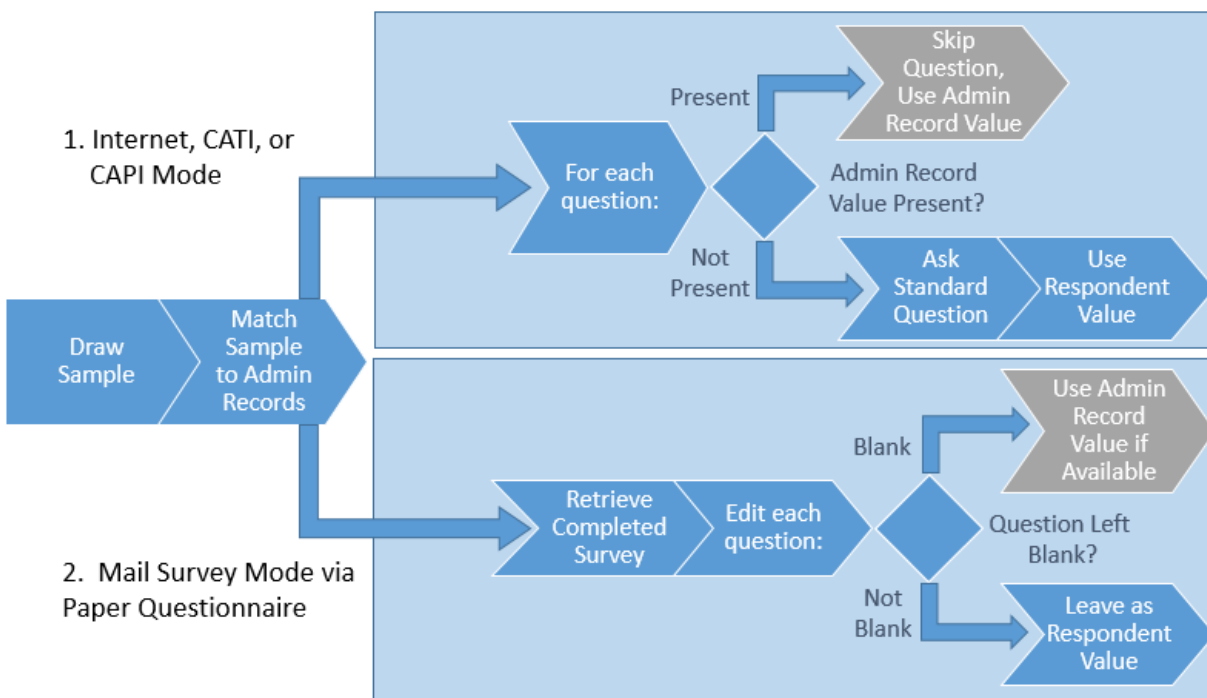
This study answers the following research questions:

1. How do the simulated estimates supplemented with administrative records differ from the published estimates (i.e., published 2015 ACS data)?
2. How much does the use of administrative records reduce ACS respondent burden? Specifically, what is the number and percentage of housing units for which administrative records could replace ACS?
3. How much does the use of administrative data reduce item allocation rates? What effect does the edit process have on ACS response values and administrative record values? What effect does the selected value (ACS or Administrative) have on other edits?

4.2. Methodology

The adaptive design that the research is designed to simulate is shown in Figure 1. The design assumes that the ACS automated modes would skip survey questions for households for which we had administrative records. To adapt our mail mode is more difficult. The administrative records do not always include data for all four of the tested items, therefore to adapt the mail mode we would need several versions of the mail questionnaire. This would be difficult to implement and manage, as well as expensive. For this reason, the design does not include adapting the mail mode. While the mail questionnaire would not skip questions, our design would use administrative data for households for which the question was applicable but left the ACS survey question blank.

Figure 1. Adaptive Design



The research required two comparable datasets, which we designate as our published and simulated. The published used response data from the 2015 ACS. These data have already been captured, edited, and used to create the official 2015 ACS estimates, which were released in 2016.

The simulated dataset was designed to simulate the adaptive design shown above in Figure 1. It used a mixture of 2015 ACS response data and administrative data from county and municipal tax records obtained from the commercial vendor, CoreLogic. The simulated dataset was created by matching the 2015 ACS unedited, unswapped data to the 2014 CoreLogic file containing the tax data using the Master Address File Identification Number (MAFID) – a number associated with each record in the Master Address File (MAF).^{3,4}

³ The 2015 ACS unedited, unswapped data only includes addresses for which we received a valid ACS return. We did not match administrative records to 2015 ACS sample addresses that did not respond to the ACS. We account for survey nonresponse in our weighting adjustments and having only administrative data for four housing items (at most) would not meet the criteria needed for the ACS return to be acceptable. For these reasons, we only matched the administrative records to the final file of addresses that were ultimately used for the 2015 ACS.

⁴ The 2014 CoreLogic file contains administrative data available through the first half of 2014. The Census Bureau received this file from the supplier in fall of 2015. We chose this file because the date it was available lined up with our production schedule (and nearest reference periods) and therefore best mimicked what would happen in a production environment.

The MAF is a Census Bureau database containing the address or location description of all known living quarters in the United States and Puerto Rico.⁵ The ACS sample is drawn from the MAF and therefore all ACS records have a unique MAFID. Census Bureau staff added MAFIDs to our supplier's datasets using a technique in which records must match exactly on "blocking" keys (Wagner and Lane, 2014). Their method first exactly matches the vendor's data to the MAF by Zip Code, then for records with missing or inaccurate Zip Codes, it exactly matches by Census Tract. Within Zip Code or Census Tract, the process attempts to match between the records on the two files with the criteria for a match becoming less restrictive in each successive pass. Since the system is a one-to-many matching system, for some records (i.e., some multi-unit structures), the process assigns the same MAFID to more than one record. The ACS dataset, however, does not have multiple records with the same MAFID. Therefore, linking the two data sources occasionally creates a many-to-one match. This project did not investigate ways to determine which is the best administrative record to use, and since less than one percent of the ACS sample are affected by this issue, administrative records that have the same MAFID as another administrative record were not used. More research on this linkage may be performed in the future.

Once the data sources were matched, ACS response data were replaced with administrative data for the following survey items: year structure built, acreage, property value, and property taxes. Since not all sample records had administrative data available, for various reasons, those records kept their original ACS value.

The ACS value was not replaced if:

- An administrative record for the ACS sample address was not found.
- An administrative record existed, but there was no value for the specific variable or the value was determined unusable because the administrative record has the same MAFID as another administrative record.
- The ACS value was from a mail questionnaire.

Additional item-specific details pertaining to the creation of the simulated file are noted below.

Property Tax: The administrative records for property tax data come from county and municipal property tax records, which often are not collected annually. The vendor's data reflect the most current record year available at the time the ACS would need it, which for many cases is not 2015. To help account for this limitation, the administrative data values for the property tax item were adjusted using an inflation adjustment factor. To inflation-adjust property tax from previous years, the dollar values were inflated to the latest year's dollar

⁵ For more information on the MAF see: <https://www.census.gov/did/www/snacc/publications/MAF-Description.pdf>.

values by multiplying by a factor equal to the average annual Consumer Price Index (CPI-U-RS-All items) factor for the current year, in this case 2015, and dividing by the average annual CPI-U-RS-All items factor for the earlier/earliest year. This converts the tax amount into 2015 dollars.

Property Value: CoreLogic offers a few different estimates of property value from administrative tax records. They have three property values pulled directly from the tax records -- assessed value, appraised value, and market value. The supplier also provides a value they create using their Automative Valuation Model (AVM). To calculate the value, the AVM uses values modeled from the tax records, along with other factors such as recent sales prices, current market conditions, geographic location, etc. The AVM value is designed to be a globally consistent current valuation of the parcel and home. The other estimates of property value are not globally consistent and do not always reflect the most up-to-date value. Therefore, this research used the AVM value. The reference year differences described for the property tax item may also affect the property value item; however, the AVM helps address this issue.

Year Built: CoreLogic offers two year built variables -- year built and effective year built. Effective year built captures more current years, after accounting for things such as teardowns and re-builds. The definition of the effective year built variable is a little ambiguous, so this research will only use the year built variable. The supplier's year built values are continuous number values, while the ACS values before 2000 are organized in categorical number ranges. Therefore, the vendor data was recoded into the analogous ACS categories.

Acreage: The vendor's acreage values are also continuous values that were recoded to match the ACS acreage categories.

After matching the ACS addresses and the administrative records, recoding the administrative data as necessary, and replacing applicable ACS values, the simulated file was processed using the 2015 ACS edit and imputation procedures. The edit and imputation procedures catch edit inconsistencies that may be caused by the administrative data as well as the ACS response data. The edits also used both data sources as donors for imputations of other missing (or inconsistent) data.

To make our simulation as realistic as possible, after running the edit and imputation procedures, the simulated file was subject to the Census Bureau's full data disclosure avoidance and weighting processes. To ensure confidentiality, we implemented disclosure avoidance steps (such as data swapping). We also created weights for each responding sample address to account for sampling, nonresponse, and population control adjustments. Chapters 10, 11, 12, and 13 in The ACS Design and Methodology Report (U.S. Census Bureau, 2014) include detailed information on these procedures.

To answer the first research question, we derived overall summary metrics and key estimates for each of the four housing items from the simulated and published datasets and compared with a focus on geographic variations.

Annually, the ACS produces 1-year estimates for certain geographic populations of 65,000 or more. The estimates are released in the form of several different types of data tables focusing on a particular topic. This research studied 38 detailed tables, comprising a total of 575 individual estimates for the nation. The number of estimates broken out by topic was: 292 property value, 238 year structure built, 21 acres, and 24 property tax.⁶

Estimates were generated for the following geographies:

- United States (nation)
- State
- Counties⁷
- Places⁷

Combining the 575 U.S. level estimates with those for all lower level geographies, the initial research dataset consisted of over two million estimates for which we had both a published and simulated version to compare. The overall summary metrics combine all the differences between the simulated and published estimates and show the distribution of the differences by geographic area. We organized the percent differences into percent range categories to display the distributions. We found that we could not calculate meaningful differences for some estimates due to one or more of the following conditions:

- We did not have both a published and simulated estimate, which made comparison impossible.
- The estimate was a median displayed as a range because it was too high or too low to meet our disclosure threshold.
- The estimate failed our data filtering and/or disclosure review checks.

Before finalizing ACS data products for public release, they undergo disclosure review to assure that confidentiality of respondents has been protected. Additionally, we employ data release

⁶ At this time, the Census Bureau does not publically publish any ACS tables for the Acreage topic. These estimates were created using the three acreage answer choices and tabulating them separately for all single-family households and mobile homes, owner-occupied households and mobile homes, and renter-occupied households and mobile homes.

⁷ Counties and places with populations under 65,000 were excluded.

rules, which check the tables for statistical reliability. For example, if more than half of the estimates in a table are not statistically different from zero, then the table fails.⁸

Fifty percent of the over two million estimates in our initial research file met one or more of the criteria in the bulleted list above. Most of those, 92 percent, failed the data filtering/disclosure review check. We do not release 1-year estimates failing these criterion, therefore we did not consider this a major limitation of the analysis. When calculating the overall summary metrics, these estimates were included in a “not calculated” category.

While the summary metrics provided an overall glimpse of the total effect on the topics studied, we chose a few key measures to look at individually.

The universes, or those required to answer a particular question, are different for some of the variables. Here are the universe definitions by variable(s):

- Year Built: asked of all responding addresses
- Acreage: asked of responding addresses that are mobile homes, detached one-family houses, or attached one-family houses
- Property Tax: asked of responding addresses that are owner-occupied with or without a mortgage
- Property Value: asked of responding addresses that are all owner-occupied, with or without a mortgage, or vacant units that are for sale only or sold, but not yet occupied

We weighted the estimates using the final survey weights, which account for sampling, nonresponse, and population control adjustments. The margins of error (MOE) were created using the replicate weights.⁹ MOEs were calculated for each estimate as well as the differences between the published and simulated estimates. These measures were used in our t-tests to determine if the differences were statistically significant (using a 90 percent confidence level).

To answer research question two, we calculated match rates and burden reduction estimates to learn how many ACS households link to an administrative record and would potentially not be asked the survey questions. Using administrative data allows the automated instruments to skip questions for internet, CATI, and CAPI responders for the year built (YBL), acreage (ACR), property value (VAL), and property tax (TAX) questions if data for the address from administrative records already exist. Our adaptive design does not remove questions from the paper questionnaire, thus the burden reduction (and the numerator for the burden reduction

⁸ See Census 2014 for more details on disclosure and filtering rules.

⁹ The ACS uses successive difference replication to produce the margins of error. For more information, see U.S. Census Bureau (2014).

calculation) was limited to addresses where someone self-responds to the internet mode and to addresses for follow-up via CATI and CAPI.

The formula for determining the burden reduction is:

Burden Reduction for item X = (Number of ACS internet, CATI, or CAPI returns in universe for item X that had a value from administrative data / Number of all ACS returns in universe for item X) * 100

Universes are the same as described above and were calculated with the ACS edited values. Therefore, burden reduction estimates calculated for this experiment will likely be slightly underestimated compared to production estimates because they will not include sample addresses that would have been asked the question (but did not since administrative data were available) and then determined later (via the edit process) that the address was not in universe for the question. To be included in the numerator for the burden reduction estimate, the address had to be in universe for the question, have responded to an ACS automated mode, and have an administrative record value.

We did not weight the burden reduction rates that we calculated to answer this research question. The Office of Management and Budget (OMB) standard estimate of burden is equal to the number of addresses in universe who must respond to the questions. It is a property of the sample only. In this case we are reducing the number of addresses who need to respond by the proportion of addresses for which we have administrative data. Additionally, these estimates were not used to make any comparisons.

Using administrative data has the potential to decrease the amount of missing data in the ACS for addresses whose responders left the ACS item missing, or provided a “Don’t know” or “Refuse” response. Research question three examined this by reviewing item allocation rates. An item allocation rate is the ratio of the number of households with allocated values for an item over the total number of households in universe for the item, multiplied by 100 and rounded to the nearest tenth. An item is considered allocated if the case is in universe and the value is blank going into the editing process, or if the editing process itself blanks the item that is in universe based on information provided for the household.¹⁰ For these calculations, the final edited variable values were used for the numerators and denominators. The household was required to be in universe for the item to be included in the table.

Alternatively, using administrative data could be ineffective if we use administrative data that end up getting blanked or changed in the edit process. Research question three, which is related to allocation rates, also studied the effect the edit process had on the value selected (ACS response or administrative data). We calculated cross-tabulations of available data (ACS

¹⁰ A small number (less than one percent) of in universe blank responses for the YBL variable are filled in using a deterministic assignment condition rather than allocation.

response only, administrative data only, or both) by selected data (which source we used) by allocation flag for each topic. An allocation flag documents the path the data took through the edit process. The edit process first checks if the item is reported and can be kept. If there are no good reported values, then the process tries to assign a value based on other reported data. As a last resort, the edit process allocates a value using a donor value from a hot-deck, which could be either an ACS response value or an administrative record value. The main focus of this analysis was to see how often we selected the administrative data only to assign or allocate another value.

Our edit and allocation procedures use reported data for one or more survey items to assign and allocate data for other items. For example, one dimension of the tenure allocation matrix is building type; so, missing values for tenure are allocated from a case with reported tenure of the same type of building. For the simulated version, administrative data were considered reported and therefore were used as donors to allocate missing data. The impact of this on the four test items is measured in research question one. The four test items play a small (if any) role in the edits for other survey items, however for research question three we also examined differences in the simulated and published datasets to see if other items not included in the study were impacted. For example, year built and year moved in are edited/allocated jointly; thus the administrative record for year built could affect self-reported move year. Similarly, administrative record for hazard insurance could affect a reported value for real estate taxes.

The results section of this report frequently references the American Community Survey - Administrative Records Experiment Results Visualization. This is a data visualization tool that provides most of the data covered in the report, along with additional data not shown. The hyperlink is available where we mention the [Administrative Records Experiment](#) tool. Holding the Ctrl key while clicking on the hyperlink will take you to the visualization tool. Once there, you select a category, choose a topic within the category selected, and the tool will display data at the state, county, or place level (depending on which geography level radio button you select). The visualization covers the following five categories: Acreage, Property Tax, Value, Year Structure Built, and Administrative Record Statistics. The topics for the first four categories include key measure estimates for the four main research topics included in the study. The remaining category, Administrative Records Statistics, covers overall match rates and burden reduction estimates. For your reference, the complete URL to the visualization is: <https://census.gov/library/visualizations/interactive/admin-record.html>.

5. ASSUMPTIONS AND LIMITATIONS

The published dataset included data collected through responses to the ACS during 2015. The simulated dataset included administrative data from the most recent tax records available to CoreLogic at the time we would need to use the data for the ACS. There is a slight lag in the time the records are completed and the time we can reasonably obtain them from CoreLogic.

Additionally, tax record collection varies by jurisdiction. Therefore, the majority of the administrative data did not align with the ACS reference period and the misalignment was not random by geography. Approximately 5 percent of the data are from tax records from 2015, 41 percent from 2014, 49 percent from 2013, 4 percent from 2012, and the remaining 1 percent from 2003 to 2011 tax years. An adjustment factor to inflate the administrative record property tax values to 2015 dollars was used. Additionally, for the property value item, CoreLogic attempts to account for some of this time lag in the methodology they use to create the AVM value (which is the property value we are using for this research). This time reference difference is a limitation of this research and will also be present if we decide to use administrative data for production ACS.¹¹

Property tax records from taxing jurisdictions often focus on aggregate measures of property structures rather than individual units within the structure. For example, for multi-unit structures, tax records often show the values for all units combined rather than individual units in sample, while the ACS generally asks for the value of the individual unit. However, the owner of a noncondominium multi-unit building, in which the owner lives in one of the units, is asked to report for the value of the entire building, the land on which it sits, and any additional structures on the property. This is consistent with the way that the values of multi-unit structures are often recorded on tax records. This may result in differences by type of housing unit in the match of administrative records to the ACS sample. This research did attempt to measure this variability or incorporate it into the results.

Tax records for some areas are not available. Additionally, there are no standards for how local governments collect tax records so the quality and meaning may vary by jurisdiction. Some jurisdictions, including but not limited to small jurisdictions, may not be well represented.

There are differences between the ACS and administrative data that may contribute to differences found in the research, such as the lack of standardization in the objectives, data collection modes and models, and questions used to obtain the data. Further differences may be attributed to inconsistencies in objectives between sources of the administrative data, such as local governmental units.

For the purposes of this research, respondent burden refers to asking the question to ACS respondents. If we have data from administrative records and do not need to ask the survey question, then the burden of asking the question is lifted from the ACS respondent (as shown in the burden reduction estimates in the results section). While there are several other measures of respondent burden, such as time to complete survey or survey items, this research measures burden only in terms of asking versus not asking the survey question.

¹¹ The impact of the time reference limitation may vary from question to question depending on the fluidity of administrative data to time reference.

We acknowledge that there is some error associated with linking administrative data and ACS data. However, we did not perform a linkage bias analysis as part of this research. Therefore, our results do not include an adjustment for linkage error.

The ultimate goal of this research was to see how ACS estimates with administrative data compare to those without administrative data. If we decide to use administrative data in production, we will do more than what is being tested in this research. For example, we could also use administrative data in our edits. This research is a first step to see what happens with a straight substitution method.

6. RESULTS

6.1. How do the simulated estimates supplemented with administrative records differ from the published estimates (published 2015 ACS data)?

There are many reasons why administrative data can differ from response data, therefore we expected differences between the simulated and published estimates. The primary goal of this project was to examine the differences and see how using administrative data could impact the ACS data products that we release to the public. To measure these differences, we compared overall summary metrics and key measures.

6.1.1. Summary Metrics

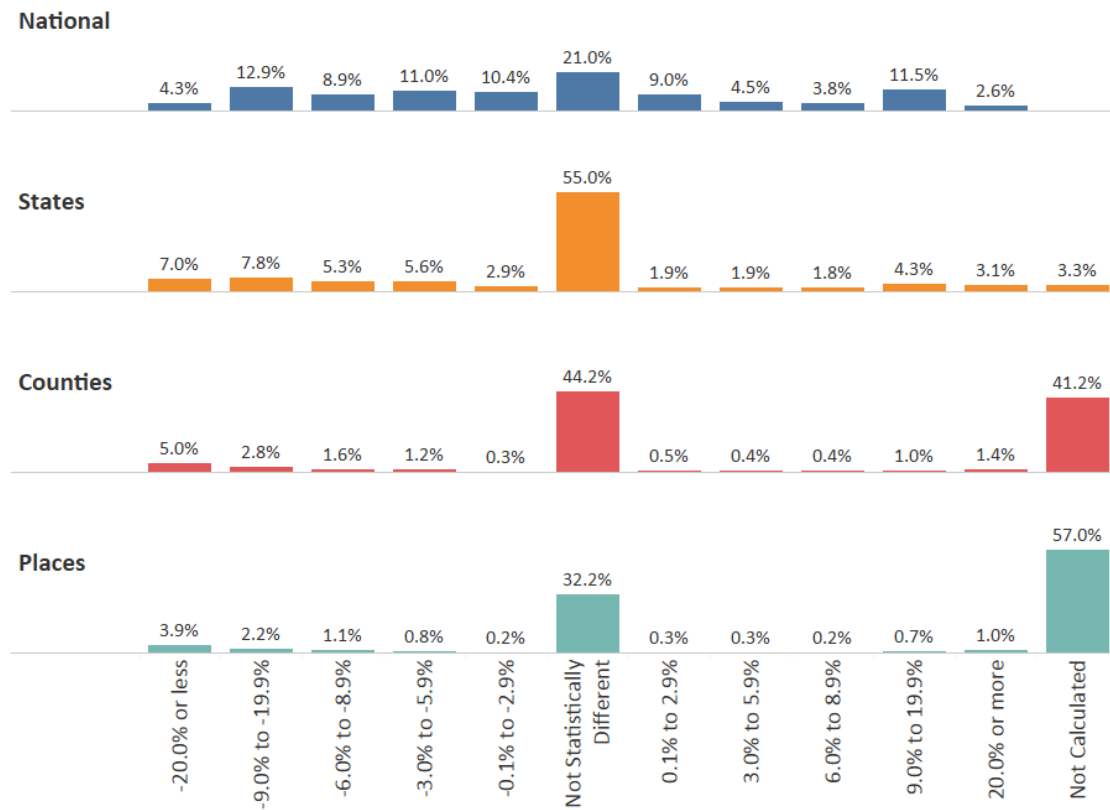
First, to gauge the overall impact to the four topics, we analyzed summary metrics by calculating and combining estimates for all topics using the simulated dataset and comparing them to equivalent estimates from the published dataset. Table 1 shows the distribution of the estimate differences (simulated minus published) by geographic area. Figure 2 shows the data graphically. The negative categories indicate that the simulated estimate (with administrative data) was lower than the published, while the positive categories indicate that the simulated estimate was higher than the published. The category in the middle of the distribution includes estimates that were not statistically different between treatments. Nearly 21 percent of the 575 U.S. level estimates and 55.6 percent of the state level estimates fell in the “not statistically different” category. This category and the “not calculated” category (described in the methodology section) accounted for a large share of the other smaller geographic area estimates.

There are some measures for which the simulated and published estimates were statistically different. It appears that the differences were mildly skewed towards the simulated estimate being lower than the published; however, there were estimates for which the simulated was higher than the published.

Figure 2 combines the estimates for all topics together, however we did review the data broken out separately by the 4 topics and the data tell a similar story. Data showing the topics combined and individually are located in Appendix A.

While our research suggests that many of the estimates we publish for lower geographic areas using 1-year ACS data would not be statistically different, using administrative records would result in some differences particularly for larger geographic areas, such as states and counties, which have enough sample cases to be included in the ACS 1-year estimates (and include aggregated data for which we cannot publish for smaller geographic areas). We were not surprised to learn that there were differences, however it is important to understand these differences better, which leads us into the second part of the research question about key measures.

Figure 2. Distribution of Simulated vs. Published Estimate Differences by Geographic Area



Source: 2015 ACS Housing Administrative Record Simulation

6.1.2. Key Measures

The previous section studied summary metrics by combining all estimates and key measures together to study them as a whole. In this section, we analyzed fifteen specific key measures. Table 1 shows the fifteen measures along with national level data. The simulated and published estimates for all but one of the fifteen key estimates were statistically different. The simulated file had fewer housing units in the more recent year built categories than the published and more older housing units. The median property value for the simulated was 6.3 percent lower than the published and the simulated had 26.6 percent fewer housing units with property values of \$2,000,000 or more. There are several other interesting differences shown in Table 1, which are discussed later, as the remainder of section 6.1.2 dives into each of the key measures separately and includes results for smaller geographic areas.

Table 1. Simulated vs. Published Key Measures for the United States

Key Measure	Simulated	Published	Percent Difference	MOE
<u>Year Built</u>				
Total housing units built 2014 or later	550,430	644,434	-14.6	0.9
Total housing units built 2010 to 2013	3,086,830	3,349,216	-7.8	0.5
Total housing units built 1939 or earlier	18,861,870	17,484,363	7.9	0.2
<u>Acres</u>				
Total single-family homes and mobile homes on less than 1 acre	78,289,605	77,690,886	0.8	0.1
Total single-family homes and mobile homes on 1 to 9.9 acres	16,260,665	16,871,685	-3.6	0.2
Total single-family homes and mobile homes on 10 or more acres	4,466,915	4,453,534	0.3*	0.4
Renter-occupied single-family homes and mobile homes on less than 1 acre	14,837,120	14,819,108	0.1	0.1
Renter-occupied single-family homes and mobile homes on 1 to 9.9 acres	1,768,310	1,833,401	-3.6	0.8
<u>Property Value</u>				
Median property value	\$182,300	\$194,500	-6.3	0.1
Property value less than \$10,000	875,020	1,045,716	-16.3	0.6
Property value \$2,000,000 or more	407,895	555,865	-26.6	1.0
<u>Real Estate Taxes</u>				
Median real estate taxes paid	\$2,190	\$2,259	-3.0	0.2
No real estate taxes paid for owner-occupied households with a mortgage	576,755	1,017,718	-43.3	0.7
Real estate taxes less than \$800 for owner-occupied households with a mortgage ¹	6,101,935	5,526,412	10.4	0.5
Real estate taxes of \$3,000 or more for owner-occupied households with a mortgage	18,926,700	19,417,936	-2.5	0.2

Source: 2015 ACS Housing Administrative Record Simulation

*Not statistically significant at alpha of 0.10

¹ Does not include no real estate taxes paid

6.1.2.1. Year Structure Built

The key measures chosen for year structure built were:

- Total housing units built in 2014 or later – a measure of recent construction
- Total housing units built between 2010 and 2013 – newer units
- Total housing units built before 1940 – Pre-1940 housing units, sometimes considered sub-standard

In addition to studying total housing units, we looked at the three measures above separately for owner-occupied and renter-occupied housing units. Owners often receive documentation at the time of the home purchase that discloses the year that the home was built; therefore, they may be more knowledgeable than renters about the precise year the unit was built. Renters move more frequently and sometimes do not know the date, or even the decade, the home or building was constructed. Despite these differences, we found that the results for owners and renters told a story similar to the results when the types were combined, so for this report, we just focus on total housing units.

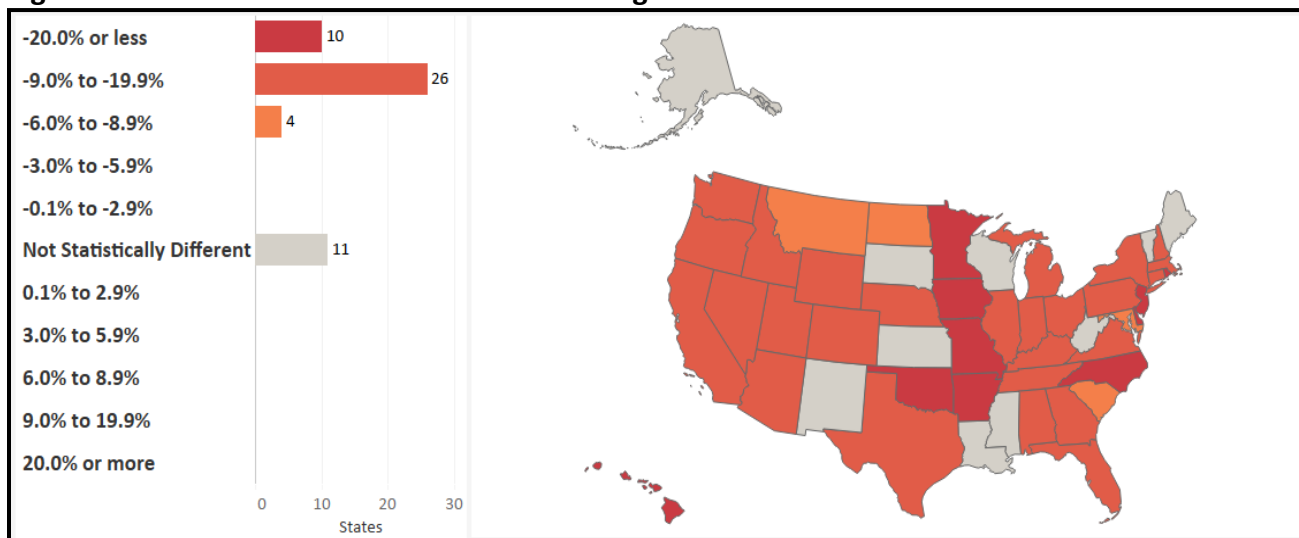
Total housing units built 2014 or later

Using the methodology described in section 4.2, all states, 734 counties, and 436 places had a difference between the published and simulated estimates that could be calculated. Figure 3 shows the differences at the state level. County and place level differences are available at the [Administrative Records Experiment](#) tool.

Of the differences calculated, 78 percent of states, 19 percent of counties, and 12 percent of places had a statistically significant difference. The statistically significant differences were heavily skewed towards the simulated estimate being lower than the published for all geographies. At the state level, the majority with differences fell into the -9.0 percent to -19.9 percent category (26 states). At the county and place level, the majority with differences fell into the -20.0 percent or less category, 106 and 43 respectively.

For all three geographic areas, there was not a single area where the simulated estimate was higher than the published estimate. In fact, the lowest category with estimates was the -8.9 percent to -6.0 percent category.

Figure 3. Simulated vs. Published - Total housing units built 2014 or later



Source: 2015 ACS Housing Administrative Record Simulation

Note: 0.10 alpha used for significance testing

Total housing units built 2010 to 2013

A difference between the published and simulated estimates could be calculated for all states, 819 counties, and 581 places. Statistically significant differences accounted for 84 percent of states, 28 percent of counties, and 18 percent of places.

Similar to total housing units built 2014 or later, the differences were skewed towards the simulated estimate being lower than the published for states, counties, and places; although, not as pronounced. At the state level, the statistically significant differences were spread out throughout the negative categories. The county and place level followed a similar pattern, but there were three county estimates and one place estimate that fell into positive categories. These data can be found at the [Administrative Records Experiment](#) tool.

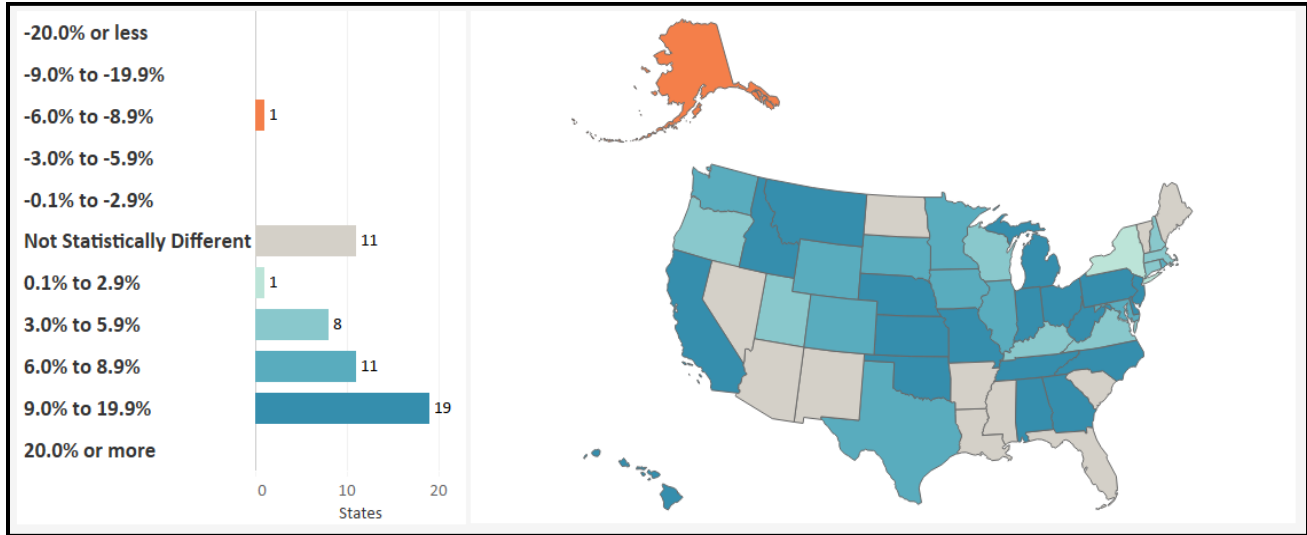
Total housing units built 1939 or earlier

A difference between the published and simulated estimates could be calculated for all states, 815 counties, and 565 places. Statistically significant differences were found for 78 percent of states, 39 percent of counties, and 32 percent of places.

Opposite of the first key measure on total housing units built in 2014 or later, the differences for total housing units built 1939 or earlier were skewed towards the simulated estimate being higher than the published (see Figure 4 showing state level data). This was true for states, counties, and places. At the state level, all statistically significant cases fell into positive categories except for one, Alaska (-8.9% to -6.0% category). No category held a majority as the data was fairly evenly spread out throughout the positive categories.

The county and place level ([Administrative Records Experiment](#)) followed a similar pattern as the state level, but there were 31 estimates for counties and 30 estimates for places that fell into negative categories. Though the county and place data was more slanted towards the higher end positive categories than states, with 60 percent of their estimates falling into categories higher than 8.9 percent.

Figure 4. Simulated vs. Published - Total housing units built 1939 or earlier



Source: 2015 ACS Housing Administrative Record Simulation

Note: 0.10 alpha used for significance testing

For the year built topic in general, the simulated estimates (with administrative data) were lower than the published estimates for more recently built categories and were higher than the published estimates for older built categories. There are several things that may have contributed to these results. Survey respondents may consider renovations and remodeling when answering the question, which may not be accounted for in the administrative data. The administrative tax record data source used in this research may not be the most up-to-date data source. For example, it may not include newly constructed homes that fell in the ACS sample. Depending on the jurisdiction, administrative data may be based on when construction on the housing unit was started, while others may be based on when construction was completed. Additionally, the reference period limitation creates lags in available administrative data and inherent differences in the two data sources.

6.1.2.2. Acreage

The following five key measures were analyzed for the topic of acreage:

- Total single-family homes and mobile homes on less than 1 acre
- Total single-family homes and mobile homes on 1 to 9.9 acres
- Total single-family homes and mobile homes on 10 or more acres

- Renter-occupied single-family homes and mobile homes on less than 1 acre
- Renter-occupied single-family homes and mobile homes on 1 to 9.9 acres

The first three measures mirror the universe and the three response options for the acreage question on the ACS and were analyzed to better understand the acreage distribution as a whole when comparing the published and simulated estimates. All housing units, vacant and occupied are included in this analysis. These three measures include both occupied and vacant housing units. Multi-family households, such as apartments are not asked the acreage question. The final two key measures include only renter-occupied housing units on less than ten acres. These estimates were deemed key measures and analyzed because the U.S. Department of Housing and Urban Development (HUD) defines renter-occupied single-family housing units on less than ten acres as an attribute of standard-quality rental housing when calculating Fair Market Rents.¹²

Total single-family homes and mobile homes on less than 1 acre

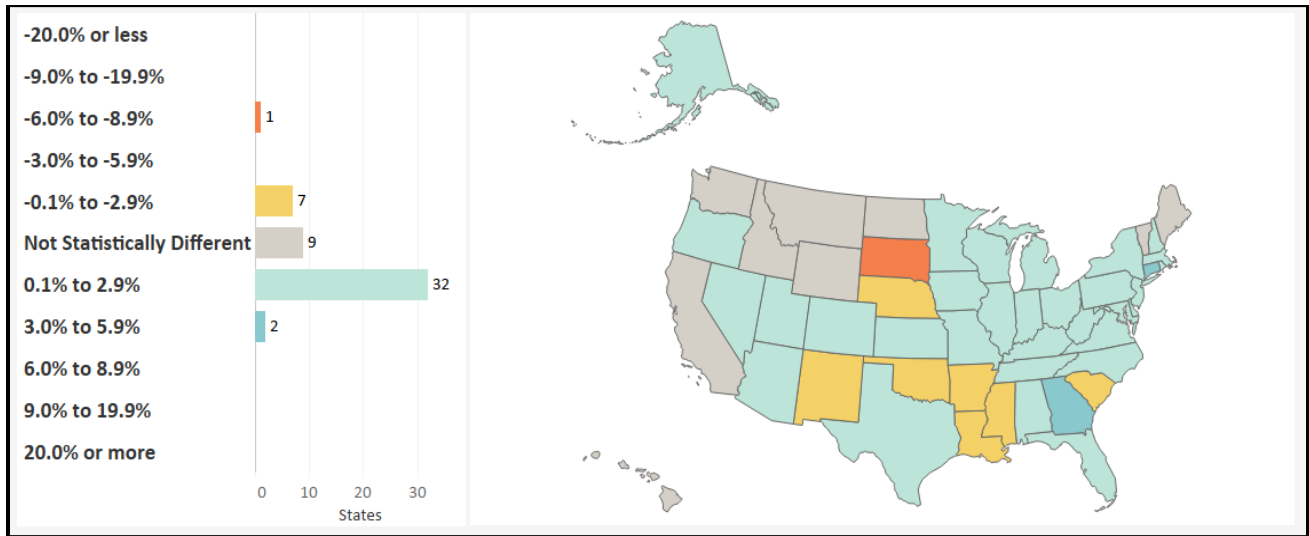
Using the methodology described in section 4.2, all states, 819 counties, and 552 places had a difference between the published and simulated estimates that could be calculated. Figure 5 shows the differences at the state level. County and place level differences are displayed at the [Administrative Record Experiment](#) tool.

Among the differences, 82.4 percent of states, 47.3 percent of counties, and 46.7 percent of places were statistically different. Of the statistically significant estimate differences, the majority of differences at all geography levels fell into the 0.1 percent to 2.9 percent category.

At the state level, estimate differences fell into the -0.1 percent to -2.9 percent and 3.0 percent to 5.9 percent categories at the second and third highest rates, respectively, and accounted for all but one of the remaining states. The only state that fell outside this range was South Dakota, which fell in the -6.0 to -8.9 percent difference category. At the county and place levels, differences fell into the 3.0 percent to 5.9 percent and 6.0 percent to 8.9 percent categories at the second and third highest rates, respectively. Very few counties (47) and places (21) fell into a negative category, which implied that at these geography levels, along with the state level, the simulated contained a mildly greater amount of single-family homes and mobile homes on less than one acre.

¹² For an overview of Fair Market Rents and the calculation process, see the document located at: https://www.huduser.gov/portal/datasets/fmr/fmrover_071707R2.doc

Figure 5. Simulated vs. Published - Total single-family homes and mobile homes on less than 1 acre



Source: 2015 ACS Housing Administrative Record Simulation

Note: 0.10 alpha used for significance testing

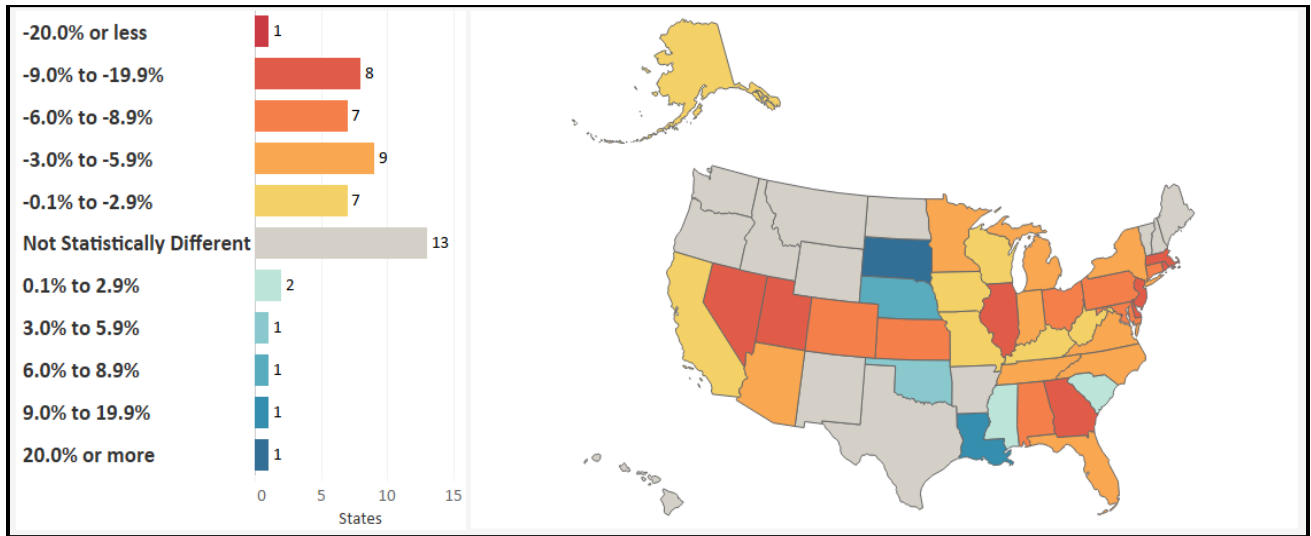
Total single-family homes and mobile homes on 1 to 9.9 acres

A difference between the published and simulated estimates could be calculated for all states, 819 counties, and 552 places. Figure 6 shows the differences at the state level. Information about county and place level differences are displayed at the [Administrative Record Experiment tool](#).

Similar to total single-family homes and mobile homes on less than once acre, statistically significant differences accounted for 74.5 percent of states, 46.4 percent of counties, and 56.2 percent of places. Overall, the majority of all statistically significant differences fell into a negative difference category. Only six states, thirty-six counties, and four places had a difference that fell into a positive category. This pattern implied that the simulated estimates resulted in less geographies with single-family homes and mobile homes on one to 9.9 acres.

The highest rate of the estimate differences at each geography level fell into a different negative category than the others. Nine states had differences in the -3.0 percent to -5.9 percent category, 168 counties fell into the -9.0 percent to -19.9 percent category, and 289 places fell into the -20.0 percent or less category. Furthermore, the 289 places with differences in the -20.0 percent or less category accounted for all but twenty-one of the statistically significant differences at this geography level.

Figure 6. Simulated vs. Published - Total single-family homes and mobile homes on 1 to 9.9 acres



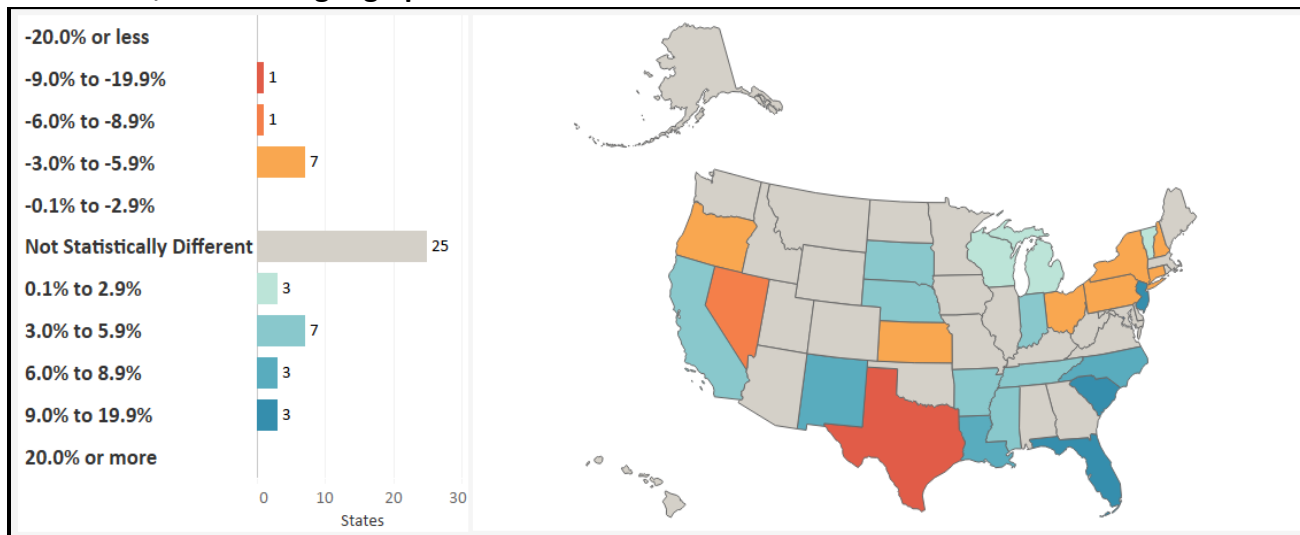
Source: 2015 ACS Housing Administrative Record Simulation
 Note: 0.10 alpha used for significance testing

Total single-family homes and mobile homes on 10 or more acres

A difference between the published and simulated estimates could be calculated for all states, 819 counties, and 552 places. Figure 7 shows the differences at the state level.

Statistically significant differences accounted for 49.0 percent of states, 21.6 percent of counties, and 12.0 percent of places. At the state level, seventeen of the twenty-five statistically significant differences fell within the range of -5.9 percent to 5.9 percent, with seven states falling in a negative category and ten falling in a positive category. No states had differences in the most extreme categories of -20.0 percent or less or 20.0 percent or more.

Figure 7. Simulated vs. Published - Total single-family homes and mobile homes on 10 or more acres, state level geographies

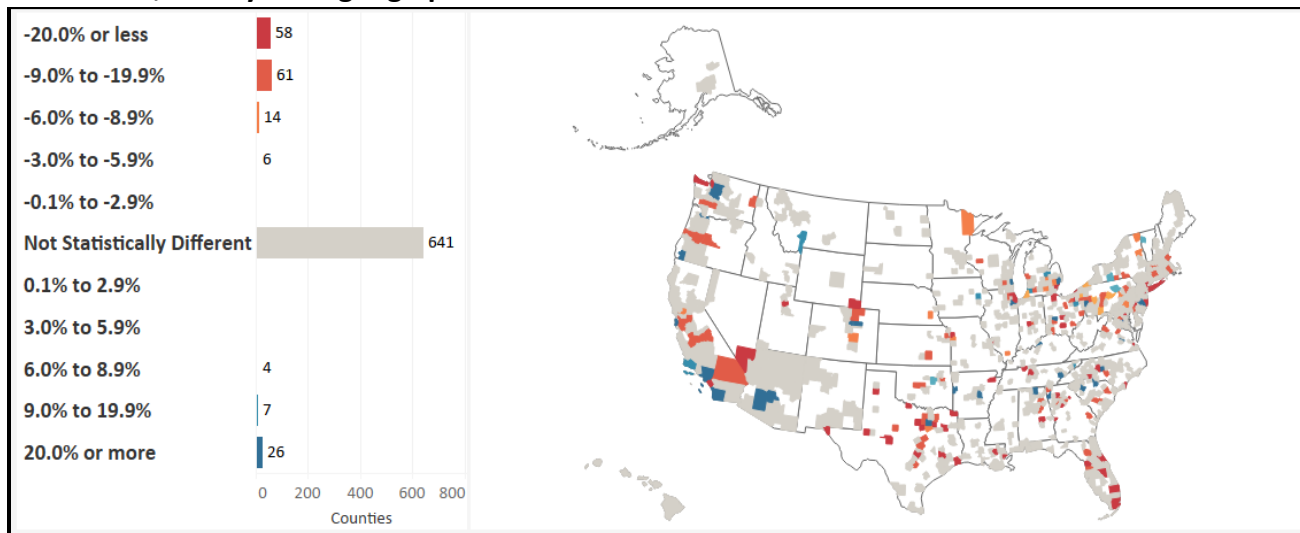


Source: 2015 ACS Housing Administrative Record Simulation

Note: 0.10 alpha used for significance testing

Unlike the state level, a differing pattern of the distribution of estimate differences emerged at the county and place level. Figure 8 displays the county level data, while place level data are displayed at the [Administrative Record Experiment](#) tool. Nearly four of every five statistically significant differences were in a negative category at the county level. Of these negative significant differences, 41.7 percent of the differences were in the most extreme category of -20.0 percent or less, while 43.9 percent were in the -9.0 percent to -19.9 percent category. Only twenty-four counties had a statistically significant difference that fell within the range of -8.9 percent to 8.9 percent. A similar pattern was also recognized at the place level, though at a far greater degree. Forty-five places had a statistically significant difference, and forty-two of them fell in the -20.0 percent or less category, while the remaining three places fell in the 20.0 percent or more category.

Figure 8. Simulated vs. Published - Total single-family homes and mobile homes on 10 or more acres, county level geographies



Source: 2015 ACS Housing Administrative Record Simulation

Note: 0.10 alpha used for significance testing

Renter-occupied single-family homes and mobile homes on less than 1 acre

A difference between the published and simulated estimates could be calculated for all states, 811 counties, and 235 places. Only 38.0 percent of states, 18.9 percent of counties, and 19.1 percent of places had a statistically significant difference. The data for this key measure and the following key measure are displayed at the [Administrative Record Experiment](#) tool.

Similar to the statistically significant differences at all geography levels for total single-family homes and mobile homes on less than one acre, the majority of differences for renter-occupied single-family homes and mobiles on less than one acre fell in a positive category, with the highest rate falling in the 0.1 percent to 2.9 percent category. Only eight states, fifty-two counties, and one place had a difference that fell in a negative category. In general, the simulated estimates contained mildly more renter-occupied single-family homes and mobile homes on less than one acre at all geography levels.

Renter-occupied single-family homes and mobile homes on 1 to 9.9 acres

A difference between the published and simulated estimates could be calculated for all states, 810 counties, and 232 places. Statistically significant differences accounted for 48.0 percent of states, 26.0 percent of counties, and 42.2 percent of places.

Like total single-family homes and mobile homes on one to 9.9 acres, the majority of the statistically significant differences fell in a negative category. Only three states, twenty counties, and one place had a difference that fell in a positive category. Of the estimate differences that fell in a negative category, the -20.0 percent or less category contained 9.5

percent of states, 86.4 percent of counties, and 99.0 percent of places. These patterns implied that the simulated estimates resulted in less estimates with renter-occupied single-family homes and mobile homes on one to 9.9 acres, and if there was a difference at the county or place level, it was highly likely the difference was negative and -20 percent or less.

For the topic of acreage in general, our findings suggest that administrative data tend to report smaller acreage than response data. Respondents may overestimate their lot's size due to rounding up to whole numbers. Administrative data is captured in fractions and anything under exactly 1.0 would be captured in the less than one acre category and anything under 10.0 would not be included in the 10 or more acres category.

6.1.2.3. Property Value

The key measures chosen for Property Value were:

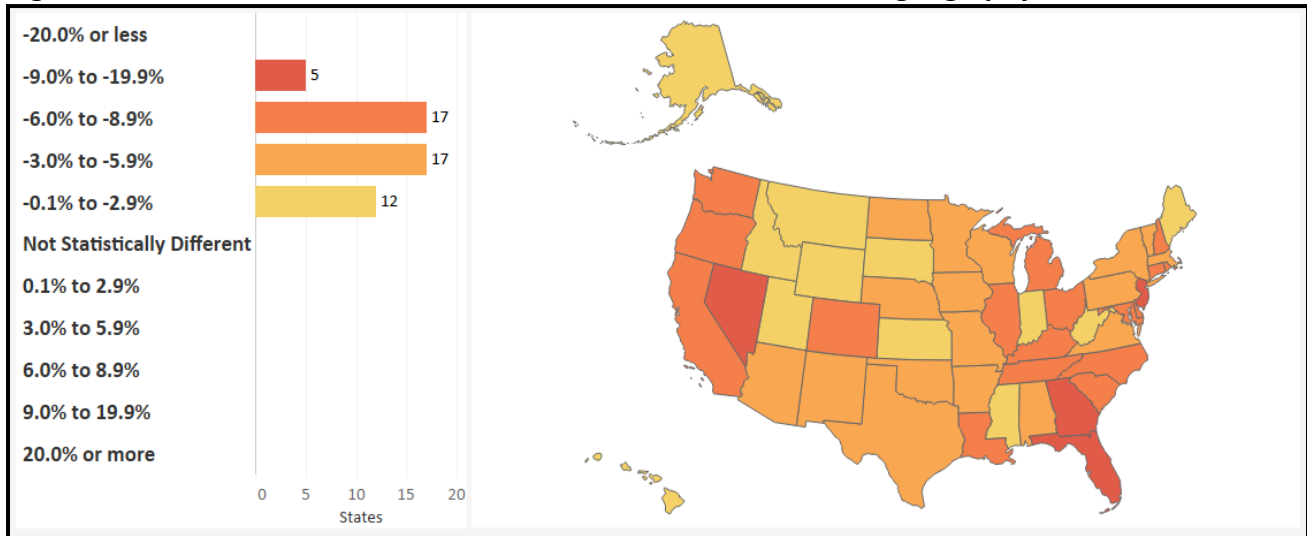
- Median home value (dollars) for housing units
- Total housing units with home value less than \$10,000
- Total housing units with home value \$2,000,000 or more

Median home value (dollars)

Figures 9 and 10 show the percent differences between simulated and published median home value (dollars) for all housing units for the states and counties.

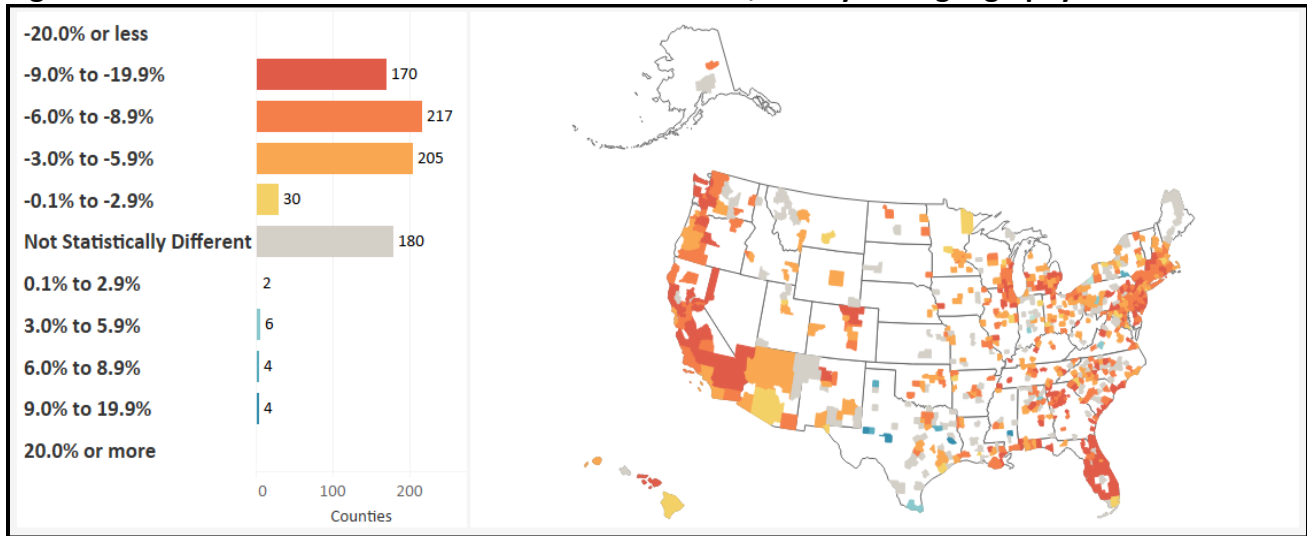
The percent differences between the simulated and published estimates were concentrated in the negative categories, implying that the administrative data medians tend to be lower than the 2015 ACS production medians. At the state level, there were no positive differences between the simulated and published estimates. Every state and the District of Columbia had a lower simulated estimate for median value, with approximately 43 percent of those differences occurring in the -9.0 percent or less categories. This same distribution of differences was largely similar for all geographic areas. There were some cases in which the published was higher, but these were primarily in smaller geographies.

Figure 9. Simulated vs. Published - Median home value, state level geography



Source: 2015 ACS Housing Administrative Record Simulation
 Note: 0.10 alpha used for significance testing

Figure 10. Simulated vs. Published - Median home value, county level geography



Source: 2015 ACS Housing Administrative Record Simulation
 Note: 0.10 alpha used for significance testing

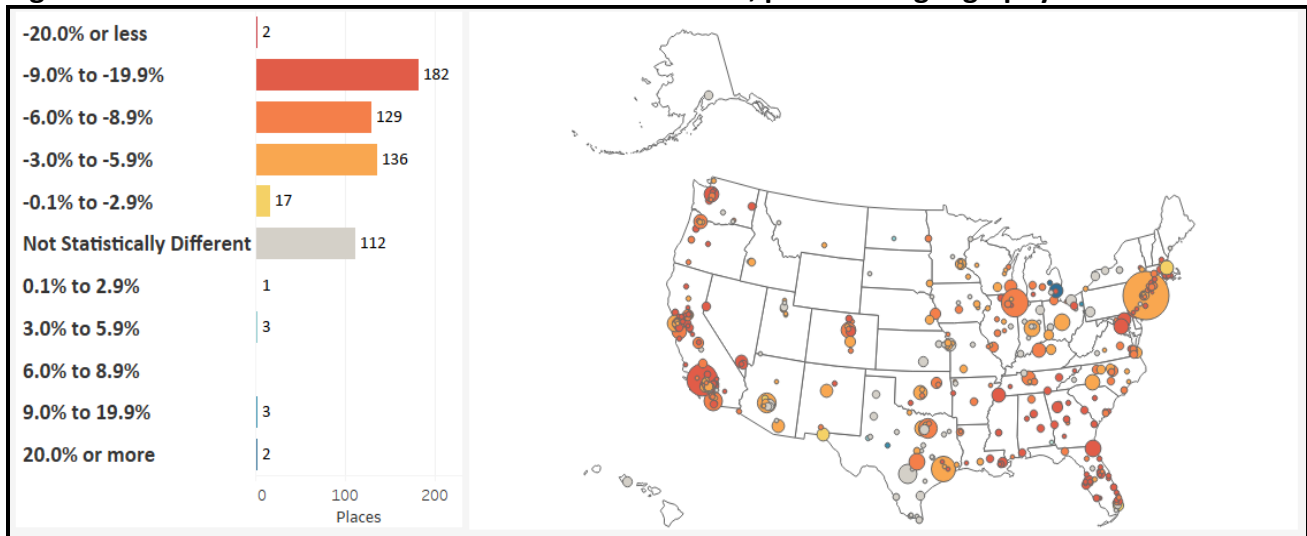
When analyzing our results for the place level geography, there was a result that conflicts with the overall trend. Figure 11 displays place level differences. In our experiment, we found that Flint, Michigan’s estimate in the dataset that contained administrative data was 66.8 percent higher than the published ACS estimates. Further research would need to be conducted to determine an exact reason for this outlier in our results.

We believe there are two possible explanations for this. Respondents in Flint, Michigan may be under valuing their homes because of the recent public health crisis. It is also possible that the

administrative data are unable to provide a timely estimate due to the rapid devaluation of homes in the area.

The ACS is often used to create statistics that can measure changes due to emergencies throughout the nation. Therefore, this result is a finding that the Census Bureau needs to be cognizant of as it considers its use of administrative records in the ACS.

Figure 11. Simulated vs. Published - Median home value, place level geography



Source: 2015 ACS Housing Administrative Record Simulation

Note: 0.10 alpha used for significance testing

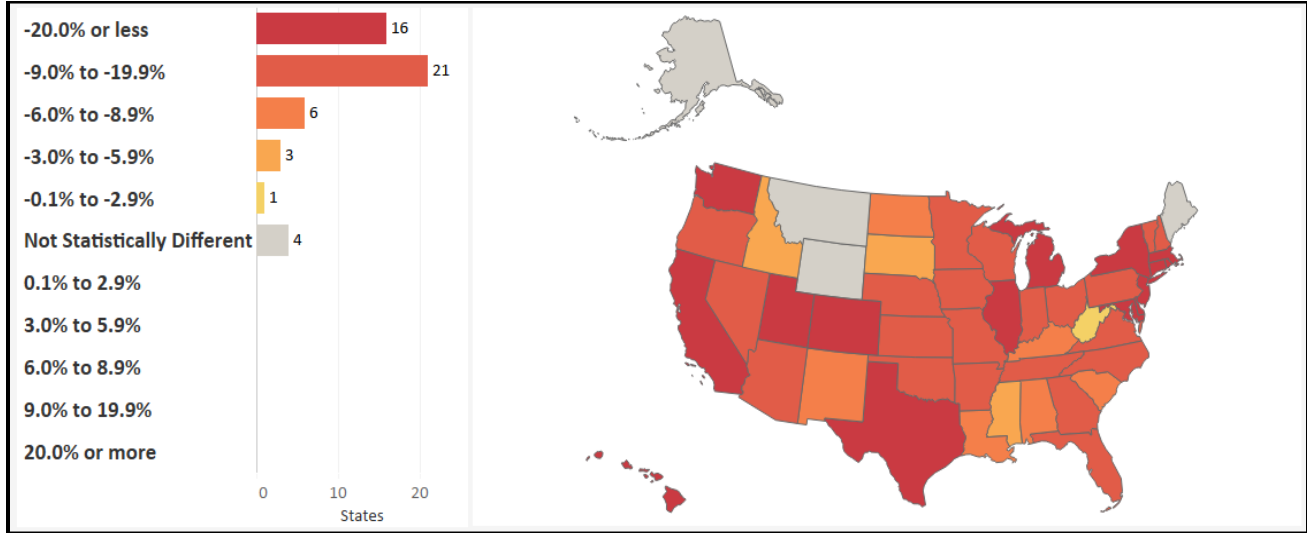
Home value less than \$10,000

Figure 12 shows the percent difference between simulated and published estimates for homes valued at less than \$10,000 for all housing units for the states. More information and geographic areas are available for review at the [Administrative Record Experiment](#) tool.

The differences in median property values imply that the simulated estimates were generally lower than the published estimates; therefore, we could assume that the simulated estimates force more values into the lower extremes and therefore there would be more homes valued at less than \$10,000 in the simulated estimate than in the published estimates (indicated by positive difference categories). However, this does not seem to be the case. At every level of geography, the majority of the differences in published and simulated had differences in the negative ranges for the key estimate of home values less than \$10,000. This means that for each geography, fewer administrative data values ended up in this category than ACS production values. At each level of geography, a large number of the differences between the published and simulated estimates occurred at the -20.0 percent or less level. The median results discussed above suggest lower simulated property values compared to published; however, it appears this is not true for homes valued at under \$10,000 by ACS respondents. It is worth mentioning that at these extremes a high proportion of estimates were not calculated in

the results because they failed the filter or disclosure rules mentioned in the methodology section. This could mean that for the cases of values in the lower bounds it was more difficult to get administrative data, which could be reliably included given the criteria in this study.

Figure 12. Simulated vs. Published - Home value less than \$10,000



Source: 2015 ACS Housing Administrative Record Simulation

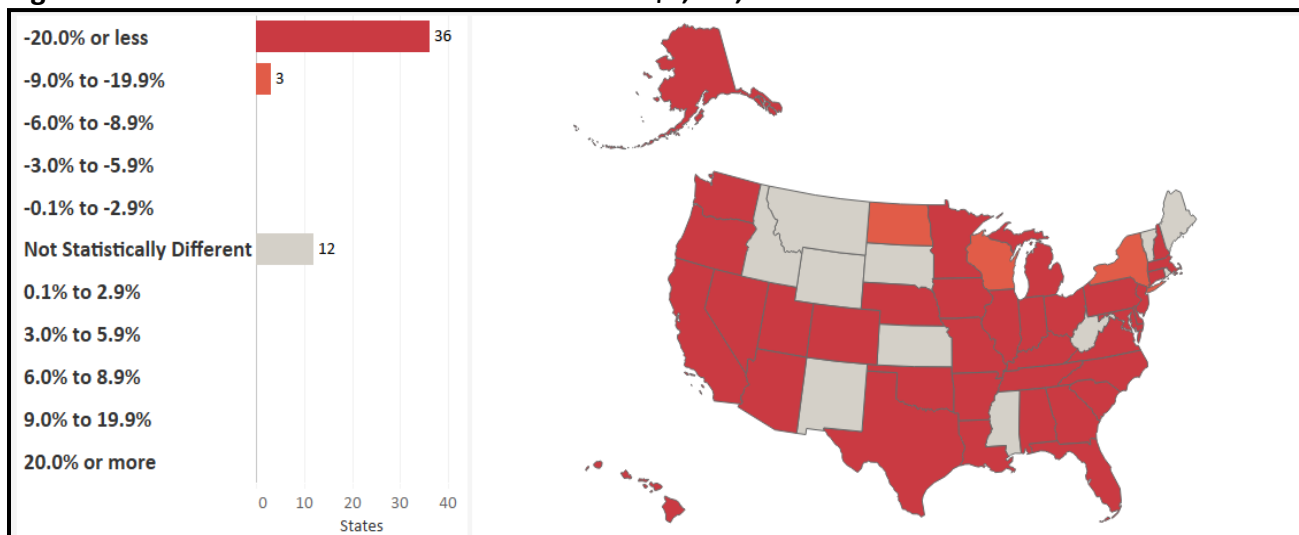
Note: 0.10 alpha used for significance testing

Home value \$2,000,000 or more

Figure 13 shows the percent differences for the key measure of home valued at \$2,000,000 or more for all housing units for the states. County and place level differences are reviewable at the [Administrative Record Experiment](#) tool.

The published versus simulated differences for homes valued in the \$2,000,000 or more category (at each level of geography) were heavily concentrated toward the negative ranges. Additionally, the majority of the differences also ended up in the -20.0 percent or less category. Again, at these extreme bounds a high proportion of these estimates failed the criteria for analysis and were not included in the results. The geographies that were included tend to be larger geographies. This means that administrative data would have a profound impact on the published estimates primarily for larger geographies.

Figure 13. Simulated vs. Published - Home value \$2,000,000 or more



Source: 2015 ACS Housing Administrative Record Simulation

Note: 0.10 alpha used for significance testing

Overall, we found that the estimates that included administrative data had lower measures of property value than the 2015 ACS estimates. This is not all that surprising, since the ACS asks respondents to report what they think their property would sell for, while the administrative data used in the simulated is modeled based on tax records and other administrative data such as recent home sales. We believe that respondents likely over-estimate their property values as compared to administrative data. Our findings were similar to those found by Kingkade (2013).

6.1.2.4. Real Estate Tax

The key measures chosen for Property Tax were:

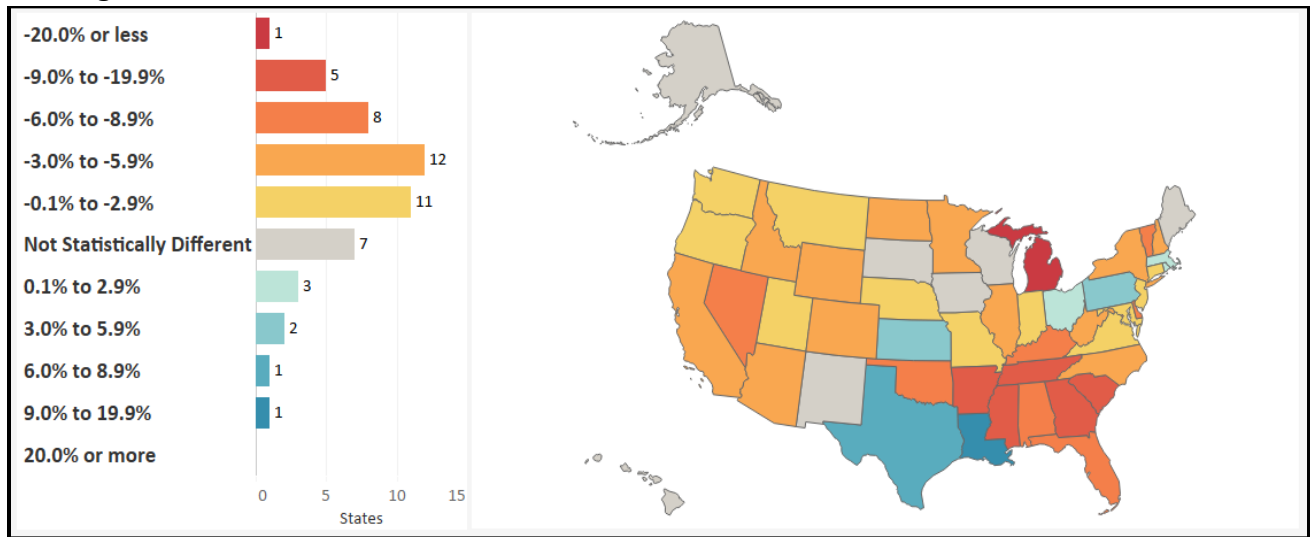
- Median real estate taxes paid for owner-occupied housing units
- Real estate taxes less than \$800 for owner-occupied household with a mortgage (does not include no real estate taxes paid)
- Real estate taxes of \$3,000 or more for owner-occupied households with a mortgage
- No real estate taxes paid for owner-occupied households with a mortgage

Median real estate taxes paid for owner-occupied housing units

Figure 14 shows the state level percent differences in median real estate taxes paid for owner-occupied housing units. County and place level differences for this measure along with all of the other real estate tax key measures (discussed below) are available for review on the [Administrative Record Experiment](#) tool.

The distribution of percent differences in the published versus simulated estimates for median real estate taxes paid tend to be negative at each level of geography. Thus, the estimates for median real estate taxes paid were generally lower for the dataset containing administrative data than the 2015 ACS production estimates. The geographies at which the differences in the published versus simulated were positive tend to be clustered together in regions around Texas and Pennsylvania, respectively, indicating the possibility that either the data from those areas or perhaps the regional real estate markets surrounding them could be affecting these differences.

Figure 14. Simulated vs. Published - Median real estate taxes paid for owner-occupied housing units



Source: 2015 ACS Housing Administrative Record Simulation

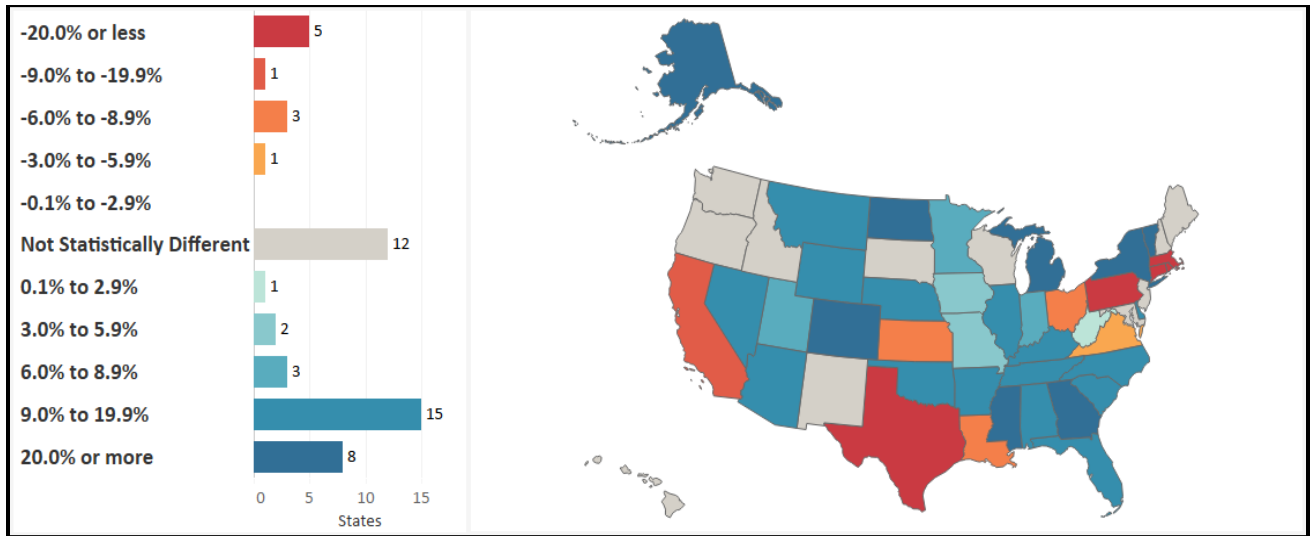
Note: 0.10 alpha used for significance testing

Real estate taxes less than \$800 for owner-occupied household with a mortgage

Figure 15 shows the state level percent differences in the estimates of real estate taxes less than \$800 for owner-occupied household with a mortgage.

At these lower bounds of the distribution of real estate taxes paid, the differences in the published versus simulated estimates were distributed amongst positive and negative categories; however, these differences tend to concentrate toward the outer extremes. These occurrences appear to be regional, with the simulated estimate being lower in the coastal Western and Northeastern regions as well as in Texas and being higher everywhere else in general. It is unclear without further analysis whether this could be due to the regional real estate markets or due to the differences in population and therefore sample size in these areas. This trend of concentration on or toward the outer bounds holds true for all levels of geography.

Figure 15. Simulated vs. Published - Real estate taxes less than \$800 for owner-occupied household with a mortgage



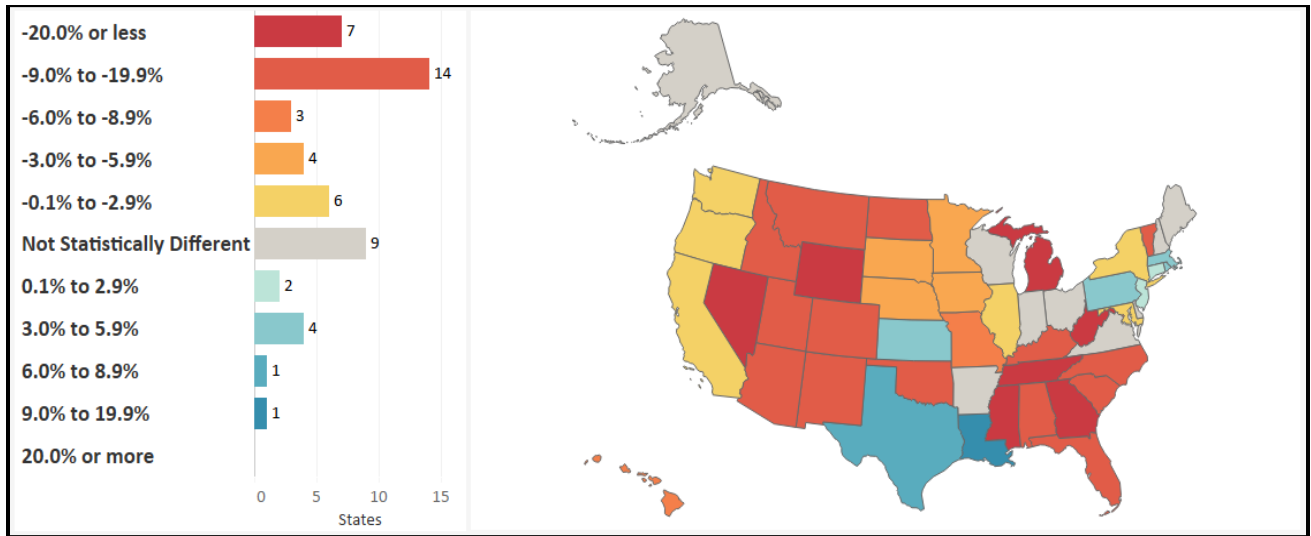
Source: 2015 ACS Housing Administrative Record Simulation
 Note: 0.10 alpha used for significance testing

Real estate taxes of \$3,000 or more for owner-occupied households with a mortgage

Figure 16 shows the state level percent difference for the estimates of real estate taxes of \$3,000 or more for owner-occupied households with a mortgage.

The upper bounds of the distribution of real estate taxes paid show a mildly different result than the lower bounds. In the \$3,000 or more range, there is a heavier concentration towards the simulated estimate being lower than the published.

Figure 16. Simulated vs. Published - Real estate taxes of \$3,000 or more for owner-occupied households with a mortgage



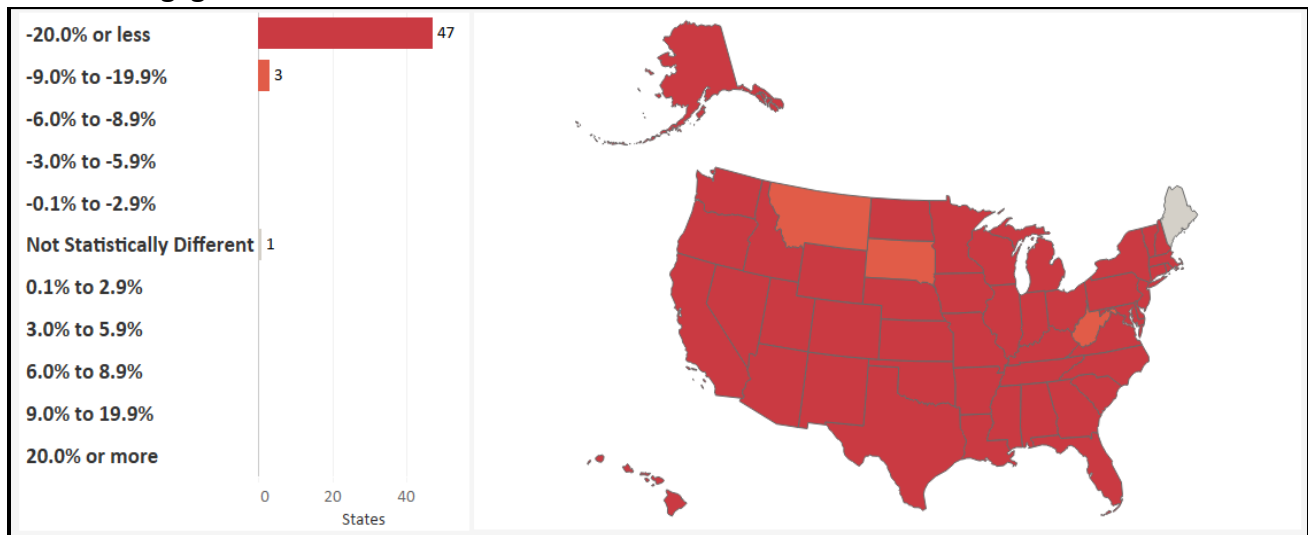
Source: 2015 ACS Housing Administrative Record Simulation
 Note: 0.10 alpha used for significance testing

No real estate taxes paid for owner-occupied households with a mortgage

Figure 17 shows the percent differences for the estimates of no real estate taxes paid for owner-occupied households with a mortgage at the state level.

The majority of differences in the published versus simulated estimates for housing units with a mortgage with no real estate taxes paid fell into the -20.0% or less category at all levels of geography. Among all geographies, 75 percent of the differences fell into the -20.0% or less category. In this analysis, the dataset with administrative data report far fewer cases of “no real estate taxes paid for housing units with a mortgage” than the 2015 ACS production estimates report at most geographies. The ACS question has a checkbox indicating “No real estate taxes paid”, which based on this finding could mean that that checkbox is incorrectly being overused in the ACS. Some people might be using this checkbox when their real estate taxes are included in their mortgage, meaning they themselves are not directly paying them.

Figure 17. Simulated vs. Published - No real estate taxes paid for owner-occupied households with a mortgage



Source: 2015 ACS Housing Administrative Record Simulation

Note: 0.10 alpha used for significance testing

The simulated estimates, in general, seemed to be lower than the published estimates which suggests that administrative data captures lower real estate tax values than survey response data. Similarly to property value estimates, it is possible that survey respondents over-report their taxes as compared to the value on their tax records. ACS respondents are more likely to give us rounded estimates, while administrative data report exact amounts. A brief review of the data shows more rounded estimates (ending in 0) in the published dataset than the simulated, however, more research is necessary to conclude that rounding contributed to estimate differences.

6.2. How much does the use of administrative records reduce ACS respondent burden? Specifically, what is the number and percentage of housing units for which administrative records could replace ACS?

To begin looking at the reduction in respondent burden, it was important to determine what percentage of households had administrative records available for the experiment. To do this, we looked at what the match rate was nationally and in subnational geographies. We found that 65.3 percent of households that responded to the 2015 American Community Survey had a match to an administrative record that had a value present for at least one of the four variables tested in the experiment.

The match rates showed some variation at subnational geographic levels. At the state level, the largest nominal match rate was for Iowa at 77.6 percent and Maine had the smallest nominal

match rate at 8.5 percent.¹³ At the county level, there was even more variation with the maximum match rate being 96.2 percent and some counties not having a single match to administrative records for any of the four variables. Looking at the match rates individually by question shows that some questions have higher match rates than others. Overall match rates are displayed at the [Administrative Record Experiment](#) tool.

To calculate the match rate for each question, the number of households that had an administrative record available for a given item was divided by the total number of households that responded to the ACS and would have been asked the question. Table 2 shows how the match rates varied by question. Figure 18 shows how match rates varied geographically at the county level. Some areas had very high match rates while others had no matches. In some instances, we were unable to match ACS data and administrative data because of duplicate MAFIDs in the administrative data. We suspect that the records with duplicates are concentrated among certain types of housing units, such as multi-unit structures, trailer parks, etc. Linkage issues, such as these, may contribute to some of the variation in match rates by geographic area and item.

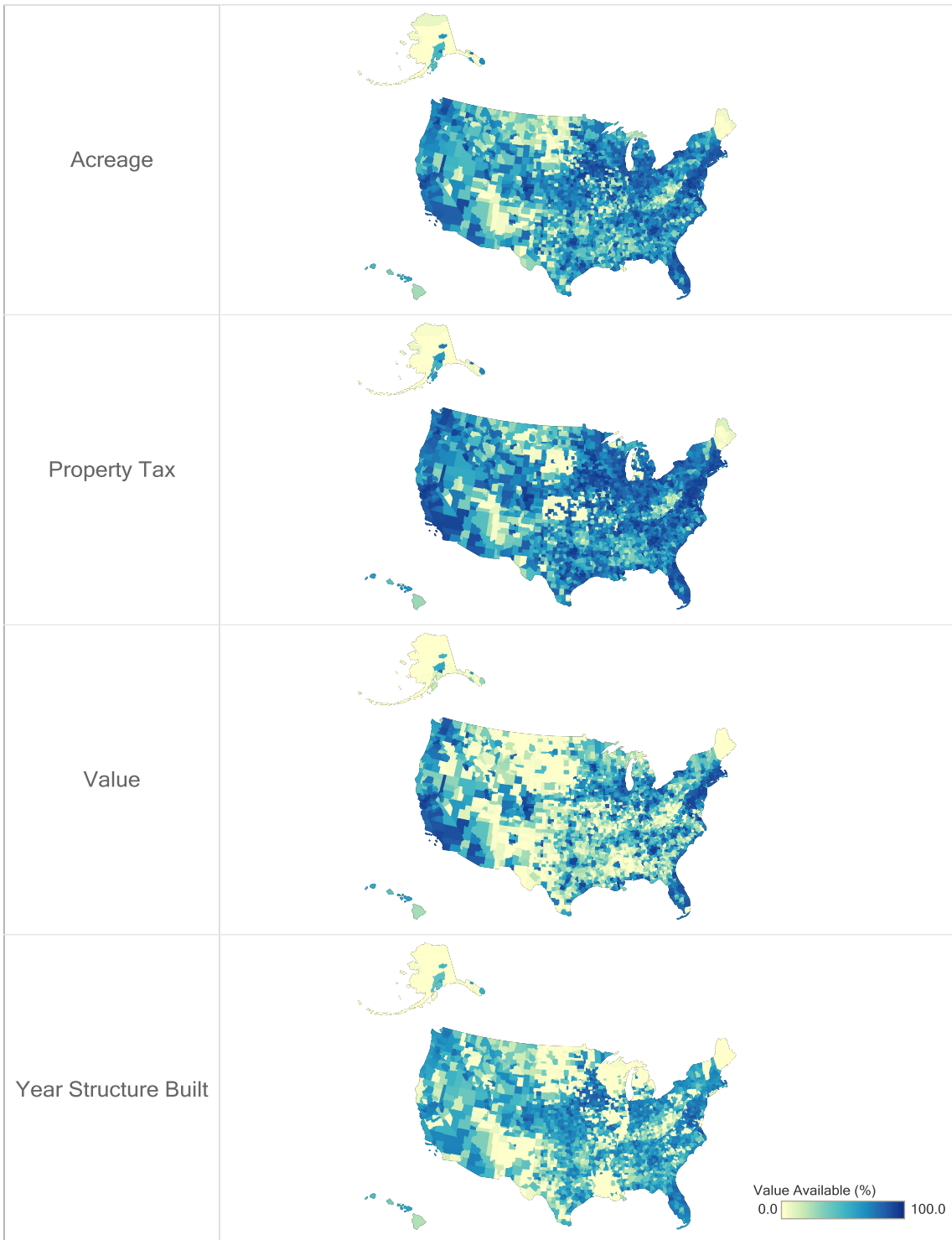
Table 2. Match Rate by Item

Item	Addresses in Universe	CoreLogic Value Present	Match Rate
Overall	2,305,707	1,505,713	65.3%
Year Built	2,305,707	1,230,287	53.4%
Acreage	1,831,641	1,374,428	75.0%
Property Value	1,500,362	1,032,867	68.8%
Property Tax	1,474,640	1,178,190	79.9%

Source: 2015 ACS Housing Administrative Record Simulation

¹³ We did not weight the match rates, create margins of errors, or test differences/comparisons for statistical significance.

Figure 18. Match Rate by Item, County Geographic Level, Contiguous U.S.



Source: 2015 ACS Housing Administrative Record Simulation

Despite matching to an administrative record, an ACS sample household may not have a value that is useable in the adaptive design used in our test (i.e., the household responds via paper questionnaire) or the record may not include data for all survey items.

Not asking the full set of survey questions when administrative data are present would reduce respondent burden for some households. This study did not include modifications to the mail questionnaire. Therefore, if the ACS were to adapt the methodology tested in this research, respondents choosing to complete the mail survey would be subject to the four survey questions even if administrative data for the household exist. However, the simulated design would modify the survey questions for respondents choosing the internet, CATI, or CAPI modes, by not asking the question if data from administrative records were available for the household, therefore reducing burden for these households.

Table 3 includes burden reduction estimates, showing the proportion of responding ACS sample addresses for which administrative records could have been used in lieu of asking the survey questions. The table includes estimates for each survey item. Note that the estimates include households responding by mail in the denominator even though it is not possible for these households to be in the numerator.

Using administrative records would alleviate the burden of asking the property tax item for 53.5 percent of responding in-universe households. The estimated burden reduction is 52.0 percent for acreage, 47.2 percent for property value and 37.5 percent for year structure built.

Table 3. Respondent Burden Reduction by Item

Item	Addresses in Universe	CoreLogic Value Present¹	Burden Reduction
Year Built	2,305,707	864,760	37.5%
Acreage	1,831,641	952,429	52.0%
Property Value	1,500,362	708,041	47.2%
Property Tax	1,474,640	789,501	53.5%

Source: 2015 ACS Housing Administrative Record Simulation

¹CoreLogic values for sample addresses responding via paper questionnaire are not included in these estimates since mail mode responders will receive the questions in our adaptive design and therefore not experience a reduction in burden

The estimates in Table 3 show the promise of using administrative records to reduce respondent burden. However, looking at burden reduction estimates by state shows that the benefit is greater for some states than others. Figure 18 shows burden reduction estimates (along with other estimates) by state for the four items. More detailed information about burden reduction can be found at the [Administrative Record Experiment](#) tool.

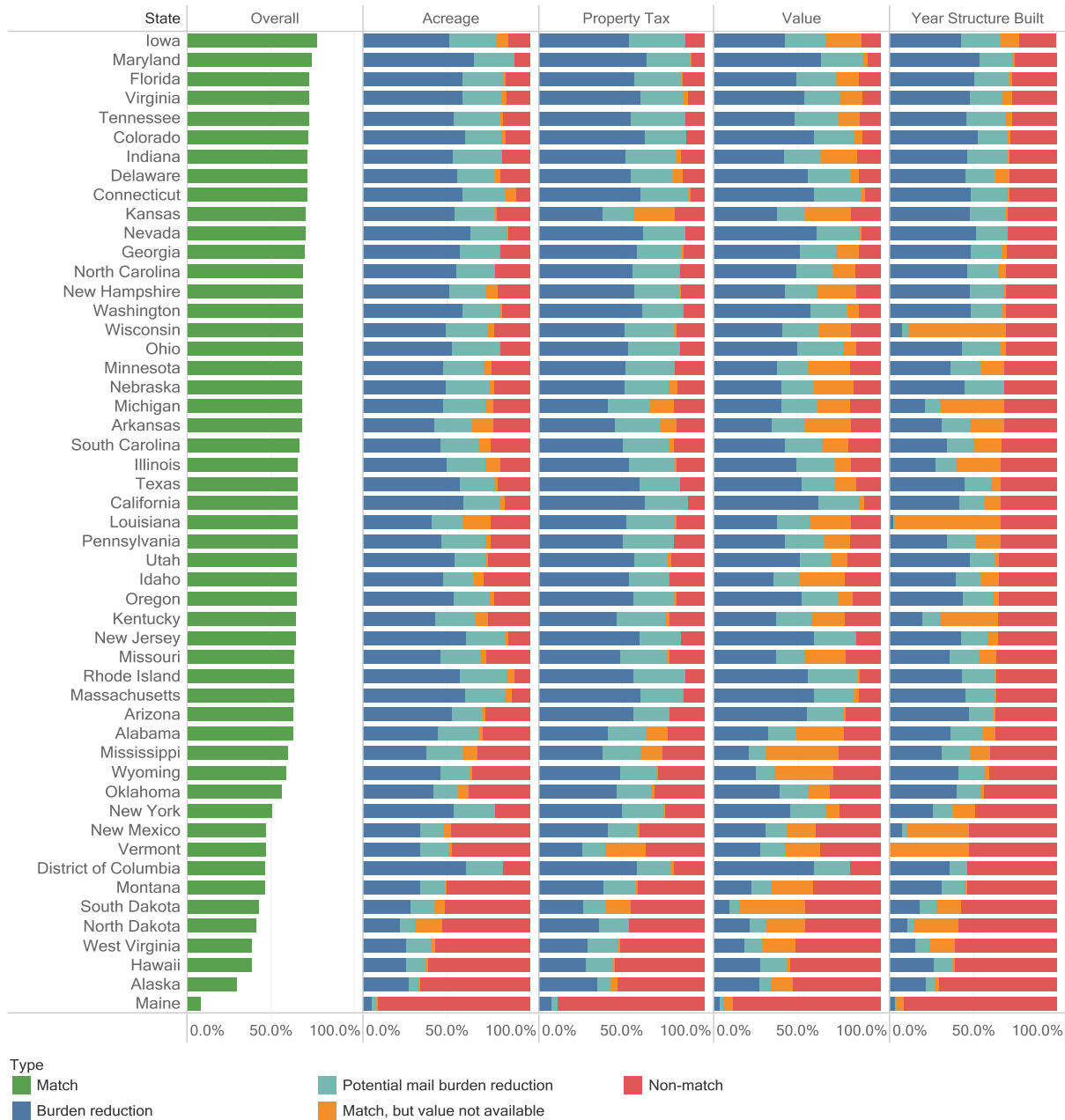
First, the green bars in Figure 19 represent the overall match rates, while the red bars are the non-match rates. Next, the figure breaks matches down further into type of match. The dark blue bars indicate the burden reduction.

Nine states (Kentucky, Louisiana, Maine, New Mexico, North Dakota, South Dakota, Vermont, West Virginia, and Wisconsin) had burden reduction estimates for the year structure built item that were less than 20 percent. Vermont did not have any administrative data for the year structure built item. For the acreage item and property tax item, only one state (Maine) had an estimate less than 20 percent; there were three states (Maine, South Dakota, and West Virginia) for the property value item. Maine stands out from the other states in that it had very low burden reductions for all four items (3.4 percent for year structure built, 5.6 percent for acreage, 4.2 percent for property value, and 7.5 percent for property tax). There are no national standards for how states collect or keep property records; therefore, low burden reduction estimates likely indicate that an area does not have the ability to offer the records or our vendor has not been able to acquire them at this time.

The light blue bars show the potential additional burden we could relieve if we are able to find a way to adapt our mail mode. As you can see, the benefit is remarkable.

Finally, the orange bars indicate that we had an administrative record match, but administrative data are not available for that particular survey question. This is interesting because it shows that just because we can match to administrative data does not mean that we will have complete data available for all survey questions. For example, the orange bars in the last column show that for a handful of states we have great match rates yet virtually no data for the year built question.

Figure 19. Type of Match by Item by State



Source: 2015 ACS Housing Administrative Record Simulation

6.3. How much does the use of administrative data reduce item allocation rates? What effect does the edit process have on ACS response values and administrative record values? What effect does the selected value (ACS or Administrative record) have on other edits?

Using administrative data allowed us to obtain significantly more data than we get from survey responses alone, which resulted in lower item allocation rates for all four tested items. Table 4 displays the simulated and published allocation rates for each item. The simulated item

allocation rates ranged from 2.3 percentage points lower for the acreage item to 12.4 percentage points lower for the real estate tax item.

Table 4. Simulated vs. Published Allocation Rates by Item

Item	Simulated	Published	Difference	MOE
Year Built	12.7	17.8	-5.2	0.2
Acreage	1.4	3.7	-2.3	0.1
Property Value	5.0	12.0	-7.0	0.1
Property Tax	4.5	16.9	-12.4	0.1

Source: 2015 ACS Housing Administrative Record Simulation

Note: 0.10 alpha used for significance testing

While administrative data can help improve allocation rates, we will still have missing data. As described in section 4.2, ACS data are subject to an edit and allocation process to validate reported data and allocate for missing data. Reported data are also used as input for hot deck matrices. The published file includes reported response data, while the simulated file includes as reported both response data and administrative data. We found that very little of the reported data, including the administrative data, was changed as a result of the edit process (less than one percent). This finding alleviates our concern that we would use administrative data only to blank or change it in our edit process.

Unlike other survey items, the items included in this test do not play a large role in our edit process and are not often used in the imputation of other items. Despite not playing a large role, we found that using administrative data as hot deck donors had an impact on some other survey items that were not included in the test. For example, property value is used to validate and allocate income. The ACS has tables on several income measures, related to both personal income and household income. Additionally, income is used to calculate poverty measures. We did not conduct a thorough investigation into the full impact to the other items, however we looked at median household income. Between the published and simulated estimates, we found mostly nominal differences in median household income. The difference for the U.S. was only 0.02 percent (though statistically significant). For smaller geographic areas, there were only a small number of statistical differences. We did, however, find some areas with large nominal differences. Federal funding is likely based on point estimates and therefore even nominal differences should be acknowledged.

The research demonstrated how relatively small changes to a few items can have a ripple effect and impact other items. While this may have a positive outcome and further improve data

quality, any changes must be thoroughly researched before implementing to avoid unexpected results.

7. CONCLUSIONS

This research confirmed that using administrative data improves item allocation rates and could potentially reduce respondent burden. While the simulated allocation rates were lower for all four survey items, the property value item experienced the largest difference with the simulated rate 12.4 percentage points lower than the published. Using administrative data would allow us to ask less questions to a significant proportion of ACS responding households. Via our tested method, the reduction is approximately 38 percent to 54 percent depending on the survey question.

Of the households that responded to the 2015 ACS, 65.3 percent had a match to an administrative record and had a value present for one of the four variables tested in the experiment. However, the match rates varied at subnational geographic levels and by survey item. The match rate for Maine was only 8.5 percent and some counties did not have an administrative data match for a single household. The match rates for the acreage and property tax items were higher than the rates for the property value and year built items.

As expected, many of the simulated estimates with the administrative data were significantly different from the published estimates. This research analyzed 15 key measures related to the four items included in the study and all but one of the national level measures were significantly different between the simulated and published. There are many reasons for differences between survey responses and administrative data sources, some of which we outline in our limitations. Not having complete coverage of administrative data for all geographic areas and housing types means that data for some areas would contain mostly ACS response data, others mostly administrative data, and others with varying combinations of the two sources. Given that we found differences between the survey and administrative data, and that the ACS needs to represent all areas and housing units as equally as possible, a direct replacement method cannot be recommended.

While we do not believe that a strictly direct replacement method (such as the one studied in this research) would work for the ACS, we learned a great deal that we believe will help with future research. Throughout the research we encountered several challenges that must be dealt with. First and foremost is the decision on how the Census Bureau obtains administrative data. For the items in this study, we currently get administrative data from an outside vendor. However, in the future this vendor could go out of business or another vendor could be awarded the contract bid. Additionally, we need to know if this vendor provides the most comprehensive administrative data or if there are other vendors that do a better job. We will need to keep up-to-date on this over time.

In addition to keeping up-to-date with the quality of the data available and received by vendors, we need to make sure the sources of our administrative data are as comprehensive as possible. During the final stages of our research, we learned of an additional data source for year structure built data. Mule (forthcoming) found that using the United States Postal Service's Delivery Sequence File addresses some of the coverage issues with administrative data for more recently built addresses.

At the time we started this research, our housing experts believed that the best source of data for the property value item was from a model created by our vendor. We were told that the model used a combination of administrative data sources to predict the most accurate and current property value estimate. However, the model is owned by the vendor and their methodology is confidential. It would be in the best interest of the Census Bureau to develop our own models so the methodology is in-house, known, and controlled by the Census Bureau. We do not want an outside vendor to have any control over our data.

We also had administrative data that we were unable to match to ACS data because of duplicate MAFIDs in the administrative data. We suspect that the records with duplicates are concentrated among certain types of housing units, such as multi-unit structures, trailer parks, etc. Therefore, it is necessary to improve how we are able to link data sources.

When developing data models, we must account for things such as: differing amounts of administrative data among geographies and types of housing units, reference period and time lag differences between administrative and response data, and current events/natural disasters that could impact estimates.

We also must be aware that using administrative data would impact the entire survey life cycle, require significant resources, and without thorough testing would be a high risk to the program. A single change has a ripple effect and could result in unexpected consequences if not thought through completely.

Administrative records provide a vast amount of data that the ACS program is eager to tap into. With the right resources, we can overcome the many challenges associated with using it. However, to use administrative data we must be cautious and well prepared in our approach. This research has opened our eyes to some of the challenges and provides input to help pave the path towards our ultimate goal of making use of administrative data to improve data quality and to reduce some of the burden placed on our respondents.

8. REFERENCES

- Brick, J. M., and Williams, D. (2013), "Explaining Rising Nonresponse Rates in Cross-Sectional Surveys," *The ANNALS of the American Academy of Political and Social Science*, 645, 36-59.
- Chappell, G., and Obenski, S. (2014). "Fiscal Year 2014 Content Review Results," *American Community Survey Fiscal Year 2014 Content Review*. Washington, DC: U.S. Census Bureau.
- Dillon, M. (forthcoming). "Preliminary Research for Replacing or Supplementing the Acreage, Number of Rooms and Bedrooms, Tenure, Property Value, & Real Estate Taxes Questions on the American Community Survey with Administrative Records," Washington, DC: U.S. Census Bureau.
- Kingkade, W. (2013). "Self-Assessed Housing Values in the American Community Survey: An Exploratory Evaluation Using Linked Real Estate Records," *Journal of Survey Statistics and Methodology*.
- Moore, B. (2016). "Preliminary Research for Replacing or Supplementing the Year Built Question on the American Community Survey with Administrative Records," Washington, DC: U.S. Census Bureau.
https://www.census.gov/content/dam/Census/library/working-papers/2015/acs/2015_Moore_02.pdf.
- Mule, V. (forthcoming). "Year Structure Built Question Removal Analysis," Washington, DC: U.S. Census Bureau.
- Ruggles, P. (2015). "Review of Administrative Data Sources Relevant to the American Community Survey," Prepared for the U.S. Census Bureau, January 31, 2015.
- U.S. Census Bureau, American Community Survey Office. (2017). "Agility in Action 2.0: A Snapshot of Enhancements to the American Community Survey," available at: <https://www.census.gov/content/dam/Census/programs-surveys/acs/operations-and-administration/2015-16-survey-enhancements/Agility%20in%20Action%20v2.0.pdf>.
- U.S. Census Bureau. (2014). "American Community Survey Design and Methodology," available at: <http://www.census.gov/programs-surveys/acs/methodology/design-and-methodology.html>.
- Wagner, D., and Layne, M. (2014). "The Person Identification Validation System (PVS): Applying the Center for Administrative Records Research and Applications' (CARRA) Record Linkage System," Washington, DC: U.S. Census Bureau.

Chun, A., Schouten, B., and Wagner, J. (2017). "JOS Special Issue on Responsive and Adaptive Survey Design: Looking Back to See Forward – Editorial," *Journal of Official Statistics*, Vol. 33, No. 3, 2017, pp. 571–577, available at: <http://dx.doi.org/10.1515/JOS-2017-0027>.

Appendix A. Distribution of Simulated vs. Published Estimate Differences by Geographic Area and Survey Item

<i>All Topics</i>	US	States	Counties	Places	PUMAs
Total Estimates	575	29,900	477,250	342,700	1,367,350
-20.01% or less	4.3	6.0	3.5	2.8	2.7
-20.00% to -9.01%	12.7	7.4	2.6	2.0	1.9
-9.00% to -6.01%	8.9	5.2	1.5	1.0	1.0
-6.00% to -3.01%	11.1	5.5	1.2	0.8	0.8
-3.00% to -0.01%	10.6	2.9	0.3	0.2	0.2
Not Statistically Significant	20.9	55.6	45.2	32.9	38.3
0.01% to 3.00%	9.0	1.9	0.5	0.3	0.3
3.01 to 6.00%	4.7	1.9	0.4	0.3	0.3
6.01 to 9.00%	3.8	1.9	0.4	0.3	0.3
9.01% to 20.00%	11.3	4.4	1.1	0.8	0.8
20.01% or more	2.6	3.8	2.1	1.6	1.6
Not Calculated	0.0	3.3	41.2	57.0	51.9

Source: 2015 ACS Housing Administrative Record Simulation

Note: 0.10 alpha used for significance testing

<i>Year Structure Built</i>	US	States	Counties	Places
Total Estimates	238	12,376	197,540	141,848
-20.0% or less	0.4	2.1	2.3	2.3
-19.9% to -9.0%	5.5	4.3	1.3	1.1
-8.9% to -6.0%	7.6	3.4	0.5	0.3
-6.00% to -3.01%	8.4	4.5	0.4	0.2
-3.00% to -0.01%	18.9	3.4	0.2	0.2
Not Statistically Significant	29.8	68.8	49.3	38.7
0.01% to 2.9%	16.0	2.7	0.4	0.3
3.0% to 5.9%	5.0	2.4	0.4	0.3
6.0% to 8.9%	3.8	2.0	0.4	0.2
9.0% to 19.9%	3.8	3.1	0.8	0.7
20.0% or more	0.8	1.3	0.7	0.6
Not Calculated	0	2.2	43.3	55.1

Source: 2015 ACS Housing Administrative Record Simulation

Note: 0.10 alpha used for significance testing

<i>Acreage</i>	US	States	Counties	Places
Total Estimates	21	1,092	17,430	12,516
-20.0% or less	4.8	4.9	8.1	9.0
-19.9% to -9.0%	0.0	2.7	3.3	0.2
-8.9% to -6.0%	0.0	3.7	1.2	0.0
-6.00% to -3.01%	14.3	4.6	0.6	0.1
-3.00% to -0.01%	9.5	4.5	0.6	0.3
Not Statistically Significant	33.3	56.7	72.4	34.4
0.01% to 2.9%	23.8	12.2	6.5	4.9
3.0% to 5.9%	0.0	2.8	2.5	2.2
6.0% to 8.9%	9.5	1.6	0.6	0.3
9.0% to 19.9%	0.0	2.9	0.5	0.1
20.0% or more	0.0	2.0	0.7	0.1
Not Calculated	0	1.4	2.9	48.3

Source: 2015 ACS Housing Administrative Record Simulation

Note: 0.10 alpha used for significance testing

<i>Property Value</i>	US	States	Counties	Places
Total Estimates	292	15,184	242,360	174,032
-20.0% or less	6.8	10.6	6.2	4.2
-19.9% to -9.0%	20.9	11.0	3.9	3.0
-8.9% to -6.0%	11.3	7.0	2.5	1.8
-6.00% to -3.01%	13.4	6.4	1.9	1.4
-3.00% to -0.01%	1.7	2.0	0.3	0.2
Not Statistically Significant	13.4	45.0	36.7	24.3
0.01% to 2.9%	2.1	0.4	0.1	0.0
3.0% to 5.9%	4.1	1.3	0.2	0.1
6.0% to 8.9%	3.4	1.7	0.3	0.2
9.0% to 19.9%	18.5	5.0	0.7	0.4
20.0% or more	4.5	4.8	1.8	1.1
Not Calculated	0.0	4.6	45.6	63.4

Source: 2015 ACS Housing Administrative Record Simulation

Note: 0.10 alpha used for significance testing

<i>Real Estate Tax</i>	US	States	Counties	Places
Total Estimates	24	1,248	19,920	14,304
-20.0% or less	12.5	12.6	14.5	13.0
-19.9% to -9.0%	0.0	8.5	4.8	3.4
-8.9% to -6.0%	0.0	4.9	2.0	1.5
-6.00% to -3.01%	4.2	8.2	2.2	1.4
-3.00% to -0.01%	33.3	6.3	1.0	0.3
Not Statistically Significant	16.7	38.3	60.2	62.2
0.01% to 2.9%	12.5	3.0	0.8	0.2
3.0% to 5.9%	8.3	4.0	1.3	1.1
6.0% to 8.9%	4.2	2.9	1.4	1.2
9.0% to 19.9%	8.3	9.0	5.8	4.3
20.0% or more	0.0	2.5	4.5	5.0
Not Calculated	0	0	1.4	6.3

Source: 2015 ACS Housing Administrative Record Simulation

Note: 0.10 alpha used for significance testing