

# Data Collection, Ratio Estimation, and Discrepancies: Analyzing ACS Respondent Counts with IPUMS USA\*

Yi Tang

Jin Zhang

Siyuan Lu

October 3, 2024

The number of respondents in each state with a doctorate as their highest level of education is examined in this document using data from the 2022 ACS IPUMS. Using data from California, we use the ratio estimators approach to estimate the total number of respondents in each state.

## 1 Introduction

We use R packages (R Core Team 2023) to clean and analyze the dataset, including libraries from `haven` (Wickham, Miller, and Smith 2023), `tidyverse` (Wickham et al. 2019), `knitr` (Xie 2014) and `labelled` (Larmarange 2024). The data we used is from IPUMS (Ruggles et al. 2021).

## 2 A brief overview of the ratio estimators approach

The ratio estimator is a method used in survey sampling to improve estimation accuracy by leveraging a known relationship between two variables. This method calculates the ratio of a particular attribute to the total population for a known subset. The ratio is then applied to other subsets to approximate totals, assuming similar correlations exist across the entire population. It is especially helpful when the precise population size is unknown but a sample yields proportional connections.

---

\*Code and data are available at: <https://github.com/YiTang2/IPUMS-USA-data.git>

### 3 Estimates and the Actual Number of Respondents

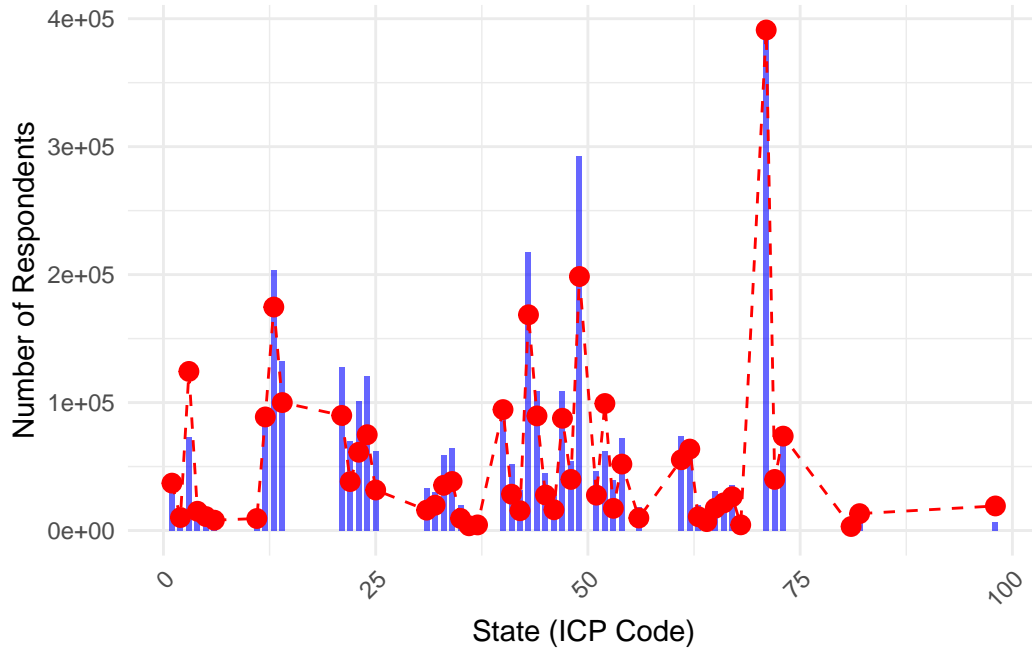


Figure 1: Comparison of Actual and Estimated Respondents by State

From Figure 1, we calculate and compare the actual and estimated total respondent counts for each state using a ratio estimation method. It assumes the ratio of doctoral degree holders to total respondents in California applies to other states. The actual counts are derived from the dataset, while the estimated counts are based on this ratio. The final figure highlights discrepancies between actual and estimated counts, helping evaluate the accuracy and limitations of the ratio estimation approach.

### 4 Explanation of why they are different

The estimated total number of respondents in each state using the ratio estimator can differ from the actual count due to several factors:

1. **Assumption of Similarity:** The ratio estimator makes the assumption that the percentage of Californians with doctorates is typical of other states; however, due to a variety of reasons, including economic, demographic, and educational infrastructure, there are considerable variations in educational attainment.
2. **Sampling Variability:** If based on a sample, random variability can impact the ratio and estimation accuracy.

3. Non-Uniform Distribution: Educational attainment isn't evenly distributed across the U.S., so any areas' ratio may not apply to other states.
4. Bias: When relationships are constant over all domains, the ratio technique performs well. The estimations will be skewed if the ratio is impacted by unobserved factors.

These factors explain why using the ratio estimator across diverse states often leads to differences from actual numbers.

## Appendix

### 5 Instructions on obtaining the data

To collect data from [IPUMS USA](#), we first navigated to the IPUMS website and selected “IPUMS USA.” We then clicked on “Get Data” and chose the “2022 ACS” sample under the “SELECT SAMPLE” section. To gather state-level data, we selected “HOUSEHOLD” followed by “GEOGRAPHIC” and added “STATEICP” to the cart. For individual-level data, we went to the “PERSON” section and added “EDUC” to the cart. After reviewing our selections by clicking “VIEW CART,” we proceeded to “CREATE DATA EXTRACT.” We changed the “DATA FORMAT” to “.csv” and clicked “SUBMIT EXTRACT.” After logging in with the account, we received an email when the extract was ready for download. Finally, we downloaded the file to use in RStudio.

## References

- Larmarange, Joseph. 2024. *labelled: Manipulating Labelled Data*. <https://CRAN.R-project.org/package=labelled>.
- R Core Team. 2023. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Ruggles, Steven, Sarah Flood, Sophia Foster, Ronald Goeken, Jose Pacas, Megan Schouweiler, and Matthew Sobek. 2021. “IPUMS USA: Version 11.0.” Minneapolis, MN: IPUMS. <https://doi.org/10.18128/d010.v11.0>.
- Wickham, Hadley, Mara Averick, Jennifer Bryan, Winston Chang, Lucy D’Agostino McGowan, Romain François, Garrett Grolemund, et al. 2019. “Welcome to the tidyverse.” *Journal of Open Source Software* 4 (43): 1686. <https://doi.org/10.21105/joss.01686>.
- Wickham, Hadley, Evan Miller, and Danny Smith. 2023. *haven: Import and Export ‘SPSS’, ‘Stata’ and ‘SAS’ Files*. <https://CRAN.R-project.org/package=haven>.
- Xie, Yihui. 2014. “knitr: A Comprehensive Tool for Reproducible Research in R.” In *Implementing Reproducible Computational Research*, edited by Victoria Stodden, Friedrich Leisch, and Roger D. Peng. Chapman; Hall/CRC.