

Risk assessment procedures for the 2020 U.S. Census

Steven Ruggles
Lara Cleveland
David Van Riper

Institute for Social Research and Data Innovation
University of Minnesota

UNECE/Eurostat Expert Meeting on Statistical Data Confidentiality (1 December 2021)

Problem statement

- US Census Bureau justified its transition to a formally private disclosure avoidance system based on the results of a database reconstruction and re-identification experiment.

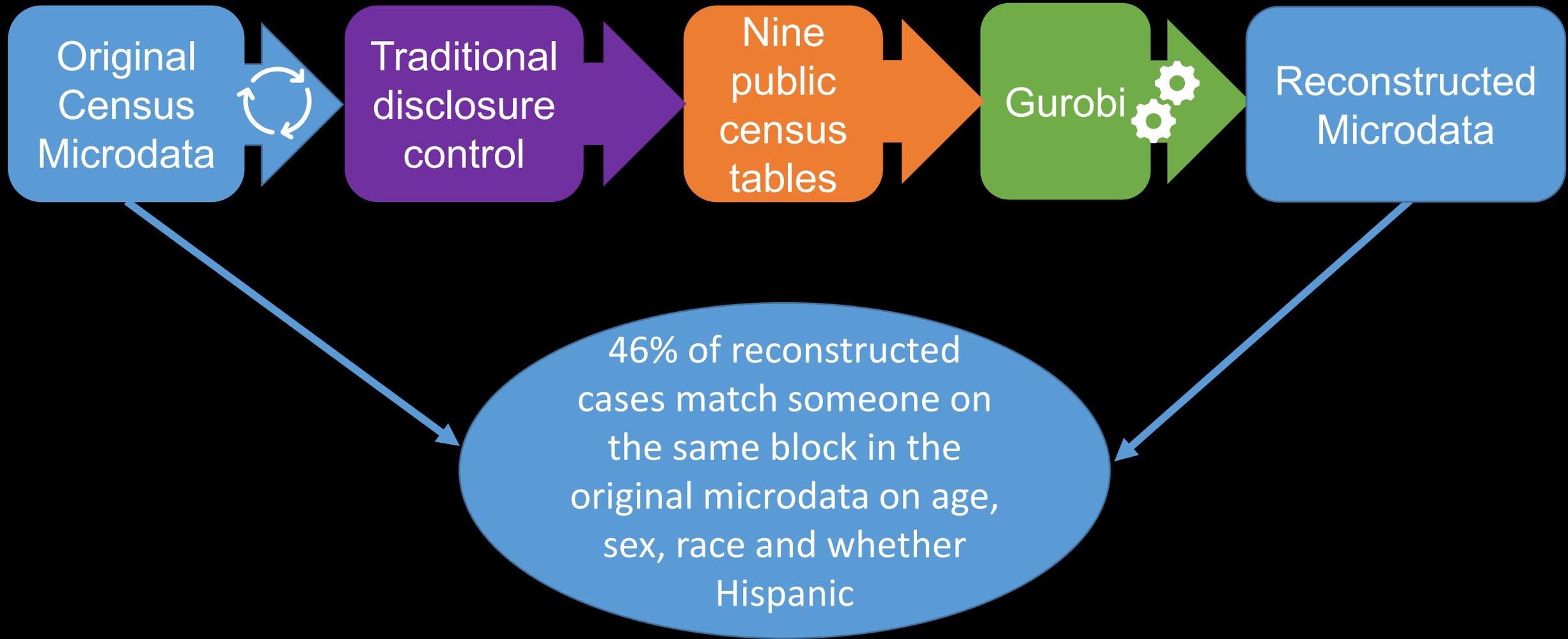
Problem statement

- US Census Bureau justified its transition to a formally private disclosure avoidance system based on the results of a database reconstruction and re-identification experiment.
- How does the Bureau's results compare with a similar (reconstruction) experiment based on chance?

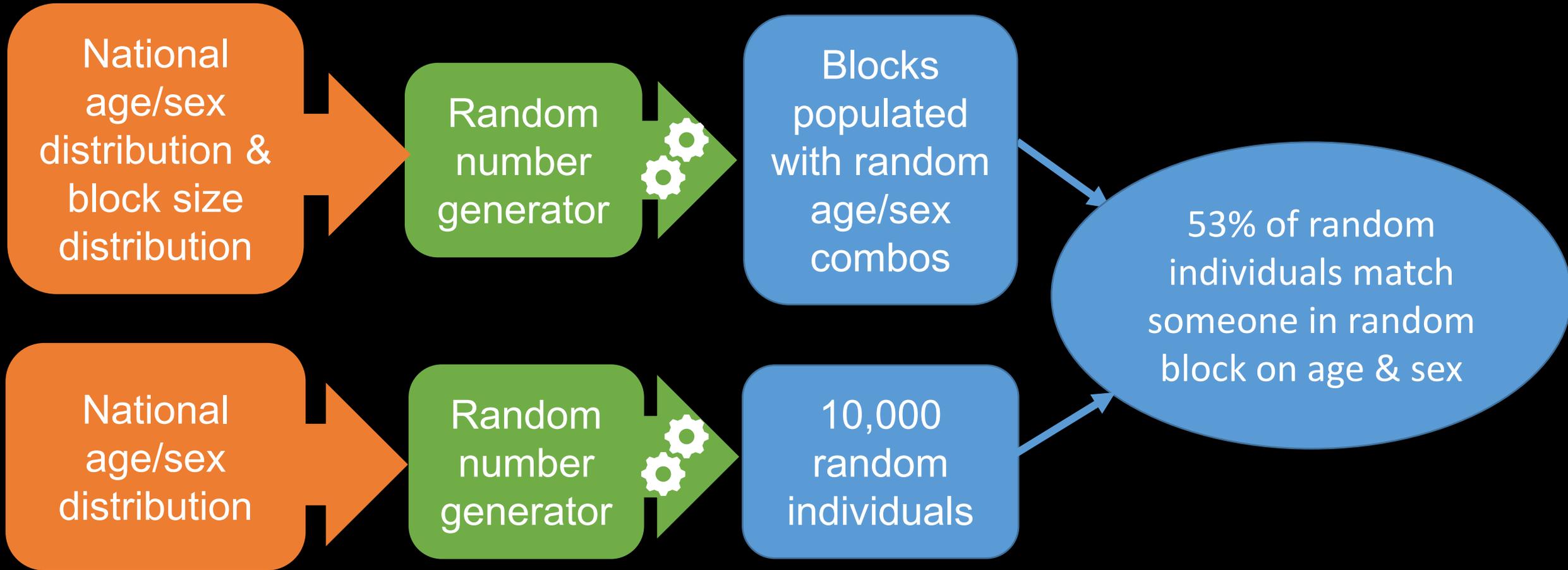
Reconstruction Results

- 46% of Bureau's reconstructed cases match someone in confidential data on census block, age, sex, race, Hispanic
- 41% of our random individuals match someone in random block on age, sex, race, Hispanic ethnicity

Census Bureau Database Reconstruction Experiment



How many matches would be expected by chance?



Adding Race and Ethnicity

53% of random individuals match someone in random block on age & sex

Assign modal race and ethnicity to all persons in each block

41% of random individuals match someone in random block on age, sex, race and whether Hispanic

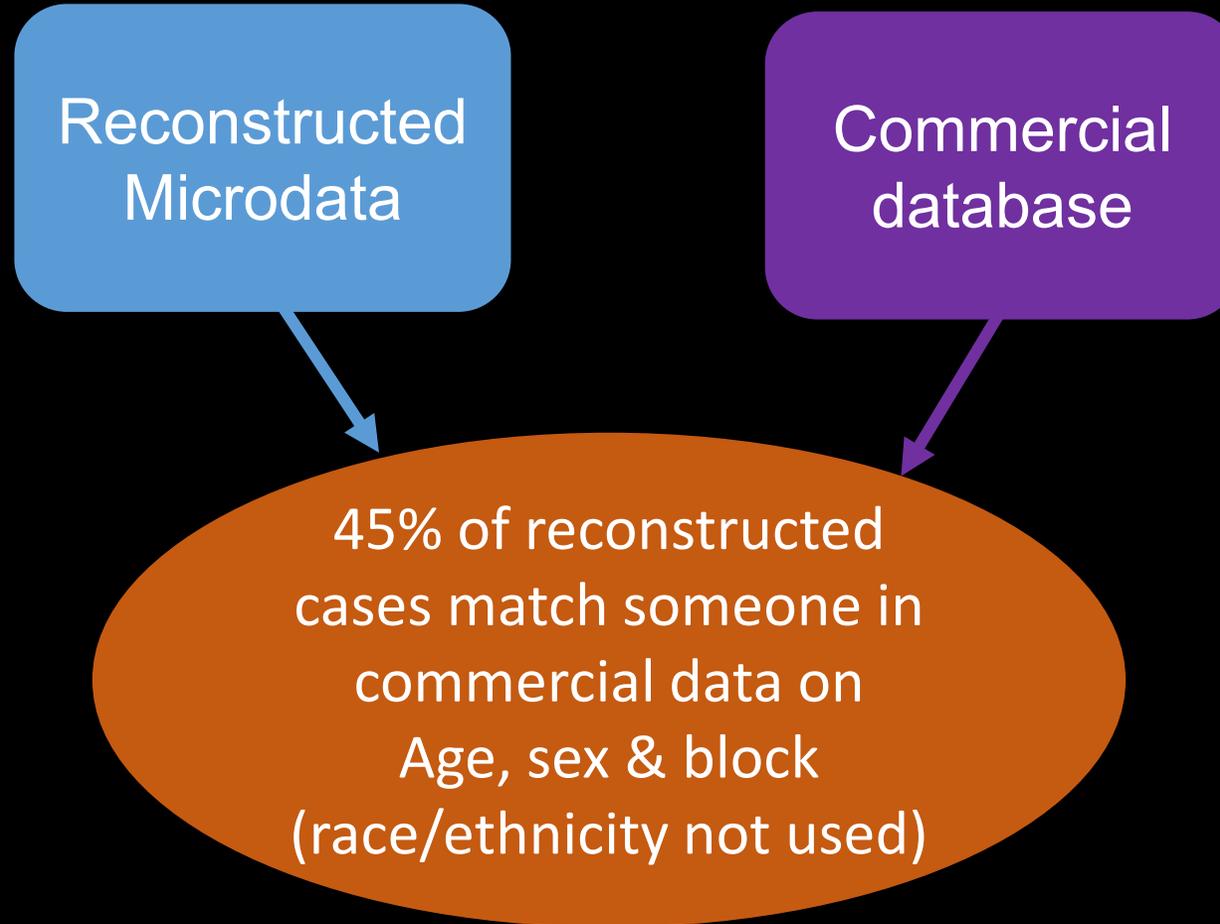
Database reconstruction experiment

- The Census Bureau “Reconstruction” matched on block, age, sex, and race/ethnicity in **46%** of cases
- A random number generator combined with a simple assignment rule for race and ethnicity yields a comparable match rate of **41%**
- Despite the Census Bureau’s massive investment of resources and computing power, their database reconstruction does not perform much better than a roll of the dice
- This is analogous to a clinical trial in which the treatment and the placebo produce virtually the same outcome.

Re-identification Results

- 45% of Bureau's reconstructed cases match someone in commercial data on census block, age, and sex
- 38% of names harvested from commercial database match a name on same block in confidential data
- 17% (38% of 45%) of Census respondents "re-identified"

Census Bureau rei-dentification experiment

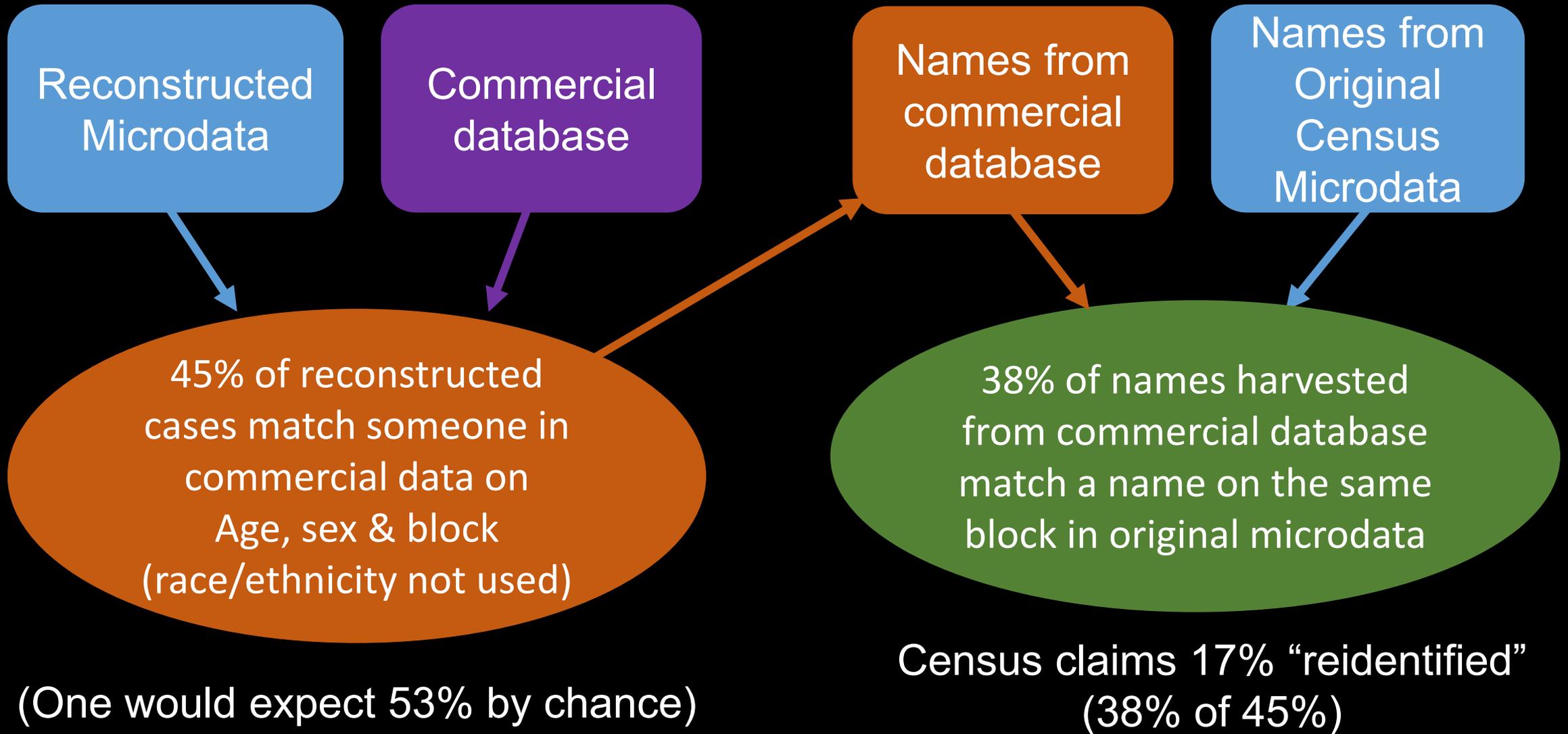


(One would expect 53% by chance)

The re-identification experiment

- Within each block the Census Bureau searched a commercial database for people who matched the age and sex of each case of their reconstructed database.
- In just 45% of cases was there at least one person in the commercial data who matched the age, sex and block number of a row of the hypothetical database.
- This is a lower match rate than the 53% one would expect by chance, probably because of the poor quality of the commercial data.

Census Bureau re-identification experiment



The re-identification experiment

- The Census Bureau then harvested the names of the 45% of cases in the commercial database that matched on age and sex and searched for them in the original census data.
- They found that 38% of the 45% were actually present on the block.
- They claimed to have “reidentified” 16.85% of the population (38% of 45%).

Census Bureau re-identification experiment

Among the 45% of cases where the reconstructed data match someone in the commercial data on age, sex, and block, 38% match a name on the same block in original microdata

- There is no null comparison
- One would expect that people residing on a particular block in the commercial data would also be enumerated in the census.
- Is 38% high or low?

Evaluating the re-identification experiment

- The Census Bureau could easily evaluate the re-identification by matching the names of people randomly selected from the commercial database to persons in the 2010 census living on the same census block, without any reference to the Census Bureau's database reconstruction.
- If the 38% match rate on names for the reconstructed population is no higher than the match rate for a randomly selected subset of the commercial data, it would mean the database reconstruction had no effect on re-identification risk.
- Because they made no null comparison, the re-identification results cannot be interpreted.

Summary: Database reconstruction experiment

- “Correctly” guessed age, sex, race and Hispanic ethnicity for 46% of persons on each block
- Great majority of those correct guesses would be expected by chance
- An outside attacker would have no means of determining which “reconstructed” records were true

Summary: Re-identification experiment

- The re-identification experiment shows only that some people found in a commercial database are also listed on the same block in the census.
- It fails to show that using the published tabulations had any impact whatsoever on the match rate between the commercial database and the original enumeration.

Consequences of disclosure panic

The disclosure avoidance system introduces unacceptable levels of error for many applications of the census.

SOCIAL SCIENCES

The use of differential privacy for census data and its impact on redistricting: The case of the 2020 U.S. Census

Christopher T. Kenny¹, Shiro Kuriwaki², Cory McCartan³, Evan T. R. Rosenman⁴, Tyler Simko¹, Kosuke Imai^{1,3*}

Census statistics play a key role in public revealing individual information, many differential privacy, which add noise to census statistics must be postprocessed Bureau's latest disclosure avoidance system electoral districts. We find that the DAS system precincts, yielding unpredictable racial "One Person, One Vote" standard as current race and ethnicity. Our findings underscore

INTRODUCTION

In preparation for the official release of U.S. Census Bureau has developed a disclosure avoidance system (DAS) to prevent Census responses from individuals (1). The DAS is based on differential privacy, which adds a certain amount of random noise to the data. The Bureau has been required by law to disclose information about Census participants implemented disclosure avoidance methods their decision to incorporate differential privacy subsequent postprocessing steps in the 2020 Census in the DAS, has been controversial. Some concerns about the potential negative impact on policy and social science research, which data (2–6).

The U.S. decennial census serves as a case study on the impact of differential privacy on the drawing of legislative districts, deterring federal funds for more than a hundred years are extensively analyzed by social scientists and international organizations, including the United Kingdom, and Australia, have adopted the adoption of differential privacy technology to its decennial census, the U.S. Census Bureau has adopted differential privacy as their "privacy definition" for data release on commuting patterns (12) considering adopting a similar approach for the American Community Survey (13).

It is a common misconception that a disclosure system only involves injecting random noise (14)

Copyright © 2021
The Authors, some
rights reserved;
exclusive licensee
American Association
for the Advancement
of Science. No claim to
original U.S. Government

- Science Advances
- Pop. Res. Policy Rev
- Socius
- PNAS

Population Research and Policy Review
<https://doi.org/10.1007/s11113-021-09664-5>

RESEARCH BRIEFS



Differential Privacy and the Accuracy of County-Level Net Migration Estimates

Richelle L. Winkler¹ · Jaclyn
David Egan-Robertson³

Received: 5 March 2021 / Accepted: 15 March 2021
© The Author(s) 2021

Abstract

Each decade since the 1950s, net migration estimates by age, sex, and race as starting and ending points for drawing thousands of times and widely used. The 2020 Census should allow for greater accuracy of new estimates using differential privacy (DP) disclosure avoidance system. This brief estimates the impact of DP on net migration estimates. Using a simulation, we construct a hypothetical set of estimates to compare them to published estimates. Findings show that bias



SOCIUS

Original Article

Differential Privacy in the 2020 Census Will Distort COVID-19 Rates

Mathew E. Hauer¹

Abstract

Scholars rely on accurate data to understand the (COVID-19) pandemic, with deaths at older ages. Population estimates subject to noise infusion from differential privacy. Using a simulation, we introduce substantial distortions, hindering our ability to understand with fewer than 1,000 persons. The U.S. Census Bureau should consider differential privacy's distortions

ASA
American Sociological Association

Socius: Sociological Research for a Dynamic World
Volume 7: 1–6
© The Author(s) 2021
Article reuse guidelines:
sagepub.com/journals-permissions
DOI: 10.1177/2378023121994014
srd.sagepub.com



How differential privacy will affect our understanding of health disparities in the United States

Alexis R. Santos-Lozada^{a,1} , Jeffrey T. Howard^b , and Ashton M. Verdery^c

^aDepartment of Human Development and Family Studies, The Pennsylvania State University, University Park, PA 16802; ^bDepartment of Public Health, University of Texas at San Antonio, San Antonio, TX 78249; and ^cDepartment of Sociology and Criminology, The Pennsylvania State University, University Park, PA 16802

Edited by Kenneth W. Wachter, University of California, Berkeley, CA, and approved April 20, 2020 (received for review February 27, 2020)

The application of a currently proposed differential privacy algorithm to the 2020 United States Census data and additional data products may affect the usefulness of these data, the accuracy of estimates and rates derived from them, and critical knowledge about social phenomena such as health disparities. We

more is known about the impact of DP on important uses of census data, particularly in light of concerns discussed at the workshop on December 2019 (4).

Why are there concerns about privacy? Title 13 of the United States Code imposes heavy obligations on the US Census Bureau

Findings of recent publications include:

- Systematically undercounts the population in mixed-race and mixed-partisan precincts.
- Introduces substantial distortion in COVID-19 mortality rates, sometimes causing mortality rates to exceed 100 percent.
- Net migration estimates by five-year age groups would only be accurate enough for use in about half of counties.

Conclusions

- The Census Bureau is now in the midst of a radical escalation of disclosure control justified by an unsuccessful database reconstruction experiment.
- The actual threat posed by the Census Bureau's reconstruction is similar to the threat posed by randomly guessing respondent characteristics.

Conclusions

- The database reconstruction experiment failed to demonstrate a disclosure threat and cannot justify degradation of the nation's statistical infrastructure.
- Reducing data quality is likely to backfire, as it will erode public trust and weaken the core justification for census data collection

Thank you.

Questions?