

US Census data vulnerable to attack without enhanced privacy measures, shows study

February 21 2023, by Devorah Fischler



Credit: CC0 Public Domain

Computer scientists at the University of Pennsylvania School of Engineering and Applied Science have designed a "reconstruction attack" that proves U.S. Census data is vulnerable to exposure and theft.

Aaron Roth, Henry Salvatori Professor of Computer & Cognitive Science in Computer and Information Science (CIS), and Michael Kearns, National Center Professor of Management & Technology in CIS, led a recent *PNAS* study demonstrating that statistics released by the

U.S. Census Bureau can be reverse engineered to reveal protected information about individual respondents.

With computing power no stronger than that of a commercial laptop and algorithm design drawn from machine learning fundamentals, the research team established risks to the [privacy](#) of the U.S. population.

The study stands out for being the first of its kind to determine a baseline for unacceptable susceptibility to exposure. In addition, it proves that an attack has the means to ascertain the likelihood that a reconstructed record corresponds to the data of a real person, making it even more probable that this kind of attack could render respondents vulnerable to identity theft or discrimination.

The findings sharpen the stakes of one of the digital era's most significant debates in [public policy](#).

"Over the last two decades it has become clear that practices in widespread use for [data privacy](#)—anonymizing or masking records, coarsening granular responses or aggregating individual data into large-scale statistics—do not work," says Kearns. "In response, [computer scientists](#) have created techniques to provably guarantee privacy."

"The [private sector](#)," adds Roth, "has been applying these techniques for years. But the Census' long-running statistical programs and policies have additional complications attached."

For example, the Census is constitutionally mandated to carry out a full population survey every ten years. This data is used for key political, economic and social functions: apportioning House seats, drawing district boundaries, determining federal funding amounts for state and local uses, financing disaster relief, welfare programs, infrastructural expansion and more. The data also provides vital tools for demographic

researchers in government and academia.

While Census information is public, strict laws govern the privacy of individual data. To this end, publicly available statistics aggregate each respondent's survey answers, reflecting the population with mathematical precision without directly revealing individuals' personal information.

The problem is that these aggregated statistics are a lock that can be picked, and all it takes are the right tools. Attackers can use these aggregates to reverse engineer sets of records consistent with confirmed statistics, a process known as "reconstruction."

In response to these risks, the Census ran its own internal reconstruction attack between the 2010 and 2020 surveys to gauge the need for a change in reporting. The findings merited a Census overhaul of confidentiality measures, and a decision to implement a provable protection technique known as "differential privacy."

Differential privacy conceals individual data while maintaining the integrity of the larger data set. Cynthia Dwork, Gordon McKay Professor of Computer Science at Harvard University and Roth and Kearns' collaborator on the study, co-invented the technique in 2006. Dwork's work is significant for being the first to provide "privacy" with a mathematically rigorous definition.

Rather than report statistics that transparently reflect true responses, differential privacy introduces strategic amounts of false data, known as "noise," which consists of randomly generated positive or negative numbers averaging out to roughly zero. At large scales, the noise's interference in statistical correctness is negligible. But complications do arise in demographic statistics describing small populations, where noise has a relatively larger effect on reporting.

The trade-off between accuracy and privacy is complex.

Certain social scientists have argued that the Census practice of publishing aggregate statistics poses no inherent risk. While acknowledging that individual records are susceptible to reconstruction through educated guessing or comparison with public documentation, this camp maintains that the Census' decision to implement differential privacy is a poor one, claiming the success rate for reconstructing individual records is no better than random chance.

But Roth and Kearns' work has proven otherwise, running queries that function like Venn diagrams with hundreds of thousands of overlapping ovals. These overlaps signal the likelihood of accuracy in possible data configurations that match publicly available statistics, allowing for attackers to outperform any possible baseline for random chance.

"What's novel about our approach is that we show that it's possible to identify which reconstructed records are most likely to match the answers of a real person," says Kearns. "Others have already demonstrated it's possible generate real records, but we are the first to establish a hierarchy that would allow attackers to, for example, prioritize candidates for identity theft by the likelihood their records are correct."

On the matter of complications posed by adding error to statistics that play such a significant role in the lives of the U.S. population, the researchers are realistic.

"The Census is still working out how much noise will be useful and fair in order to balance the trade-off between accuracy and privacy. And, in the long run, it may be that public policymakers decide that the risks posed by non-noisy statistics are worth the transparency," says Roth.

But when it comes to absolute guarantees for individual data protection, Roth and Kearns both affirm beyond a doubt: "Differential privacy is the only game in town."

More information: Travis Dick et al, Confidence-ranked reconstruction of census microdata from published statistics, *Proceedings of the National Academy of Sciences* (2023). [DOI: 10.1073/pnas.2218605120](https://doi.org/10.1073/pnas.2218605120)

Provided by University of Pennsylvania

Citation: US Census data vulnerable to attack without enhanced privacy measures, shows study (2023, February 21) retrieved 5 May 2024 from <https://techxplore.com/news/2023-02-census-vulnerable-privacy.html>

This document is subject to copyright. Apart from any fair dealing for the purpose of private study or research, no part may be reproduced without the written permission. The content is provided for information purposes only.