



Confidence-ranked reconstruction of census records from aggregate statistics fails to capture privacy risks and reidentifiability

David Sánchez^{a,1} , Josep Domingo-Ferrer^a, and Krishnamurthy Muralidhar^b

Dick *et al.* (1) describe a method to reconstruct record prototypes from aggregate statistics of the US Decennial Census data, where we call prototype a record present with some multiplicity (i.e., number of repetitions) in the original data D . The proposed method ranks the reconstructed prototypes by how frequently they appear in multiple reconstructions.

The authors evaluate the effectiveness of their method by measuring the rate of the top k -ranked reconstructed prototypes that were present in a synthetically generated version of D . The reported results show that this rate is near to 1 for small values of k , which the authors interpret as a privacy risk on the subjects to whom those records refer. However, since the rate considers the number of prototypes in D rather than the actual number of records, the authors are neglecting the multiplicity of each prototype in D , which is key to assessing privacy risks: A prototype appearing, say, 10 times in D means that 10 individuals share the same record values for the considered attributes and, therefore, are intrinsically protected against reidentification (2, 3) because “our privacy is protected to the extent that we blend in with the crowd” (4). In terms of the well-known k -anonymity model (2), this means that unequivocal reidentification of any of those 10 individuals is not possible and that the probability of correct reidentification is just $1/10$. It can be noted that the authors’ method would still (wrongly) claim high rates if D was already protected via k -anonymity, i.e., all prototypes in it had multiplicity k . Therefore, the reported rates do not capture the privacy risks in D .

The authors’ method is incapable of ascertaining the multiplicity of the reconstructed prototypes and, therefore, cannot assess whether they appear only once in D and can therefore lead to reidentification if linked with external identified sources. Besides, since their ranking is based on the frequency of prototypes in multiple reconstructions, the top k -ranked prototypes can be expected to be the k most common prototypes in D ; that is, those with the greatest number of repetitions and, thus, those intrinsically more private. Conversely, outlying (i.e., less frequent) records in D , which

actually correspond to individuals at highest risk, are unlikely to be among the top k . Hence, outlying records are paradoxically beyond the reach of the proposed method. Therefore, contradictory to the authors’ claims, their method is *ineffective* to i) assess the privacy risks in D , ii) detect the most privacy-sensitive records in D , and iii) guide targeted reidentification attacks on D .

The authors refer to their method as a reconstruction attack, which implies that it may compromise D , and suggest that the produced ranking can be used by an adversary for “identity theft.” As a conclusion, they raise “sober warnings on the privacy risks of releasing precise aggregate statistics of a dataset” and suggest using differential privacy, a privacy-enhancing technique that can severely damage the utility of the published data (5). Given the above, this conclusion is unwarranted and inappropriate.

ACKNOWLEDGMENTS. Partial support to this work has been received from the European Union’s Horizon 2020 research and innovation program under grant agreement no. 871042 (“SoBigData++”), the Spanish Ministry for Science and Innovation (MCIN)/Spanish National Research Agency (AEI)/10.13039/501100011033/European Fund for Economic and Regional Development (FEDER), European Union (EU) (project PID2021-123637NB-I00, “Co-utile decentralized computing (CURLING)”), and the Government of Catalonia (Catalan Institution for Research and Advanced Studies (ICREA) Acadèmia Prizes to J.D.-F. and D.S. and grant 2021SGR-00115). We are with the United Nations Educational, Scientific and Cultural Organization (UNESCO) Chair in Data Privacy, but the views in this paper are their own and are not necessarily shared by UNESCO.

Author affiliations: ^aUniversitat Rovira i Virgili, Department of Computer Engineering and Mathematics, UNESCO Chair in Data Privacy, Center for Cybersecurity Research of Catalonia, Tarragona, Catalonia 43007, Spain; and ^bDepartment of Marketing and Supply Chain Management, University of Oklahoma, Norman, OK 73019

Author contributions: D.S., J.D.-F., and K.M. analyzed data; and wrote the paper.

The authors declare no competing interest.

Copyright © 2023 the Author(s). Published by PNAS. This article is distributed under Creative Commons Attribution-NonCommercial-NoDerivatives License 4.0 (CC BY-NC-ND).

¹To whom correspondence may be addressed. Email: david.sanchez@urv.cat.

Published April 24, 2023.

1. T. Dick *et al.*, Confidence-ranked reconstruction of census microdata from published statistics. *Proc. Natl. Acad. Sci. U.S.A.* **120**, e2218605120 (2023).
2. P. Samarati, Protecting respondents identities in microdata release. *IEEE Trans. Knowl. Data Eng.* **13**, 1010–1027 (2001).
3. S. Chawla, C. Dwork, F. McSherry, A. Smith, H. Wee, “Toward Privacy in Public Databases” in *Theory of Cryptography. TCC 2005. Lecture Notes in Computer Science*, J. Kilian, Ed. (Springer, Berlin, Heidelberg, 2005), vol. 3378.
4. J. Gehrke, M. Hay, E. Lui, R. Pass, “Crowd-Blending Privacy” in *Advances in Cryptology – CRYPTO 2012. Lecture Notes in Computer Science*, R. Safavi-Naini, R. Canetti, Eds. (Springer, Berlin, Heidelberg, 2012), vol. 7417.
5. V. J. Hotz *et al.*, Balancing data privacy and usability in the federal statistical system. *Proc. Natl. Acad. Sci. U.S.A.* **119**, e2104906119 (2022).