Analyzing the impact of differential privacy on the accuracy of decennial census data

David Van Riper

vanriper@umn.edu

Brookings Institution September 26, 2019

IPUMS.ORG

Outline

- What is differential privacy?
- Applying differential privacy to data
- Implementing differential privacy for census
- Analyzing impact of differential privacy

WHAT IS DIFFERENTIAL PRIVACY?

IPUMS.ORG

Differential privacy is...

• A formal (mathematical) definition of privacy

 $\frac{\Pr[M(D) \in S]}{\Pr[M(D') \in S]} \le e^{\varepsilon}$

Differential privacy is not...

• An algorithm for disclosure control



Differential privacy is not...

- An algorithm for disclosure control
- An absolute guarantee against disclosure risk



APPLYING DIFFERENTIAL PRIVACY

IPUMS.ORG

"True" microdata

<u>Sex</u>	<u>School</u>	<u>Sex</u>	<u>School</u>
Male	Never	(Female	Never
Male	Never	x4 🖌 🗄	
Male	Never	Female	Never
🖌 Male	Attending	Female	Attending
v12 Male	Attending	×17 ₹	
	:	Female	Attending
Male	Attending	Female	Past
Male	Past	x31	
x33 〈	:	Female	Past
🗸 Male	Past		

Construct cross-tabs from "true" data

	School Attendance		
	Never	Attending	Past
Male	3	12	33
Female	4	17	31

Population = 100

Draw noise from Laplace distribution



Add noise to cross-tab

	S	chool Attendanc	е
	Never	Attending	Past
Male	3 - 1 = 2	12 + 0 = 12	33 + 1 = 34
Female	4 + 8 = 12	17 + 2 = 19	31 - 2 = 29

Sum = 108

Construct synthetic microdata



DIFFERENTIAL PRIVACY AND CENSUS

IPUMS.ORG

Differential privacy and census

POLICY DECISIONS

IPUMS.ORG

Policy decisions

- Global privacy loss budget (ε)
- Geographic levels
- Tables
- Invariants and constraints

- 1940 Geographic levels
 - Nation
 - State
 - County
 - Enumeration district

Privacy Loss Budget



• 1940 tables

Geographic Levels/Tables

- Voting age [2] x Hispanic [2] x Race [6]
- Households/group quarters type [8]

Privacy Loss Budget

- Detailed [192]
 - Voting age [2] x Hispanic [2] x Race [6] x GQ Type [8]

Noise Injection

IPUMS ORG

ANALYZING DIFFERENTIALLY PRIVATE 1940 CENSUS DATA

IPUMS ORG

- Census Disclosure Avoidance System (DAS) source code published in April 2019
 - 2020 Census DAS Development Team, 2019



- Census Disclosure Avoidance System (DAS) source code published in April 2019
 - 2020 Census DAS Development Team, 2019
- Implements their TopDown algorithm
 - Abowd et al, 2019



IPUMS.ORG

Fixed parameters

• Four geographic levels

– Nation, state, county, enumeration district



Fixed parameters

• Four geographic levels

– Nation, state, county, enumeration district

- Three queries / tables
 - Voting age Hispanic Race
 - Houshold group quarters
 - Detailed

IPUMS

Modifiable parameters

• Global privacy loss budget (ε)



Modifiable parameters

- Global privacy loss budget (ε)
- Fractional allocation to
 - Geographic levels
 - Tables

Modifiable parameters

- Global privacy loss budget (ε)
- Fractional allocation to
 - Geographic levels
 - Tables
- Number of runs

 Comparisons between "true" data (IPUMS 1940 complete-count) and differentially private data



- Differences in total population for counties and enumeration districts
- County-level African American population
- ED-level proportion of total population who identify as African American
- Index of dissimilarity (D)
- Multigroup entropy (H)

Key takeaways

- Geographic units with smaller populations are less accurate
- Small sub-populations are less accurate
- Bias for segregation metrics concerning

Differentially private datasets

CENSUS DAS



Global privacy loss budget (ε)

- 8 values: [0.25, 0.50, 0.75, 1.0, 2.0, 4.0, 6.0, 8.0]

• Four runs for each value of ε

Geographic levels	Fraction
Nation	0.25
State	0.25
County	0.25
Enumeration district	0.25

Tables	Fraction
Voting age—Hispanic – Race	0.675
Household – Group quarters	0.225
Detailed	0.1

Difference between IPUMS and Census DAS total population counts US counties (orange) and EDs (teal)



African American population under different levels of noise injection US counties



Percentage of population who is African American US enumeration districts



Source: Ruggles et al. (2018); US Census Bureau (2019)

Index of dissimilarity (D) under different levels of noise injection US counties



Source: Ruggles et al. (2018); US Census Bureau (2019)

Differentially private datasets

GEOGRAPHIC LEVELS

IPUMS.ORG

- Global privacy loss budget (ε)
 One value: 1.0
- One run

Geographic levels	Fraction*
Nation	0.85
State	0.05
County	0.05
Enumeration district	0.05

Tables	Fraction
Voting age – Hispanic – Race	0.675
Household – Group quarters	0.225
Detailed	0.1

Difference between IPUMS and Census DAS total population counts US counties (orange) and EDs (teal)

African American population - noise injection varies by geolevel US counties

Percentage of population who is African American - noise injection varies by geolevel US enumeration districts

Index of dissimilarity (D) - noise injection varies by geolevel US counties

Source: Ruggles et al. (2018); US Census Bureau (2019)

Multigroup entropy (H) - noise injection varies by geolevel US counties

Differentially private datasets

TABLES

- Global privacy loss budget (ε)
 One value: 1.0
- One run

Geographic levels	Fraction
Nation	0.25
State	0.25
County	0.25
Enumeration district	0.25

Tables	Fraction*
Voting age – Hispanic – Race	0.9
Household – Group quarters	0.05
Detailed	0.05

Difference between IPUMS and Census DAS total population counts US counties (orange) and EDs (teal)

African American population - noise injection varies by table US counties

Percentage of population who is African American - noise injection varies by table US enumeration districts

Source: Ruggles et al. (2018); US Census Bureau (2019)

Index of dissimilarity (D) - noise injection varies by query US counties

Source: Ruggles et al. (2018); US Census Bureau (2019)

Multigroup entropy (H) - noise injection varies by query US counties

• Diff. privacy less complicated than expected

- Diff. privacy less complicated than expected
- Fundamental importance of policy decisions

- Diff. privacy less complicated than expected
- Fundamental importance of policy decisions
- Largest impact on accuracy of small areas and small sub-populations

TPUMS

- Diff. privacy less complicated than expected
- Fundamental importance of policy decisions
- Largest impact on accuracy of small areas and small sub-populations
- Bias for segregation metrics concerning

TPUMS

References

2020 Census DAS Development Team. (2019) 2019. *Disclosure Avoidance System for the 2020 Census, End-to-End Release: Uscensusbureau/Census2020-Das-E2e*. Python. US Census Bureau. <u>https://github.com/uscensusbureau/census2020-das-e2e</u>.

Abowd, John, Daniel Kifer, Brett Moran, Robert Ashmead, Philip Leclerc, William Sexton, Simson Garfinkel, and Ashwin Machanavajjhala. 2019. "Census TopDown: Differentially Private Data, Incremental Schemas, and Consistency with Public Knowledge." US Census Bureau. <u>https://github.com/uscensusbureau/census2020-das-</u> <u>e2e/blob/master/doc/20190711 0945 Consistency for Large Scale Differentially Private Histograms.pdf</u>.