UNCERTAINTY IN DEMOGRAPHIC AND SOCIOECONOMIC DATA THE USE OF DIFFERENTIAL PRIVACY FOR DISCLOSURE CONTROL, AND ITS POTENTIAL IMPACT ON AGE AND **RACE/ETHNICITY COUNT**

DAVID VAN RIPER, UNIVERSITY OF MINNESOTA SETH SPIELMAN, UNIVERSITY OF COLORADO

SOURCES OF RACE/ETHNICITY DATA

American Community Survey

- ACS is a sample of housing units.
- 3.5M housing units are sampled, about 2M responses are collected (~60% response rate).
- ACS population counts are estimates, from a sample, and and as a result carry uncertainty.

Decennial Census

- The Decennial Census is complete enumeration of the US population.
- Published data have always obfuscated responses to prevent reidentification.
- For 2020 Census is adopting a formal privacy framework which injects noise into the data.

ACS ESTIMATE QUALITY: HISPANIC POP BY TRACT



Margin of error 100% of the estimate

15% of all Census Tracts have a margin of error that is 100% of more of the estimate.

Margin of error 50% of the estimate 53% of all Census Tracts have a margin error that is 50% of more of the estimate.

Margin of error 10% of the estimate 99% of all Census Tracts have a margin of error that is 10% of more of the estimate.

Note: Zero estimates included in figure but excluded from % of tracts > than 100%, 50%, 10% thresholds Figures truncate outliers, excluding the largest 1% of areas

ACS ESTIMATE QUALITY: TOTAL POP BY TRACT



Margin of error 100% of the estimate

Less than 1% of Census Tracts have a margin of error that is 100% or more of the estimate.

Margin of error 50% of the estimate

Less than 1% of all US Census Tracts have a margin error that is 50% of more of the estimate.

Margin of error 10% of the estimate

35% of all US Census Tracts have a margin of error that is 10% of more of the estimate.



African-American Population

Margin of Error and Estimate

Margin of Error and Estimate Native Hawaiian/Pacific Islander: All US Census Tracts 2019 ACS



American Indian/Alaska Native: All US Census Tracts 2019 ACS

ACS ESTIMATE QUALITY: HISPANIC POP BY COUNTY



Margin of error 100% of the estimate

3% of all counties have a margin of error that is 100% of more of the estimate.

Margin of error 50% of the estimate

10% of counties have a margin error that is 50% of more of the estimate.

Margin of error 10% of the estimate

13% of all counties have a margin of error that is 10% of more of the estimate.



A SOLUTION?

- Census tracts can be too small to provide reliable racial/ethnic group estimates. Counties can be too large, meaningful intra metropolitan sociospatial variation is lost.
- NYC Approach: Define custom "Neighborhood Tabulations Areas" bespoke groups of tracts.
- We've developed software and a website to generate "optimal" geographies based on user defined constraints (margin of error, population size).
- There are methods to refine census/ACS estimates using ancillary data sources.

| 10002000 | a https://reducinguncertainty.org/1858 |)/ger/occupied | | 0 ··· 0 | 2 ⊈ N CD ≰ Ф ≦ |
|---|--|--|--|--|---|
| ertainty | | | | 🖽 🛛 Deta Que | Ity Margin of Error Regionalization Credits |
| | | | Data quality and the Am Community Survey (ACS) is the household s(3,5 mlion homes contacted ead bourse for neighborhood scale information at head to social scientific research in the US the ACS can be high writeliable. For example, in the estimated number of children under S in Uncertainy of this magnitude complicates the policy making, research, and governance. Our project presents away to reduce the max technical details of this paper and example in described in the ISOSone Faguer. This website 1885 metropolican statistical areas, before and process, in order to explain, demonstrate, and the data. | erican largest survey of US nyear) and is the princip sout the US population. US population. Seeding and is a critical seever, 28% of census tract soverity +/- 150, use of Social data in gins of error in survey d s called regionalization. righementations are gresens the data from after the regionalization after the regionalization | al he s. |
| ••• 🗉 | | | € github.com | Ċ | @ ₾ |
| <> Code | e 🕕 Issues 3 🛛 👔 | Pull requests ⓒ Actions 삔 F | Projects 🖽 Wiki 🕐 Security | └─ Insights (| 3 Settings |
| ្លា mas | ter 👻 🖓 2 branches 🛇 | > 0 tags | Go to file Add file - | ± Code → | About |
| 양 mas | ter - P 2 branches 🤆 | >0 tags | Go to file Add file - 1494fec on Mar 8, 2018 | <u>↓</u> Code ↓○ 27 commits | About A tool to improve the usability of census data via "good" |
| 2º mas | ter - ¹ / ₂ 2 branches ^c | > 0 tags Updating line 356 | Go to file Add file - f494fec on Mar 8, 2018 | ✓ Code → ✓ 27 commits 5 years ago | About A tool to improve the usability of census data via "good" gerrymandering |
| 2º mas | ter • 1/2 branches C ass Update README.md le ENSE.bd | > 0 tags Updating line 356 Update LICENSE.txt | Go to file Add file - f494fec on Mar 8, 2018 | ✓ Code → ✓ 27 commits 5 years ago 5 years ago | About A tool to improve the usability of census data via "good" gerrymandering ${\mathcal O}$ journals.plos.org/plosone/article |
| * mass • gen • con • LIC • RE | ter - 1/2 branches C ass Update README.md de ENSE.txt ADME.md | O tags Updating line 356 Update LICENSE.txt Update README.md | Go to file Add file - f494fec on Mar 8, 2018 | ★ Code ▼ ♦ 27 commits 5 years ago 5 years ago 3 years ago | About A tool to improve the usability of census data via "good" gerrymandering |
| 2º mas | ter v 1/2 branches C ass Update README.md de ENSE.txt ADME.md README.md | O tags Updating line 356 Update LICENSE.txt Update README.md | Go to file Add file - | Code Code | About A tool to improve the usability of census data via "good" gerrymandering |
| P mas | ter • 1/2 branches © pass Update README.md je ENSE.txt ADME.md README.md FOOI for Reduction Immunity Sur | Updating line 356 Update LICENSE.txt Update README.md | Go to file Add file - r494fec on Mar 8, 2018 | Code Code Code Code Code Code Code Code | About A tool to improve the usability of census data via "good" gerrymandering |

DECENNIAL 2020: A NEW APPROACH TO DISCLOSURE AVOIDANCE

Swapping (2010 and earlier)

Noise infusion (2020)



| | School Attendance | | | | |
|--------|--------------------------|---------------------------|---------------------------|--|--|
| | Never | Attending | Past | | |
| Male | 3 - 1 = 2 | 12 + 0 = 12 | 33 + 1 = 34 | | |
| Female | 4 + 8 = 12 | 17 + 2 = 19 | 31 - 2 = 29 | | |
| | | | | | |

N = 100 N = 108

DECENNIAL 2020: POLICY DECISIONS



| Query | Allocation (%) |
|--|-------------------|
| Voting age * Hispanic * Race * Citizen | 50 |
| Household – Group quarters | 20 |
| Detailed | 10 |
| Sex * Age (single year of age) | 5 |
| Sex * Age (4-year age bins) | 5 |
| Sex * Age (16-year age bins) | 5 |
| Sex * Age (64-year age bins) | 5 |

RESULTS



Sex by Single Year of Age: Census Tract 303

Sex by Age: G2701230030300 200-Count 0-200-400-10-14 15-19 20-24 25-29 30-34 35-39 40-44 45-49 50-54 55-59 60-64 65-69 70-74 75-79 80-84 85+ 5-9 0-4 Age 16 Source: US Census Bureau 2011; US Census Bureau 2019; Van Riper et al. 2020

REAL-WORLD PUBLIC HEALTH EXAMPLE

- Asthma ED visit rates
 - Asthma ED visits in 2010 for Massachusetts towns
 - 0-4, 5-14, 15-34, 35-64, 65+ age bins
 - Age counts (denominators)
 - 2010 Summary File 1
 - Gold standard
 - Vintage 1 (October 2019)
 - Vintage 2 (May 2020)
 - How do rates based on diff. private different denominators compare to rates based on 2010 SF1 data?

Percent Difference in Age-Adjusted Asthma ED Visits in 2010 (MA towns)



Source: Van Riper et al. 2020; US Census Bureau 2019; Massachusetts Department of Health 2020

2020 DECENNIAL STILL A MOVING TARGET

- Additional demonstration data release scheduled for April 30
 - Only data on race, ethnicity, and voting age
- Demonstration data on sex, age, race/ethnicity forthcoming
 - No firm timeline, though, from Census Bureau
- Scientists should study guidance on handling uncertainty in decennial counts
 - Handbook on differential privacy will be available
 - Unlikely to get measure of uncertainty for decennial counts

CONCLUSION

- From ACS and Decennial one can expect good city-level rates/population estimates
 - But within-city or county is harder to understand
- Geographic and demographic resolution matter
 - Units with larger counts will be more accurate
 - Demographic groups with larger counts will be more accurate
- It is possible to process publicly released data to improve estimates.
 - Particularly for ACS data
 - Less certain about decennial data

QUESTIONS OR FEEDBACK:

vanriper@umn.edu

seth.spielman@colorado.edu

Code and data:

https://github.com/geoss/cdc-ncvhs-covid-2021