



Five Ethical Risks to Consider before Filling Missing Race and Ethnicity Data

Workshop Findings on the Ethics of Data Imputation and Related Methods

Megan Randall, Alena Stern, and Yipeng Su

March 2021

A growing set of methods from data science and statistics could fill critical gaps in race and ethnicity data by matching, imputing, or otherwise adding demographic and locational characteristics to existing datasets. As the potential for appending race and ethnicity variables grows, however, so does the risk of ethical violations and potential harm to Black, Indigenous, and other people of color (BIPOC).*

Disaggregating data by race and ethnicity can help shine a light on racialized systems of privilege and oppression.¹ However, many high-value datasets either don't report race or ethnicity or have missing race and ethnicity data. For example, the lack of race and ethnicity information in credit bureau data has inhibited efforts to examine how credit scores affect racial homeownership gaps or to challenge the use of credit screens in hiring, while the lack of racial identifiers in federal income tax data prevents efforts to assess the racially disparate impact of federal tax policies.² Given the absence of race and ethnicity information, researchers are often forced to choose between using imprecise methods to estimate race (such as the predominant race of an individual's zip code) or forgoing data disaggregation altogether.

As the availability of disaggregated data has grown, and as policymakers increasingly recognize its value for identifying and addressing racial disparities, researchers and advocates have called attention to associated challenges, such as accurately capturing small racial and ethnic groups; collecting data on

* Throughout this brief, the authors use "BIPOC" when generally referring to nonwhite people and communities. When discussing a specific demographic group or identity, we opt to name them as specifically as possible and note that questions about ethically representing the specific personhood and experiences of diverse individuals and communities are central practices in equitable data use. We recognize that "BIPOC" groups racial and ethnic groups together, obscuring the specific experiences of each, and we remain committed to using inclusive language whenever possible.

intersections with class, gender and other identities at smaller disaggregation levels; shifting the focus from individual outcomes to structural indicators; and appropriately visualizing racial disparities.³

Sophisticated methods for generating or appending racial and ethnic identifiers (box 1) share many of these broader challenges but also pose distinct risks. While the body of work examining harms and harm-mitigation approaches for algorithms and big data in general is considerable,⁴ it lacks specific guidance on the ethical risk areas that data scientists, statisticians, and researchers encounter when using imputation, matching, or related methods to fill missing race and ethnicity data.

BOX 1

Methods for Generating and Appending Missing Race and Ethnicity Data

Imputation. *Probabilistic methods* can help generate new data that maintain statistical properties of the “real” data. One widely used technique for imputing race and ethnicity on administrative data is the Bayesian Improved Surname Geocoding tool, developed by RAND for the US Department of Health and Human Services and also used by the Equal Opportunity Employment Commission (Harris 2020).^a Using *multiple imputation*, researchers can analyze variation resulting from the imputation process and, drawing from sets of probabilities, determine whether results are robust to the randomness inherent in the imputation process.^b

Machine-learning methods. Machine-learning methods can be especially useful when using text data or modeling complex, non-linear relationships. For example, the Urban Institute applied machine-learning methods to impute property-level zoning density limits (Nechamkin and MacDonald 2019) and impute sentiment toward police from tweets (Oglesby-Neal, Tiry, and Kim 2019). Several groups are already working on ways to mitigate biases in machine learning, such as the [Algorithmic Justice League](#) and researchers from the University of Chicago who developed the [Aequitas tool](#).^c

Data linkage. *Probabilistic data linkage*, popularly known as fuzzy matching, connects information from separate sources based on the probability of two records representing the same person or entity, using multiple and/or non-unique keys. Urban probabilistically linked names and addresses of Community Development Financial Institutions to estimate community development financial flows and linked mothers’ and infants’ records to evaluate the Strong Start for Mothers and Newborns initiative (Hill et al. 2014).^d *Data fusion* integrates multiple data sources to achieve more accuracy than a single data source would (e.g., combining multiple administrative data sources, such as the American Community Survey, city surveys, and United States Postal Service data, to understand change in a neighborhood).

Notes

^a Known as the Medicare Bayesian Improved Surname Geocoding (MBISG) tool (LeRoy et al. 2013), the latest MBISG 2.0 method combines name, administrative data, and Census data based on address in a calibrated Bayesian framework (multinomial logistic regression model) to estimate probabilities by race and ethnicity for each record in the dataset (Haas et al. 2019).

^b Many public data products use multiple imputation, including the Survey of Income and Program Participation Synthetic Beta, the National Survey of Children’s Health (Benedetto, Stinson, and Abowd 2013), and the Survey of Consumer Finances (Lindamood, Hanna, and Bi 2007).

^c See also Ziyuan Zhong, “A Tutorial on Fairness in Machine Learning,” *towards data science*, October 21, 2018, <https://towardsdatascience.com/a-tutorial-on-fairness-in-machine-learning-3ff8ba1040cb>.

^d See also “Community Development Financial Flows: How US Counties Compare,” Urban Institute, June 26, 2018, <https://apps.urban.org/features/community-development-financing/>.

In November 2020, the Urban Institute's [Racial Equity Analytics Lab](#) and Office of Technology and Data Science convened experts from the data science, government, racial justice, and data privacy fields to discuss the ethics of using advanced statistical methods to fill gaps in race and ethnicity data (box 2). Workshop participants affirmed the demand for disaggregated data, and the strong appetite for tools to fill data gaps, while identifying ethical risk areas that data scientists, statisticians, and researchers must grapple with when employing advanced techniques for generating these data.

BOX 2

Design Thinking Workshop on the Ethics of Imputation and Related Methods

On November 9, 2020, the Urban Institute held a virtual workshop on the ethics of using imputation and related methods to fill missing race and ethnicity data for savings and wealth datasets. We invited representatives from local, state, and federal government; justice and privacy advocates; imputation experts; consumer financial protection researchers; and organizations representing impacted communities, ensuring racial and ethnic diversity among the invitees. Twenty-nine people attended the workshop, including representatives from the Urban Institute, MetroLab Network, the New York City Department of Consumer and Worker Protection, and the Federal Reserve Bank of Atlanta.

We used a design thinking workshop format, a collaborative model for convening, engaging, and harnessing diverse perspectives and experiences to accelerate problem-solving.^a For most of the two-and-a-half-hour session, participants worked in four breakout groups to identify ethical risks associated with using these methods and surface potential risk-mitigation measures. Though the workshop prompt asked participants to consider applications and risks related to savings and wealth datasets in detail, participants also brainstormed and discussed applications to other high-value datasets that lack race and ethnicity identifiers. At the end of the workshop, breakout groups presented their key findings, and the whole group discussed cross-cutting themes. Immediately following the workshop, the authors synthesized key findings from each breakout group and shared them with the participants for review and comment.

The authors thank the following individuals who generously participated in the workshop and shared their insights with us. Workshop participants also provided invaluable feedback on early versions of this brief. Views expressed do not reflect the position of their organizations listed below.

Zayne Abdessalam from the New York City Department of Consumer and Worker Protection
Jasmine Burnett and Donta Council from the Federal Reserve Bank of Atlanta
Hector Dominguez from the City of Portland
Kimberly D. Lucas from MetroLab Network
Kevin Moore from the Federal Reserve
Shena Ashley, Claire Bowen, Breno Braga, Steven Brown, Ilham Dehry, LesLeigh Ford, Graham MacDonald, Signe-Mary McKernan, Michael Neal, Khuloud Odeh, Jamila Patterson, Kathryn L.S. Pettit, and Aaron Williams from the Urban Institute

We are also thankful to those workshop participants who chose to remain anonymous and are thus not named above, but who contributed valuable time and insights to our process.

Note

^a See Katrina Ballard, "How Can Human-Centered Design Uncover Policy Solutions?" *Data@Urban* (blog), Urban Institute, January 14, 2020, <https://urban-institute.medium.com/how-can-human-centered-design-uncover-policy-solutions-4ebf2ec0d89b>.

This brief summarizes five ethical risk areas that surfaced during the workshop. These risks pose serious harms to BIPOC and erode trust in government and other data stewards. While not exhaustive, this list provides researchers with a glimpse into where they should proceed with caution—as well as where the field must collaborate to develop equity-centered tools and resources. Workshop participants highlighted the need for a code of conduct that builds on these identified risks and stipulates ethical norms and guidelines on the use of emerging techniques for appending race and ethnicity to data.

Excluding People and Communities of Color from Ownership of Their Data and from Decisions on Research Process and Methods

Power dynamics between individuals whose data are being collected and the organizations who are funding, collecting, and using the data prevent people from exercising authority over their own data.⁵ Failing to provide channels for critical individual and community-level input increases the likelihood that researchers will overlook how individuals and communities connect to and identify with the research process.

Although this risk is present in *all* research featuring people and communities, workshop participants reported that researchers working primarily with advanced statistical techniques and secondary data (rather than directly with people in community) need more guidance on how to build community engagement principles,⁶ stakeholder vetting opportunities,⁷ and intentional participatory mechanisms into their work.⁸ In addition to best practices kits and self-assessment tools that help research teams incorporate engagement principles into advanced analytics and “big data” projects,⁹ workshop participants said that the field needs to develop models for data ownership and governance that give BIPOC more authority over their data. Direct engagement is likely not possible nor desirable for every researcher analyzing quantitative data, but researchers should pursue opportunities to incorporate the perspectives of affected communities within the project, for instance through guidance provided by the data collector or with the assistance of other researchers more proximate to the affected communities. Planning for community engagement that is appropriate to the type and context of each analytic project can help researchers equitably navigate the risks and trade-offs throughout the research process.

Violating Individual Informed Consent

Separate from community engagement, informed *individual* consent is traditionally a requirement and expectation in any research involving people.¹⁰ The collection and dissemination of multiple sources of anonymized secondary data, however, can dilute or circumvent traditional informed consent processes.¹¹ Historically, BIPOC have been systematically deprived of opportunities for informed consent in research.¹² Even today, researchers overrepresent Black patients in US Food and Drug

Administration-approved clinical trials that do not require informed consent.¹³ As identified in this brief, imputation and related methods carry material risks for individuals and communities. Someone who consented to provide sensitive financial or health data may not have done so if plans to append race or ethnicity to their data were fully disclosed. Additionally, “informed refusal” to participate or provide data is a meaningful personal and political choice.¹⁴ If someone declines to report their race or ethnicity in a survey, later generating that value through advanced analytical methods overrides that initial refusal.

Obtaining and enforcing more robust individual consent practices while preserving the availability of some secondary data for research comes with practical challenges. An individual’s risk of exposure and the potential uses of secondary data may not be known ahead of time. Moreover, even if risks are more explicitly acknowledged in the “fine print” of an informed consent procedure, this does not guarantee that people will take full notice, understand the risks, or feel empowered to say no (e.g., if the consent process is connected to an application for a benefit program or other needed service).

At minimum, workshop participants expressed a need for clearer language and individual notice.¹⁵ More ambitiously, data collectors should establish data governance and sharing practices that ensure data are not used beyond the consented-to purposes. For researchers who use secondary data, workshop participants expressed a need for more guidance on how to incorporate individual voice and choice into their analytics projects.

Compromising Individual Privacy or Confidentiality

When researchers append racial or ethnic identifiers to *other* identifiers at small units of geography, individuals are at increased risk of re-identification,¹⁶ even if those datasets have previously been anonymized.¹⁷ To address these concerns, researchers can aggregate data, reporting only on larger geographies or on combined racial and ethnic groups. They can also use methods for generating synthetic data that preserve statistical properties of the original data while adding enough “noise” to preserve privacy. In both cases, researchers must make important trade-offs between privacy and accuracy.

Workshop participants said that the possibility of re-identification, either through linking of multiple sources or by imputation,¹⁸ needs to be more explicitly acknowledged in informed consent procedures.¹⁹ For research teams who rely on secondary data, workshop participants identified a need for additional tools, like privacy impact assessments,²⁰ to help researchers assess and mitigate privacy risks specific to their projects.

Producing Inaccurate Estimates and Misleading Conclusions

Imputation and related methods often come with a degree of statistical uncertainty. But if imputed race and ethnicity variables do not meaningfully predict *actual* race and ethnicity, the conclusions policymakers draw from the imputed data could lead to misinformed policy choices that harm BIPOC. Workshop participants said that researchers need to produce and share estimates of the variation in their imputation process and analyze whether their results are robust to that variation.

The level of variability is likely higher for smaller race and ethnicity subgroups because fewer observations are being imputed and the data used to perform the imputation are more variable. For example, the Bayesian Improved Surname Geocoding tool is less accurate for Native American and multiracial people than for Black or Asian people (LeRoy et al. 2013). Researchers will encounter a trade-off between fully representing specific race and ethnicity subgroups in the data (while tolerating higher levels of variation and uncertainty) and aggregating those subgroups (producing less variable results, but also concealing important heterogeneity across subgroups).

Workshop participants identified a need for guidelines on defining acceptable ranges of uncertainty for different use cases. Whether uncertain estimates can be used to responsibly inform policy will vary depending on the objective of the research and limits of the data in question. Additionally, participants emphasized the importance of communicating and reflecting this uncertainty in final research products so policymaking audiences understand the risks related to uncertainty in the dataset.

Additionally, methods like imputation only produce results as accurate as the underlying data, which often reflect structural disparities and racial biases.²¹ Linking biased datasets together, using them to power data-driven decision systems, or training predictive algorithms with them can magnify erroneous results. For example, a lack of diversity in publicly available image datasets has contributed to racial bias in facial recognition systems, which use those data to train and evaluate their underlying algorithms.²² Similarly, a significant undercount of Black and Latinx populations in the 2020 Census could lead to bias in algorithms that use Census data to produce and append racial identifiers.²³ In a different example, in 2020, the California Attorney General revoked law enforcement departments' access to a database of suspected gang members because of pervasive errors. Gang affiliations had been assigned using largely unsubstantiated (and, in some cases, demonstrably falsified) reports from individual law enforcement officers, reflecting significant racial bias.²⁴ Using any such data to inform algorithmic approaches, as Chicago similarly did for predicting gun crime,²⁵ will produce racially biased results.²⁶

Before incorporating data into any imputation, matching, or machine-learning process, researchers need to understand *how* those data are collected and for what purpose. Workshop participants identified a need for more routine and robust quality assurance processes, stakeholder input mechanisms, and continuous review of underlying data and their analytic outputs to help identify biases. Developing more verification tools like the Urban Institute's Spatial Equity Data Tool may, for example,

help researchers identify which neighborhoods and demographic groups are underrepresented and overrepresented in certain datasets.²⁷

Generating Data for Purposes That Harm People or Communities of Color

Datasets that exclude race and ethnicity may do so for good reason. For example, the federal government purposefully prohibits credit bureaus from collecting data on race and ethnicity to protect against discriminatory lending.²⁸ Workshop participants expressed a strong concern that imputed, or otherwise appended, racial and ethnic identifiers could be weaponized against BIPOC.

This concern drives at the heart of a larger debate about the responsible use and presentation of racially disaggregated data. Any tools that empower disaggregation, including imputation and related methods, can contribute to racist narratives if the data reinforce harmful stereotypes about BIPOC that lead to discrimination against groups and neighborhoods. This can happen either through the selection and visual presentation of the data or through framing racially disparate outcomes as the result of individual choices and behaviors rather than structural forces.²⁹

Workshop participants also expressed specific misgivings about imputation, matching, and other methods for linking highly sensitive, personally identifiable financial or health data (for example) to an individual's race and ethnicity. Linked credit bureau data, for example, could be used punitively to reinforce racially discriminatory lending practices or target predatory products.³⁰ Meanwhile, calls to “personalize law” based on an individual's data footprint, or to personalize medicine based on race and ethnicity, provide opportunities for imputed datasets to advance racial stereotyping and discrimination in criminal justice, medicine, and financial services.³¹ In efforts to comply with President Trump's executive order on citizenship data, after federal courts blocked the administration's attempt to add a citizenship question to the 2020 Census questionnaire, the US Census Bureau considered whether it could determine individual citizenship status through linked government datasets.³² Even if individual citizenship estimates were kept confidential,³³ many worried that public block-level estimates could help target deportation raids or inform redistricting efforts that would deprive noncitizens of political representation.³⁴ The Census Bureau is no longer pursuing this effort.³⁵

Workshop participants agreed that sharing and vetting analytic methods with community representatives who have a stake in the research process is a critical component of an equity-centered approach. They also identified a tension between the need for openness in data and methods (which enables different community stakeholders to identify biases and, when necessary, intervene to prevent ethical violations) and the need to protect against potential misuse of BIPOC data (which could occur if providing unfettered access those data or methods). Workshop participants expressed a need for thoughtful principles on—or an independent entity that can help govern—when, how, and to whom sensitive data containing individual race and ethnicity variables are released.

Conclusion and Next Steps

The risks detailed in this brief provide a small window into the ethical questions that researchers and data scientists need to consider before embarking on efforts to attach racial and ethnic identifiers to data. While this brief has focused on the many roadblocks to ethically generating and appending missing race and ethnicity data, there are also material and ethical costs of *not* filling these gaps, including obscuring ongoing patterns of racial discrimination and perpetuating invisibility for communities of color who do not see themselves represented in data.³⁶

Researchers seeking to ethically generate and append missing race and ethnicity data need further guidance and resources to mitigate the ethical risks as best as possible and navigate the numerous trade-offs between and among different ethical priorities surfaced in this document—trade-offs that are rarely clear-cut. As Urban’s Racial Equity Analytics Lab develops its capacity to generate and deploy sophisticated data on race and ethnicity, it will build on these workshop learnings in our forthcoming work to explore approaches and best practices that mitigate ethical risk, center equity, and build empathy for people and communities.

Notes

- ¹ See “Making the Case for Data Disaggregation to Advance a Culture of Health,” PolicyLink, <https://www.policylink.org/our-work/community/health-equity/data-disaggregation>; and “Why Disaggregating Data by Race Is Important for Racial Equity,” Annie E. Casey Foundation blog, August 18, 2020, <https://www.aecf.org/blog/taking-data-apart-why-a-data-driven-approach-matters-to-race-equity/>.
- ² For more on the racial homeownership gap, see Choi et al. (2019). For a discussion on the lack of race and ethnicity in tax data, see Bearer-Friend (2019). For more problems with employer credit checks, see Traub and McElwee (2016).
- ³ See “The Essentials of Disaggregated Data for Advancing Racial Equity,” Race Matters Institute, 2019, <https://viablefuturescenter.org/racemattersinstitute/resources/disaggregated-data/>; “Workshop Series: Addressing Health Equity through Data Disaggregation,” UCLA Center for Health Policy Research, <http://events.r20.constantcontact.com/register/event?oeidk=a07ehc2hrvo920ec420&llr=6yxgr6cab>; Andrews, Parekh, and Peckoo (2019); Annie E. Casey Foundation (2008); National Forum on Education Statistics (2016); Rubin et al. (2018); Jonathan Schwabish and Alice Feng, “Applying Racial Equity Awareness in Data Visualization,” *Data@Urban* (blog), Urban Institute, September 3, 2020, <https://urban-institute.medium.com/applying-racial-equity-awareness-in-data-visualization-bd359bf7a7ff>; Natalie Spievack and Cameron Okeke, “How We Should Talk about Racial Disparities,” *Urban Wire*, Urban Institute, February 26, 2020, <https://www.urban.org/urban-wire/how-we-should-talk-about-racial-disparities>; and Sonia Torres Rodriguez, “Three Steps to Improving Data to Help Combat the Public Health Emergency of Structural Racism,” *Urban Wire*, Urban Institute, January 11, 2021, <https://www.urban.org/urban-wire/three-steps-improving-data-help-combat-public-health-emergency-structural-racism>.
- ⁴ See Alice Feng and Shuyan Wu, “The Myth of the Impartial Machine,” *Data@Urban* (blog), Urban Institute, July 9, 2019, <https://urban-institute.medium.com/the-myth-of-the-impartial-machine-9fecb291abe0>; Caroline Lair, “12 Black Women in AI paving the way for a better world,” LinkedIn, June 15, 2020, <https://www.linkedin.com/pulse/12-black-women-ai-paving-way-better-world-caroline-lair/>; Ziyuan Zhong, “A Tutorial on Fairness in Machine Learning,” *towards data science*, October 21, 2018, <https://towardsdatascience.com/a-tutorial-on-fairness-in-machine-learning-3ff8ba1040cb>; and Zook et al. (2017).

- ⁵ See Chicago Beyond (2018) and Gangadharan et al. (2019).
- ⁶ For example, see “Community Engagement Methods at Urban,” Urban Institute, <https://www.urban.org/research/data-methods/community-engagement-methods-urban>; and Alena Stern, Graham MacDonald, and Khuloud Odeh, “How Can Cities Equitably Engage Marginalized Communities?” September 2020, <https://apps.urban.org/features/how-to-create-equitable-technology-programs/#marginalized>.
- ⁷ As discussed in Annie E. Casey Foundation (2020); or through mechanisms like a community advisory board, as described in Kubicek and Robles (2016).
- ⁸ For additional resources, see “Equitable Data Practice,” Urban Institute, <https://www.urban.org/elevate-data-equity/resources-elevate-data-equity/equitable-data-practice>; Actionable Intelligence for Social Policy (2020); and Gaddy and Scott (2020).
- ⁹ For a discussion on data literacy, community engagement, and “big data,” see Bhargava et al. (2015). Loukides, Mason, and Patil (2018), meanwhile, highlight (among other data ethics considerations) how a lack of racial and ethnic diversity on a research team can make ethical breaches more likely.
- ¹⁰ For a discussion of racial equity and the Institutional Review Board process, which requires informed consent, see Jenita Parekh and Manica F. Ramos, “Racial Equity Considerations and the Institutional Review Board,” *Child Trends*, March 11, 2020, <https://www.childtrends.org/publications/racial-equity-considerations-and-the-institutional-review-board>.
- ¹¹ See, for example, Froomkin (2019), and Petroni and Long (2016).
- ¹² See, for example, Nuriddin, Mooney, and White (2020).
- ¹³ See Ike Swetlitz, “African-Americans are disproportionately enrolled in studies that don’t require informed consent,” *STAT*, October 1, 2018, <https://www.statnews.com/2018/10/01/african-americans-clinical-trials/>.
- ¹⁴ For a description and conception of “informed refusal” in bioethics, see Benjamin (2016).
- ¹⁵ See, for example, “plainlanguage.gov,” Plain Language Action and Information Network, <https://www.plainlanguage.gov/>.
- ¹⁶ See Claire Bowen, “Will the Census’s Data Privacy Efforts Erase Rural America?” *Urban Wire*, Urban Institute, March 4, 2020, <https://www.urban.org/urban-wire/will-census-data-privacy-efforts-erase-rural-america>.
- ¹⁷ See, for example, Gina Kolata, “Your Data Were ‘Anonymized’? These Scientists Can Still Identify You,” *New York Times*, July 23, 2019, <https://www.nytimes.com/2019/07/23/health/data-privacy-protection.html>.
- ¹⁸ As data sharing becomes more common among public and private institutions, some have begun developing templates for informed consent that incorporate the possibility of data sharing, although such mechanisms do not always directly address risks with a racial equity lens. See, for example, “Recommended Informed Consent Language for Data Sharing,” Inter-university Consortium for Political and Social Research, <https://www.icpsr.umich.edu/web/pages/datamanagement/confidentiality/conf-language.html>.
- ¹⁹ See also “Time to discuss consent in digital-data studies” *Nature*, July 31, 2019, <https://www.nature.com/articles/d41586-019-02322-z>.
- ²⁰ For examples of Privacy Impact Assessments, see “Privacy Impact Assessments,” Federal Trade Commission, <https://www.ftc.gov/site-information/privacy-policy/privacy-impact-assessments>; “Privacy Impact Assessments (PIA),” General Services Administration, <https://www.gsa.gov/reference/gsa-privacy-program/privacy-impact-assessments-pia>; and “Privacy Impact Assessments,” Department of Health and Human Services, <https://www.hhs.gov/pia/index.html>.
- ²¹ See, for example, Alena Stern, Graham MacDonald, and Khuloud Odeh, “What evidence should cities assess when designing equitable technology programs?” Urban Institute, September 2020, <https://apps.urban.org/features/how-to-create-equitable-technology-programs/#evidence>.

- ²² For information on racial bias in facial recognition training data see Joy Buolamwini and Timnit Gebru, “Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification,” *Proceedings of Machine Learning Research*, 2018, <http://proceedings.mlr.press/v81/buolamwini18a/buolamwini18a.pdf>.
- ²³ See “2020 Census: Who’s At Risk of Being Miscalculated?” Urban Institute, June 4, 2019, <http://apps.urban.org/features/2020-census/>; and Elliott et al. (2019).
- ²⁴ As discussed in Benjamin (2019). See also Anita Chabria and Leila Miller, “Reformers Want California Police to Stop Using a Gang Database Seen as Racially Biased,” June 24, 2020, <https://www.latimes.com/california/story/2020-06-24/california-police-urged-to-stop-using-gang-database-deemed-biased>.
- ²⁵ See Sam Charles, “Chicago Police Department ends use of ‘Strategic Subject List,’” *Chicago Sun-Times*, January 27, 2020, <https://chicago.suntimes.com/city-hall/2020/1/27/21084030/chicago-police-strategic-subject-list-party-to-violence-inspector-general-joe-ferguson>; and Maryam Saleh, “Chicago’s Promise: Caught in a Gang Dragnet and Detained by ICE, an Immigrant Tests the Limits of a Sanctuary City,” *Intercept*, January 28, 2018, <https://theintercept.com/2018/01/28/chicago-gangs-immigration-ice/>.
- ²⁶ See Will Douglas Heaven, “Predictive Policing Algorithms Are Racist. They Need to Be Dismantled.” *MIT Technology Review*, July 17, 2020, <https://www.technologyreview.com/2020/07/17/1005396/predictive-policing-algorithms-racist-dismantled-machine-learning-bias-criminal-justice/>.
- ²⁷ See Ajjit Narayanan and Alena Stern, “Introducing a Spatial Equity Data Tool,” *Data@Urban* (blog), Urban Institute, September 23, 2020, <https://urban-institute.medium.com/introducing-a-spatial-equity-data-tool-b959c40298cf>.
- ²⁸ See Cooper and Getter (2015).
- ²⁹ See Ibram X. Kendi, “Stop Blaming Black People for Dying of the Coronavirus,” April 14, 2020, <https://www.theatlantic.com/ideas/archive/2020/04/race-and-blame/609946/>; and Spievack and Okeke, “How We Should Talk about Racial Disparities.”
- ³⁰ As Harcourt (2014, 7) observed, the emerging ability to link datasets across commercial and public surveillance platforms is producing “a trove of data to identify individuals, nudge them politically, manipulate them slightly, encourage and stimulate their consumption and disclosures, watch and surveil them, detect, predict, and punish.” See also Gillian B. White, “Why Blacks and Hispanics Have Such Expensive Mortgages,” *Atlantic*, February 25, 2016, <https://www.theatlantic.com/business/archive/2016/02/blacks-hispanics-mortgages/471024/>.
- ³¹ While such calls for personalization of laws and services based on personal data introduce potential benefits, experts have also identified risks and potential misuses. See Porat and Strahilevitz (2013) for a proposal to personalize default rules based on data, while Devins et al. (2017) discuss the possible misuse of big data for personalizing law in a variety of domains. For a discussion of the merits and risks of using race as a variable in precision medicine, see Geneviève et al. (2020); and Emily Singer, “Race and Personalized Medicine,” December 20, 2006, <https://www.technologyreview.com/2006/12/20/227301/race-and-personalized-medicine/>.
- ³² See Abowd et al. (2020); Aaron Boyd, “How Census Is Building a Citizenship Database Covering Everyone Living in the U.S.” *Nextgov*, April 1, 2020, <https://www.nextgov.com/analytics-data/2020/04/how-census-building-citizenship-database-covering-everyone-living-us/164275/>; and Hansi Lo Wang, “To Figure Out Who’s A Citizen, Trump Administration Is Using These Records,” *NPR*, May 20, 2020, <https://www.npr.org/2020/05/20/855062093/to-figure-out-whos-a-citizen-trump-administration-is-using-these-records>.
- ³³ By law, the US Census Bureau is required to keep individuals’ data confidential. However, historians have documented how confidentiality rules were previously suspended to allow the US government to make use of Census data for identifying and interning Japanese Americans. See Lori Aratani, “Secret use of census info helped send Japanese Americans to internment camps in WWII,” *Washington Post*, April 6, 2018, <https://www.washingtonpost.com/news/retropolis/wp/2018/04/03/secret-use-of-census-info-helped-send-japanese-americans-to-internment-camps-in-wwii/>.
- ³⁴ See Shira Mitchell, “Worries Grow Census Data Could Be Used to Target the Undocumented,” May 3, 2018, <https://independent.org/2018/05/census-data-could-be-used-to-target-the-undocumented/>.

³⁵ See Mike Schneider, “Citizenship Data Is Latest Rollback of Trump Census Efforts,” Associated Press, January 23, 2021, <https://apnews.com/article/joe-biden-census-2020-aa774e5d530354767c712ec44d7cfa04>.

³⁶ See Williams, Brooks, and Shmargad (2018) on how data lacking race and ethnicity identifiers can produce algorithmic bias; Bearer-Friend (2018) on the balance between racial privacy and protection against discrimination as it relates to federal income tax data; and Namratha Kandula and Nilay Shah, “Asian Americans Invisible in COVID-19 Data and in Public Health Response,” *Chicago Reporter*, June 16, 2020, <https://www.chicagoreporter.com/asian-americans-invisible-in-covid-19-data-and-in-public-health-response/>.

References

Abowd, John M., William R. Bell, J. David Brown, Michael B. Hawes, Misty L. Heggeness, Andrew D. Keller, Vincent T. Mule, Jr., et al. 2020. “Determination of the 2020 U.S. Citizen Voting Age Population (CVAP) Using Administrative Records and Statistical Methodology.” Center for Economic Studies Working Paper 20-33. Washington, DC: US Census Bureau.

Actionable Intelligence for Social Policy. 2020. *A Toolkit for Centering Racial Equity Throughout Data Integration*. Philadelphia: Actionable Intelligence for Social Policy.

Andrews, Kristine, Jenita Parekh, and Shantai Peckoo. 2019. “How to Embed a Racial and Ethnic Equity Perspective in Research: Practical Guidance for the Research Process.” Washington, DC: Child Trends.

Annie E. Casey Foundation. 2008. “Advancing Better Outcomes for All Children: Reporting Data Using a Racial Equity Lens.” MORE Race Matters Occasional Update 3. Baltimore: Annie E. Casey Foundation.

———. 2020. *Four Principles to Make Advanced Data Analytics Work for Children and Families*. Baltimore: Annie E. Casey Foundation.

Bearer-Friend, Jeremy. 2018. “Should the IRS Know Your Race? The Challenge of Colorblind Tax Data.” *Tax Law Review* 73 (1): 1–68.

Benedetto, Gary, Martha H. Stinson, and John M. Abowd. 2013. “The Creation and Use of the SIPP Synthetic Beta.” Washington, DC: Census Bureau.

Benjamin, Ruha. 2016. “Informed Refusal: Toward a Justice-Based Bioethics.” *Science, Technology, & Human Values* 41 (6): 967–90.

———. 2019. *Race After Technology: Abolitionist Tools for the New Jim Code*. Cambridge: Polity Press.

Bhargava, Rahul, Erica Deahl, Emmanuel Letouzé, Amanda Noonan, David Sangokoya, and Natalie Shoup. 2015. “Beyond Data Literacy: Reinventing Community Engagement and Empowerment in the Age of Data.” Cambridge, MA: Data-Pop Alliance.

Chicago Beyond. 2018. *Why Am I Always Being Researched?* Chicago: Chicago Beyond.

Choi, Jung Hyun, Caitlin Young, Alanna McCargo, Michael Neal, and Laurie Goodman. 2019. *Explaining the Black-White Homeownership Gap*. Washington, DC: Urban Institute.

Cooper, Cheryl R., and Darryl E. Getter. 2020. “Consumer Credit Reporting, Credit Bureaus, Credit Scoring, and Related Policy Issues.” Washington, DC: Congressional Research Service.

Devins, Caryn, Teppo Felin, Stuart Kauffman, and Roger Koppl. 2017. “The Law and Big Data.” *Cornell Journal of Law and Public Policy* 27: 357–413.

Elliott, Diana, Robert Santos, Steven Martin, and Charmaine Runes. 2019. *Assessing Miscounts in the 2020 Census*. Washington, DC: Urban Institute.

Froomkin, A. Michael. 2019. “Big Data: Destroyer of Informed Consent.” *Yale Journal of Health Policy, Law, and Ethics* 18 (23).

Gaddy, Marcus, and Kassie Scott. 2020. “Principles for Advancing Equitable Data Practice.” Washington, DC: Urban Institute.

- Gangadharan, Seeta Peña, Tawana Petty, Tamika Lewis, and Mariella Saba. 2018. "Digital Defense Playbook: Community Power Tools for Reclaiming Data." Detroit: Our Data Bodies.
- Geneviève, Lester Darryl, Andrea Martani, David Shaw, Bernice Simone Elger, and Tenzin Wangmo. 2020. "Structural Racism in Precision Medicine: Leaving No One Behind." *BMC Medical Ethics* 21 (17).
- Haas, Ann, Marc N. Elliott, Jacob W. Dembosky, John L. Adams, Shondelle M. Wilson-Frederick, Joshua S. Mallett, Sarah Gaillot, Samuel C. Haffer, and Amelia M. Haviland. 2019. "Imputation of Race/Ethnicity to Enable Measurement of HEDIS Performance by Race/Ethnicity." *Health Services Research* 54 (1): 13–23.
- Harcourt, Bernard E. 2014. "Governing, Exchanging, Securing: Big Data and the Production of Digital Knowledge." Public Law and Legal Theory Working Paper Group 14–390. New York: Columbia Law School.
- Harris, Ada. 2020. "Using Bayesian Improved Surname Geocoding (BISG) to Classify Race and Ethnicity in Administrative Employment Data by Industry: A Validation Study." Paper presented at the JSM Virtual Conference, American Statistical Association, August 6.
- Hill, Ian, Sarah Benatar, Brigitte Courtot, Fredric Blavin, Embry M. Howell, Lisa Dubay, Bowen Garrett, et al. 2014. *Strong Start for Mothers and Newborns Evaluation: Year 1 Annual Report*. Washington, DC: Urban Institute.
- Kubicek, Katrina, and Marisela Robles. 2016. *Resource for Integrating Community Voices into a Research Study: Community Advisory Board Toolkit*. Los Angeles: Southern California Clinical and Translational Science Institute.
- LeRoy, Lisa, Melanie Wasserman, Michael Rezaee, and Alan White. 2013. "Understanding Disparities in Persons with Multiple Chronic Conditions: Research Approaches and Datasets." Cambridge, MA: Abt Associates.
- Lindamood, Suzanne, Sherman D. Hanna, and Dan Bi. 2007. "Using the Survey of Consumer Finances: Some Methodological Considerations and Issues." *Journal of Consumer Affairs* 41 (2): 195–214.
- Loukides, Mike, Hilary Mason, and DJ Patil. 2018. *Ethics and Data Science*. Sebastopol, CA: O'Reilly Media.
- National Forum on Education Statistics. 2016. "Forum Guide to Collecting and Using Disaggregated Data on Racial/Ethnic Subgroups." NFES 2017-017. Washington, DC: National Center for Education Statistics.
- Nechamkin, Emma and Graham MacDonald. 2019. "Predicting Zoned Density Using Property Records." Washington, DC: Urban Institute.
- Nuriddin, Ayah, Graham Mooney, and Alexandre I. R. White. 2020. "Reckoning with Histories of Medical Racism and Violence in the USA." *Lancet* 396 (10256): 949–51.
- Oglesby-Neal, Ashlin, Emily Tiry, and KiDeuk Kim. 2019. "Public Perceptions of Police on Social Media." Washington, DC: Urban Institute.
- Petroni, MJ, and Jessica Long. 2016. "Data Ethics: Informed Consent and Data in Motion." Mountain View, CA: Accenture.
- Porat, Ariel, and Lior Strahilevitz. 2013. "Personalizing Default Rules and Disclosure with Big Data." *Michigan Law Review* 112 (8): 1417–78.
- Rubin, Victor, Danielle Ngo, Ángel Ross, Dalila Butler, and Nisha Balaram. 2018. *Counting a Diverse Nation: Disaggregating Data on Race and Ethnicity to Advance a Culture of Health*. Oakland, CA: PolicyLink.
- Traub, Amy, and Sean McElwee. 2016. *Bad Credit Shouldn't Block Employment: How to Make State Bans on Employment Credit Checks More Effective*. New York: Demos.
- Williams, Betsy Anne, Catherine F. Brooks, and Yotam Shmargad. 2018. "How Algorithms Discriminate Based on Data They Lack: Challenges, Solutions, and Policy Implications." *Journal of Information Policy* 8:78–115.
- Zook, Matthew, Solon Barocas, Danah Boyd, Kate Crawford, Emily Keller, Seeta Peña Gangadharan, Alyssa Goodman, et al. 2017. "Ten Simple Rules for Responsible Big Data Research." *PLOS Computational Biology* 13 (3): e1005399.

About the Authors

Megan Randall is a research associate in the Urban Institute's Research to Action Lab where she brings expertise on state and local public finance to the Lab's work on federal place-based programs, state and local economic development policy, and inclusive economic recovery. She also provides research and project management support to Urban's Racial Equity Analytics Lab. Randall graduated summa cum laude with a bachelor's degree in political science from the University of California, Berkeley, and earned master's degrees in public affairs and in community and regional planning from the University of Texas at Austin.

Alena Stern is a senior data scientist at the Urban Institute studying policy solutions to advance equity and inclusion in cities. Before joining Urban, she worked as a senior program manager with AidData, an Open Cities fellow at the Sunlight Foundation, and a graduate research assistant at the Center for Data Science and Public Policy, where she used machine learning, natural language processing, statistical analysis, and geospatial data to inform the design of government policies and international development programs. Alena holds a BA in economics and international relations from the College of William and Mary and an MS in computational analysis and public policy from the University of Chicago.

Yipeng Su is a research associate in the Metropolitan Housing and Communities Policy Center at the Urban Institute. Her research interests include housing, economic development and the intersection between urban planning and technology. She is skilled in quantitative research methods and recently joined Urban's Racial Equity Analytics Lab to provide research and technical support to researchers across Urban. She holds a master's degree in public administration from New York University's Wagner Graduate School of Public Service

All authors contributed equally to this work.

Acknowledgments

This brief is a product of the Urban Institute’s Racial Equity Analytics Lab, which operates with the generous support of Salesforce Foundation, the Robert Wood Johnson Foundation, and Urban’s general support donors. We are grateful to them and to all our funders, who make it possible for Urban to advance its mission.

The views expressed are those of the authors and should not be attributed to the Urban Institute, its trustees, or its funders. Funders do not determine research findings or the insights and recommendations of Urban experts. Further information on the Urban Institute’s funding principles is available at urban.org/fundingprinciples.

We would like to thank our Urban Institute colleagues Jonathan Schwabish and Jessica Kelly for excellent feedback on early versions of this brief.



500 L’Enfant Plaza SW
Washington, DC 20024

www.urban.org

ABOUT THE URBAN INSTITUTE

The nonprofit Urban Institute is a leading research organization dedicated to developing evidence-based insights that improve people’s lives and strengthen communities. For 50 years, Urban has been the trusted source for rigorous analysis of complex social and economic issues; strategic advice to policymakers, philanthropists, and practitioners; and new, promising ideas that expand opportunities for all. Our work inspires effective decisions that advance fairness and enhance the well-being of people and places.

Copyright © March 2021. Urban Institute. Permission is granted for reproduction of this file, with attribution to the Urban Institute.