# Data Ethics

## Investing Wisely in Data at Scale

September 2016

**Prepared for the MacArthur and Ford Foundations**

David Robinson & Miranda Bogen
**www.dataethics.io**

Upturn

# Executive Summary

"Data at scale" — digital information collected, stored and used in ways that are newly feasible — opens new avenues for philanthropic investment. At the same time, projects that leverage data at scale create new risks that are not addressed by existing regulatory, legal and best practice frameworks. Data-oriented projects funded by major foundations are a natural proving ground for the ethical principles and controls that should guide the ethical treatment of data in the social sector and beyond.

This project is an initial effort to map the ways that data at scale may pose risks to philanthropic priorities and beneficiaries, for grantmakers at major foundations, and draws from desk research and unstructured interviews with key individuals involved in the grantmaking enterprise at major U.S. foundations. The resulting report was prepared at the joint request of the MacArthur and Ford Foundations.

Grantmakers are exploring data at scale, but currently have poor visibility into its benefits and risks. Rapid technological change, the scarcity of data science expertise, limited training and resources, and a lack of clear guideposts around emergent risks all contribute to this problem.

Funders have important opportunities to invest in, learn from, and innovate around data-intensive projects, in concert with their grantees. Grantmakers should not treat the new ethical risks of data at scale as a barrier to investment, but these risks also must not become a blind spot that threatens the success and effectiveness of philanthropic projects. Those working with data at scale in the philanthropic context have much to learn: throughout our conversations with stakeholders, we heard consistently that grantmakers and grantees lack baseline knowledge on using data at scale, and many said that they are unsure how to make better informed decisions, both about data's benefits and about its risks. Existing frameworks address many risks introduced by data-intensive grantmaking, but leave some major gaps. In particular, we found that:

- **Some new data-intensive research projects involve meaningful risk to vulnerable populations, but are not covered by existing human subjects regimes, and lack a structured way to consider these risks.** In the philanthropic and public sector, human subject review is not always required and program officers, researchers, and implementers do not yet have a shared standard by which to evaluate ethical implications of using public or existing data, which is often exempt from human subjects review.

- **Social sector projects often depend on data that reflects patterns of bias or discrimination against vulnerable groups, and face a challenge of how to avoid reinforcing existing disparities.** Automated decisions can absorb and sanitize bias from input data, and responsibly funding or evaluating statistical models in data-intensive projects increasingly demands advanced mathematical literacy which foundations lack.

- **Both data and the capacity to analyze it are being concentrated in the private sector, which could marginalize academic and civil society actors.** Some individuals and organizations have begun to call attention to these issues and create their own trainings, guidelines, and policies — but ad hoc solutions can only accomplish so much.

To address these and other challenges, we've identified eight key questions that program staff and grantees need to consider in data-intensive work:

1. For a given project, what data should be collected, and who should have access to it?

2. How can projects decide when more data will help — and when it won't?

3. How can grantmakers best manage the reputational risk of data-oriented projects that may be at a frontier of social acceptability?

4. When concerns are recognized with respect to a data-intensive grant, how will those concerns get aired and addressed?

5. How can funders and grantees gain the insight they need in order to critique other institutions' use of data at scale?

6. How can the social sector respond to the unique leverage and power that large technology companies are developing through their accumulation of data and data-related expertise?

7. How should foundations and nonprofits handle *their own* data?

8. How can foundations begin to make the needed long-term investments in training and capacity?

Newly emergent ethical issues inherent in using data at scale point to the need for both a broader understanding of the possibilities and challenges of using data in the philanthropic context as well as conscientious treatment of data ethics issues. Major foundations can play a meaningful role in building a broader understanding of these possibilities and challenges, and they can set a positive example in creating space for open and candid reflection on these issues. To those ends, we recommend that funders:

○ **Include data ethics as an element of larger efforts to build data literacy among grantmakers and grantees.** Create spaces for conversation and reflection for funders in order to promote data literacy and sensitivity, and invest in education on data-related topics for current and future staff.

○ **Incorporate data ethics in the grantmaking process.** Create an internal "data ethics point of contact" who can facilitate access to relevant expertise and keep an eye out for latent data ethics risk in projects, and consider changing grant applicant procedures to encourage applicants to prospectively consider data-related issues. Larger foundations should support a central resource to address data ethics concerns for the philanthropic community.

○ **Create a data ethics checklist for grantees and program staff.** Even without introducing any new requirements or policies, equipping staff with guiding questions they can ask about new, data-oriented projects can help funders identify and address areas of ethical concern.

# Table of Contents

Privacy

Open access and creative commons

Some risks of data at scale are not addressed by existing frameworks

"Public" data can now be used in unethical ways.

Automated decisions can absorb and sanitize bias from input data.

Data at scale may allow powerful institutions to marginalize academic and civil society actors.

## 5  Managing Data Ethics Risk

Data science practitioners and researchers are stepping up.

Companies are creating their own, internal ethical review processes.

## 6  Practical Steps for Funders

Include data ethics as an element of larger efforts to build data literacy among grantmakers and grantees.

Create spaces for conversation and reflection.

Invest in education on data-related topics for current and future staff.

Incorporate data ethics in the grantmaking process.

Create an internal "data ethics point of contact" who can facilitate access to relevant expertise and keep an eye out for latent data ethics risk in projects.

Change grant applicant procedure, to encourage applicants to prospectively consider data-related issues.

Support a central resource to address data ethics concerns for the philanthropic community.

Create a data ethics checklist for grantees and program staff.

## 7  Conclusion

## 8  Reading List on Managing Data-Related Risks

General Background

Data Privacy and Fairness

Human Capital and The Role of Firms

# Introduction

Across society, historically unprecedented quantities of digital information are being gathered, stored, and analyzed. This phenomenon — the growth of "data at scale" — is opening new avenues for philanthropic investment, even as it also transforms key governmental and private institutions. Across nearly every area of their work, funders today enjoy growing opportunities to support projects that gather, analyze, and apply large amounts of digital information to advance charitable aims.

At the same time, because many grantees seek to understand, inform, and shape the behavior of government institutions and large companies, funders are supporting work that brings critical scrutiny to the social impact of data at scale. That work is beginning to highlight ways in which data at scale can harm vulnerable groups and threaten philanthropic priorities.

Philanthropic projects that leverage data at scale are a natural proving ground not only for the benefits of data at scale, but also for the ethical principles and controls that should guide data-oriented projects in the social sector and beyond.

This project is an initial effort to map the ways that data at scale may pose risks to philanthropic priorities and beneficiaries, for US-based grantmakers at major foundations. It reflects desk research and unstructured interviews with 15 key individuals involved in the grantmaking enterprise at major U.S. foundations. We describe our findings, provide a menu of potential practical steps to address these questions, and offer an annotated reading list of relevant resources. Upturn prepared this report at the joint request of the MacArthur and Ford Foundations; both MacArthur and Ford are charter members of the NetGain philanthropic partnership, which aims to help philanthropy address the opportunities and challenges of the Internet age.

This report is particularly focused on areas of risk that may not be well covered by existing regulatory, legal and best practice frameworks. Such risks, where they do arise, are not presently easy for stakeholders to recognize, classify, and address. "Data ethics" is our deliberately capacious term for the risks that arise in connection with data at scale.

It is vitally important for funders to invest in, learn from, and innovate around data-intensive projects, in concert with their grantees. Funders should not treat the new ethical risks of data at scale as a barrier to investment, but neither should such risks be allowed to become a blind spot that can threaten the success and effectiveness of philanthropic projects. We offer this report to support well informed and effective grantmaking in a new and rapidly evolving area.
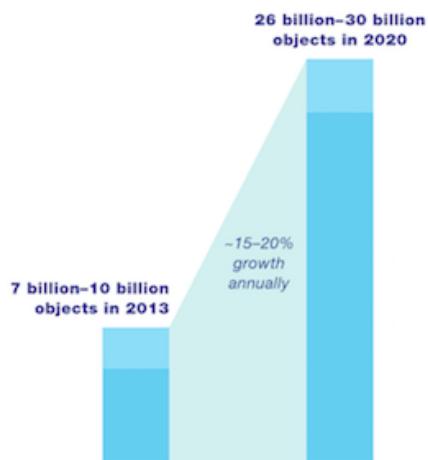
Funders should not treat the new ethical risks of data at scale as a barrier to investment, but neither should such risks be allowed to become a blind spot that can threaten the success and effectiveness of philanthropic projects.

As we detail below, existing ethical frameworks provide strong and effective guidance and risk management for many applications of data at scale. On the other hand, there are some important new risks that do fall between the cracks of existing ethical frameworks applicable to the grantmaking enterprise, and we believe it would be worthwhile for foundations to take simple, low-cost steps to mitigate these risks. At the outset of this study — and throughout our interviews with grantmakers and grantees who consider or execute data-oriented projects — we heard consistently that grantmakers and grantees lack baseline knowledge on using data at scale. Many said that they are unsure how to make better informed decisions, both about data's benefits and about its risks. We assess that these needs will grow in the future as new applications of data at scale become ever more pervasive, including through the growth of digitally enhanced physical objects known as the "Internet of Things" and the intensified use of machine learning in critical contexts that philanthropies aim to shape. Philanthropies will need to be well-informed, networked with experts, and responsive to new challenges as they arise.

## "Data at scale" means what is newly possible with computers.

In this report, "data at scale" refers to the collection, storage, analysis or use of digital information in ways that have only recently become technically and economically feasible. As a rough, admittedly imprecise rule of thumb, therefore, we are referring to activities and grantmaking opportunities that were *infeasible* as recently as the start of 2007.

For context, ten years ago there was no such thing as an iPhone or an Android phone; there are now more than 200 million smartphone users in the U.S. and an expected two billion worldwide.[1] In the same period, the field of machine learning — techniques that allow computers to automatically infer rules and patterns from large quantities of data — has experienced almost exponential growth.[2] As of 2014, the world was generating more data in ten minutes than all of humanity had from its inception through the year 2003.[3] Industry analysts anticipate that worldwide data production will be 44 times greater in 2020 than it was in 2009,[4] thanks in part to 25-30 billion (or more) new connected devices.[5]



*Source: McKinsey & Company*

Worldwide data production will be 44 times greater in 2020 than it was in 2009.

Data collection, storage, and automated analysis — which have long been exponentially plummeting in cost — are now so inexpensive that a growing range of institutions (including some nonprofits) can afford to pervasively monitor and analyze not only their own activities, but also those of the people and organizations around them. For example, rather than classifying music song by song to introduce listeners to new artists or genres, streaming service Spotify monitors and compares the listening patterns of its 100 million users[6] to deliver algorithmically personalized soundtracks to every user each week.[7] Politics too has been transformed: the 2008 Obama presidential campaign ran 66,000 computer simulations each day as part of its voter registration, persuasion and turnout efforts.[8] In the nonprofit context, Direct Relief International used real-time resource tracking and GPS data to augment disaster relief efforts following Hurricane Sandy.[9]
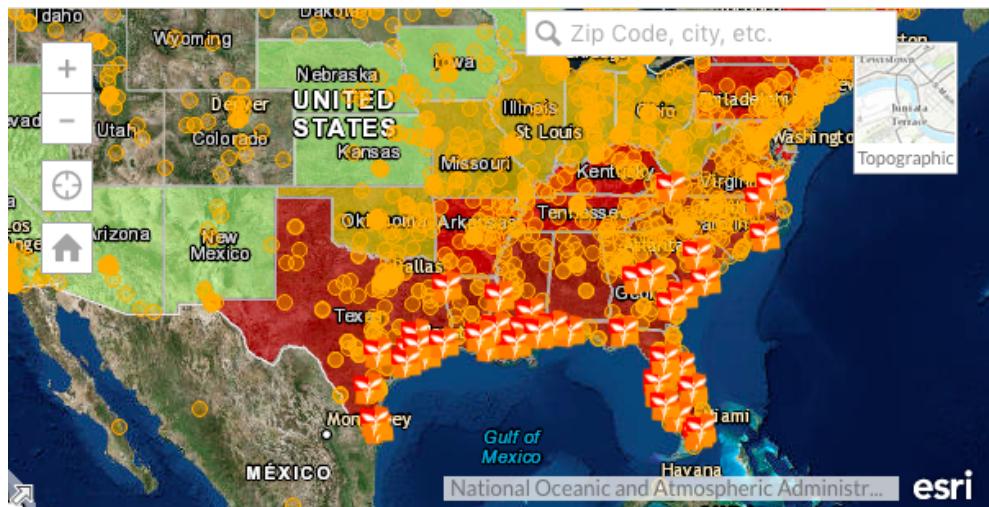
*Image: Direct Relief*

Large technology companies, armed with the data from Internet browsing behavior, smartphones, and other sensors, now automatically construct detailed and cumulative individual portraits of the nearly two billion people (including 65 percent of Americans) who use smartphones in their daily lives, as well as of others whose SMS messages, government forms, and other indicia of daily life are automatically and digitally tracked. Government offices, from intelligence agencies to local law enforcement, increasingly deploy powerful tools to "gather the haystack," such as by recognizing and tracking the license plates and facial features of everyone, not just people subject to suspicion. Even a startup with a few dozen employees can today collect, manage and analyze the daily communications of hundreds of millions of users. For example, at the time WhatsApp was acquired by Facebook, it was a company of 55 employees, processing an estimated 50 billion messages per day, for 450 million users.[10] In the private sector and in government, organizations are gathering vast amounts and varied types of data, and analyzing it (often extensively and automatically) to inform their decisions. As the supply of data expands, so too does the range of important decisions that can be automated or guided by large-scale data analysis.

## The risks of data at scale are complex.

Ensuring data privacy by removing people's names from a dataset, or by storing data in a safe place, is no longer a complete recipe for addressing the risks of data at scale. Today, even novice analysts can *cross-reference data sources and re-identify* what, until recently, would have been a completely anonymous set of data points. For example, researchers have discovered that most anonymized credit card transactions could be re-identified with as few as four rough pieces of information (e.g. purchase location, time of day, purchase amount, and random transaction IDs).[11] As more data is released to the public — often with positive civic aims — these risks become more salient.

Moreover, people are generating large volumes of information that they may consider to be personal but that is technically public, like connections with friends on social media or detailed histories of their whereabouts. While relatively innocuous at a small scale, analyzing this **"public" data in aggregate can reveal sensitive details** like financial status or sexual preference. Analysis of this type of data also reduces some of the privacy that has traditionally been protected by obscurity. Some **semi-private data** that seems benign — like health statistics gathered by fitness tracking devices — **can become hazardous when shared with**

**third parties** like health insurance companies who could preemptively raise rates. The "data exhaust" we generate by using the web, mobile devices and apps feels private and personal, but is increasingly analyzed, shared, and sold.

Data intensive projects also can bring an imprimatur of objectivity or neutrality to processes that actually reflect or reinforce social biases. **Subgroups and communities can be over- or underrepresented**, depending on what method of data collection is used. Collecting information via smartphones, for example, means that populations with low smartphone penetration or poor connectivity may not be fairly represented in the resulting data, leading to skewed results. Indicators like zip code or gender, when included in data-driven decision making tools, may **perpetuate longstanding biases** since models would be based on the status quo, not the desired state.

## Grantmakers are exploring data at scale, but currently have poor visibility into its benefits and risks.

Both the new opportunities and the new risks created by data at scale are significant for foundations of every size. Funders are exploring ways that data at scale can advance the public interest, ranging from more affordable medical treatment to increased economic opportunity,[12] to documenting and addressing disparities in the criminal justice system,[13] and elsewhere.[14] At the same time, funders are also helping to highlight areas where corporate and government applications of data at scale may pose underappreciated risks for marginalized groups.[15]

## Foundations face structural challenges when attempting to invest wisely in data at scale.

**The technologies and practices of data at scale are new and fast-changing.** Cutting edge capabilities to collect and analyze data are being developed in the private sector, in government, and at universities, and the flow of technology into the nonprofit sector has always lagged. This leaves foundations and grantees at a distinct disadvantage in using, but also in following and participating in the conversations about, these new technologies. At the same time, governments and private companies are embracing new applications of data to address societal issues, and foundations are right to want to get involved — whether to collaborate or critique — but lack organizational flexibility to quickly develop institutional expertise.

**Relevant expertise is scarce.** Not only the specialized skills of data science, but also baseline familiarity with the field and its activities, are currently rare among both grantmakers and grantees. As a result, many grantmakers and many grantees currently lack the knowledge necessary to find and judge data or projects involving it.

> Many grantmakers and grantees currently lack the knowledge necessary to find and judge data or projects involving it.

**Program staff lack a clear path for learning more about data at scale.** Data science itself is a nascent field, with no standard curriculum or industry-wide certifications. Some training resources exist for "responsible data" in social contexts, but topics related to data at scale are more complex and technical than other aspects of the "responsible data" discussion, and have not yet been translated for less technical program staff who nevertheless are presented with issues that demand robust knowledge of statistics and large-scale data analysis.

**Guideposts and controls are lacking.** While many issues deriving from uses of data at scale are addressed by existing ethical, legal, and governance frameworks, new risks have emerged that current frameworks fail to sufficiently address. Moreover, program staff and grantees attempting to leverage data at scale have a structural disincentive to highlight data-related risks, since this could prevent potentially impactful projects from being broadly implemented, or could invite criticism that could undermine what might otherwise be positive impact of projects.

# Data at Scale

When working on projects involving data at scale, program officers and grantees sometimes find themselves in new ethical territory, without clear guideposts and, often, without guides. This is not only an issue for philanthropies in their grantmaking but also for nonprofits and funders in their own internal operations as they increasingly adopt data collection, database and tracking tools for communications and development purposes. Our investigation focused on the moments stakeholders encountered ethical uncertainty or hesitation, and our discussions frequently centered on questions that grantees and program staff are unsure how to answer — and, at times, unsure how to ask.

We heard and saw eight key challenges repeatedly across organizations and projects.

## For a given project, what data should be collected, and who should have access to it?

The types of data that are possible to collect are rapidly changing and expanding, but program staff do not always realize why certain new types of data may be particularly sensitive, even when collected from public sources or anonymized. Moreover, while practitioners understand that in theory data ought to be protected, practical questions often remain. These include how to structure data sharing agreements, whether access might be given to partners or law enforcement, whether study results from a sensitive dataset should be made public, and even how to more inclusively involve the subjects of the data in its collection and use.

The Array of Things project, a network of digital sensors across Chicago to measure factors including air quality, noise, and pedestrian and vehicular traffic levels, has faced these concerns with the sensor and camera data it plans to collect. Advocates publicly noted that the project's proposed governing policies do not actually include concrete assurances about when data would be deleted or which partners would have access to it.[16] As sensors continue to get cheaper and easier to deploy, and promise new insights into areas including population movement, public



*Image: Smart Chicago Collaborative*

health and urban life, we anticipate these same question will frequently arise in donor-funded projects.

Call data records (CDRs) are another form of data that have been recognized as a potential source of both benefit and risk in the humanitarian context, since they offer the potential for important insights into populations that might otherwise be difficult to measure. Mobile company Orange confronted this tension when it launched a corporate social responsibility initiative, the Data 4 Development (D4D) challenge, which involved releasing a year's worth of CDRs from
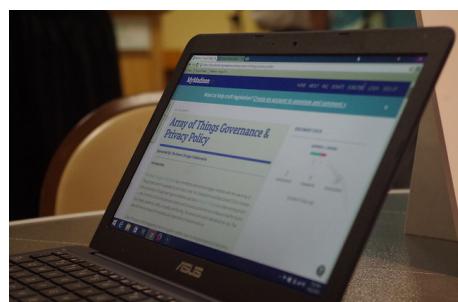
Côte d'Ivoire for population-level analysis. The dataset included call and SMS records, low and high resolution spatial trajectories, and communications graphs of subsets of the company's mobile subscribers. Given the sensitivity of the data, and high risk of re-identification, the company had to decide who would be able to access the data, and what from the resulting research would be published. The company was remarkably conscientious in its treatment of ethical considerations (which we will return to in a following section) — but unfortunately the philanthropic community has not been quite so thoughtful in similar circumstances.



Mapping Ebola. *Image: Mapbox*

During the West African Ebola crisis that began in 2014, nearly every country affected (with the sole exception of Liberia) approved the release of CDR datasets. While most discussions related to the use of CDRs recognized that privacy should be a key concern as organizations attempted to extract value from the information, in the context of this crisis few appeared to consider the highly revealing nature of CDRs as a cause for hesitation in releasing non-anonymized data (de-identification would have defeated the purposed of the epidemiological "contact tracing" process in this case). The intense time pressure and the desire to take action overpowered all these factors, which in another case could have led to a more robust conversation about responsible data use. This experience, in which legal data-sharing regimes were ignored, suggests any sort of conscientious ethical review beyond basic privacy considerations was rare, if it happened at all.[17]

## How can projects decide when more data will help — and when it won't?

Funders and organizations frequently clamor for data before it is clear how this will serve a project's goals. And they often face difficult questions about responsible use only after obtaining the data.



*Image: DFID - UK Department for International Development*

One humanitarian technologist pointed out that there has been "a feeding frenzy of new methodology," particularly technology like sensors and drones in the humanitarian context. For example, several organizations have launched or invested in crisis mapping projects, using GIS and location data to visualize natural and humanitarian disasters or vulnerable populations. However, the maps produced have on occasion been coopted by exactly the actors who are causing danger in the first place (e.g. militias, abusers) — putting subjects at even greater risk. There are also important questions about the validity and accuracy of some of these new methods.[18]

During the Ebola crisis, multiple humanitarian organizations and funders pressed mobile operators to release their CDRs to expedite the contact tracing process meant to reveal peoples' associations and contain the spread of the disease. But several experts have reflected that it was not clear whether call data records would aid in contact tracing, whether contact tracing would work to contain Ebola, or whether location data from the local cell phone towers was sufficiently reliable. This didn't stop major foundations from pushing countries to break their

own laws to release the data. The same organizations that had previously articulated the need for privacy policies and frameworks in other contexts bulldozed through laws and international agreements that governed the use and sharing of this data, under the mistaken impression that more data would solve problems that may have been caused by more straightforward challenges like too few doctors and health workers on the ground.

## How can grantmakers best manage the reputational risk of data-oriented projects that may be at a frontier of social acceptability?

Advocates, the media, and the public are particularly sensitive to what might appear to be intrusive or unfair uses of data. Several projects have faced unexpected negative public reaction related to data ethics issues, requiring resources outside of project scope and budget. Program managers and funders have either needed to scramble to address, or simply were not able to sufficiently respond to, external concern. These experiences suggest dual opportunities for improvement in the grantmaking process for data-intensive projects: on the one hand, an opportunity to better anticipate potential concerns with data-intensive projects during grantwriting and project design, potentially altering course in order to address concerns, and on the other hand, the potential to strengthen capacity to engage on these questions as they arise during project execution.

## When concerns are recognized with respect to a grant, what is the best currently available way to address them?

In several projects we reviewed, grantees and program staff recognized early on that the project's particular methodology of data collection and use would brush up against legitimate privacy concerns. Awareness of such risks, though, did not necessarily or automatically lead to better outcomes.

Projects dealing with data that is already public do not generally receive substantive review by institutional review boards (IRBs).

For example, we were repeatedly informed that projects dealing with already public data — such as social media postings, or sensor recordings in public places — do not generally receive substantive review by institutional review boards (commonly referred to as IRBs, these boards are institutional bodies that evaluate research on several ethical guidelines, which are discussed in a subsequent section of this report).

When they do review projects, IRBs sometimes provide flawed remedies that are meant to address ethical concerns but instead can lead projects to make superficial changes that either do not truly address harms, or even unnecessarily undermine project objectives. When one human rights research organization wanted to help communities in a developing nation to advocate for water rights in the face of mining companies whose operations threatened local water supply, program staff decided to "crowdsource" data about community water needs, tapping local researchers to conduct door-to-door surveys (with the concurrent aim of engaging the community members in a more inclusive research process instead of treating them as passive subjects). The organization hoped to return the aggregated data to local activists to support organizing and advocacy efforts. Since the researchers were collecting information about vulnerable populations, they were sensitive to the need for IRB review, but were concerned that the IRB might not understand the crowdsourcing or data-sharing methods the project hoped to use and would impose restrictions that would unduly hurt the project.

In order to forestall that possibility, the team opted to strip meaningful GPS information from the collected data. This allowed the project to receive an exemption from IRB review, but also removed data that would have been vitally helpful to the community. Moreover, a project manager also recognized that the local researchers would easily be able to re-identify data, since they were the ones who collected it in the first place. In other words, while the IRB did not directly mandate that particular step, the team's understanding of IRB requirements and limits nudged it in a direction



*Image: Surrey County Council*
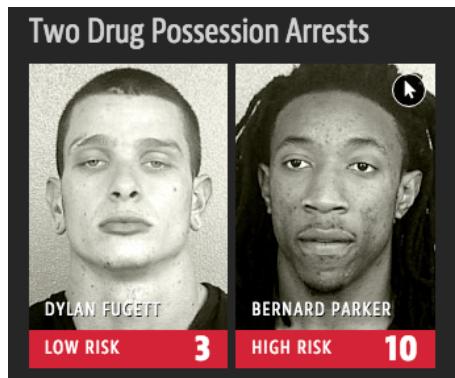
that was neither productive in preventing harm, nor beneficial to the project.

## How can funders and grantees gain the insight they need in order to critique *other* institutions' use of data at scale?

Foundations and the nonprofits they support not only want to use data to improve their own work, but also to evaluate and critique applications of data at scale by other organizations deploying interventions. Staff see other people, organizations, and government bodies using data analysis, profiling or prediction in questionable ways, but cannot access the underlying data in question or do not have the technical know-how or resources themselves to critique these uses of data.

Take, for example, the Arnold Foundation's pre-trial risk assessment tool that aims to reduce bias by excluding certain factors in its risk calculation. The design process for this tool involved thoughtful consideration of key issues by program staff. But the tool is currently used in a growing range of settings and public validation studies have yet to be released, leaving key stakeholders from other institutions working on similar issues with an inability to comment on whether it truly achieves the outcomes for which it was built. (Validation studies are underway now, and planned for public release.) Likewise, in ProPublica's recent investigation into bias in sentencing algorithms,[19] a lack of access to key



*Source: ProPublica*

information (and the resulting need to piece together data through public records requests and other challenging methods) complicated journalists' efforts to assess the system.[20] Several funders of criminal justice work felt that the publication's analysis might be statistically flawed, but did not personally have the skillset to articulate how the study might be interrogated.

## How can the social sector respond to the unique leverage and power that large technology companies are developing through their accumulation of data and data-related expertise?

As companies accumulate growing stores of data about people, as well as trends across communities and society (not to mention having the computing power necessary to analyze it), funders and nonprofits recognize that this data and capacity for analysis would be hugely

valuable for social science research and interventions. But many are uncertain about both how to arrange access to the data and whether using it would be ethical, given that it may have been collected without what researchers would consider fully informed consent.

One program officer described the paradox of knowing that gig economy platforms like Uber hold critical data that would be valuable to labor advocates, but also being unsure about how to protect that data if access was granted, or how to best use it to empower people without violating their privacy.



*Image: CA Dept of Insurance*

## How should foundations and nonprofits handle *their own* data?

At the same time as foundations and nonprofits are grappling with questions of using data-intensive methods to positively impact society, they also struggle with internal data-related challenges. Development and marketing staff are experimenting with A/B testing, tracking, and data mining to reach and mobilize supporters, raise money, and communicate impact, all areas that present their own privacy and data security concerns. Organizations are under acute pressure to use inexpensive and readily available tools, even if in doing so they may undermine their own values with respect to how information should be handled, and may miss the opportunity to help build a better global baseline on privacy. Groups like the Nonprofit Technology Network exist to improve technical competency in the social sector, but conversations about data responsibility and ethics in internal use of data are well behind those happening in the private sector.

## How can foundations begin to make the needed long-term investments in training and capacity?

While it is not happening yet at a meaningful level, we feel it is crucial to point out that many interviewees identified a strong need for long-term investment in training on data ethics issues. Some progress has been made at a project-by-project level, but more sustained investment above and beyond the piecemeal efforts could go a long way.

Many interviewees identified a strong need for long-term investment in training on data ethics issues.

More than half of our interviewees cited organizations like the Responsible Data Forum, which facilitates in-person and online discussions to work through data-related challenges, as a resource they found useful when working through data and ethics issues. Specifically, training webinars, collaborative listservs, and intensive, intimate in-person gatherings have proven valuable to practitioners in integrating a more thoughtful approach into data-intensive projects.

Several grantee staff described how they have organized their own trainings and meetups, but articulated that more engagement and follow-up after these trainings would increase the impact of these one-off efforts. Many also noted that finding funding for trainings and conferences is relatively easy, but getting support for sustained training and long-term follow up has proven nearly impossible.

# Emerging Practices for Addressing Data Ethics Risk

While many of the ethical challenges presented by data at scale are not addressed by existing guideposts, grantmakers are improvising some methods to navigate decisions and to promote thoughtful discussion about these questions with colleagues and peers. Many of these solutions have been moderately successful, but are limited in their effectiveness for the reasons described below.



*Image: Book Sprints*

## Person-to-person, informal consultations

Many interviewees described processes of informally consulting with colleagues or emailing discussion lists (such as the Responsible Data Forum list set up by The Engine Room) for advice on complicated ethical issues. Several Slack channels (open, topic-specific group chat forums) have been introduced to create more spaces for discussion, including one set up by the Humanitarian Operations Mobile Acquisition of Data (NOMAD) group. Data scientists outside of philanthropy also consider related questions on platforms like Stack Overflow, Reddit, and Quora. They most often discover these resources through word of mouth.

While some practitioners have become go-to resources for questions about data ethics, privacy and security, this ad hoc method doesn't catch every issue crossing funders' desks. The people consulted don't always have the tools or time to provide support, and the people requesting help are not certain what kind of due diligence is necessary. Some don't recognize the need for help at all.

## Third-party evaluations of data-related ethical risks

Some program officers are treating programs involving data-induced risk like any other program in their portfolio, and pushing for rigorous testing, monitoring, and evaluation just as they deal with other projects. For example, several foundations have pressed for external evaluation of all pre-trial risk assessment tools, noting that while each may be addressing ethical concern in different ways, the most important goal is that they accurately evaluate recidivism risk and improve outcomes, with ethical risk managed as part of that broader calculus.

Although this approach may address some elements of data ethics risk, foundations don't necessarily have expertise to review, monitor, or evaluate data-intensive programs in-house, and may not be able to identify when a more extensive analysis or ethical review is most needed.

## Ethical reflection on data issues as a required part of the grant application process

One way to ensure grantees are considering ethical questions in their projects from the beginning is to require that they explain their thought process around data-related issues when applying for funding even if the project does not necessarily require formal IRB review. The Sloan Foundation, for example, added an "information products" appendix section to its proposal where grantseekers must articulate their thought processes on data-related questions.[21] While Sloan believes that blanket rules do not provide helpful guidance, they preface the appendix with a statement of the foundation's principles and ask grantees to describe their proposed work through that lens. The process also gives program officers a way to flag grant applications that need extra review.

Such an approach seems to be a promising method of identifying programs that stand to present complex ethical issues, and appears limited only in that it is not yet a broader practice among funders.

## Practitioner-developed resources and sector-specific codes of conduct

Confronted with minimal guidance and inherent discomfort over certain uses of data, some people involved in data-intensive projects have taken the initiative to create their own frameworks to guide responsible collection and use of technology and data in certain contexts.

Participants of the Data Science for Social Good Fellowship at the University of Chicago developed "An Ethical Checklist for Data Science" that presents a series of questions to prompt reflection prior to and during any given project.[22] The Engine Room developed a series of case studies they hope organizations will look to when thinking through challenges in responsible data use, and human rights technologists led the formulation of a code of conduct around SMS use in humanitarian contexts, eventually working with mobile operator interest group GSMA, to formalize and circulate the guidance to both telecommunication and humanitarian organizations that might consider using mobile data-gathering and communication technology in disaster settings.

So far, these efforts have depended on the initiative of motivated individuals who drive the creation of and mobilization around new resources and guidance. But formalizing and promoting widespread adoption of them requires additional investment and central coordination, like the GSMA provided for the SMS code of conduct. Even with institutional support, the effectiveness of self-policing may be limited.

Some broader efforts are underway to develop greater clarity into issues related to data at scale. Scholars and researchers from anthropology, philosophy, law, and economics work under the auspices of the Council for Big Data, Ethics, and Society to address similar questions to those raised in this report. The Council's output has included research and regulatory filings about big data and human subjects research (including a comment on the government's Notice of Proposed Rulemaking to update the Common Rule), case studies on data use, and reviews of ethical codes of conduct in a variety of fields.[23]



*Image: GSMA*

# Policies for Data at Scale

As data becomes easier to collect, store, and analyze, understanding of ethical implications of using data in the context of both social research and program implementation is evolving more quickly than the frameworks that have traditionally defined ethical issues. Some resources provide partial coverage for issues related to data-intensive work, which we outline below, but none were formulated to specifically address data at scale and the patchwork of resources leaves several glaring holes. Many practitioners generally recognize the need to update existing policies, but the conversation is ongoing and in the meantime practitioners are left to navigate the ones that are currently in use, imperfect though they may be.
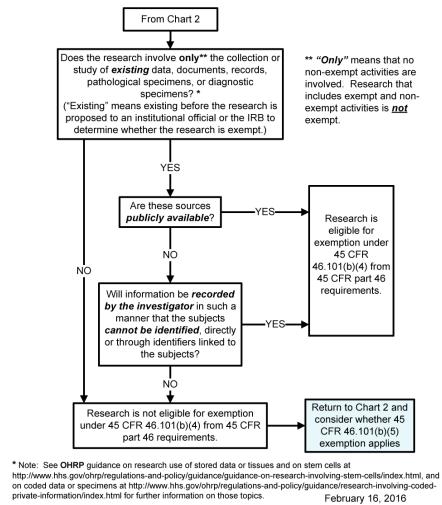
## Existing frameworks address many risks in data-intensive grantmaking — but leave some major gaps.

Foundation staff are broadly familiar with traditional human subject review processes, and incorporate ethical review in this context when relevant risks become apparent. Responsible data practices such as data security and basic privacy protections have also generated some attention across the philanthropic and nonprofit communities, to a limited extent. These existing frameworks provide some key guideposts for projects and initiatives involving the use of data at scale.

### Human subjects review

Initially introduced to address risks in medical experiments, the rules known as "human subjects" protections have since expanded to cover social science research and some social interventions. The Common Rule,[24] the U.S.'s federal policy for the protection of human subjects, requires any university-based and federally-funded research to have research reviewed by Institutional Review Boards (IRB). These boards evaluate projects on several ethical principles that must be protected through the research methods including informed consent, risk-benefit balancing, and considerations of fairness and equity. Researchers dealing with data-intensive inquiries — at institutions requiring or recommending ethical review, at least — can submit projects to IRBs for exemption or approval, a process which can help refine approaches to ethical considerations.

Not all projects that raise ethical issues need to go through human subject review (particularly those outside of the university context) and even those that do can fall between the cracks if they rely on certain types of public or preexisting data that are presumed to not add marginal additional risk to subjects and therefore fall outside of the scope of IRB review, or if future harms are underestimated or not easily predicted.

From Chart 2

Does the research involve **only**\*\* the collection or study of *existing* data, documents, records, pathological specimens, or diagnostic specimens? \*
("Existing" means existing before the research is proposed to an institutional official or the IRB to determine whether the research is exempt.)

\*\* *"Only"* means that no non-exempt activities are involved. Research that includes exempt and non-exempt activities is *not* exempt.

YES

Are these sources *publicly available*?

YES

Research is eligible for exemption under 45 CFR 46.101(b)(4) from 45 CFR part 46 requirements.

NO

NO

Will information be *recorded by the investigator* in such a manner that the subjects *cannot be identified*, directly or through identifiers linked to the subjects?

YES

NO

Research is not eligible for exemption under 45 CFR 46.101(b)(4) from 45 CFR part 46 requirements.

Return to Chart 2 and consider whether 45 CFR 46.101(b)(5) exemption applies

\* Note: See **OHRP** guidance on research use of stored data or tissues and on stem cells at http://www.hhs.gov/ohrp/regulations-and-policy/guidance/guidance-on-research-involving-stem-cells/index.html, and on coded data or specimens at http://www.hhs.gov/ohrp/regulations-and-policy/guidance/research-involving-coded-private-information/index.html for further information on those topics.

February 16, 2016

Human Subject Regulations Decision Chart
*Image: U.S. Department of Health and Human Services*

## Data security best practices and controls

Data at scale introduces significant risk that sensitive information newly digitized or stored as part of comprehensive individual profiles might be accessed by someone besides the intended researcher or program implementer. Based partly on prior experience funding projects involving technology tools, funders are beginning to take action to ensure the security of collected and stored data is incorporated into grantmaking by developing guides for digital security and grantcraft. Particularly when working in sensitive situations or with at-risk populations, simple data protection guidelines such as encouraging the encryption of data at rest and in transit (including by using the cloud-based services of reputable companies) are often adopted to mitigate much of the risk associated with government surveillance or malicious hackers who might attempt to access and use compromising information. But foundations recognize that their ability to provide information security guidance within major areas of grantmaking like nonprofit journalism, climate change and environmental defense, documentary filmmaking, and human rights and humanitarian assistance is still underdeveloped.

## Privacy

Data privacy laws and principles of deidentification provide practical and ethical protection for many risks involving data at scale. Compliance with local laws around protection of privacy, while complicated by the patchwork nature of privacy laws for geographically diverse data sets, requires researchers and project implementers to consider risks to the privacy of the people included in the data in question.

Advances in data analysis tools and computing power mean that basic privacy protections are no longer sufficient.

Even outside of the context of legal requirements, the practice of removing or obscuring identifying data is a cornerstone of responsible treatment of data, and is increasingly recognized as such in the philanthropic community. Data-intensive projects carry markedly less risk when individuals cannot be readily identified from the data they contributed (whether voluntarily or indirectly). Nevertheless, advances in data analysis tools and computing power mean that basic privacy protections like deidentification and data minimization are probably not sufficient to truly protect

data that can, when layered with other private or public information, paint fine-grained and revealing individual profiles of research subjects or program beneficiaries.

### Open access and creative commons

Digital technologies have caused the cost of distributing knowledge and culture to plummet, sparking a revolution in how content is shared. In light of these developments, funders increasingly recognize an ethical imperative to ensure that the research and other outputs they underwrite are broadly shared with the world. There is now a consensus among major foundations including Ford, MacArthur, and OSF that open licensing, such as Creative Commons, should be the default licensing approach for grant-funded work.[25]

These policies, like the other existing guideposts mentioned above, were not developed specifically for projects that involve "data at scale." But these guidelines may well apply to some "data at scale" projects, and they are an illustration of funders' capacity to align information handling expectations with programmatic priorities. Funders' move toward open licenses followed a deliberate, incremental exploration of the issue, and responded in part to the input of the open licensing community.[26]

## Some risks of data at scale are not addressed by existing frameworks

The increased popularity of ever more data-intensive research and analysis techniques has raised three main categories of ethical risk that are both not covered by existing frameworks and less familiar to the philanthropic community.

### "Public" data can now be used in unethical ways.

Even in the most traditional of research contexts, the Common Rule, governing the standard ethical review process for institutions that receive federal funding, explicitly exempts research that relies on preexisting or "public" datasets, under the assumption that reusing data already collected does not add any significant risk for people included in that data.[27] Institutional Review Boards frequently dismiss or wave through research covered under this exemption, without considering downstream effects that using existing data or combining it with other data might bring about.

While this approach may be appropriate for biomedical research protocols, where data collection itself is an intrusive and potentially harmful process, or social research in which disclosure of private data would clearly be harmful, potential harms from data collected from sources like public sensors, social media, or web traffic are less obvious. Datasets can be merged or overlaid with others to uncover previously inconceivable correlations, or triangulated to reidentify data points in ways that would previously have been impossible.[28] Big data analysis, for example, enabled Harvard computer science researcher Latanya Sweeny to re-identify 43 percent of anonymized hospital records in the state of Washington by cross-referencing a publicly available dataset with newspaper stories that mentioned the word "hospitalized."[29] The dataset in this case was not particularly large and the exercise did not require any significant amount of computing power, but it is easy to imagine how a similar tactic could be applied in many other contexts.

In another instance, a full, deidentified dataset from the New York Taxi Commission that included all trip records let analysts determine the religion of some drivers as well as identify certain individual passengers and their riding behavior.[30] While the dataset was "public" in the sense that it was classified as a public record by a municipal office, the people who populated that dataset did not likely consider their taxi trips to be benign public data — particularly given
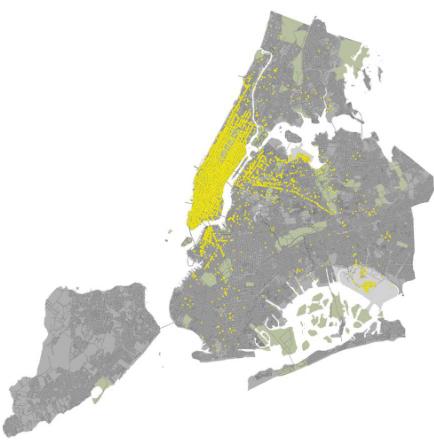
Image: NYC Taxi & Limousine Commission

that the data was later used to reveal private preferences and personally identifying information.

Common understanding of whether certain behavior is public or private has also evolved rapidly: By posting on social media sites, or walking through an urban center equipped with sensors, can we assume that people intended, or even knew, that their data or conduct could be reused for research purposes? This question raised flags in 2006 when a Harvard research project used Facebook network data of 1,700 college students to study the evolution of interests and friendships over time, and then released the data for others to study. The data released turned out to be relatively easy to re-identify, and the students included in the research had never given explicit permission for their data to be used in the first place — but the students had connected with friends and shared interests "publicly." What is the proper ethical process for this sort of data?[31] What about using data that is released in a hack or leak?

Using public data now poses more risk than is recognized by traditional ethical review processes. Asking whether data is public or already exists is no longer a reliable threshold question for judging whether a proposed project, relying upon that data, would be ethical. It is exactly this loophole that leads research using various types of public information to fall through the cracks of existing ethical review procedures.
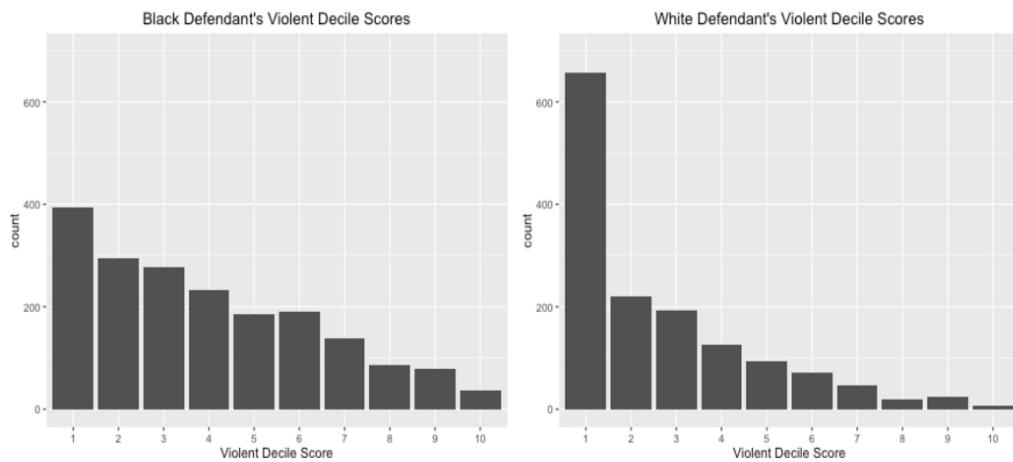
In the philanthropic and public sector, program officers, researchers, and implementers do not yet have a shared schema by which to evaluate ethical implications of using public or existing data, leading to uncertainty and apprehension. This exception is quickly becoming a major ethical liability, and practitioners, with insufficient tools to consider these dilemmas, are increasingly asked to design research and social programs, provide critique, and even make funding decisions based on underdeveloped expertise in the ethical implications of data-intensive work. While some efforts exist in pockets of the community to better understand data-related issues, we did not find any evidence of collective efforts by key players to advance shared understandings across the philanthropic community, even though coalitions like the NetGain partnership have proven effective in advancing conversations around technology in the philanthropic sector.

## Automated decisions can absorb and sanitize bias from input data.

As automated, data-driven decision making systems become more pervasive and more technically advanced, there is a growing appetite to apply automated methods in the service of philanthropic priorities. But along with the consistency that characterizes data-driven decisions, there is a growing recognition that data-driven decisions can also encode social bias, particularly when the decisions are based on patterns reflected in existing data.

*ProPublica analysis of COMPAS's violent recidivism scores. Image: ProPublica*

For example, an online system designed to screen out users who aren't providing their real names may inadvertently stigmatize Native American users, because a name like "Elaine Yellow Horse" is similar to fictitious names provided by non-Native users.[32] And if search engine users subconsciously associate black-sounding names with criminal conduct, then click patterns may lead to socially stigmatizing messages (such as ads suggestive of a prior arrest record) appearing alongside ethnically black names even of people who have never been arrested.[33]

Even when sensitive categories like race or gender are not used in a decision making system, other metrics that may seem benign at first glance (like zip code) can easily reflect the same biases and become accidental proxies for data that might otherwise be excluded.[34]

When data reflects bias, there is a risk of unfair outcomes; as the saying goes, "garbage in, garbage out." For example, if prior arrests are used to target police enforcement activities that themselves *lead* to future arrests, the result may be a self-reinforcing pattern of excessive focus on heavily policed communities.[35] Anyone designing or reviewing new automated tools must be alert to the possibility of reinforcing discriminatory patterns that are reflected in existing data, but program officers and project implementers don't necessarily have the technical or subject-matter expertise to evaluate the full extent of ethical implications of such trade-offs.

Even in seemingly benign cases, the mode of data collection itself can lead to biases in how people are treated. This reality is a costly downside to the "streetlight effect," where the data that is easiest to find is studied, and conclusions are drawn from an incomplete data set. As Dr. Kate Crawford has pointed out, "not all data is created or collected equally," creating shadows in data sets that lead to underrepresentation of some people or communities.[36]

Responsibly funding or evaluating statistical models in data-intensive projects increasingly demands either advanced mathematical literacy, or access to technical experts who recognize these sorts of data-related risks.

## Data at scale may allow powerful institutions to marginalize academic and civil society actors.

As private firms consolidate ownership of key online platforms and succeed in attracting field-leading talent away from academia, these companies are taking on a newly central role not only in shaping daily life, but also in creating and disseminating knowledge. Some of the most

interesting social data is now collected and owned by these large, private firms (and intelligence agencies) with billions of data points, few externally imposed limits on how the data can be analyzed, and near-limitless resources to hire the best data scientists to play with the data.

The movement of data science research from universities to private and government settings further complicates ethical evaluation of research proposals, since companies are not subject to IRB review and the work occurring in the intelligence community is protected by national security exceptions.[37] In-house researchers may be unfamiliar with or resentful of restrictions imposed by ethics review boards on analysis of unprecedented amounts and types of data at the researcher's fingertips.

As foundations and civic actors explore partnerships with private firms to analyze data for social purposes, the appropriate ethical process to take advantage of this rich data for public benefit remains hazy.

Some have raised concerns about the possibility of IRB "laundering," in which researchers bound by ethical review processes broker access to privately-owned, preexisting data to fuel their inquiries, effectively outsourcing ethically nebulous portions of the research to be completed outside the IRB's auspices.[38] While social media platforms and other digital services often include clauses in their End User Licensing Agreements that require users to agree that their data may be used for research, many disagree about whether these check-box agreements provide enough meaningful notice to the users included in research about which they are not informed.[39] As foundations and civic actors explore partnerships with private firms to analyze data for social purposes, the appropriate ethical process to take advantage of this rich data for public benefit remains hazy.

One well-known collaboration between a company and a university team leveraged these circumstances in 2014. Facebook data scientists experimentally modified the platform's News Feed algorithm for a small subset of users in order to test whether the proportion of positive or negative posts people saw affected the emotional tone of their own posts. With the availability of such a large set of previously unavailable data, the team was able to determine that there was indeed an "emotional contagion" effect, and they collaborated with a Cornell-based researcher to publish the findings in a scientific journal.[40] The study may have constituted human subject research in the traditional sense, since the activity could be described as "a systematic investigation, including research development, testing and evaluation, designed to develop or contribute to generalizable knowledge," involving "manipulations of the subject or the subject's environment … for research purposes."[41] But because the data was collected by and only accessible to the Facebook team and not the university affiliate, Cornell's IRB concluded that the project did not fall under its purview and declined to review it.[42] Following public backlash against the research, Facebook voluntarily instituted its own, internal ethics review committee. Not all companies will choose the same course, and similar issues are arising with other actors as more and more data is accumulated in the private sector.

Online dating service OKCupid has no qualms about performing experiments on its users, going so far as to publish a blog post proclaiming "We Experiment on Human Beings!" and describing how the platform has manipulated its users — by indicating to users that were ostensibly "bad" matches according to OKCupid's own matching algorithm that they were actually "good" matches — to learn to what extent perception of romantic compatibility influences behavior.[43]

## We Experiment On Human Beings!

July 28th, 2014 by Christian Rudder
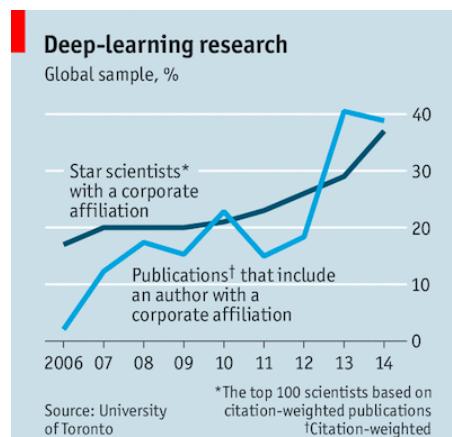
Tweet   Share 0

*Image: OkCupid*

OKCupid is far from alone within the private sector in believing that various forms of A/B testing are standard business practices and that businesses should not hesitate to run tests on their users to improve product performance or gain insight into their user base. Even when ethical concerns are raised in situations like these, some data scientists resist the imposition of ethical constraints, skeptical that anyone without technical expertise can deal with ethical issues any better than data scientists.

Not only is data being concentrated in private hands, but so too is the capacity to analyze it. Data analysis and machine learning talent that formerly gravitated to research universities now flocks to the private sector, which promises big salaries and even bigger datasets to analyze.[44] Such a trend means that, in all likelihood, public-facing research that develops and catalyzes important conversations about ethical uses of data may be the exception to the rule, with most research outputs and technical knowledge staying in-house, in private (or secret governmental) hands.



*Advanced research is moving to the private sector.*
*Image: The Economist*

In 2016, a data-sharing agreement between Google subsidiary company DeepMind and the UK's National Health Service became public, outlining how NHS would give DeepMind access to data from over 1.5 million patient records to be analyzed using the company's artificial intelligence technology.[45] While the agreement covered issues like data security and information governance, the arrangement left no way for patients to opt out of analysis, relying on a broad definition of "implied consent" for direct patient care — but applying the exception on an unprecedented and enormous scale. Traditionally, research and analysis of medical data would be classified as "indirect care," but the promise of personalized medicine through artificial intelligence and big data analysis blurred the line.[46] The arrangement also raised concerns about lack of transparency into Google's work, which was to take place outside of a traditional research context and would arguably directly serve the company's commercial goals.

Ethical review for traditional academic research in the era of big data has lagged behind the evolving nature of data and its related risks, and even more so in the private sector. Industry oversight remains lax.

## Managing Data Ethics Risk

In 2015, the U.S. government published a notice of proposed rulemaking outlining an intent to modernize policies related to human subject protection in such a way that would exempt a large portion of data analysis from review, including exempting novel and potentially revealing analyses of data that is already considered public. Under the current system, review boards are not equipped or prepared to review ethical implications of most cutting-edge data-heavy research,[49] and path-breaking researchers may try to avoid the IRB process altogether.[50] Funding or collaborating with an academic partner or turning to external review boards does not ensure sufficient ethical review, so researchers are turning to other, ad hoc solutions.

### Data science practitioners and researchers are stepping up.

Outside the academic context, a growing group of data scientists and practitioners are using blog posts, articles and conference talks to call attention to ethical issues in data analysis, and are beginning to encourage peers to consider some key questions as they pursue projects. Some organize community events, collect and distribute case studies, conduct internal trainings, and participate actively on mailing lists to pose questions about data ethics in upcoming projects.

Academic institutions are also taking some steps to amplify underrepresented voices in technology-intensive research, albeit in ways that are not specifically focused on data at scale. For example, the Tech Policy Lab at the University of Washington adopted a strategy of "Diverse Voices" to include voices of non-mainstream populations who are nevertheless affected by technology research and policy.[51]

Both individuals and institutions have also introduced internal checklists and methodologies to guide data intensive research.[52] Some guiding questions include:

- Is data science the right tool for this job?

- Does the team include and consider people and institutions who will ultimately be affected by the tool?

- Were the systems and processes used to collect the data biased against any groups?

- Should the team include features that could be discriminatory?

- What are the consequences of acting on false positives or of failing to act on false negatives?

- Can the scope of data collection be limited while still achieving the project aim?

- We want this data, but do we need it? Can other data be ethically collected and used?

- Are particular stakeholders empowered or disempowered as a result of this project?

- Could the collection of the data in this study be reasonably expected to cause tangible harm to any person's well-being?

- Does the project potentially set a precedent for unethical methodologies that could be misused by others in the future?

Others report that they turn to external ethical guidelines, such as the Association of Internet Researchers Ethics Statements.[53] However, many data scientists and practitioners we came across admitted that even these recommendations, which are more applicable to research in the digital context, do not fully capture the needs of social science studies or practical projects.

## Companies are creating their own, internal ethical review processes.

Despite (and in some cases, because of) public controversies over ethics of private research, some institutions have established internal ethical review processes for their own research and for research based on data released externally.

In 2012, mobile provider Orange launched its "Data 4 Development" challenge, in which it released call data records from entire African countries so researchers could try to identify new solutions to development challenges. While the dataset was extremely sensitive, and its release beyond the researchers presented significant risk to the privacy of its subscribers if it were to be de-anonymized, it was not subject to any type of ethical constraint apart from national privacy laws of the target countries. Unsurprisingly in retrospect, the resulting papers submitted by researchers who had accessed the data revealed a mess of previously unpredicted risks. In the second iteration of contest, Orange introduced a dual ethical review process in which both an internal and external panel would review research on ethical and methodological grounds — for not only the research itself but the implications of the projects — and classify research into one of several categories: *Publish without restriction*, *Ask to consider adjustments*, and *Do not publish/Keep for restricted audience.*[54]

Though only 47 percent of surveyed statisticians characterized Facebook's emotional contagion study as explicitly unethical, Facebook too introduced its own internal research review group in response to the backlash the company experienced following the study's publication. The company describes the group as considering implications of using historical data, whether research deals with sensitive populations, and whether the proposal should be flagged for additional expert consultation. Facebook has also said that is has added research-specific employee training, and consults external experts for feedback when presented with research that might have outsized impact on a certain group of individuals.[55]

It is hard to say what impact industry self-policing will have on the state of ethics in big data research, or in data-intensive projects. The movement of major companies toward incorporating heightened consideration of ethical implications is encouraging in part, adding a measure of reflection that was previously not required. But the decision to do so is still voluntary, and committees across companies lack standardization provided by imposed, external rules. Updated federal regulations with broader scope to cover data research (for example, mandating that scientific journals require submitted research to undergo ethical review) could lead to more protection for research subjects, but also runs the risk that cutting-edge research conducted independently by or in partnership with industry would remain unpublished and immune to peer review or public critique.

## Practical Steps for Funders

The challenges described in this report are new and fast changing, and the topic of "data ethics" overlaps with territory covered by other well-established ethical guideposts. We believe that articulating principles or imposing rules on data ethics at this point would be premature. Nonetheless, newly emergent ethical issues inherent in using data at scale points to the need for both a broader understanding of the possibilities and challenges of using data in the philanthropic context as well as conscientious treatment of data ethics issues. Major foundations can play a meaningful role in setting an example in both of these areas, creating space for open and candid reflection on these challenges.

### Include data ethics as an element of larger efforts to build data literacy among grantmakers and grantees.

The data ethics conversation does not, and should not, exist in a silo. General data literacy as well as familiarity with basic responsible data practices are a prerequisite to functional comfort with ethical issues related to data at scale, and these skills should be promoted in concert, including for monitoring and evaluation, development, and communications professionals who are engaging in increasingly data-intensive efforts of their own as they work to measure the impact of grantmaking and interventions and increase capacity of their organizations. Increased data literacy across and outside the philanthropic sector will enhance the ability to implement effective, targeted programs as well as to insightfully critique less responsible uses of data within the public and private sector. Discussions on data ethics specifically will find more fertile ground in the public and social sector as the baseline of fundamental understanding is enhanced.

#### Create spaces for conversation and reflection.

Foundations can promote data literacy at a number of different levels. A medium impact first step would be to organize a cross-institution cohort of program staff who are already frequently exposed to data-intensive projects that would meet on a semi-regular basis to increase their baseline level of data literacy and also discuss data-related challenges they face.

A higher-impact but more challenging proposition would be to reach all program officers at major foundations, even those not yet dealing with data-intensive projects, with the understanding that these staff-people can serve as gatekeepers and initial screeners for ethical issues in proposals and programs. A series of workshops, or inclusion of conference sessions on data literacy at strategically selected major gatherings in the philanthropic and nonprofit sectors would raise the profile of the topics within the community, and prompt more people to seek out additional resources.

#### Invest in education on data-related topics for current and future staff.

To create basic comfort in engaging with these emerging issues, we recommend avoiding rote or factual trainings such as those required to complete IRB protocols. In-person, collaborative workshops may be a better place to begin, particularly given that early experiences with such trainings are varied and it is not yet clear what wisdom ought to be formalized.

A longer term strategy would include promoting data literacy among future foundation and nonprofit staff. Specific steps might be funding university-based initiatives to promote coursework and seminars incorporating data-related topics within graduate programs that traditionally feed talent into foundations, NGOs and large nonprofits, as well as encouraging monitoring and evaluation training programs to include more about data use and ethics within their curricula. Building ethics into professional training has been largely successful in the field of journalism (albeit also challenged by technology and new media) and medicine, and data scientists are similarly beginning to call for ethical training at the university level. As greater awareness of data applications in philanthropy spreads across the sector, conversations about ethical implications will be significantly more concrete and productive.

## Incorporate data ethics in the grantmaking process.

Create an internal "data ethics point of contact" who can facilitate access to relevant expertise and keep an eye out for latent data ethics risk in projects.

A recurrent theme in our research was that, for those data ethics issues not covered by IRB review, program staff lack clear sources of expert advice. That is, when a program officer feels uncertainty or concern in this area, they may not know whom to consult. We believe many program staff would welcome the creation of a central resource that they could turn to for help in evaluating grant proposals that require deep statistical, software engineering or other digital technology related expertise — or background in ethical reasoning specific to those fields — to evaluate.

Thus, we recommend designating a data ethics point of contact who can serve as a voluntary, supportive internal coordinator for expertise and institutional knowledge, to assist program staff in evaluating data-oriented grants. This person's role would be not so much to *be* an expert as to know who to reach out to as questions come up, though the more literate he or she is in the relevant areas, the more helpful this position will prove.

In creating such a role, it will be important to consider the structural incentives of program staff, who may feel self-conscious revealing technological uncertainty and may judge that they have little to gain from highlighting hesitations or ethical questions pertinent to the grants they are considering. As a result, a confidential, internal and informal process — perhaps based on the model of an ombudsperson — could increase the likelihood that candid conversations about these issues will take place where they are needed.

An alternative approach — if one were pessimistic about whether program staff could be enticed, or would be equipped, to flag these issues on their own — would be to build an opportunity for data ethics review into all grants, by allowing the point of contact to read all or certain classes of grants and to initiate conversations where potential concerns are recognized. We do not suggest that this review be required as part of the grantmaking process; at this early stage, it does not seem appropriate for any projects to formally require a data ethics approval. Even without formal power, though, such a review opportunity could help ensure that needed conversations happen.

Beyond its immediate utility, a "point of contact" approach would also help to collect institutional knowledge, surface patterns, and build capacity to engage these issues in the

future. We would not suggest that each discussion of a data ethics issue be formally documented, since such documentation might chill the environment for candid internal discussions. But the point of contact might periodically describe overall experiences or topics that arise.

### Change grant applicant procedure, to encourage applicants to prospectively consider data-related issues.

Adding sections to grant applications where applicants are asked to articulate their intended uses of data as well as to what extent they have thought through ethical challenges (described in more detail below) would help flag projects for further review, and also signal to grantees that the funder values thoughtful consideration of responsible data use in the projects they opt to support.

### Support a central resource to address data ethics concerns for the philanthropic community.

A central organization supported by several foundations could develop a deep expertise in relevant ethical considerations, data security, and data science methodology as relevant to social impact goals, and could serve as an informal or formal resource for grantees to consult prior to implementing their projects. Grantees' engagement with such a resource might be voluntary, at the recommendation of a data ethics point of contact, or might be required, as is the case with IRBs the academic context. Given the early state of awareness and conversation, though, we would recommend that such a body remain an opt-in resource for a meaningful period of time before becoming part of any required process.

## Create a data ethics checklist for grantees and program staff.

Even without any new requirements or policies, it could be useful to equip program staff with questions they can ask about new, data-oriented projects. We have included some examples of questions here, inspired by those identified by practitioners, but caution that these should be considered illustrative rather than exhaustive.[56] As grantmakers develop greater fluency in data-related issues and initial checklists are tested in the field, more effective ways to ask and integrate these questions will likely emerge.

- **Knowing that it is *possible* to do this with data, is it the right thing to do?** Will this use of data achieve the intended goals, or might a different method be more powerful? Does the grantee actually have the capacity to use the proposed tool or evaluate the data collected?

- **How will the data be protected?** What is the risk of re-identification or potentially harmful reuse, and can it be reduced?

- **How will data be obtained?** What, if any, risk of bias is introduced in the method of data collection or the type of data collected or included?

- **Is there a risk of bias here?** If there is risk of bias in using particular classes of data, should the research include this potentially discriminatory data? What if it appears that including such data, even at the cost of harming a few people, may help achieve a better outcome for many more people? How might the project correct for, or otherwise address, this risk of bias?

- **What will mistakes in this data mean for the people involved?** What are the risks or consequences of acting on false positives or not acting on false negatives identified?

- **Do data subjects have meaningful choice where appropriate?** Are the subjects of the project given sufficient notice and opportunity to opt-out, if feasible? If not, are we doing everything we can to mitigate any future harm our project might cause to people or communities by nature of their inclusion in our project?

- **Could this data or its analysis reasonably lead to concrete harm to anyone, including but not limited to the people we are trying to help?** If so, have those risks been thought through?

- **Might this project set a precedent for methods that could be misused by others in the future?**

- **Can the project more effectively involve affected populations in research design and implementation process?** If certain groups are likely to be affected by the research or project outcome, how are we ensuring perspectives from that community are included from the outset?

# Conclusion

Data at scale creates powerful new opportunities for impact and effectiveness for major U.S. philanthropies, and at the same time generates new risks that require thoughtful attention from foundations, program staff, and grantees.

Our report identifies three basic challenges in data-oriented grantmaking that are not well addressed by existing practices: First, public data can now be used in a growing variety of potentially harmful ways, even though its collection and use are not carefully regulated. Second, decisions driven by data at scale offer profound benefits in many areas of work, but also pose a risk of reinforcing longstanding social biases — and the use of automation may give decisions an unearned patina of social neutrality. Third, human capital and institutional expertise for leveraging data at scale are concentrated in certain large companies and government organizations, creating a long-term challenge for the nonprofit sector and academic researchers' capacity to harness these powerful new methods and shape how they are used across society.

We believe that major foundations have a unique role to play in leveraging data at scale for social good, as well as in shaping norms around how this data is treated. We hope that this report helps enable foundation leaders and program staff to act with confidence, embracing change even as they mitigate emergent risks.

# Reading List on Managing Data-Related Risks

## General Background

**Big Data: Seizing Opportunities, Preserving Values (Executive Office of the President), May 2014, https://www.whitehouse.gov.**
This White House report finds that big data analytics "have the potential to eclipse longstanding civil rights protections in how personal information is used in housing, credit, employment, health, education, and the marketplace."

**danah boyd, *Where Do We Find Ethics? — Data & Society*: Points, MEDIUM (Apr. 150, 2016), https://points.datasociety.net/where-do-we-find-ethics-d0b9e8a7f4e6#.iljn0f2gv.**
"How do we enable ethics in the complex big data systems that are situated within organizations, influenced by diverse intentions and motivations, shaped by politics and organizational logics, complicated by issues of power and control?"

**danah boyd & Kate Crawford, *Critical Questions for Big Data*, 15 INFORMATION, COMMUNICATION & SOCIETY 662 (2012), http://dx.doi.org/10.1080/1369118X.2012.678878.**
From the abstract: "Diverse groups argue about the potential benefits and costs of analyzing genetic sequences, social media interactions, health records, phone logs, government records, and other digital traces left by people. Significant questions emerge. Will large-scale search data help us create better tools, services, and public goods? Or will it usher in a new wave of privacy incursions and invasive marketing? Will data analytics help us understand online communities and political movements? Or will it be used to track protesters and suppress speech? Will it transform how we study human communication and culture, or narrow the palette of research options and alter what 'research' means? Given the rise of Big Data as a socio-technical phenomenon, we argue that it is necessary to critically interrogate its assumptions and biases."

**Personal Democracy Forum, Kate Crawford | Know Your Terrorist Credit Score (Jun. 2016), https://www.youtube.com/watch?v=xXs3c42WP8E.**
A short and worthwhile talk in which Crawford explains why people ought to be concerned about inequality that is being embedded in algorithmic systems.

**Jacob Metcalf et al., *Perspectives on Big Data, Ethics, and Society*, COUNCIL FOR BIG DATA, ETHICS AND SOCIETY (Council for Big Data, Ethics and Society), May 23, 2016.**
An informative white paper drawn from an NSF-funded scoping process.

**Tim Harford, *Big Data: Are We Making a Big Mistake?*, FINANCIAL TIMES, Mar. 28, 2014, http://www.ft.com/cms/s/2/21a6e7d8-b479-11e3-a09a-00144feabdc0.html.**
Harford questions some of the assumption driving hype around big data, pointing out basic statistical weaknesses that have plagued prominent applications of large-scale data analysis.

**Benjamin Wittes, *Databuse: Digital Privacy and the Mosaic* (Brookings), Apr. 1, 2011, https://www.brookings.edu/research/databuse-digital-privacy-and-the-mosaic.**

The author "explore[s] the possibility that technology's advance and the proliferation of personal data in the hands of third parties has left us with a conceptually outmoded debate, whose reliance on the concept of privacy does not usefully guide the public policy questions we face" and proposes that the "relevant concept is not…protecting some elusive positive right of user privacy but, rather, protecting a negative right—a right against the unjustified deployment of user data in a fashion adverse to the user's interests, a right, we might say, against database."

## Data Privacy and Fairness

**Moritz Hardt, *How big data is unfair: Understanding sources of unfairness in data driven decision making*, MEDIUM (Sept. 26, 2014), https://medium.com/@mrtz/how-big-data-is-unfair-9aa544d739de.**
The author argues there are "powerful forces that can render decision making that depends on learning algorithms unfair" and explains technical and statistical reasons why this is the case.

**Kate Crawford, *The Hidden Biases in Big Data*, HARVARD BUSINESS REVIEW (Apr. 1, 2013), https://hbr.org/2013/04/the-hidden-biases-in-big-data.**
"Data and data sets are not objective; they are creations of human design. We give numbers their voice, draw inferences from them, and define their meaning through our interpretations. Hidden biases in both the collection and analysis stages present considerable risks, and are as important to the big-data equation as the numbers themselves.

**Woodrow Hartzog, *There Is No Such Thing as "Public" Data*, SLATE, May 19, 2016. http://www.slate.com/articles/technology/future_tense/2016/05/okcupid_s_data_leak_shows_there_s_no_such_thing_as_public_data.html.**
Hartzog points out that there is no universal definition for "public" and that this has led to logical fallacies in how the use of such data is justified

## Human Capital and The Role of Firms

**_Million-Dollar Babies_, THE ECONOMIST, Apr. 2, 2016, http://www.economist.com/news/business/21695908-silicon-valley-fights-talent-universities-struggle-hold-their.**
This piece describes the exodus of qualified machine learning talent from academia into private settings, and warns this could concentrate expertise in disproportionately few firms whose research might not be shared publicly and stunt the growth the field.

# Endnotes

**1**    Number of smartphone users in the U.S., Statista,
http://www.statista.com/statistics/201182/forecast-of-smartphone-users-in-the-us/.

**2**    Yaroslav Bulatov, *Machine Learning, etc: Trends in Machine Learning according to Google Scholar*, Machine Learning, etc (Dec. 16, 2005),
https://yaroslavvb.blogspot.com/2005/12/trends-in-machine-learning-according.html.

**3**    *Big Data: Changing the Way Businesses Compete and Operate* (EY), Apr. 2014.

**4**    *Big Data Universe Beginning to Explode*, Computer Sciences Corporation,
http://www.csc.com/insights/flxwd/78931-big_data_universe_beginning_to_explode.

**5**    Harald Bauer et al., *The Internet of Things: Sizing up the opportunity*, McKinsey &
Company, http://www.mckinsey.com/industries/high-tech/our-insights/the-internet-of-things-
sizing-up-the-opportunity.

**6**    *Monthly active Spotify users worldwide 2016*, Statista,
http://www.statista.com/statistics/367739/spotify-global-mau/.

**7**    Adam Pasick, "The Magic That Makes Spotify's Discover Weekly Playlists so Damn Good,"
*Quartz*, Dec. 21, 2015, http://qz.com/571007/the-magic-that-makes-spotifys-discover-weekly-
playlists-so-damn-good/.

**8**    Zeynep Tufekci, *Beware the Big Data Campaign*, The New York Times, Nov. 16, 2012,
http://www.nytimes.com/2012/11/17/opinion/beware-the-big-data-campaign.html.

**9**    "Big Data vs. Big Storm: New Technology Informs Hurricane Sandy Preparedness,
Response," *Direct Relief*, October 29, 2012, https://www.directrelief.org/2012/10/press-big-data-
vs-big-storm-new-technology-informs-hurricane-sandy-preparedness-response/.

**10**    *Big Data Universe Beginning to Explode, supra* note 4.

**11**    Y.A. de Montjoye et al., *Unique in the Shopping Mall: On the Reidentifiability of Credit
Card Metadata*, 347 Science 536 (2015).

**12**    *See e.g.* Arjuna Costa et al., *Big Data, Small Credit: The Digital Revolution and Its Impact
on Emerging Market Consumers* (Omidyar Network), Feb. 24, 2016.

**13**    *See e.g.* Shoshannah Sayers, *Activists Wield Search Data to Challenge and Change
Police Policy*, Southern Coalition for Social Justice (Nov. 20, 2014),
http://www.southerncoalition.org/activists-wield-search-data-challenge-change-police-policy/.

**14**    *See e.g. Fairness by Design*, Ford Foundation, https://www.fordfoundation.org/the-
latest/ford-live-events/fairness-by-design/.

**15**    *See e.g. The Color of Surveillance: Georgetown Law Conference to Explore Racial Bias of Government Monitoring*, Georgetown Law, http://www.law.georgetown.edu/news/press-releases/the-color-of-surveillance-georgetown-law-conference-to-explore-racial-bias-of-government-monitoring.cfm; Latanya Sweeney, Discrimination in Online Ad Delivery, 56 Communications of the ACM 44 (2013).

**16**    Amina Elahi, *City needs more detail in Array of Things privacy policy, experts say*, Chicago Tribune (Jun. 20, 2016), http://www.chicagotribune.com/bluesky/originals/ct-expert-array-of-things-privacy-policy-bsi-20160621-story.html.

**17**    Sean Martin McDonald, "Ebola: A Big Data Disaster: Privacy, Property, and the Law of Disaster Experimentation," *CIS Papers 2016.01*, accessed July 13, 2016, http://cis-india.org/papers/ebola-a-big-data-disaster.

**18**    *See, e.g.,* Patrick Ball, Jeff Klingner, and Kristian Lum, "Crowdsourced Data Is Not a Substitute for Real Statistics," Beneblog (March 17, 2011), http://benetech.blogspot.com/2011/03/crowdsourced-data-is-not-substitute-for.html.

**19**    Julia Angwin et al., *Machine Bias: There's Software Used Across the Country to Predict Future Criminals. And it's Biased Against Blacks.*, ProPublica (May 23, 2016), https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing.

**20**    Jeff Larson et al., *How We Analyzed the COMPAS Recidivism Algorithm*, ProPublica, May 23, 2016, https://www.propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm.

**21**    *See* Alfred P. Sloan Foundation, Grant Application Guidelines, http://www.sloan.org/fileadmin/media/files/application_documents/proposal_guidelines_research_officer_grants.pdf.

**22**    Alan Fritzler, *An Ethical Checklist for Data Science*, Data Science for Social Good (Sept. 18, 2015), http://dssg.uchicago.edu/2015/09/18/an-ethical-checklist-for-data-science/.

**23**    *See e.g.* Jacob Metcalf, Letter on Proposed Changes to the Common Rule Council for Big Data, Ethics, and Society (2016), http://bdes.datasociety.net/council-output/letter-on-proposed-changes-to-the-common-rule; Jacob Metcalf & Kate Crawford, Where are Human Subjects in Big Data Research? The Emerging Ethics Divide (2016), http://papers.ssrn.com/abstract=2779647.

**24**    *The Common Rule*, 45 CFR. § 46.

**25**    *See e.g.* Katie Shilton, *Emerging Ethics Norms in Social Media Research*, https://bigdata.fpf.org/wp-content/uploads/2015/12/Shilton-Emerging-Ethics-Norms-in-Social-Media-Research1.pdf; Jane Park, *The Open Society Foundations encourage grantees to use CC licenses*, Creative Commons (Jun. 30–2011), https://creativecommons.org/2011/06/30/the-open-society-foundations-encourage-grantees-to-use-cc-licenses/; *Commitment to Open Licensing*, The William and Flora Hewlett Foundation, http://www.hewlett.org/about-us/values-policies/commitment-open-licensing.

**26**    *See* Phil Malone, *Foundation Funding: Open Licenses, Greater Impact* (The Berkman Center for Internet & Society), Feb. 2011.

**27**    Jacob Metcalf et al., *Perspectives on Big Data, Ethics, and Society*, Council for Big Data, Ethics and Society, May 23, 2016.

**28**     Jacob Metcalf & Kate Crawford, *Where are Human Subjects in Big Data Research? The Emerging Ethics Divide* (2016), http://papers.ssrn.com/abstract=2779647.

**29**     Latanya Sweeney, *Matching Known Patients to Health Records in Washington State Data*, SSRN Electronic Journal (2013).

**30**     Metcalf, *supra* note 29.

**31**     danah boyd & Kate Crawford, *Critical Questions for Big Data*, 15 Information, Communication & Society 662 (2012).

**32**     Moritz Hardt, *How big data is unfair: Understanding sources of unfairness in data driven decision making*, Medium (Sept. 26, 2014), https://medium.com/@mrtz/how-big-data-is-unfair-9aa544d739de.

**33**     *See* Sweeney, *supra* note 30.

**34**     Cathy O'Neil, *The Ethical Data Scientist*, Slate, Feb. 4, 2016.

**35**     *See* Jack Smith, *"Minority Report" Is Real — And It's Really Reporting Minorities*, Mic.com (Nov. 9, 2015), https://mic.com/articles/127739/minority-reports-predictive-policing-technology-is-really-reporting-minorities.

**36**     Kate Crawford, Think Again: Big Data, FOREIGN POL'Y (May 10, 2013), http://www.foreignpolicy.com/articles/2013/05/09/think_again_big_data.

**37**     Omar Tene & Jules Polentsky, *Beyond IRBs: Ethical Guidelines for Data Research.*

**38**     *Id.* at 5

**39**     Metcalf, *supra* note 29.

**40**     Metcalf et al., *supra* note 28.

**41**     *The Common Rule*, 45 C.F.R. § 46.

**42**     *Cornell ethics board did not pre-approve Facebook mood manipulation study*, Washington Post, https://www.washingtonpost.com/news/morning-mix/wp/2014/07/01/facebooks-emotional-manipulation-study-was-even-worse-than-you-thought/.

**43**     Christian Rudder, *We Experiment On Human Beings!*, OkTrends (Jul. 28, 2014), http://blog.okcupid.com/index.php/we-experiment-on-human-beings/.

**44**     *Million-Dollar Babies*, The Economist, Apr. 2, 2016.

**45**     Hal Hodson, *Revealed: Google AI Has Access to Huge Haul of NHS Patient Data*, New Scientist, https://www.newscientist.com/article/2086454-revealed-google-ai-has-access-to-huge-haul-of-nhs-patient-data/.

**46**     Subhajit Basu, "Should the NHS Share Patient Data with Google's DeepMind?," WIRED UK, accessed July 28, 2016, http://www.wired.co.uk/article/nhs-deepmind-google-data-sharing.

**49**     Even among academics working with relevant data-heavy projects, studies have found that none of the researchers interviewed were challenged on ethical considerations by IRBs or institutional equivalents, but rather by peers and colleagues or funders. Katie Shilton, Emerging

Ethics Norms in Social Media Research, https://bigdata.fpf.org/wp-content/uploads/2015/12/Shilton-Emerging-Ethics-Norms-in-Social-Media-Research1.pdf.

**50**     One study found anecdotal evidence that younger university faculty members viewed ethics review boards as obstacles to research and tended not to submit work to IRBs for approval. Kalev Leetaru, *Are Research Ethics Obsolete In The Era Of Big Data?*, Forbes, http://www.forbes.com/sites/kalevleetaru/2016/06/17/are-research-ethics-obsolete-in-the-era-of-big-data/. Considering the fact that some evidence from the field of human resources shows that the majority of data scientists have fewer than 10 years of experience, with a median of 7 years, the sentiments from academia are likely similar, if not amplified, outside of the academic context. The Burtch Works Study: Salaries of Data Scientists (Apr. 2016), http://www.burtchworks.com/files/2016/04/Burtch-Works-Study_DS-2016-final.pdf.

**51**     Diverse Voices - Tech Policy Lab, http://techpolicylab.org/diverse-voices.

**52**     *See for example* Patrick Meier, *Launching: SMS Code of Conduct for Disaster Response*, iRevolutions (Feb. 25, 2013), https://irevolutions.org/2013/02/25/launching-sms-code-of-conduct; Alan Fritzler, *An Ethical Checklist for Data Science*, Data Science for Social Good (Sept. 18, 2015), http://dssg.uchicago.edu/2015/09/18/an-ethical-checklist-for-data-science; Networked Systems Ethics, http://networkedsystemsethics.net/index.php?title=Networked_Systems_Ethics#Summary_questions_.28TL.3BDR.29; and Nathaniel Poor & Roei Davidson, *The Ethics of Using Hacked Data: Patreon's Data Hack and Academic Data Standards* (Data & Society Research Institute), Mar. 17, 2016.

**53**     *Ethics*, Association of Internet Researchers, http://aoir.org/ethics/.

**54**     Linnet Taylor, *No Place to Hide? The Ethics and Analytics of Tracking Mobility Using Mobile Phone Data,* Forthcoming in Environment & Planning D: Society & Space (2015); *Challenge 4 Development*, Orange, http://www.d4d.orange.com.

**55**     Molly Jackman & Lauri Kanerva, *Evolving the IRB: Building Robust Review for Industry Research*, 72 Wash. & Lee L. Rev. 442 (2016).

**56**     Fritzler, *supra* note 52. Nathaniel Poor & Roei Davidson, *supra* note 52.