
FORUM ON HEALTH AND NATIONAL SECURITY

ETHICAL USE OF BIG DATA FOR
HEALTHY COMMUNITIES AND
A STRONG NATION

CONFERENCE REPORT

Center for the Study of Traumatic Stress
Department of Psychiatry
Uniformed Services University of the Health Sciences



FORUM ON HEALTH AND NATIONAL SECURITY

ETHICAL USE OF BIG DATA FOR
HEALTHY COMMUNITIES AND
A STRONG NATION

A Workshop of the
Health Services Research Program and the
Center for the Study of Traumatic Stress
Uniformed Services University of the Health Sciences
and
South Big Data Innovation Hub

EDITED BY

Tracey Pérez Koehlmoos, Ph.D., M.H.A.
Robert J. Ursano, M.D.
Carol S. Fullerton, Ph.D.
Paul E. Hurwitz, M.P.H.
Robert K. Gifford, Ph.D.

December 10, 2018

Uniformed Services University of the Health Sciences
Bethesda, Maryland

CSTS



From the Conference Series:

FORUM ON HEALTH AND NATIONAL SECURITY

ETHICAL USE OF BIG DATA FOR HEALTHY COMMUNITIES
AND A STRONG NATION

Editor's Note: This transcript has been edited, however, as in most transcripts some errors may have been missed. The editors are responsible for any errors of content or editing that remain.

IPD 2018 by Center for the Study of Traumatic Stress
Department of Psychiatry
Uniformed Services University of the Health Sciences
4301 Jones Bridge Road
Bethesda, MD 20814-4799

First Edition

CONTENTS

	Page Number
ATTENDEES	<i>i</i>
EXECUTIVE SUMMARY	I
INTRODUCTION	3
BIG DATA AND DATA ETHICS	5
POTENTIAL SOLUTIONS	11
CASE STUDIES	17
CONCLUSION	21
APPENDIX: REFERENCES, FORUM READINGS, VIDEO AND CASE STUDIES	23

ATTENDEES

Cecilia Aragon, Ph.D. (via teleconference)
Professor, Human Centered Design &
Engineering
University of Washington
Seattle, WA

Gari Clifford, D.Phil.
Interim Chair, Associate Professor
Biomedical Informatics
Emory University and Georgia Institute of
Technology
Atlanta, GA

Adam Davis, Ph.D.
Director, Registries and Research Data
Banks
Uniformed Services University of the
Health Sciences
Bethesda, MD

Jared Elzey, M.S., C.R.A.
Director, Services & Support Unit
Office for Research
Meharry Medical College
Nashville, TN

Robert Gifford, Ph.D.
Assistant Professor
Department of Psychiatry
Executive Officer, Center for the Study of
Traumatic Stress
Uniformed Services University of the
Health Sciences
Bethesda, MD

Valery Gordon, Ph.D., M.P.H.
Senior Advisor for Human Subjects
Protection
National Center for Advancing
Translational Sciences
National Institutes of Health
Bethesda, MD

Khara Grieger, Ph.D.
Senior Environmental Research Scientist
RTI International
Research Triangle Park, NC 27708

Edmund Howe, M.D., J.D.
Professor, Department of Psychiatry
Scientist, Center for the Study of
Traumatic Stress
Uniformed Services University of the
Health Sciences
Bethesda, MD

Angela Icaza, M.S., M.B.A.
Program Director
Clinical Informatics Policy
Office of the Assistant Secretary of
Defense for Health Affairs
Arlington, VA

Mary Kelleher, M.S.
Office of Research and Development
Department of Veterans Affairs
Arlington, VA

Tracey Pérez Koehlmoos, Ph.D. M.H.A.
Associate Professor of Preventive
Medicine & Biostatistics
Uniformed Services University of the
Health Sciences
Bethesda, MD

Continued

Brandeis Marshall, Ph.D.

Professor, Computer Science
Spelman College
Atlanta, GA

Joshua Morganstein, M.D.

CAPT, USPHS
Assistant Chair and Associate Professor
Department of Psychiatry
Assistant Director, Center for the Study
of Traumatic Stress
Uniformed Services University of the
Health Sciences
Bethesda, MD

Elizabeth Newbury, Ph.D.

Director of the Serious Games Initiative
The Wilson Center
Washington, DC

Wendy Nilsen, Ph.D.

Program Director
National Science Foundation
Alexandria, VA

COL John Scott, U.S. Army

Data Manager
Defense Health Agency
Falls Church, VA

Lea Shanley, Ph.D.

Fellow, Nelson Institute for
Environmental Studies
University of Wisconsin-Madison
(former) Co-Executive Director
South Big Data Innovation Hub
Madison, WI

Robert Ursano, M.D.

Professor, Department of Psychiatry
Director, Center for the Study of
Traumatic Stress
Uniformed Services University of the
Health Sciences
Bethesda, MD

Mimi Whitehouse, M.B.A., M.I.S.

Manager
Accenture Federal Services
Washington, DC

John Wilbanks

Chief Commons Officer
Sage Bionetworks
Washington, DC

Gary Wynn, M.D.

LTC, MC, USA
Assistant Chair and Professor
Department of Psychiatry
Senior Scientist, Center for the Study of
Traumatic Stress
Uniformed Services University of the
Health Sciences
Bethesda, MD

CONFERENCE PLANNING COMMITTEE

Tracey Pérez Koehlmoos, Ph.D., M.H.A.
Robert J. Ursano, M.D.
Brandeis Marshall, Ph.D.
Jared Elzey, M.S., C.R.A.
Karl Gustafson
Lea Shanley, Ph.D.

EDITING COMMITTEE

Tracey Pérez Koehlmoos, Ph.D., M.H.A.
Robert J. Ursano, M.D.
Carol S. Fullerton, Ph.D.
Paul E. Hurwitz, M.P.H.
Robert K. Gifford, Ph.D.
Alexander G. Liu, M.P.H.
Hieu M. Dinh, B.S.
Madeline E. Crissman, B.A.

EXECUTIVE SUMMARY

Recent advances in artificial intelligence (AI) have made AI a powerful tool for research, particularly for extracting meaningful insights from extremely large data sets. These developments dramatically increase both the research benefits of big data and the risks posed to individual privacy, forcing a new examination of ethics in research. To advance discussion of research ethics in this context, the *Forum on Health and National Security: Ethical Use of Big Data for Healthy Communities and a Strong Nation* was held on December 10, 2018 at the Uniformed Services University of the Health Sciences (USUHS) in Bethesda, MD. The workshop was sponsored by the Health Services Research Program and the Center for the Study of Traumatic Stress (www.cstsonline.org) of the Uniformed Services University of the Health Sciences, along with the South Big Data Innovation Hub (southbigdatahub.org). The workshop was designed to identify ethical questions relevant to military health research studies using big data.

Using two case studies as a focal point, participants outlined the ethical problems that can arise during a research project. Through the case studies and group discussions, they explored researchers' ethical obligations to research subjects, particularly in the area of privacy, trust, and consent, as well as potential methods to improve researchers' ability to collect, access, and share data while protecting privacy. The discussions yielded several solutions that endeavor to embed ethics more thoroughly into each step of a research project. These include creating risk management frameworks and data governance policies, improving education and workplace training, and increasing community involvement in research design and practice.

Through the case studies and group discussions, they explored researchers' ethical obligations to research subjects, particularly in the area of privacy, trust, and consent, as well as potential methods to improve researchers' ability to collect, access, and share data while protecting privacy.

INTRODUCTION

The *Forum on Health and National Security: Ethical Use of Big Data for Healthy Communities and a Strong Nation* was held at the Uniformed Services University of the Health Sciences (USUHS) in Bethesda, Maryland on December 10, 2018. The workshop was developed as a partnership between the Health Services Research Program and the Center for the Study of Traumatic Stress of the Uniformed Services University of the Health Sciences (the federal university of the health sciences), and the South Big Data Innovation Hub.

The workshop evolved from ongoing discussions at the intersection of ethics and big data in the military health field. To explore these areas more deeply and solicit diverse ideas, all the organizations partnered to host a forum on the ethics of combining big data, artificial intelligence (AI) tools, and military health and performance information. The workshop was designed as an interactive discussion to encourage participants to share their experiences, ideas, and expertise while learning from others. Pre-forum readings, and video and case study links were distributed (see Appendix).

The forum began with welcoming remarks outlining the workshop objectives and context and participant introductions. A discussion of ethical questions in three dominant themes then began: what is big data, what are the sources of big data, and how can big data be handled ethically? In the afternoon, participants discussed the ethical quandaries raised by two case studies. A final discussion summarized highlights from the day's discussions.

The workshop evolved from ongoing discussions at the intersection of ethics and big data in the military health field.

BIG DATA AND DATA ETHICS

Researchers are increasingly applying AI to large data sets to advance health research. While this is being done with the worthy goal of improving the health of communities, it is essential to ensure that ethical considerations are woven throughout the entire process to mitigate potential negative consequences to individuals and the nation. Unfortunately, the path to doing this is not always clear. In medical research, data is coming in from multiple sources with multiple stakeholders and undergoes multiple iterations in increasingly complex, confusing, and vulnerable systems. In the context of military health data, the use and application of AI can carry threats to national security as well as implications for individual privacy.

The *Forum on Health and National Security: Ethical Use of Big Data for Healthy Communities and a Strong Nation* sought to gain clarity on several key issues surrounding the best ethical approaches to health research with large data. Its stated objectives were to identify key ethical issues, determine mechanisms to mitigate harms, identify gaps in research systems, and identify possible solutions.

In the context of military health data, the use and application of AI can carry threats to national security as well as implications for individual privacy.

What is Big Data?

While the world is awash in data, the term “big data” refers to data from many sources, is merged with other data, and contains multiple data marks, data points, servers, and metadata (information about data that scientists use for interpretation and correlation). The scale varies, for example, big data can mean a large amount of data about one person, or one data point about each member of a large population. It can refer to one very large data set, or thousands of small data sets analyzed together. To illustrate the latter, an example of disaster-response data crowdsourced from thousands of smartphones was discussed.

A series of parameters were discussed, colloquially called the “Five V’s,” that are commonly considered when discussing big data issues. The V’s include velocity (How fast does the data come in?), variety (What form is the data in and where does it come from?), veracity (How accurate is the data?), volume (How much data is there?), and value (Is the data benign or dangerous? What value will your findings or outcomes have?).

Wearable devices, which rapidly log multiple types of health data, represent one source of big data relevant to health research. In one study, patient data from 22,000 real-time streaming sensors was so large, varied, and fast that it essentially broke the algorithms intended to analyze the data. Big data is typically too large and complex to handle without computational assistance. As a result, analyzing big data increasingly requires advanced AI algorithms. Just as each disease requires its own treatment plan, each data type requires its own algorithms. In addition, big data usually requires strong security, although its type determines its requirements, because, like in medicine, a “one-size-fits-all” security model is not possible.

All of the characteristics of big data, from collection to security, access, and analysis, have ethical challenges.

What are the Ethical Challenges?

All of the characteristics of big data, from collection to security, access, and analysis, have ethical challenges. For example, the velocity of data coming in can put stress on researchers to speed up their experiments (Hayden, 2018). In addition, researchers frequently conduct multi-site, multi-data collaborations, bringing in a variety of data types whose interaction, AI implementation, and management present many complexities which can result in active data repositories, data integrations, data analyses, and publications with potential ethical quagmires.

Another ethical challenge is that businesses and research institutions approach data very differently. Government agencies or government-funded research institutions must follow regulations which businesses do not. In a company, working quickly with data can fuel innovation, but without ethical or regulatory constraints, companies might bypass ethical considerations.

Military health data has an extra ethical challenge. Primary health data is readily accessible for patient care, but secondary uses, such as for research that can improve health delivery requires special considerations. In addition, it can be a challenge to find analysts with appropriate data and security expertise who will choose to work in government rather than industry.

Potential secondary use of data is another important issue that impacts the privacy of the data holders. Consideration of individuals' ability to own their data is important. Data selling has been seen to be a potentially exploitative practice. The complexities of data ownership and use are many.

A thoughtful deliberation of ethics takes time to balance benefits and harms to a participant, time to reach data sharing agreements, time to consider unintended consequences, and time to implement them.

Why are Ethics Important?

Ethical use of big data is particularly important for health-relevant big data. The area is new and law and practice are still evolving. Big data reflects and represents actual people. Therefore big data studies require consideration of data as conscientiously as one would for human subjects or whole populations. For further elaboration on these issues, reference was made to *Ten Simple Rules for Responsible Big Data Research* (Zook, Barocas, Boyd, et al., 2017).

Ethics are especially important in military health research because scientists have a dual obligation in this context: to protect people and to protect national security. For example, health data of individual armed service members or their families could indicate troop deployment movements that may have national security implications.

Current laws that address privacy or security are limited and require reconsideration with big data in mind. While laws do govern the collection or use of certain data on individuals, the process by which individuals can actually recover their personal data is often cumbersome, slow, or non-existent. Violating such laws may require actions, but often people are more likely to have a frustrating experience at a cost to their personal time and expense.

For researchers and health care providers, the Health Insurance Portability and Accountability Act (HIPAA) creates certain restrictions, but private companies have been known to infringe on these areas. In addition, while HIPAA regulations may prevent unethical or illegal uses of protected health information (PHI), personally identifiable information (PII) has few regulatory protections. Safeguarding PII is important, in part because PHI and PII are so interrelated. Some have stated that

all data is ultimately health data because data collected by our credit cards, internet search engines, and location apps can point to health issues, including information that is otherwise protected by HIPAA. Consumers get multiple warnings about protecting PHI or financial information, however, PII is readily available to commercial companies. Those companies may take strides to protect data, but unlike PHI which is protected by standardized reviews, documentation, and institutional oversight, consent and the potential uses of PII are much less regulated/standardized.

Some private sector agencies have proposed four specific perspectives on ethical issues: fairness, accountability, transparency, and explainability. Identifying problems in these areas helps the organization make ethical recommendations. For example, in the financial services industry, it has been reported that some have been denied personal loans based on data points such as their zip code, a possible proxy for race. This practice fails both the fairness and explainability tests. Ethical considerations may also be misused as rationalizations for decisions.

Can Data be Protected?

Data protection is often central to efforts to promote the ethical use of data. Once data is collected, shared, or made accessible it becomes vulnerable to exposure, capture, and unethical re-use or repurposing without user consent and possibly in violation of user privacy. Where data travels and how it may subsequently be used is difficult to determine in today's interconnected world. Cybersecurity protections are an important and complex area of big data ethical collection, storage and use.

Data is often voluntarily given, such as in the case of store rewards programs, however data can still be stolen, whether federally or commercially protected. Published data can be used in ways that the original researchers or subjects never intended. For example, data points unrelated to study outcomes have been used to predict individuals' sexual orientation (Wang & Kosinski, 2018). Users often voluntarily post data or opinions in one context that can unintentionally spread online. This suggests that people have a poor understanding of privacy and context on the internet and that resilient systems are necessary to account for these changing contexts and consequences.

Health data is particularly sensitive. PII can be used to infer health data about a person. Health data itself can also uncover information a user wants private.

Discussion addressed if it was possible to use AI to reverse the barrage of privacy intrusions and alert a user if their data is tampered with, flag abnormalities, or counterattack companies or actors handling data irresponsibly. Could such applications help to counter AI being used against people?

It is necessary to balance scientific inquiry with data use consequences and support the ethical principles of equality, dignity, and justice for research subjects.

Data Sources, Research Practices, and Community Engagement

Workshop participants next addressed how various data sources, research practices, and ways of engaging with communities offer opportunities for research, as well as the potential for ethical conflicts. They explored how these facets may be addressed through different mechanisms for advancing data protection and ethical research.

Data Sources, Engagement, and Transparency

Health data comes from many sources—from GPS units and accelerometers

Ethics are especially important in military health research because scientists have a dual obligation in this context: to protect people and to protect national security.

Many believe that users should be able to access their own data.

to surveys to proteomics, to name just a few. Researchers seeking health data must engage ethically and transparently throughout their interactions with both the data and the participants from which it is sourced in ways that encourage trust, respect individuals' comfort with risk, and enable truly informed consent. It is also important for researchers to have diversity in their teams in order to have broad perspective. The public engages with data, and researchers engage with data, and each of these relationships has ethical dimensions. Ethical researcher engagement with data requires security mechanisms that balance data privacy with data access.

Personal data can be shared knowingly or unknowingly by the public. People willingly give out personal data to banking as well as fitness apps. The transactions and data use in such contexts is mostly transparent. However, data is also unknowingly given out and used without consent or notification. In this case, data use may be opaque because users do not know when, or for what purpose, their data is being used.

Many believe that users should be able to access their own data. For example, some heart monitors only provide data to an approved doctor, but a patient might want to get a second opinion, see the data for themselves, or even see the algorithm itself. Researchers who ask for data should anticipate such requests. Programs are working to balance meaningful, open data for individuals with the potential to use the data for business opportunities.

Trust

Ethical considerations can help researchers understand the reasons why people may not trust data collectors. Individuals may view unsanctioned use of their data, especially by the government or large corporations, as threatening. Recent large-scale data breaches were discussed. How to build and ensure trust is a task for researchers and systems that collect data.

Younger people are more familiar with the digital world and while they may not necessarily trust it, they appear to often be more tolerant and accepting of data collection, especially if it adds value or convenience to their lives. Others are often less familiar with digital data and more suspicious when asked to provide information about themselves. Communities can also influence trust and a person's willingness to share data. Native Americans and other minority communities in particular, have been noted to be protective of data ownership and interpretation because of historic experiences of violations of trust.

How organizations treat users can influence how much trust they hold. Organizations that act transparently and are responsive to complaints appear to increase a user's sense of agency and trust. Being unresponsive and opaque appears to reduce trust.

For citizens of some countries, trust in health care is not an issue since the government is the only health care provider. The complexity of health care in the United States may create more opportunities for reduced trust and acceptance. In the U.S., individuals often trust that the law will protect their data privacy. Data breaches and their consequences are a new legal area.

Including electronic devices and their data in research studies adds another layer of requiring trust. Data agreements with device manufacturers are rare. For example, subjects might wear a data collection unit on their body in a study, and while there is an agreement between the researchers and the subject, there is not an agreement with the device company, which may also have access to the data. It is also possible

to make inferences from a device worn by a spouse about potential behaviors of the spouse not wearing a device.

Risk Perception and Risk Tolerance

In addition to varying levels of trust, the public has varying levels of risk perception (how serious is the risk?) and risk tolerance (what level of risk can I accept?). For example, driverless cars have a high risk perception—the downsides can be fatal—and a low risk tolerance—fatal consequences are often assessed as not worth the risk. These risk perceptions can prevent widespread adoption of alternatives that are overall safer. People often react to a perceived loss of control rather than statistical evidence.

Ethical data systems need to consider, respect and respond to users' risk perceptions and risk tolerances. Risk perception is important. High risk perception can cause health care workers to stay home during an outbreak to avoid contracting a disease, however, if administrators present risks accurately and educate workers about preventive measures, risk perception decreases and workers report for duty. A target population's risk perception and risk tolerance can also be influenced by their specific health issues. Someone with a terminal diagnosis has a very different risk profile than someone who is healthy.

Familiarity with devices, ease of navigating the digital world, and a sense of personal agency affect a person's risk perception and risk tolerance and, like trust, vary widely across cultures and demographics. Self-driving cars score poorly in these three areas, which is one reason they make people uneasy. Scientists must consider risk perceptions which are often dramatically understated gaps between real and perceived risk. People can also underappreciate future risk. For example, genetic testing technologies that use saliva samples of DNA are increasingly popular and have even been used to solve crimes. They collect very personal data on non-criminals that are vulnerable to exposure. Regardless of the level of risk tolerance, ethics requires that data be adequately protected.

Consent

Consent and trust are closely linked. When users consent to provide data, they are placing trust in the researchers, perhaps akin to a patient's trust in a doctor's care. Ethical consent goes beyond a signature on a form. Ethical consent agreements must be understandable to the general public, communicate risks and benefits, and accurately convey how the data will be used. Ethical consent also may need consideration of non-participants. With today's interconnected digital platforms, data about one person can expose data about family or community members, who have not given consent, to unintended risks.

Participants should understand the research goal, and should also understand who is expected to benefit from the research. Usually, participants do not directly benefit, as most studies are testing a hypothesis. However, some studies can tell subjects whether they have a gene for certain diseases, such as breast cancer, allowing participants to then act on that knowledge if they wish. The Department of Defense requires that research on children conducted under its purview has to show not just an overall benefit to all children but a specific and individual benefit to each participating child. While this is a difficult hurdle to overcome because research is unpredictable, researchers have adapted their methodology accordingly.

Participants explored several shortcomings of existing consent mechanisms. First,

In addition to varying levels of trust, the public has varying levels of risk perception (how serious is the risk?) and risk tolerance (what level of risk can I accept?).

A second problem is that consent for one study rarely covers secondary uses of data after the initial research is concluded.

obtaining consent from individuals is a time-consuming task. Participants suggested a mechanism for obtaining consent from large populations, akin to businesses' user agreement policies, would dramatically speed the process. However, research participation has fewer immediate or apparent benefits than, for example, joining a social network and can put users' data at risk. Such hurdles would need to be addressed before such a process could work for large populations.

A second problem is that consent for one study rarely covers secondary uses of data after the initial research is concluded. To address this, the U.S. Department of Health and Human Services' (HHS) Revised Common Rule includes a new section on "Broad Consent" to allow such use. However, if people refuse to grant that broad consent, their data becomes unusable for any purpose, which can make Broad Consent impractical for health researchers.

Consent protocols have been developed that specifically address future data use. If participants do not consent to future use, they are excluded from studies where outcomes depend on secondary uses. In other studies, participants can opt out of specific types of future data use, such as by allowing future use of their clinical data but disallowing future use of their genomic data.

Another thorny problem is the use of employee data. An employer's decisions about whether and how to use this data can affect employees directly.

The perceived benefits of a research endeavor often play an important role in the process of obtaining consent for use of data. Health studies using big data can sometimes offer valuable benefits to a community in exchange for data access. For example, in a malaria vaccine research study conducted in low-income countries, researchers made an agreement with the involved communities that if the vaccine worked, it would be made available *and* affordable there.

One action area may be for researchers to pivot from thinking about getting "informed consent" to a framework of "advise and consent," in which a researcher provides information to participants, advises them of the risks, and then asks for consent. Adding the "advise" component creates a true interaction with separate responsibilities, i.e., the researcher must properly inform a user about the research question, data collection, and how it will impact the participant, and the participant decides whether to give consent or not.

The Importance of Diversity

Participants agreed that ethical consideration requires a diversity of backgrounds, opinions, and ideas on the research team. Technology is created by people whose personal perspectives, biases, and limitations can end up in their work (Buolamwini, 2016). In one example, facial recognition software was unable to "see" the person because her skin was darker than the faces on which it was trained. Ethical researchers must identify these limitations.

Participants also agreed that inviting outsiders and experts in different discipline areas into the research process garners surprising insights and improvements. For example in creating educational games, standards require that several outside experts are consulted, such as educational psychologists and teachers of the targeted grade range. Non-experts can also provide added benefits, for example, including representatives from the communities whose data is being collected.

Diversity can also be important in the data management and research communities, requiring an expanded "pipeline" to bring underrepresented groups into AI research. Science requires diversity for many elements.

POTENTIAL SOLUTIONS

Ethics is a holistic endeavor. Big data health research requires ethical considerations in each step of the research process from designing the research question, to determining effective data collection methods, to creating algorithms for analysis (Loukides, Patil, & Mason, 2018). Participants suggested multiple ways to embed ethics into the broader research endeavor creating effective risk management frameworks and data governance policies, improving education, trust, and diversity, learning from existing systems, rethinking the approval process, and reframing the human-AI relationship as a collaboration instead of a competition.

Ethics is a holistic endeavor. Big data health research requires ethical considerations in each step of the research process from designing the research question, to determining effective data collection methods, to creating algorithms for analysis.

Risk Management Frameworks

The risks of using AI on health data are complex and often unfamiliar. The creation of an ethical risk management framework that identifies problems, assesses risk, makes mitigation plans, communicates risk, seeks feedback and considers the community and reassesses risk can facilitate ethical research. Having a comprehensive framework in place can facilitate management.

While big data is diverse it is still possible to create a framework that facilitates researchers acting ethically, communicating risk, and encouraging innovation while being malleable enough to adapt to the range of projects. Risk-prevention mechanisms can be designed and incorporated enabling researchers to add resilience to a system or strengthen security features. For example, it would be helpful to build in “reverse mechanisms” to retrieve or destroy sensitive data in the case of a security breach. Despite the best intentions, however, problems are inevitable. Data cannot be “recalled” like a flawed consumer product. It can be endlessly copied or transferred and become untraceable. Ethical and effective risk management frameworks can mitigate problems and ensure consideration of ethical approaches to actions.

Following a risk management framework can create extra work. Developers or analysts may need a tangible incentive to take it on, in addition to a better awareness of the risks of not handling data ethically. A code of conduct or checklist can also nudge workers to “do the right thing.”

Data Governance

Data governance is an essential element for ethical big data research. Current data management systems have limited ability to handle today’s challenges. Ethical data governance ensures that data is findable, accessible, interoperable, and reusable. It also requires knowledge of risk management tools and mechanisms for predicting and mitigating risks. Organizations are increasingly appointing a responsible steward to oversee this process.

Ethical data governance creates guidelines for issues such as limiting data collection to only that necessary to satisfy the research question, safely sharing raw data between researchers, and restricting large data set transfers. There is a growing

Cybersecurity is an essential piece of the ethical puzzle. Data governance plans must include proper security for data, and concerns about cloud-based services, expense, risk, and failure mitigation must also be considered.

awareness that sharing data can increase risk. Analyzing large data sets on-site can, at times, decrease the space and security needed to transfer large files but has its own challenges.

Different research questions require different data, multiple algorithms, or separate analyses. Secure data repositories are a core part of data governance, providing researchers with different levels of access (e.g. raw data, metadata only, or results from algorithm deployment only) depending on their associations. Such data repositories can add cost and require specialized training.

Sharing data is not a new concept. Accomplishing sharing in large governmental organizations can be particularly challenging and time consuming. There are few current incentives to improve data sharing and governance. Biospecimens have relatively well-established governance protocols that data scientists may be able to learn from.

Cybersecurity is an essential piece of the ethical puzzle. Data governance plans must include proper security for data, and concerns about cloud-based services, expense, risk, and failure mitigation must also be considered. It is also important to consider who bears the computational costs of securing, accessing, or analyzing data sets. Large organizations are constantly balancing the ability to have useful health data with its cybersecurity considerations, national security implications, and the ability to improve health care.

Data literacy is the flip side of data accessibility, and requires tools to aid data interpretation. Tools are available or in progress to improve health literacy among the public and to encourage researchers to consider health literacy throughout the entire research process.

Data governance plans must also ensure that data is transparent and usable. Many organizations are working on improving the usability of data. The VA and DoD, for example, are engaged in developing such platforms for health quality surveillance.

Education

Multiple participants stressed the importance of education across the spectrum of stakeholders, including ethical big data research for students, scientists, developers, and community members. The general public can also benefit from a better understanding of the risks and benefits of sharing data.

While K-12 students may receive some informal teaching in this area, particularly reference to usual identifying information in commercial uses, big data and research is more likely at the undergraduate or graduate level when students become meaningfully exposed to data concepts. Ethical use of big data should be integrated into the overall data and analytic education.

Several organizations are working to create an ecosystem of responsible big data use. Industries can be encouraged to adopt and publicize their practices to establish transparency and foster trust. The California Consumer Privacy Act, modeled on the European Union's General Data Protection Regulation (GDPR), enables individuals to decide how their data can be used, including being removed completely from a system or collection. Corporations have also developed codes of data ethics with practical applications offered to employees who work with AI and big data (Accenture Labs, 2016). Some organizations have also created free ethics curricula packets for college students and expressed a commitment to diverse hiring practices as a part of being a responsible and ethical business.

Workforce training for data scientists can be in person or via webinars, and can

be offered on a routine basis. A training portfolio that spans multiple groups—K-12, university level, “boot camp” style courses, and professional—may also increase ethics awareness and ethical practices. In addition to students and workers, communities can also benefit from education outreach. A better understanding of the overall research process and how data is used can enhance community trust. Community education and outreach efforts that are available and convenient can both foster trust and encourage participation.

Developing an ethics curriculum for students who plan to work as data analysts can be a targeted goal. Such a core curriculum on ethics would increase a familiarity with ethical thinking and demonstrate different ways it can be applied to big data research.

Teaching ethical issues related to big data also highlights the question of what do we mean by “ethics”? Who is qualified to teach it? What priorities are emphasized? How can impact be measured? Ethical behavior is not necessarily arriving at a “right answer” but the ability to consider multiple avenues that arrive at multiple outcomes and the ethical implications of each. Partnerships across agencies and private groups can support implementation of ethical approaches for collecting or using big data.

Diversity

Improved diversity will enhance a data ethics strategy in order to foster implicit biases in both research and leadership teams. Bias is a well-known problem in AI illustrating the need for broad considerations when developing AI based apps. A debating bot fared better when its intake was curated instead of learned in real-time, because its developers were able to control the level of bias (Lee, 2018). It may not be feasible to curate data in every situation, but in public-facing applications, ethics may require it.

Ethical discussions need teams that include a diversity of backgrounds, experience, opinions, and expertise to best tackle the complex problems. Research is enhanced when scientists seek out different opinions in order to move research forward and find solutions. To the same end, initial data governance plans would be best determined by a diverse committee of experts and stakeholders who also define the role and responsibilities of a chief data strategist.

Community Involvement

Understanding the community whose members are taking part in a research project increases the community’s level of trust in a research endeavor and promotes ethical research behavior. An effective data ethics system takes a community’s culture and perspectives into consideration throughout the research process. In addition, participants should be able to see how their data is being used, what findings emerge from the research, and whether they will be impacted by the data use. A community advisory panel whose contributions are valued can improve the research process, flag potential abuses and approve secondary uses of data when appropriate.

As discussed earlier in the workshop, there are multiple ways to define a community. A community of rare disease sufferers may disregard privacy in the search for a cure, whereas chronic pain sufferers may prefer to keep their health data private to avoid denial of care. Poor treatment in the past is a strong barrier to gaining trust, especially if research leaders are perceived as outsiders. A number of health care panels are required to have at least one community member on an advisory panel to represent the patient population. However, researchers must be careful not to burden

Teaching ethical issues related to big data also highlights the question of what do we mean by “ethics”?

While existing IRB guidelines for big data use can be helpful, most IRBs are more experienced in HIPAA compliance and may not have the data, privacy, or cybersecurity expertise that ethical big data health research requires.

one individual with the role of representing an entire community. For example, “veterans” is a community, but within it, there are veterans of different ages, male and female veterans, and urban and rural veterans, who all have different perspectives.

Learning from Other Models

Ethical data practices can be borrowed from other organizations and countries that are confronting these same issues. For example, in some countries biological repositories (“biobanks”) must adhere to strict security rules, and individuals can report data concerns to a governmental ombudsman. The separation between government and industry, and related data sharing also varies by country. Industry and academia are much more integrated in some countries. Countries also vary in their concerns about large private or commercial groups collecting private information.

The publishing industry can also offer lessons, for example, the requirement that researchers state that they obtained informed consent before their research can be published. Data enclaves as in DoD and VA also offer protections for big data, limiting who can have access and how. The data protection companies may also provide valuable information and informative examples. Security precautions in this industry are significant and they tend to encrypt data, have notification when third parties access personal data, control over personal data access, and the ability for data owners to charge for data use. This process can be made transparent, private, and gives agency and financial incentives to the data owner at a time when those options are unavailable in nearly every other sphere.

Researchers can also learn from the National Science Foundation (NSF). NSF’s grant application review process often includes outside expertise, community input, and second opinions to ensure appropriate attention to ethical implications. Similarly, NSF’s Smart and Connected Communities program integrates community involvement into research programs in a way that acknowledges inequality, asks for community acceptance, and aims to improve a community’s overall quality of life.

On the other hand, there are so many commercial AI projects currently underway that it is difficult to learn from them. For example, predictive analytics systems are now used to approve customers for mortgages. This is rarely open for review or analysis to determine if a particular bias skewed the approval process in favor of one demographic or to the detriment of another.

Institutional Review Boards

Institutional review boards (IRBs) include safeguards to protect subjects. However, they also have multiple shortcomings that can leave data or subjects vulnerable. While existing IRB guidelines for big data use can be helpful, most IRBs are more experienced in HIPAA compliance and may not have the data, privacy, or cybersecurity expertise that ethical big data health research requires. In addition, IRBs often do not cover every aspect of data collection. For example, some organizations may want to own the intellectual property that is the research outcome and license it for research, whereas IRBs rarely handle intellectual property issues.

In new research, health data is often collected from wearable devices. When manufacturers are not associated with a research institution or a federal agency, they are not included in IRB oversight, and ethical practices can be neglected. The IRB review process can also be very slow with significant impact on the research. Another research challenge is the Paperwork Reduction Act, which can require substantial procedural requirements on data from the public.

Some organizations use other layers of oversight in addition to, or instead of, IRBs, such as information security officers to review research proposals more quickly. Other federal agencies in particular may require cybersecurity measures and be approved by the Chief Information Officer. Unfortunately, these extra layers can delay projects and frustrate researchers. In some countries IRBs are not always mandated. Some communities in the U.S. do not rely on IRBs to protect them but instead, set up separate, representative committees to review projects from the community's perspective. This is also a common practice in crowdsourced or citizen science projects.

The Relationship Between Humans and AI

AI is often viewed as in opposition to human control but in reality, collaboration between humans and AI is the key to success. There are things that machines can do better than humans, and there are things that humans can do better than machines. For example, medical practice is increasingly reliant on human-AI collaboration in diagnostics, where doctors and tailored AI programs both participate in a consultation. The doctor can draw on the AI's curated body of knowledge while making the final decision. This collaboration can also work well in imaging, for example, an algorithm can be used to screen images and prompt a human radiologist to take a closer look at certain cases. AI and human collaboration in imaging may work especially well because machines are much better at repetitive tasks—they don't get bored or tired—but humans have a better understanding of the nuances required to make correct diagnoses.

In big data, AI's fast computations can give researchers more time to interpret the results, another nuanced task where humans outperform machines. On the other hand, AI is not always the answer. For example, an algorithm to test for tuberculosis in x-rays could not surpass humans despite years of work. Collaboration with AI should be encouraged where it is efficient, but not overly relied upon where it does not add value.

AI can also make the general public uncomfortable. In aviation there is a push for more automation, when crashes are determined to be caused by "pilot error." However, research has shown that "pilot error" usually stems from poor interface design between humans and AI systems. In a recent crash, pilots tried to save the plane, but their actions were overridden by the AI system and all 189 people died. Better human-AI collaboration can improve safety without relinquishing too much human control.

In some cases, the potential for AI in real-world scenarios has been oversold. Not every field will benefit from AI or human-AI collaboration. In most cases, human creativity is needed to design an AI system, fine-tune it, and analyze the outcomes. In addition, it is humans who will know when to break the rules in order to achieve justice, and when we are merely automating inequality.

AI is often viewed as in opposition to human control but in reality, collaboration between humans and AI is the key to success.

CASE STUDIES

Workshop participants were asked to read through two case studies: that of the researcher at the center of the Cambridge Analytica scandal, and a fictional public sector-data analytics scandal. Each case provided opportunities to ask questions, determine ethical issues, and brainstorm solutions.

Case 1: Cambridge Analytica

In 2018 it was revealed that in 2016, a company leveraged psychological data from almost 50 million online user accounts to create voter profiles that it then sold to political campaigns. Part of the scandal was that while only some online users took a personality quiz, the app recorded data from quiz takers and their friends. The researcher behind the app saved that data in a private database (when he should have deleted it, according to terms of service), and then shared it with the company. Kogan admits that he failed to thoroughly consider the consequences and shares some of the responsibility for the data exposure.

Several questions came up in the discussion of this case. What was included in the documentation which must have been signed? It may have had information about what the company intended to do with the data. Whether it covered that or not, the researcher failed to imagine the harms enabled by his actions.

Paying attention early in the innovation process can help to catch and address warning signs, but that relies on making sure the right questions are being asked. In all research there are “known unknowns” that can be examined through thorough questioning. However, there are also “unknown unknowns.” Risk management frameworks can offer mechanisms to identify these and guide ethical decision-making.

Did the researcher consult with anyone before sharing his data? In his experience, while researchers *should* consult with lawyers, representatives from their institutions, and other experts from different domains, often they do not, leaving them unprotected and unaware of potential consequences. The researcher may not have been clear on the relevant intellectual property ownership details involved. Many institutions believe that they own all employees’ intellectual properties, despite this being a contentious legal area. In some cases, researchers have developed applications only to see them claimed by companies or institutions.

This legal gray area can make it difficult to make ethical decisions. An ethics and data science checklist could have helped the researcher make a different decision. It was noted that for a checklist to be useful, it must be balanced. Overly detailed checklists are overwhelming, but important aspects can fall through the cracks of checklists that are too vague. A checklist for researchers should include outside

Paying attention early in the innovation process can help to catch and address warning signs, but that relies on making sure the right questions are being asked.

Secrecy means intentionally preventing knowledge from becoming public.

advisement and explicit documentation for data use agreements to protect researchers, who often work in isolation and may not see the full implications of the data they have collected. The online company has changed its terms of service, but at the time apps were able to reach users' friends without asking for consent.

Case 2: New Leviathan and the Wales Consulting Group

The second case study describes a fictional town, New Leviathan, whose mayor partners with Wales Consulting Group (WCG) to apply AI to municipal data to reduce the crime rate by intervening when citizens may become victims of crime. While the mayor had unilateral power to make this decision, she did not notify anyone of the partnership. After conducting an internal ethics review, WCG's employees agreed to take the project on pro bono. Crime initially dropped, but it rose again and the mayor, facing public pressure and hoping to win re-election, asked WCG to go further and predict who is likely to commit a crime, a controversial practice known as "predictive policing." A short time later, the contract was exposed, many citizens were outraged, and it was revealed that WCG was selling similar technology to a dictator who was using it to silence critics.

This case study illustrates several ethical problems: the delicate balance of privacy and public safety, failure to plan adequately for problems, and how easily AI can be adapted to different desires. In discussing the case, participants focused on the ethics of privacy, secrecy, and public safety needs. Privacy is not a binary value—it exists on a continuum and balances expectations with actual information (Zook, Barocas, Boyd, et al., 2017). Privacy does not have a simple definition, because it is a bundle of multiple threads that communities and individuals may view differently. Women have different privacy needs than men and Native American communities have different privacy needs than others. But all people have a certain expectation of privacy, even when they are being actively policed. In New Leviathan, the mayor stated that individual privacy was respected, but the community's overall expectation of privacy was radically breached.

In addition, there is a discrepancy between the legal definition of privacy and what individuals or communities feel is private. Private information may not necessarily be important, but that does not mean people wish it to become public. Communities may share private information among themselves, and understand they have given up some loss of control over that information, but that is very different from it being given over to an AI tool and used for something to which they may be philosophically opposed. The data WCG used was collected for a different purpose, and the public did not consent to this use. Hence there was a feeling that their privacy was violated.

Another aspect of privacy is agency. Private data is under an individual's control. When it is no longer private, that person no longer has control over who can see the data or how it is used. Privacy and agency can be especially important in communities who fear being punished if their secrets became public or who are unfairly targeted as a high-crime population.

Secrecy is another issue at play in this case study. Secrecy means intentionally preventing knowledge from becoming public. Secrecy is very different from privacy. Data ethics has a strong emphasis on transparency. Participants agreed that the

mayor in the scenario made a mistake when she acted in secret. Of course, secrecy can also be helpful. The federal government categorizes the level of secrecy needed for different types of data based on impact levels. Levels 5 and 6 are the most secretive, but lower level data may still contain PII or PHI and require security parameters for authentication or authorization. With today's advanced AI tools making inferences more easily, even low-impact data can present a potential national security risk. Large, de-identified data sets that cover a large swath of the U.S. would be valuable to many countries.

This case study also illustrates the dangers of poorly planned projects that fail to consider negative consequences. The mayor did not consider a scenario where her project was revealed and negatively received. If she had first assessed the community's reaction and any privacy risks, she would have been better prepared for the backlash and might have acted differently. Whitehouse noted that police departments eagerly use video analytics to identify suspects, despite that fact that existing photo banks frequently misidentify faces and do not always accurately represent people of color, causing false positives that can disrupt an innocent person's life.

Several participants suggested that a representative advisory group could have collaborated with the mayor, the police, and WCG to find an alternate method to reduce crime. An advisory group also may have uncovered unrecognized biases in the mayor, the police, or the algorithm. Did the algorithm unwittingly target only poor or minority neighborhoods? Did the engineers who wrote the algorithm unknowingly incorporate biases? Humans are often not aware of the biases that influence their work.

This case study also illustrates how easily an algorithm can switch from benign to ominous. It can take only a few changes in coding and retraining to go from helping a potential victim, a public safety endeavor, to targeting a potential criminal, a violation of privacy and of the values of the community whose data was used. Although the mayor had the power to make the decision, ethical considerations about potential project drift or unethical applications of the AI may have steered her to act differently, preventing New Leviathan from being a testing ground before the AI tool is used to oppress free speech or identify political enemies. This secondary use of data is especially problematic because people's data can not only be used without their consent, but also potentially in ways that go against their values.

For anyone using data—companies, researchers, government—engaging communities can shine a light on unnoted ethical problems. For example, algorithms are not yet nuanced to understand that in communities, people can be both crime victims and crime perpetrators.

The mayor bears much of the responsibility for what happened in this scenario, largely because she did not seek community involvement or expert advice. While this situation is hypothetical, agencies may operate similarly. Police departments routinely “scrape” public data with no IRB requirements. Similarly, many research projects use only publicly available data and therefore do not feel the need to seek approvals or advisory groups, despite the potential for similar ethical problems.

Both case studies demonstrate that a failure to fully consider ethical implications and negative consequences can doom research with even the best of intentions. While the actors in these scenarios may not have violated laws, they did not act ethically. Such cases offer important learning opportunities. To act ethically, researchers must look at projects holistically and keep ethics in mind through every step of the

Both case studies demonstrate that a failure to fully consider ethical implications and negative consequences can doom research with even the best of intentions.

process. This can be done through planning, testing, implementation, execution, review, and revision.

CONCLUSION

Big data is both powerful and complex. Our understanding of how best to use big data/AI, how to interpret it, and how to keep it safe are new fields of work. While discussions of ethics in big data research are relatively new, AI applications for big data are accelerating rapidly. The application of AI to big data raises the prospect of unintended consequences. Therefore, ethical considerations must be part of big data research from formulating the question to how to answer the question and what to do with results.

The *Forum on Health and National Security: Ethical Use of Big Data for Healthy Communities and a Strong Nation* examined a wide range of challenges to ethical use of big data, including considerations of the best way to approach community members about providing their data, how to ensure that all data collection follows federal regulations, the need for extra caution when dealing with military data, and considerations for potential secondary uses of data. These challenges require addressing in many organizations and agencies. Data represents people whose privacy is protected, both for the benefit of the individual and for the benefit of our nation's security.

After discussing these multiple challenges, participants crafted a pathway to embed ethics throughout the big data research process. Ethical big data practice means treating the owners of data with dignity, earning their trust, respecting community values and fears, being transparent about data uses, and acquiring consent carefully. It also means ensuring that research teams are diverse, consider outside expertise, and have specific training in ethics and big data research.

Participants agreed that ethical research can only happen with planning. In order to accomplish this task, institutions will need to prioritize the creation of risk management frameworks, data governance plans, and advisory groups or IRBs with appropriate expertise, diversity, and ethical training. There are many models and case studies that can guide the development of these plans. Researchers should also closely examine the successes and failures of data governance, risk management, and privacy violations to ensure ethical actions.

Researchers have long worked by the maxim of “do no harm,” but the harms that could come from AI-big data pairings are new and not yet clearly delineated. Unintended risks in this space can have significant consequences. Ethical due diligence is a component of good research. In addition, despite the best intentions and even in the context of strong cybersecurity protections, data is vulnerable to accidental misuse, intentional misuse, unauthorized secondary uses, or application pivots that endanger privacy, civil liberties, or national security.

The research agenda of the nation is best served by building ethics into the entire research ecosystem. There are substantial challenges to fully realizing this goal. Technology applications and development are fueled by large funding streams that advance extremely rapidly, and it is a challenge to keep up with the ethical implica-

Big data is both powerful and complex. Our understanding of how best to use big data/AI, how to interpret it, and how to keep it safe are new fields of work.

Ethical big data practice means treating the owners of data with dignity, earning their trust, respecting community values and fears, being transparent about data uses, and acquiring consent carefully.

tions of new capabilities and vulnerabilities as they emerge. It takes commitments of time and funding to address the ethical complexities, train others to understand them, and create appropriate ethical frameworks before research begins.

The intersection of AI systems, health records, and big data is newly charted territory. Ethical science requires that researchers prioritize using big data responsibly, respecting individual privacy, and protecting the nation.

APPENDIX:

REFERENCES, FORUM READINGS, VIDEO AND CASE STUDIES

References

- Accenture Labs. (2016). *Building digital trust: The role of data ethics in the digital age*. Retrieved from https://www.accenture.com/_acnmedia/PDF-22/Accenture-Data-Ethics-POV-WEB.pdf
- Buolamwini, J. (2016, November). How I'm fighting bias in algorithms [Video file]. Retrieved from https://www.ted.com/talks/joy_buolamwini_how_i_m_fighting_bias_in_algorithms
- Hayden, E. C. (2018, May 18). Experimental drugs poised for use in Ebola outbreak. *Nature*. Retrieved from <https://www.nature.com>
- Lee, D. (2018, June 19). IBM's machine argues, pretty convincingly, with humans. *BBC News*. Retrieved from <https://www.bbc.com>
- Loukides, M., Mason, H., & Patil, D. J. (2018). *Ethics and Data Science*. Sebastopol, CA: O'Reilly Media.
- Wang, Y., & Kosinski, M. (2018). Deep neural networks are more accurate than humans at detecting sexual orientation from facial images. *Journal of Personality and Social Psychology*, 114(2), 246-257. doi:10.31234/osf.io/hv28a.
- Zook, M., Barocas, S., Boyd, D., Crawford, K., Keller, E., Gangadharan, S. P., Goodman, A., Hollander, R., Koenig, B. A., Metcalf, J., Narayanan, A., Nelson, A., & Pasquale, F. (2017). Ten simple rules for responsible big data research. *PLoS Computational Biology*, 13(3): e1005399. doi:10.1371/journal.pcbi.1005399.

Forum Readings

- Buchanan, E. (2017). Considering the ethics of big data research: A case of Twitter and ISIS/ISIL. *PloS one*, 12(12), e0187155. doi:10.1371/journal.pone.0187155
- Eubanks, V. (2018). *Automating inequality: How high-tech tools profile, police, and punish the poor*. New York, NY: St. Martin's Press.
- Loukides, M., Mason, H., & Patil, D. J. (2018). *Ethics and Data Science*. Sebastopol, CA: O'Reilly Media.
- Zook, M., Barocas, S., Boyd, D., Crawford, K., Keller, E., Gangadharan, S. P., Goodman, A., Hollander, R., Koenig, B. A., Metcalf, J., Narayanan, A., Nelson, A., & Pasquale, F. (2017). Ten simple rules for responsible big data research. *PLoS Computational Biology*, 13(3): e1005399. doi:10.1371/journal.pcbi.1005399

Video and Case Studies

Video

Bias in machine learning training sets: Joy Buolamwini TEDx talk.
https://www.ted.com/talks/joy_buolamwini_how_i_m_fighting_bias_in_algorithms?language=en#t-512023

Case Study 1

Cambridge Analytica: What Role for the University as Researcher at Centre of Scandal Admits ‘I should have questioned the ethics of the exercise’?
<http://magnacartafordata.org/cambridge-analytica-what-role-for-the-university-as-the-researcher-at-centre-of-the-scandal-admits-i-should-have-questioned-the-ethics-of-the-exercise/>

Case Study 2

Public Sector Data Analytics.
<https://aiethics.princeton.edu/wp-content/uploads/sites/587/2018/10/Princeton-AI-Ethics-Case-Study-6.pdf>

Center for the Study of Traumatic Stress
Department of Psychiatry
Uniformed Services University of the Health Sciences
4301 Jones Bridge Road
Bethesda, MD 20814
www.CSTSonline.org

