



**Selected Papers of #AoIR2020:
The 21st Annual Conference of the
Association of Internet Researchers**
Virtual Event / 27-31 October 2020

ETHICAL REVIEW BOARDS AND PERVASIVE DATA RESEARCH: GAPS AND OPPORTUNITIES

Michael Zimmer, PhD
Marquette University

Edward Chapman
Marquette University

Introduction

The growing prevalence of data-rich networked information technologies—such as social media platforms, smartphones, wearable devices, and the internet of things—brings an increase in the flow of rich, deep, and often identifiable personal information available for researchers. As the Computational Social Science group at Microsoft Research notes:

With an increasing amount of data on every aspect of our daily activities – from what we buy, to where we travel, to who we know, and beyond – we are able to measure human behavior with precision largely thought impossible just a decade ago, creating an unprecedented opportunity to address longstanding questions in the social sciences. (“Computational Social Science,” n.d.)

More than just “big data,” the datasets envisioned above are unique in that they represent people’s lives and activities, bridge multiple dimensions of a person’s life, and are often collected, aggregated, exchanged, and mined without them knowing. We call this data “pervasive data,” and the increased scale, scope, speed, and depth of pervasive data available to researchers require that we confront the ethical frameworks that guide such research activities.

Multiple stakeholders are embroiled in the challenges of research ethics in pervasive data research. For example: researchers struggle with questions of privacy and consent

Suggested Citation (APA): Zimmer, M and Chapman, E. (2020, October 28-31). *Ethical Review Boards and Pervasive Data Research: Gaps and Opportunities*. Paper presented at AoIR 2020: The 21th Annual Conference of the Association of Internet Researchers. Virtual Event: AoIR. Retrieved from <http://spir.aoir.org>.

(Shilton, 2015; Zimmer, 2016); user communities may not even be aware of the widespread harvesting of their data for scientific study (Fiesler & Proferes, 2018); platforms are increasingly restricting researcher's access to data over fears of privacy and security (Bruns, 2018); and ethical review boards face increasing difficulties in properly considering the complexities of research protocols relying on user data collected online (Buchanan & Ess, 2009; Vitak et al., 2017).

The results presented in this paper expand our understanding of how ethical review board members think about pervasive data research.¹ It provides insights into how IRB professionals make decisions about the use of pervasive data in cases not obviously covered by traditional research ethics guidelines, and points to challenges for IRBs when reviewing research protocols relying on pervasive data.

Methodology

A survey instrument was created to assess IRB members' training and attitudes around research protocols that rely on pervasive data. The survey asked respondents for anonymous reporting on whether their IRB regularly reviews research that utilizes pervasive data, and the kind of training or resources they rely on to review such submissions, and how confident they are in their assessment of research protocols using pervasive data.

Respondents were also presented with eleven hypothetical research scenarios that rely on pervasive data [See Appendix A]. Scenarios varied in the source of the data, the intended inferences to gain from the data, whether the data was publicly-available, whether specific informed consent was sought, the level of identifiability of the data, whether steps were taken to anonymize the data, and whether a platform's terms of service might have been violated to obtain the data. Table 1 in Appendix B provides a summary of these variables. Respondents were asked to predict how their IRB would view the hypothetical protocol, and to identify the key factors that contributed to their response.

The survey was available online from November 2018 through July 2019, and was restricted to ethical review board members in the U.S.

Data cleaning—ensuring respondents answered questions for at least one of the eleven scenarios—yielded 77 usable responses, of whom 64 (83%) were located in an IRB based at a college or university, with the majority (34, 53%) of these from R1 institutions while 15 (23%) identifying as liberal arts-focused institutions. Not all questions were required to be answered, and thus some responses total less than 77.

Summary of Findings

General Experience with and Preparedness for Protocols Using Pervasive Data

¹ Within the United States, university ethical review boards are typically named "Institutional Review Boards", and thus the acronym IRB will be used throughout the remainder of the paper.

We asked a series of questions to gauge an IRBs exposure to, and preparation for, research protocols that rely on pervasive data. Half of respondents (30/59) reported receiving 10 or fewer proposals that used pervasive data annually, while over one-third (20/59) reviewed more than 50 each year. The most common types of pervasive data appearing in research protocols reviewed within the past year are provided in Table 2, and the disciplines represented in those protocols are provided in Table 3 (see appendix B).

In terms of preparedness for reviewing protocols using pervasive data, less than one-third of respondents indicated their institution provided specific training sessions for IRB members that addressed the collection and use this kind of data (see Table 4, Appendix B). The vast majority of respondents (54/59) also indicated their IRB lacked any specific checklist, review tool, policy or set of guidelines for reviewing protocols that rely on pervasive data, while eleven indicated such guidelines were under development. Only four respondents indicated such materials existed, and when asked if the available materials were “excellent,” “adequate,” or “poor,” each of the four respondents indicated “adequate.” In the absence of specific internal guidelines, respondents relied on various external regulations or guidelines when reviewing protocols utilizing pervasive data (see Table 5, Appendix B).

Respondents were also asked to agree or disagree with statements about how well-versed their institutions’ IRB members were regarding both the *technical* and *ethical* aspects of the collection and use of pervasive data. Only 25% of respondents agreed that their IRB members are well-versed in the technical aspects of these type of research protocols, while nearly one-half felt they were well-versed in the ethical dimensions (see Tables 6 and 7 in Appendix B).

Responses to Hypothetical Scenarios Using Pervasive Data

Respondents were presented with eleven hypothetical research scenarios (see Appendix A) relying on pervasive data and were asked to consider how their IRB would likely review each case. A summary of results is provided in Table 8 in Appendix B.

Along with the assessment with each scenario, respondents were asked what the key factors would be in making their determination. A treemap summary of responses for each scenario are provided in Appendix C, highlighting the most common factors indicated in support of each respondent’s determination (the size of each box indicates its relative importance compared to other factors in that determination).

Discussion

Our findings suggest that IRBs are largely unprepared for addressing the unique challenges that stem from research protocols relying on pervasive data. Over two-thirds of the respondents indicated that no specific training was provided for IRB members to address the ethics of pervasive data research, and over 90 percent noted the absence any specific set of guidelines for reviewing protocols that rely on pervasive data. Further, while nearly half of our respondents felt they were well-versed in the ethical issues related to pervasive data, only 25% felt their IRB was had sufficient technical

knowledge to understand such protocols. This might present cases of misunderstanding or over-confidence in the ability to adequately assess pervasive data research protocols.

Our initial analysis of the results from the various hypothetical research scenarios supports this concern. While many of the eleven scenarios presented yielded largely consistent assessments (for example, 75% of respondents indicated scenario 10 would be “Expedited”), numerous scenarios revealed large diversity of viewpoints on how an IRB would review the protocol. In Scenario 1, for example, respondents were spread across the four possible IRB review categories, and the treemap summary reveals how the same factor (“Terms of service” in this case) might be listed as a key reason for divergent determinations. Other scenarios share this feature, suggesting confusion exists across IRB members regarding how to address pervasive data research.

Further analyses and implications for the research and IRB community will be presented.

Acknowledgements

This material is based upon work supported by the National Science Foundation under Grant No. 1947754.

References

- Bruns, A. (2018). Facebook shuts the gate after the horse has bolted, and hurts real research in the process. *Internet Policy Review*.
- Buchanan, E., & Ess, C. (2009). Internet research ethics and the institutional review board: Current practices and issues. *ACM SIGCAS Computers and Society*, 39(3), 43–49.
- Computational Social Science. (n.d.). *Microsoft Research*. Retrieved February 4, 2020, from <https://www.microsoft.com/en-us/research/group/computational-social-science/>
- Fiesler, C., & Proferes, N. (2018). “Participant” Perceptions of Twitter Research Ethics. *Social Media + Society*, 4(1), 2056305118763366. <https://doi.org/10.1177/2056305118763366>
- Shilton, K. (2015). *Emerging Ethics Norms in Social Media Research*. Workshop on Beyond IRBs: Ethical Review Processes for Big Data Research. <https://bigdata.fpf.org/papers/emerging-ethics-norms-in-social-media-research/>
- Vitak, J., Proferes, N., Shilton, K., & Ashktorab, Z. (2017). Ethics Regulation in Social Computing Research: Examining the Role of Institutional Review Boards. *Journal of Empirical Research on Human Research Ethics*, 12(5), 372–382. <https://doi.org/10.1177/1556264617725200>
- Zimmer, M. (2016, May 14). *OkCupid Study Reveals the Perils of Big-Data Science*. *Wired*. <https://www.wired.com/2016/05/okcupid-study-reveals-perils-big-data-science/>

Appendix A: Scenarios

Eleven hypothetical scenarios were presented to respondents. For each scenario, respondents were asked the following questions:

Your institution's IRB would likely consider this proposal to be:

- Not human subjects research
- Exempt
- Expedited
- Requiring full board review

In this case, what would be the key factor(s) in making that determination? (check all that apply)

- Public vs. private site
 - Public vs. private data
 - Level of analysis (group vs. individual)
 - Whether data is identifiable
 - Whether data gathering violates terms of service
 - Whether data is reused
 - Whether the project combines datasets
 - Method of obtaining data
 - Type of data
 - Whether informed consent was obtained
 - Purpose of the research
 - Impact beyond the participants
 - Other
1. Researchers plan to scrape public comments from online newspaper pages to predict election outcomes. They will aggregate their analysis to determine public sentiment. The researchers don't plan to inform commenters, and they plan to collect potentially-identifiable user names. Scraping comments violates the newspaper's terms of service.
 2. Researchers plan to scrape public Twitter feeds to predict risky drug-use behaviors. They will analyze individual behaviors. The researchers don't plan to inform Twitter users, but they will not collect any identifying information. Scraping Tweets does not violate Twitter's terms of service.
 3. Researchers plan to analyze private interaction data from a dating site to understand the sexual behavior of groups. The researchers plan to collect informed consent from dating site users, and they plan to collect identifiable information from participants. Asking users for permission to use their data does not violate the dating site's terms of service.
 4. Researchers plan to collect newspaper comments by reading articles and cutting and pasting all associated comments into spreadsheets. They will use qualitative analysis to understand individual political views. The researchers don't plan to inform

commenters, and they plan to collect potentially-identifiable user names. Cutting and pasting comments does not violate the newspaper's terms of service.

5. Researchers plan to work with a mobile phone company to collect geolocation data to understand group mobility patterns in a city. The researchers will not inform the mobile phone users, and they will not collect any additional identifying information. Partnering with the mobile phone company to collect data does not violate the company's terms of service.
6. Researchers plan to combine mental health records provided by a university and public social media activity to predict mental health conditions among students. The researchers plan to collect informed consent, and they plan to collect identifiable information from participants.
7. Researchers plan to use a database of public tweets curated and shared by another researcher to study a political event. Researchers do not plan to inform the original posters, and researchers have taken measures to de-identify the data.
8. Researchers plan to scrape data from an open health forum and combine it with scraped tweets to predict mental health conditions. The researchers will not inform forum users, and they may collect potentially identifying information. Scraping data violates neither the health forum nor Twitter's terms of service.
9. Researchers plan to scrape profile photos, which are visible to any member of the service, from a dating site to build models that predict sexual preference or behavior. Researchers will not inform the dating site users, but they will not collect any identifying information and their photograph dataset will not be released publicly. Creating a fake profile, necessary to access the photos, violates the dating site's terms of service.
10. Researchers plan to ask Apple HealthKit users to voluntarily submit their activity data to understand the general impact of exercise on a health condition. The researchers plan to obtain informed consent, and they plan to collect identifiable information from participants. Asking users to submit activity data does not violate Apple Health Kit's terms of service.
11. Researchers plan to scrape public posts and interactions from Facebook to study group-level dynamics. They plan to collect informed consent from the original poster, but not those they interacted with, and they may collect identifying information. Scraping posts with permission of the original poster does not violate Facebook's terms of service.

Appendix B: Tables

Table 1: Overview of Variables in Hypothetical Research Scenarios

Scenario	Consent	Publicness	Anonymity	Terms of Service
1. Scrapping public newspaper comments to predict elections	Low	High	Low	Low
2. Scrapping public Twitter feeds to predict risky drug-use behaviors	Low	High	High	High
3. Analyzing dating site data to infer sexual behavior	High	Low	Low	High
4. Analyzing newspaper comments to understand political views	Low	High	Low	High
5. Collect geolocation data from mobile provider to understand group mobility patterns in a city	Low	Low	High	High
6. Combine mental health data with social media activity	High	Neutral	Low	Neutral
7. Analyzing preexisting Twitter dataset to study political event	Low	High	High	Neutral
8. Scraping health forum and combining with Twitter data to predict mental health	Low	High	Low	High
9. Scraping profile photos to predict sexual behavior	Low	Neutral	High	Low
10. Analyze Apple HealthKit data to assess impact of exercise on health	High	Neutral	Low	High
11. Scrape public Facebook posts to study group-level dynamics	Neutral	High	Low	High

(High = high compliance; Neutral = neutral or unknown level; Low = low compliance)

Table 2: Types of pervasive data in protocols reviewed in past 12 months

Social media posts	58
Sensor data	54
Social media profiles	33
Locational data	33
Social media images	30
Network traffic data	26
Other	7

(n=73; multiple selections allowed)

Table 3: Disciplines submitting protocols using pervasive data in past 12 months

Social Sciences	47
Medical/Health	36
Computer Science/Engineering	24
Arts/Humanities	12
Natural Sciences	3
Other	7

(n=74; multiple selections allowed)

Table 4: Does your institution provide specific training sessions for IRB members that addresses the ethics of the collection/use of pervasive data?

Yes, it is part of required training	14	18.4%
Yes, but it is optional	7	9.2%
No	53	69.7%
I don't know	2	2.6%

(n=76)

Table 5: Regulations or guidelines are relied on when reviewing protocols relying on pervasive data

Federal Policy for the Protection of Human Subjects ("Common Rule")	71
The Belmont Report: Ethical Principles and Guidelines for the Protection of Human Subjects of Research	61
SACHRP Considerations and Recommendations Concerning Internet Research and Human Subjects Research Regulations	33
American Psychological Association (APA) Psychological Research Online: Opportunities and Challenges	9
Association of Internet Researchers (AoIR) Ethics Guidelines	7
The Menlo Report: Ethical Principles Guiding Information and Communication Technology Research	3
ACM SIGCHI Research Ethics Committee Guidelines	1
Other	7

(n=77; multiple selections allowed)

Table 6: IRB members at my institution are well-versed in the *technical* aspects of the collection/use of pervasive data

Strongly agree	5	6.6%
Somewhat agree	14	18.4%
Neither agree nor disagree	17	22.4%
Somewhat disagree	28	36.8%
Strongly disagree	12	15.8%

(n=76)

Table 7: IRB members at my institution are well-versed in the *ethical* aspects of the collection/use of pervasive data

Strongly agree	13	17.1%
Somewhat agree	23	30.3%
Neither agree nor disagree	10	13.2%
Somewhat disagree	20	26.3%
Strongly disagree	10	13.2%

(n=76)

Table 8: Your institution's ERB would likely consider this proposal to be...

<i>Scenario</i>	<i>Not Human Subjects Research</i>	<i>Exempt</i>	<i>Expedited</i>	<i>Requiring Full Board Review</i>
1. Scrapping public newspaper comments to predict elections (<i>n=40</i>)	11 27.5%	8 20.0%	8 20.0%	13 32.5%
2. Scrapping public Twitter feeds to predict risky drug-use behaviors (<i>n=45</i>)	18 40.0%	15 33.3%	9 20.0%	3 6.7%
3. Analyzing dating site data to infer sexual behavior (<i>n=43</i>)	1 2.3	7 16.3%	21 48.8%	14 32.6%
4. Analyzing newspaper comments to understand political views (<i>n=44</i>)	25 56.8%	10 22.7%	7 15.9%	2 4.5%
5. Collect geolocation data from mobile provider to understand group mobility patterns in a city (<i>n=45</i>)	18 40.0%	13 28.9%	8 17.8%	6 13.3%
6. Combine mental health data with social media activity (<i>n=44</i>)	0 0.0%	0 0.0%	19 43.2%	25 56.8%
7. Analyzing preexisting Twitter dataset to study political event (<i>n=46</i>)	26 56.5%	13 28.3%	6 13.0%	1 2.2%
8. Scraping health forum and combining with Twitter data to predict mental health (<i>n=41</i>)	11 26.8%	8 19.5%	7 17.1%	15 36.6%
9. Scraping profile photos to predict sexual behavior (<i>n=46</i>)	3 6.5%	2 4.3%	10 21.7%	31 67.4%
10. Analyze Apple HealthKit data to assess impact of exercise on health (<i>n=44</i>)	0 0.0%	7 15.9%	33 75.0%	4 9.1%
11. Scrape public Facebook posts to study group-level dynamics (<i>n=44</i>)	5 11.4%	5 11.4%	17 38.6%	17 38.6%

Appendix C: Scenario Treeplots

1. Scrapping public newspaper comments to predict elections

Full review				Expedited			
Terms of service		Impact beyond participants		Identifiable data		Terms of service	
Identifiable data		Obtains informed consent		Method of obtaining data		Public vs. private data	
Not human subjects research				Exempt			
Public vs. private data		Public vs. private site	Level of analysis	Other	Public vs. private data		Terms of service
		Identifiable data	Purpose of the research		Public vs. private site		

2. Scrapping public Twitter feeds to predict risky drug-use behaviors

Not human subjects research			Expedited			
Public vs. private data	Public vs. private site	Identifiable data	Identifiable data		Public vs. private data	
			Level of analysis		Type of data	
Exempt			Full review			
Identifiable data	Public vs. private data	Public vs. private site	Identifiable data	Level of analysis	Method of obtaining data	
			Impact beyond participants	Obtains informed consent	Purpose of the research	Type of data

3. Analyzing dating site data to infer sexual behavior

Expedited		Full review	
Identifiable data	Obtains informed consent	Identifiable data	Public vs. private data
Type of data	Public vs. private data	Type of data	Obtains informed consent
Method of obtaining data		Exempt	Obtains informed consent
		Level of analysis	Public vs. private data
		Identifiable data	Method of obtaining data
		Purpose of the research	Type of data

4. Analyzing newspaper comments to understand political views

Not human subjects research		Exempt	
Public vs. private data	Type of data	Public vs. private data	Terms of service
Public vs. private site		Purpose of the research	Identifiable data
Method of obtaining data		Public vs. private site	
		Expedited	
		Identifiable data	Purpose of the research
			Level of analysis
		Public vs. private site	Terms of service
			Type of data

5. Collect geolocation data from mobile provider to understand group mobility patterns

Exempt		Expedited		Not human subjects research	
Identifiable data		Identifiable data	Type of data		Identifiable data
Type of data	Level of analysis	Combines datasets	Level of analysis	Method of obtaining data	
Public vs. private data	Purpose of the research	Impact beyond participants	Public vs. private data	Purpose of the research	Level of analysis
		Full review		Public vs. private data	
		Identifiable data	Type of data		Public vs. private data

6. Combine mental health data with social media activity

Full review		Expedited	
Type of data	Identifiable data	Identifiable data	Type of data
Obtains informed consent	Purpose of the research	Obtains informed consent	Public vs. private data

9. Scraping profile photos to predict sexual behavior

Full review			Expedited		
Terms of service	Method of obtaining data	Obtains informed consent	Public vs. private data	Identifiable data	
			Public vs. private site	Type of data	
			Terms of service		
Identifiable data	Public vs. private site	Type of data	Exempt Impact beyond participants	Not human subjects research	
			Level of analysis	Public vs. private data	Other
			Obtains informed consent	Purpose of the research	

10. Analyze Apple HealthKit data to assess impact of exercise on health

Expedited		Full review	
Identifiable data	Method of obtaining data	Identifiable data	Type of data
		Level of analysis	Method of obtaining data
		Obtains informed consent	Purpose of the research
Type of data	Public vs. private data	Exempt	Obtains informed consent
		Other	

11. Scrape public Facebook posts to study group-level dynamics

Full review		Expedited	
Identifiable data	Obtains informed consent	Public vs. private data	Public vs. private site
Impact beyond participants	Method of obtaining data	Identifiable data	Impact beyond participants
Public vs. private data	Purpose of the research	Exempt	Not human subjects research
		Public vs. private site	Obtains informed consent
		Identifiable data	Public vs. private data
			Public vs. private site
			Level of analysis