



Privacy Concerns with Using Public Data for Suicide Risk Prediction Algorithms: A Public Opinion Survey of Contextual Appropriateness

Journal:	<i>Journal of Information, Communication & Ethics in Society</i>
Manuscript ID	JICES-08-2021-0086.R1
Manuscript Type:	Journal Paper
Keywords:	Privacy, contextual integrity, data ethics, algorithmic ethics

SCHOLARONE™
Manuscripts

Privacy Concerns with Using Public Data for Suicide Risk Prediction Algorithms: A Public Opinion Survey of Contextual Appropriateness

Abstract

Purpose: Existing algorithms for predicting suicide risk rely solely on data from electronic health records, but such models could be improved through the incorporation of publicly available socioeconomic data – such as financial, legal, life event, and sociodemographic data. This study’s purpose is to understand the complex ethical and privacy implications of incorporating sociodemographic data within the health context. We present results from a survey exploring what the general public’s knowledge and concerns are about such publicly available data and the appropriateness of using it in suicide risk prediction algorithms.

Design/methodology/approach: A survey was developed to measure public opinion about privacy concerns with using socioeconomic data across different contexts. We presented respondents with multiple vignettes that described scenarios situated in medical, private business, and social media contexts, and asked participants to rate their level of concern over the context and what factor contributed most to their level of concern. Specific to suicide prediction, we presented respondents with various data attributes that could potentially be used in the context of a suicide risk algorithm and asked participants to rate how concerned they would be if each attribute was used for this purpose.

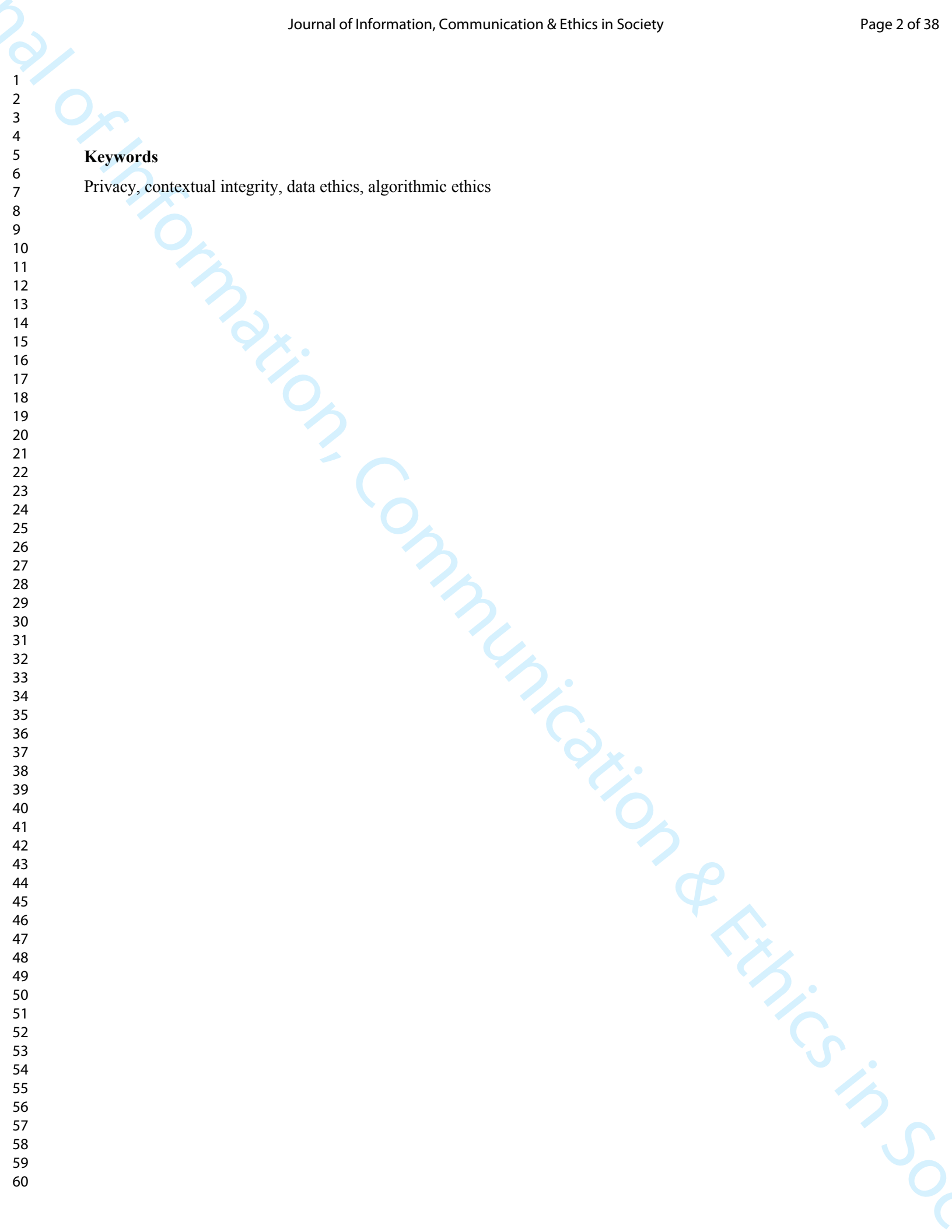
Findings: We found considerable concern across the various contexts represented in our vignettes, with greatest concern in vignettes that focused on the use of personal information within the medical context. Specific to the question of incorporating socioeconomic data within suicide risk prediction models, our results show a clear concern from all participants in data attributes related to income, crime and court records, and assets. Data about one’s household were also particularly concerns for our respondents, suggesting that even if one might be comfortable with their own being used for risk modeling, data about other household members is more problematic.

Originality: Previous studies on the privacy concerns that arise when integrating data pertaining to various contexts of people’s lives into algorithmic and related computational models have approached these questions from individual contexts. This study differs in that we captured the variation in privacy concerns across multiple contexts. We also specifically assessed the ethical concerns related to a suicide prediction model and determining people’s awareness of the publicness of select data attributes, as well as which of these data attributes generated the most concern in such a context. As far as we know, this is the first study to pursue this question.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Keywords

Privacy, contextual integrity, data ethics, algorithmic ethics



Introduction

Suicide is currently the 10th leading cause of death in the United States, and suicide rates have continued to increase in recent years (“Suicide statistics”, 2019). About half of suicide decedents have contact with the healthcare system in the month before their death (Ribeiro *et al.*, 2017), indicating that there is a significant opportunity to identify patients who are at risk for a suicide attempt when they visit their healthcare provider. Most algorithms for predicting suicide risk relies solely on data from electronic health records (Barak-Corren *et al.*, 2017), but it has been proposed that algorithmic solutions for predicting suicide risk could be improved through the incorporation of publicly available socioeconomic data – such as financial, legal, life event, derogatory, and sociodemographic data – alongside electronic health record data (Barak-Corren *et al.*, 2017; O’Connor and Portzky, 2018; Turecki and Brent, 2016).

If incorporation of publicly available socioeconomic data into the suicide prediction model – and other predictive models – is to be pursued, we must understand the complex ethical and privacy implications so as to not erode the trust between a patient and the healthcare system. To start this effort, we present results from a survey exploring what the general public’s knowledge and concerns are about such publicly available data.

Background

In an age of rapidly increasing technological advancements, scholars from all fields of expertise are investigating how big data can aid them in their work. From bioinformaticians who analyze genetic code to political scientists who analyze government records to sociologists who analyze social media data, everyone wants in on the benefits of big data and the algorithmic systems it empowers (boyd and Crawford, 2012). While the potential uses for big data and algorithm-based systems are nearly endless, it is important and necessary to consider the ethical and privacy concerns related to these endless possibilities (Metcalf *et al.*, 2016).

One critical question focuses on the ethical and privacy complications of using publicly-available data in research. Significant gaps among researchers on what constitutes “public” data that does not require explicit consent prior to harvesting (Zimmer, 2010, 2016), whether a platform’s terms of service might allow the automated scraping of public data (Fiesler *et al.*, 2020), and even at what stage does computational research become human subjects research requiring particular ethical protection (Metcalf and Crawford, 2016). Further, users are often not aware of the types of access researchers have to public data via social media platforms and their APIs (Fiesler and Proferes, 2018). Given the above, uncertainty persists among the research community on how to address the ethical and privacy complications of using ostensibly public data in computational research projects (Shilton, 2015; Vitak *et al.*, 2016).

These challenges become even more complicated within the context of medical data, where users’ knowledge and attitudes about the public availability of their health data is particularly muddled. While some

1
2
3 studies have shown users expressing little concern over privacy issues related to their personal fitness
4 information, noting such data was not inherently sensitive and expressing ambivalence over the possibility of
5 sharing that data with third parties (Zimmer *et al.*, 2020), other studies have shown users of self-tracking
6 technologies are frequently unaware of the details of external data access to which they agree in the context of
7 clicking “accept” to terms of use” (Bietz *et al.*, 2016). Further, a recent study assessing privacy concerns related to
8 the use of publicly-available health-related tweets data in research found the acceptability of harvesting such data
9 depended greatly on the nature of the health ailment, who was collecting it, and the context of use (Reuter *et al.*,
10 2019).

11
12 The potential of incorporating publicly-available socioeconomic data – typically obtained from
13 commercial data brokers – into medical research complicates things even further. The role of socioeconomic
14 factors in various aspects of healthcare we well studied (Fiscella *et al.*, 2000; Marmot *et al.*, 2008), and their
15 potential for assisting in the particular challenge of suicide prevention are promising (Barak-Corren *et al.*, 2017;
16 O’Connor and Portzky, 2018; Turecki and Brent, 2016). But integrating such a wide range of data points – such
17 as financial, legal, life event, derogatory, and sociodemographic data – into different contexts is often met with
18 resistance, irrespective of how public the information might be (Crain, 2018; Hoofnagle, 2004; Martin and
19 Nissenbaum, 2017; Tene and Polonetsky, 2013).

20
21 To help assess the ethics of incorporating public socioeconomic data in the development of algorithms to
22 predict suicide risk, we invoke Nissenbaum’s (Nissenbaum, 2010) theory of privacy as contextual integrity (CI).
23 CI rejects the traditional dichotomy of public versus private information, as well as the notion that privacy
24 preferences and decisions in one context universally apply to other contexts. Instead, CI rests on the
25 understanding that our interactions with other people, institutions, and technologies, occur within particular
26 contexts. Norms of appropriateness govern people’s expectations of how personal information should flow within
27 any given context. Therefore, responding to a data ethics question—e.g., should third-party socioeconomic data
28 be integrated with health data?—needs to start not with privacy as a static set of principles but with an
29 understanding of norms of appropriateness within the context in which the data is being collected and used, and
30 whether it is deemed appropriate to move information from one context – such as socioeconomic data – and apply
31 it within a new context – such as medical research in suicide risk prediction.

47 **Study Objective**

48
49 We rely on contextual integrity to investigate people’s privacy concerns about integrating publicly
50 available socioeconomic data within various contexts, including a suicide risk prediction algorithm. We aim to
51 measure public opinion about privacy concerns across medical-, business-, and social media-related contexts, as
52 well as specifically investigate the privacy concerns related to a suicide prediction model. To do this, we first
53
54
55
56
57
58
59
60

1
2
3 assess participants' general privacy concerns and their privacy knowledge. Then we evaluate participants'
4 opinions and concerns towards various contexts, which we will then be able to compare to concerns towards the
5 use of socioeconomic data in the suicide prediction algorithm. Overall, we aim to determine whether it is ethical
6 to combine publicly available socioeconomic data with health data to improve a suicide prediction algorithm
7 through an assessment of individuals' awareness of this very possibility as well as their comfort with the
8 appropriateness of this use of said data within this particular context.
9

10
11 We approach this through the following research questions:
12

13 RQ1: To what extent do people understand that their socioeconomic data – such as financial, legal, life
14 event, derogatory, and sociodemographic data – is publicly available?
15

16 RQ2: In which contexts do individuals find the use of publicly available information most concerning?
17

18 RQ3: Regarding suicide risk prediction algorithms specifically, which socioeconomic data points are
19 most problematic for inclusion?
20
21
22
23

24 **Methods**

25 *Survey Instrument*

26
27 A survey was developed to measure public opinion about privacy concerns with using various
28 socioeconomic data points across various contexts. The survey consisted of five sections that allowed us to collect
29 information about demographics of respondents, general privacy concerns, awareness of what types of
30 information are publicly available, privacy concerns associated with specific contexts, and concerns over the use
31 of thirty publicly available data attributes in the context of a suicide risk prediction algorithm. To inquire about
32 personal data use across a broad range of contexts, we presented respondents with ten vignettes that described
33 scenarios situated in medical, private business, and social media contexts (see Appendix A). Following each
34 vignette, we asked participants to rate their level of concern over the context, and we asked what factor
35 contributed most to their level of concern. Specific to suicide prediction, we presented respondents with various
36 data attributes¹ that could potentially be used in the context of a suicide risk algorithm and asked participants to
37 rate how concerned they would be if each attribute was used for this purpose (see Appendix B). The survey was
38 tested for clarity and consistency, and the research protocol received Institutional Review Board approval.
39
40
41
42
43
44
45
46
47

48 *Data Collection*

49 The survey was deployed on Qualtrics from July 17-24, 2020. We contracted Qualtrics to recruit
50 approximately respondents between the ages of 26-99 who live in the United States. Respondents with a response
51
52
53

54 ¹ Data attributes were based on a list of over 400 "Socioeconomic Health Attributes" marketed by LexisNexis to improve predictive modeling:
55 <https://risk.lexisnexis.com/products/socioeconomic-health-attributes>
56
57
58
59
60

1
2
3 time of fewer than 2.35 minutes, which was one-half of the median soft launch time, were excluded from the data.
4
5 Respondents who left comments that were unrelated to what we were asking and exhibited a lack of thoughtful
6
7 consideration were also removed from the data. In total, we had 420 respondents in our data set.
8
9

10 *Data Analysis*

11 Descriptive statistics were first used to analyze all questions from the survey. All descriptive statistical
12
13 procedures were done in Microsoft Excel (Version 2006) and R (Version 3.6.1). We then computed a privacy
14
15 knowledge score by giving a participant 1 point for every question in the privacy knowledge section that they
16
17 answered with “Publicly available” and ½ point for every question they answered “Could be determined based on
18
19 other publicly available information.” Their points were summed to obtain their privacy knowledge score. We
20
21 computed a privacy concern score for each participant by averaging their responses to the nine questions in the
22
23 general privacy section. A Pearson correlation coefficient was computed in IBM SPSS Statistics for Windows,
24
25 version 24 to assess the relationship between privacy concern score and privacy knowledge score.

26 To further assess privacy concern, we created two new data sets from the original, one containing
27
28 participants with a high privacy concern score (greater than 3) and one containing participants with a low privacy
29
30 concern score (less than or equal to 3) to assess how each group responded to all of the survey questions. Chi-
31
32 square tests were conducted to determine whether there was a relationship between all of the demographic factors
33
34 and general privacy concern rating (high or low) and between concern over each vignette and privacy concern
35
36 rating. All chi-square analysis was done in SPSS. A t-test was performed in R to determine if there was a
37
38 significant difference between the privacy knowledge score means of the high general privacy concern group and
39
40 low general privacy concern group. To evaluate the relationship between privacy concern and knowledge over
41
42 specific data attributes, we asked about data attributes from related categories in both the concern and knowledge
43
44 section of the survey. For each concern question, we split the data into two groups, some concern (participants
45
46 who answered extremely, moderately, or somewhat concerned) and little-to-no concern (participants who
47
48 answered slightly or not at all concerned) and analyzed how much knowledge each group had regarding that
49
50 specific data attribute. We isolated all participants an additional time based on their concern toward the suicide
51
52 prediction model vignette. Participants who selected extremely, moderately, or somewhat concerned in response
53
54 to this vignette were placed in a “some concern” group, and participants who selected slightly or not at all
55
56 concerned were placed in a “little-to-no concern” group. We then analyzed the groups' responses to their concern
57
58 over the use of various data attributes in a suicide prediction model.
59
60

52 **Results**

54 *Demographics*

1
2
3 Over half (251/420, 59.8%) of the participants were aged 26-45, while the other 169 participants (40.2%)
4 ranged from 46-66 or more years old. The mean age was 44.97 years (SD = 14.92). Exactly half of the
5 participants were female (210/420, 50.0%). Participants could select multiple ethnicities: 327 participants were
6 white, 46 were African American/Black, 17 were Asian/Pacific Islander, 25 were Hispanic/Latinx, six were
7 Middle Eastern, and seven were Native American/Indigenous. Most participants were married or in a domestic
8 partnership (260/420, 61.9%), and 100 participants (23.8%) were single and never married. Education levels
9 varied; 330 participants (78.6%) had at least some level of college experience, whereas 88 (21.0%) had either less
10 than a high school degree, a high school degree or equivalent, or had attended trade school. Most participants
11 (274/420, 65.2%) were employed, and 72 (17.1%) were retired. There was a broad range of household income;
12 69.8% of participants (293/420) had an income of less than \$100,000, whereas 28.1% of participants (118/420)
13 had an income of greater than \$100,000. Household sizes tended to be small (286/420, 68.1%), although 129
14 participants (30.7%) lived with 3-5 members and 5 participants (1.2%) lived with six or more members.
15
16
17
18
19
20
21
22
23

24 *Knowledge of Publicly Available Information*

25 We assessed the knowledge of participants regarding the publicness of various socioeconomic data
26 attributes by testing their awareness of 15 data elements that are publicly available. Summary results are provided
27 in Table 1. Overall, 4179 of the total 6300 answer responses (66.3%) to the privacy knowledge questions were
28 correct (participant selected "Could be determined based on other publicly available information" or "Publicly
29 available), and 2121 responses (33.7%) were incorrect (participant selected "Not publicly available" or "I don't
30 know"). The publicness of several data attributes was fairly common knowledge: whether you own or rent at your
31 current address (77.4% correctness), whether or not you are registered to vote (76.4% correctness), the amount of
32 time you lived at your previous address (75.2% correctness), the last recorded sale price of your current address
33 (74.3% correctness), and the amount of time since you last moved (72.4% correctness). The most missed question
34 asked about the publicness of your total number of relatives and associates that own a boat or airplane (46.0%
35 correctness). Questions asking about knowledge of the publicness of derogatory record data all had between 64-
36 70% correctness: time since your most recent arrest (64.5% correctness), the total number of misdemeanor
37 convictions (68.3% correctness), the total count of household members with felony convictions (68.6%
38 correctness), whether you have been housed in a correctional facility (69.0% correctness), and total bankruptcy
39 filings (69.3% correctness). Questions with approximately half of the correct responses were: amount of time
40 since your last car accident (53.3% correctness), the total number of relatives and associates who have attended
41 college (56.9% correctness), number of members in your household with licenses for concealed weapons (59.8%
42 correctness), and your estimated household income range (63.6% correctness).
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 [Table 1 here]
4
5

6 *Vignette Privacy Concerns*

7
8 Respondents were presented with ten vignettes to gauge their privacy concerns across different types of
9 information gathered across different contents. Overview results are provided in Table 2. The vignette that
10 generated the most concern was about a public health worker collecting GPS data to assess adherence to stay-at-
11 home orders during the COVID-19 pandemic (319/420, 76% answered that they were somewhat, moderately, or
12 extremely concerned about this vignette). The vignette that generated the least concern was about a restaurant
13 owner conducting surveys to improve the restaurant's quality of service (184/420, 43.8% answered that they were
14 not at all concerned about this vignette).
15
16
17
18

19 The vignette discussing the use of socioeconomic data in a suicide prediction algorithm showed
20 considerable concern, with two-thirds of participants answering that they were somewhat, moderately, or
21 extremely concerned about this use of data, while 32.1% answered that they were not at all concerned or slightly
22 concerned.
23
24
25

26
27 [Table 2 here]
28
29

30 *Factors Contributing to Concern*

31
32 Overall, the factors that contributed most to concern over the vignettes were the *purpose of collecting the*
33 *data* (784/3780 total responses to factor questions, 20.7%) and the *potential future use of data* (737/3780 total
34 responses to factor questions, 19.5%). The *purpose of data collection* was the greatest contributing factor for the
35 vignettes about fitness data being used to develop a weight loss product, purchase transactions being used to stock
36 products, search history being used for targeted advertising, Twitter 'Following' lists being used to identify
37 accounts of people getting information about the Black Lives Matter movement, and email tracking being used to
38 understand a company's target audience. *Future use* was the greatest contributing factor for the vignettes about
39 tracking COVID-19 stay-at-home orders and the use of genomic data to identify a cancer-causing mutation.
40
41
42
43

44 For the vignette discussing the use of socioeconomic data in a suicide prediction algorithm, a technical
45 error prevented users from seeing the options of *potential future use of the data* and *none*. Based on the choices
46 available to respondents for this vignette, the *type of data being collected* generated the most concern (139/420,
47 33.1%) for this vignette.
48
49
50

51 *Concern Over Data Attributes in Suicide Prediction Model*

52
53
54
55
56
57
58
59
60

1
2
3
4 Specific to the development of a suicide prediction algorithm, we presented respondents with various data
5 attributes that would be used for that purpose (see Appendix B). These results are summarized in Table 3.

6 Attributes with the greatest expressed concern (indicating some, moderate, or extreme concern) include those
7 about annual income (78.6%), ownership of assets (71.2%) or value of real estate (71.4%), court appearances
8 (69.0%), arrest records (68.3%), and felony records (67.1%), and whether one holds a license for concealed
9 weapons (64.8%). Concerns were also evident in data attributes about one's entire household, with many
10 exceeding the concern expressed for the individual data attribute.
11
12

13
14 Attributes with the least amount of expressed concern include possessing a hunting or fishing license
15 (50.0%), whether one attended college (51.7%), or the number of times in a car accident (54.5%).
16
17
18

19 [Table 3 here]
20
21

22 *General Concerns Translate to Specific Concerns*

23
24 Based on participants' responses to the general privacy questions, we computed a privacy concern score
25 and assigned participants to a high or low general privacy concern group. Two hundred fifteen participants
26 (51.2%) were assigned to the high general privacy concern group and 205 (48.8%) were assigned to the low
27 general privacy concern group. The mean privacy concern score for all participants was 3.06 out of 5 (SD = 0.60).
28 The mean privacy concern score for the high concern group was 3.53 (SD = 0.33), and the mean privacy concern
29 score for the low concern group was 2.56 (SD = 0.38).
30
31
32

33 Overall, and as expected, those with low general privacy concerns tended to have lower concerns with the
34 vignettes, while those with greater general privacy concerns found the vignettes more concerning. For the vignette
35 discussing the use of socioeconomic data in a suicide prediction algorithm showed considerable concern, 80.0%
36 of respondents with high overall privacy concerns found this particular vignette concerning. And of the
37 respondents with low overall privacy concerns, 55.1% found this vignette concerning.
38
39
40

41 We then isolated the responses of participants in the high privacy concern group and the low privacy
42 concern group to evaluate how each group responded to the appropriateness of specific socioeconomic data
43 attributes being used in a suicide prediction model. The overall privacy concern groupings tended to be indicative
44 of participants' responses to the specific data attributes. Of the 6450 total responses to questions regarding concern
45 over the use of various data attributes from the high privacy concern group, there were 3865 moderately or
46 extremely concerned responses (59.9%). Of the 6150 total responses to questions regarding concern over the use
47 of various data attributes from the low privacy concern group, there were only 1767 moderately or extremely
48 concerned responses (28.7%) and 1842 not at all concerned responses (30.0%).
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 Among the high overall privacy concern group, considerable concern was expressed for all data attributes,
4 with the lowest-rated factor being “Whether you have a hunting or fishing license” with only 64.2% expressing
5 some, moderate, or extreme concern. For those with low overall privacy concerns, some attributes still presented
6 considerable concern: 70.2% of this group expressed some, moderate, or extreme concern about “Your estimated
7 annual income,” and 70.7% expressed similar concern over “The total estimated annual income for your entire
8 household.” Data attributes referencing household members also tended to rate higher than other attributes for this
9 group.
10
11
12
13
14
15

16 Discussion

17 Overall Findings

18
19 While the benefits of big data are manifold, it is necessary to consider the ethical questions and privacy
20 concerns that arise when integrating data pertaining to various contexts of people’s lives into algorithmic and
21 related computational models. Previous studies have approached these questions from a variety of contexts,
22 including personal fitness data (Bietz *et al.*, 2016; Zimmer *et al.*, 2020) and social media monitoring (Reuter *et*
23 *al.*, 2019). This study differs in that we aimed to capture the variation in privacy concerns across several contexts,
24 spanning from medicine to business to social media. We were also specifically interested in assessing the ethical
25 concerns related to a suicide prediction model and determining people’s awareness of the publicness of select data
26 attributes, as well as which of these data attributes generated the most concern in such a context.
27
28
29
30
31

32 Addressing RQ1, we found that overall, two-thirds of our respondents correctly determined that the data
33 elements in the survey were publicly available (either directly available or through some sort of inference). While
34 this can be viewed positively, the fact remains that one-third did not have complete awareness of the extent of the
35 publicness of various socioeconomic data points.
36
37

38 Addressing RQ2, we found considerable concern across the various contexts represented in our vignettes.
39 With the exception of vignette 7 (a restaurant using customer satisfaction surveys to improve quality), the
40 majority of respondents expressed some level of concern about the data use proposed within the hypothetical
41 vignettes. The highest levels of concern centered on GPS tracking for social distancing compliance, and marketers
42 monitoring email and search engine activities. General concerns over the collection and use of personal data
43 during the COVID-19 pandemic might be a contributing factor to concerns over vignette 1. Overall, concern was
44 greatest (68.0% expressing some, moderate, or extreme concern) in vignettes that focused on the use of personal
45 information within the medical context (vignettes 1, 2, and 10).
46
47
48
49
50

51 Our results also show that the most common factors contributing concerns across the various vignettes
52 were the *purpose of data collection* and the *potential future use of data*, a finding supported by existing research
53
54
55
56
57
58
59
60

1
2
3 showing consumers are most concerned about how companies are and might be using their personal information
4 in the future (Hoffman *et al.*, 1999; Phelps *et al.*, 2018) [11].
5

6 By splitting the data into a high privacy concern group and a low privacy concern group, we were able to
7 identify trends within and across these groups. While we expected people with generally high levels of overall
8 privacy concerns to therefore express concerns with our vignettes, we were more curious as to whether
9 individuals who typically have low privacy concerns might suddenly express concern for a particular scenario. As
10 with the general findings, even those with low privacy concerns expressed considerable worry about GPS-
11 tracking during the COVID-19 pandemic, as well as having researchers monitor Black Lives Matter activity on
12 Twitter. Here, our low privacy concern respondents expressed similar worries from the high concern group about
13 how such data might be used for other purposes.
14
15
16
17
18
19
20

21 *Data Concerns with Suicide Risk Prediction Modelling*

22 RQ3 reflects on our specific interest in measuring individuals' comfort with incorporating socioeconomic
23 data within suicide risk prediction models. While this particular vignette ranked in the middle of overall concern,
24 various data elements stood out as particularly problematic among our respondents. Our results show a clear
25 concern from all participants in data attributes related to income, crime and court records, and assets. This is
26 consistent with other research [14,17], indicating that most consumers were unwilling to share information about
27 household income and other financial information. Data about one's household – beyond just the individual –
28 were also particularly concerns for our respondents, suggesting that even if one might be comfortable with their
29 own being used for risk modeling, data about other household members is more problematic. This held true even
30 for respondents with generally low privacy concerns, suggesting these data elements are particularly troublesome
31 when used within this context.
32
33
34
35
36
37

38 Connected to RQ1, a concerning finding is that many attributes that a majority of respondents failed to
39 recognize were publicly-available were also flagged as particularly concerning in the detailed assessment of data
40 used within suicide risk prediction algorithms. For example, 54.0% of respondents didn't recognize that the
41 "count of relatives and associates that own a boat or airplane" was publicly-available, yet 75.5% found it
42 somewhat, moderately, or extremely concerning that the data element "Whether or not anyone in your household
43 owns assets (such as a watercraft, an aircraft, or real estate property)" might be used in a suicide risk prediction
44 algorithm. Similarly, over 40% of respondents did not realize "number of members in your household with
45 licenses for concealed weapons" was publicly-available, while 66.4% found using such data concerning. This
46 suggests many respondents have concerns over the use of certain data elements while underestimating the general
47 availability of the data.
48
49
50
51
52
53
54
55
56
57
58
59
60

Study Limitations

We recognize that participants recruited through Qualtrics are likely digitally savvy individuals and of a high enough socioeconomic status to own a device on which to take the survey. We acknowledge that these characteristics likely had some impact on our results. To help mitigate the effects of these characteristics, we requested that Qualtrics provide us with a specific distribution of individuals across age, income, and gender. Nonetheless, these characteristics undoubtedly had an influence on how participants responded, in particular, to the vignettes.

We also recognize that had we framed the suicide prediction model vignette in slightly different terms, it could have elicited a different response from participants. For example, had we put a greater emphasis on the benefits and societal good of creating such an algorithm and had we clarified that all personal information would be de-identified, perhaps participants would have been less concerned over the use of data in this way. **Future work could focus more specifically on participants' concerns over a suicide risk prediction algorithm and include vignettes all with the main purpose of creating a suicide prediction model but altering more minor factors about the vignettes, such as the type of data used, how it was obtained, and whether the algorithm would be used by someone other than clinicians.**

Conclusion

In this study, we measured public opinion regarding the use of data in various contexts. In particular, we were interested in assessing opinion over the use of publicly available socioeconomic data in a suicide risk prediction algorithm. To aid in our analysis of these contexts, we also measured public knowledge of select data attributes and concern over the incorporation of these attributes into the suicide prediction model.

Combining socioeconomic data with existing medical records gives researchers the opportunity to improve suicide prediction models. It is clear that the overall goal of this initiative, minimizing suicide attempts, is good and beneficial to society. However, informed by the lens of contextual integrity, the incorporation of socioeconomic data within suicide risk prediction models threatens to violate existing norms of what information is appropriate within the medical context. We found that over two-thirds of participants have at least some concern level toward using socioeconomic data in the suicide prediction algorithm. In comparison to the response to the nine other vignettes, this suicide prediction model vignette fell approximately in the middle in terms of the level of concern. This indicates that while this case is less concerning than some popular uses of data today, such as tracking of search history or email tracking, it is undoubtedly more concerning than researchers accessing genomic data from an ancestry website or fitness data from a wearable device.

We also found that the publicness of some data attributes was well known, such as voter registration records and address records, whereas the publicness of other types of information was less well known, such as

1
2
3 asset records of relatives and accident records. We highlighted certain data attributes that were particularly
4 sensitive to individuals who exhibited both high and low privacy concern, such as data related to income, assets,
5 and criminal records. Taken together, medical patients may have a lack of awareness that their doctors have
6 access to their socioeconomic data and data about their household members which has been aggregated by a third
7 party, and some of those data elements are particularly problematic.
8
9
10

11 Ultimately, we were able to determine that the appropriateness of incorporating personal data within
12 various computational applications is contextually dependent, with the appropriateness often determined by the
13 type of use and concern over future uses of data. We found that the use of certain data attributes is more
14 concerning than others, and that individuals often lack full knowledge of the availability of public data, especially
15 certain sensitive socioeconomic data attributes about our lives and our broader households. Specifically, we
16 determined that participants were most concerned about the use of income records, asset data, and criminal
17 records in suicide risk prediction models, with asset data also being among the data elements participants were
18 least aware were publicly available. Therefore, researchers hoping to rely on such data need to take steps to fully
19 consider the broader ethical and privacy implications of relying on such data, despite their possible predictive
20 value.
21
22
23
24
25
26

27 In the broadest sense, we have shown how confronting the ethical and privacy implications of
28 incorporating publicly available socioeconomic data into algorithmic models presents a unique challenge that
29 requires more than simply relying on the public availability of such data. Researchers – and the general public –
30 are better off when we rely on robust conceptual frameworks such as contextual integrity and engage in social
31 science-based research to better understand the knowledge and expectations of the general public. Algorithmic
32 models like those to help predict suicide risk can be of great public benefit, but only if pursued in an ethically
33 informed manner.
34
35
36
37
38
39

40 **Acknowledgments**

41 This material is based upon work supported by the National Science Foundation REU site grant #IIS-
42 1950826 “Data Science Across the Disciplines.” We also thank Dr. Jordan Smoller and his colleagues at Harvard
43 Medical School and in the Psychiatric & Neurodevelopmental Genetics Unit (PNGU) at Massachusetts General
44 Hospital for their feedback and support.
45
46
47
48

49 **References**

50 Barak-Corren, Y., Castro, V.M., Javitt, S., Hoffnagle, A.G., Dai, Y., Perlis, R.H., Nock, M.K., *et al.* (2017),
51 “Predicting Suicidal Behavior From Longitudinal Electronic Health Records”, *The American Journal of*
52 *Psychiatry*, Vol. 174 No. 2, pp. 154–162.
53
54
55
56
57
58
59
60

- 1
2
3 Bietz, M.J., Bloss, C.S., Calvert, S., Godino, J.G., Gregory, J., Claffey, M.P., Sheehan, J., *et al.* (2016),
4 “Opportunities and challenges in the use of personal health data for health research”, *Journal of the*
5 *American Medical Informatics Association : JAMIA*, Vol. 23 No. e1, pp. e42–e48.
6
7
8 boyd, danah and Crawford, K. (2012), “Critical Questions for Big Data”, *Information, Communication & Society*,
9 Vol. 15 No. 5, pp. 662–679.
10
11 Crain, M. (2018), “The limits of transparency: Data brokers and commodification”, *New Media & Society*, SAGE
12 Publications, Vol. 20 No. 1, pp. 88–104.
13
14 Fiesler, C., Beard, N. and Keegan, B.C. (2020), “No Robots, Spiders, or Scrapers: Legal and Ethical Regulation
15 of Data Collection Methods in Social Media Terms of Service”, *Proceedings of the International AAAI*
16 *Conference on Web and Social Media*, Vol. 14, pp. 187–196.
17
18 Fiesler, C. and Proferes, N. (2018), “‘Participant’ Perceptions of Twitter Research Ethics”, *Social Media +*
19 *Society*, Vol. 4 No. 1, p. 2056305118763366.
20
21
22 Fiscella, K., Franks, P., Gold, M.R. and Clancy, C.M. (2000), “Inequality in Quality: Addressing Socioeconomic,
23 Racial, and Ethnic Disparities in Health Care”, *JAMA*, Vol. 283 No. 19, p. 2579.
24
25 Hoffman, D.L., Novak, T.P. and Peralta, M. (1999), “Building consumer trust online”, *Communications of the*
26 *ACM*, Vol. 42 No. 4, pp. 80–85.
27
28 Hoofnagle, C. (2004), “Big Brother’s Little Helpers: How ChoicePoint and Other Commercial Data Brokers
29 Collect and Package Your Data for Law Enforcement”, *North Carolina Journal of International Law*,
30 Vol. 29 No. 4, p. 595.
31
32 Marmot, M., Friel, S., Bell, R., Houweling, T.A. and Taylor, S. (2008), “Closing the gap in a generation: health
33 equity through action on the social determinants of health”, *The Lancet*, Vol. 372 No. 9650, pp. 1661–
34 1669.
35
36
37
38 Martin, K. and Nissenbaum, H. (2017), “Privacy Interests in Public Records: An Empirical Investigation”,
39 *Harvard Journal of Law & Technology*, Vol. 31, p. 111.
40
41 Metcalf, J. and Crawford, K. (2016), “Where are human subjects in Big Data research? The emerging ethics
42 divide”, *Big Data & Society*, Vol. 3 No. 1, p. 2053951716650211.
43
44 Metcalf, J., Keller, E.F. and boyd, danah. (2016), *Perspectives on Big Data, Ethics, and Society*, Council for Big
45 Data, Ethics, and Society, available at: [https://bdes.datasociety.net/council-output/perspectives-on-big-](https://bdes.datasociety.net/council-output/perspectives-on-big-data-ethics-and-society/)
46 [data-ethics-and-society/](https://bdes.datasociety.net/council-output/perspectives-on-big-data-ethics-and-society/) (accessed 19 April 2019).
47
48
49 Metzger, M.J. (2007), “Communication Privacy Management in Electronic Commerce”, *Journal of Computer-*
50 *Mediated Communication*, Vol. 12 No. 2, pp. 335–361.
51
52 Nissenbaum, H. (2010), *Privacy in Context: Technology, Policy, and the Integrity of Social Life*, Stanford Law
53 Books, Stanford, Calif.
54
55
56
57
58
59
60

- 1
2
3 O'Connor, R.C. and Portzky, G. (2018), "Looking to the Future: A Synthesis of New Developments and
4 Challenges in Suicide Research and Prevention", *Frontiers in Psychology*, Vol. 9, available
5 at:<https://doi.org/10.3389/fpsyg.2018.02139>.
6
7
8 Phelps, J., Nowak, G. and Ferrell, E. (2018), "Privacy Concerns and Consumer Willingness to Provide Personal
9 Information":, *Journal of Public Policy & Marketing*, SAGE PublicationsSage CA: Los Angeles, CA,
10 available at:<https://doi.org/10.1509/jppm.19.1.27.16941>.
11
12
13 Reuter, K., Zhu, Y., Angyan, P., Le, N., Merchant, A.A. and Zimmer, M. (2019), "Public Concern About
14 Monitoring Twitter Users and Their Conversations to Recruit for Clinical Trials: Survey Study", *Journal*
15 *of Medical Internet Research*, Vol. 21 No. 10, p. e15455.
16
17
18 Ribeiro, J., Gutierrez, P., Joiner, T., Kessler, R., Petukhova, M., Sampson, N., Stein, M., *et al.* (2017), "Health
19 care contact and suicide risk documentation prior to suicide death: Results from the Army Study to
20 Assess Risk and Resilience in Servicemembers", *Journal of Consulting and Clinical Psychology*, Vol.
21 85 No. 4, pp. 403–408.
22
23
24 Shilton, K. (2015), "Emerging Ethics Norms in Social Media Research", presented at the Workshop on Beyond
25 IRBs: Ethical Review Processes for Big Data Research, available at:
26 <https://bigdata.fpf.org/papers/emerging-ethics-norms-in-social-media-research/> (accessed 28 December
27 2016).
28
29
30 "Suicide statistics". (2019), *American Foundation for Suicide Prevention*, 15 November, available at:
31 <https://afsp.org/suicide-statistics/> (accessed 4 October 2020).
32
33
34 Tene, O. and Polonetsky, J. (2013), "Big Data for All: Privacy and User Control in the Age of Analytics",
35 *Northwestern Journal of Technology and Intellectual Property*, Vol. 11 No. 5, pp. 239–273.
36
37 Turecki, G. and Brent, D.A. (2016), "Suicide and suicidal behaviour", *Lancet (London, England)*, Vol. 387 No.
38 10024, pp. 1227–1239.
39
40 Vitak, J., Shilton, K. and Ashktorab, Z. (2016), "Beyond the Belmont Principles: Ethical Challenges, Practices,
41 and Beliefs in the Online Data Research Community", *Proceedings of the 19th ACM Conference on*
42 *Computer-Supported Cooperative Work & Social Computing*, ACM, New York, NY, USA, pp. 941–
43 953.
44
45
46 Zimmer, M. (2010), "'But the data is already public': On the ethics of research in Facebook", *Ethics and*
47 *Information Technology*, Vol. 12 No. 4, pp. 313–325.
48
49 Zimmer, M. (2016), "OkCupid Study Reveals the Perils of Big-Data Science", *Wired*, 14 May, available at:
50 <https://www.wired.com/2016/05/okcupid-study-reveals-perils-big-data-science/> (accessed 28 May 2016).
51
52
53
54
55
56
57
58
59
60

1
2
3 Zimmer, M., Kumar, P., Vitak, J., Liao, Y. and Chamberlain Kritikos, K. (2020), “‘There’s Nothing Really They
4 Can Do with This Information’: Unpacking How Users Manage Privacy Boundaries for Personal Fitness
5 Information”, *Information, Communication & Society*, Vol. 23 No. 7, pp. 1020–1037.
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Privacy Concerns with Using Public Data for Suicide Risk Prediction Algorithms: A Public Opinion Survey of Contextual Appropriateness

Abstract

Purpose: Existing algorithms for predicting suicide risk rely solely on data from electronic health records, but such models could be improved through the incorporation of publicly available socioeconomic data – such as financial, legal, life event, and sociodemographic data. This study’s purpose is to understand the complex ethical and privacy implications of incorporating sociodemographic data within the health context. We present results from a survey exploring what the general public’s knowledge and concerns are about such publicly available data and the appropriateness of using it in suicide risk prediction algorithms.

Design/methodology/approach: A survey was developed to measure public opinion about privacy concerns with using socioeconomic data across different contexts. We presented respondents with multiple vignettes that described scenarios situated in medical, private business, and social media contexts, and asked participants to rate their level of concern over the context and what factor contributed most to their level of concern. Specific to suicide prediction, we presented respondents with various data attributes that could potentially be used in the context of a suicide risk algorithm and asked participants to rate how concerned they would be if each attribute was used for this purpose.

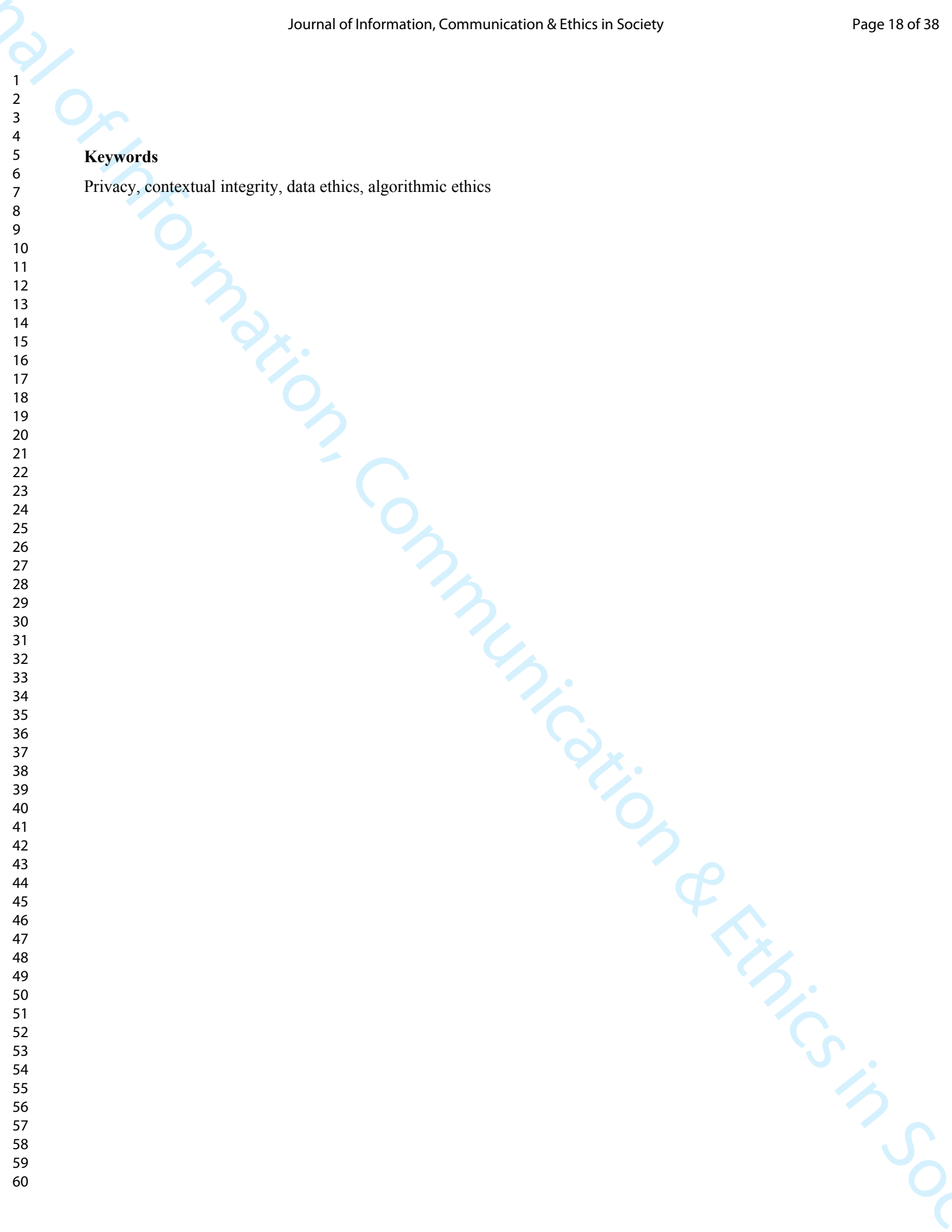
Findings: We found considerable concern across the various contexts represented in our vignettes, with greatest concern in vignettes that focused on the use of personal information within the medical context. Specific to the question of incorporating socioeconomic data within suicide risk prediction models, our results show a clear concern from all participants in data attributes related to income, crime and court records, and assets. Data about one’s household were also particularly concerns for our respondents, suggesting that even if one might be comfortable with their own being used for risk modeling, data about other household members is more problematic.

Originality: Previous studies on the privacy concerns that arise when integrating data pertaining to various contexts of people’s lives into algorithmic and related computational models have approached these questions from individual contexts. This study differs in that we captured the variation in privacy concerns across multiple contexts. We also specifically assessed the ethical concerns related to a suicide prediction model and determining people’s awareness of the publicness of select data attributes, as well as which of these data attributes generated the most concern in such a context. As far as we know, this is the first study to pursue this question.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Keywords

Privacy, contextual integrity, data ethics, algorithmic ethics



Introduction

Suicide is currently the 10th leading cause of death in the United States, and suicide rates have continued to increase in recent years (“Suicide statistics”, 2019). About half of suicide decedents have contact with the healthcare system in the month before their death (Ribeiro *et al.*, 2017), indicating that there is a significant opportunity to identify patients who are at risk for a suicide attempt when they visit their healthcare provider. Most algorithms for predicting suicide risk relies solely on data from electronic health records (Barak-Corren *et al.*, 2017), but it has been proposed that algorithmic solutions for predicting suicide risk could be improved through the incorporation of publicly available socioeconomic data – such as financial, legal, life event, derogatory, and sociodemographic data – alongside electronic health record data (Barak-Corren *et al.*, 2017; O’Connor and Portzky, 2018; Turecki and Brent, 2016).

If incorporation of publicly available socioeconomic data into the suicide prediction model – and other predictive models – is to be pursued, we must understand the complex ethical and privacy implications so as to not erode the trust between a patient and the healthcare system. To start this effort, we present results from a survey exploring what the general public’s knowledge and concerns are about such publicly available data.

Background

In an age of rapidly increasing technological advancements, scholars from all fields of expertise are investigating how big data can aid them in their work. From bioinformaticians who analyze genetic code to political scientists who analyze government records to sociologists who analyze social media data, everyone wants in on the benefits of big data and the algorithmic systems it empowers (boyd and Crawford, 2012). While the potential uses for big data and algorithm-based systems are nearly endless, it is important and necessary to consider the ethical and privacy concerns related to these endless possibilities (Metcalf *et al.*, 2016).

One critical question focuses on the ethical and privacy complications of using publicly-available data in research. Significant gaps among researchers on what constitutes “public” data that does not require explicit consent prior to harvesting (Zimmer, 2010, 2016), whether a platform’s terms of service might allow the automated scraping of public data (Fiesler *et al.*, 2020), and even at what stage does computational research become human subjects research requiring particular ethical protection (Metcalf and Crawford, 2016). Further, users are often not aware of the types of access researchers have to public data via social media platforms and their APIs (Fiesler and Proferes, 2018). Given the above, uncertainty persists among the research community on how to address the ethical and privacy complications of using ostensibly public data in computational research projects (Shilton, 2015; Vitak *et al.*, 2016).

These challenges become even more complicated within the context of medical data, where users’ knowledge and attitudes about the public availability of their health data is particularly muddled. While some

1
2
3 studies have shown users expressing little concern over privacy issues related to their personal fitness
4 information, noting such data was not inherently sensitive and expressing ambivalence over the possibility of
5 sharing that data with third parties (Zimmer *et al.*, 2020), other studies have shown users of self-tracking
6 technologies are frequently unaware of the details of external data access to which they agree in the context of
7 clicking “accept” to terms of use” (Bietz *et al.*, 2016). Further, a recent study assessing privacy concerns related to
8 the use of publicly-available health-related tweets data in research found the acceptability of harvesting such data
9 depended greatly on the nature of the health ailment, who was collecting it, and the context of use (Reuter *et al.*,
10 2019).

11
12 The potential of incorporating publicly-available socioeconomic data – typically obtained from
13 commercial data brokers – into medical research complicates things even further. The role of socioeconomic
14 factors in various aspects of healthcare we well studied (Fiscella *et al.*, 2000; Marmot *et al.*, 2008), and their
15 potential for assisting in the particular challenge of suicide prevention are promising (Barak-Corren *et al.*, 2017;
16 O’Connor and Portzky, 2018; Turecki and Brent, 2016). But integrating such a wide range of data points – such
17 as financial, legal, life event, derogatory, and sociodemographic data – into different contexts is often met with
18 resistance, irrespective of how public the information might be (Crain, 2018; Hoofnagle, 2004; Martin and
19 Nissenbaum, 2017; Tene and Polonetsky, 2013).

20
21 To help assess the ethics of incorporating public socioeconomic data in the development of algorithms to
22 predict suicide risk, we invoke Nissenbaum’s (Nissenbaum, 2010) theory of privacy as contextual integrity (CI).
23 CI rejects the traditional dichotomy of public versus private information, as well as the notion that privacy
24 preferences and decisions in one context universally apply to other contexts. Instead, CI rests on the
25 understanding that our interactions with other people, institutions, and technologies, occur within particular
26 contexts. Norms of appropriateness govern people’s expectations of how personal information should flow within
27 any given context. Therefore, responding to a data ethics question—e.g., should third-party socioeconomic data
28 be integrated with health data?—needs to start not with privacy as a static set of principles but with an
29 understanding of norms of appropriateness within the context in which the data is being collected and used, and
30 whether it is deemed appropriate to move information from one context – such as socioeconomic data – and apply
31 it within a new context – such as medical research in suicide risk prediction.

32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 **Study Objective**

48
49 We rely on contextual integrity to investigate people’s privacy concerns about integrating publicly
50 available socioeconomic data within various contexts, including a suicide risk prediction algorithm. We aim to
51 measure public opinion about privacy concerns across medical-, business-, and social media-related contexts, as
52 well as specifically investigate the privacy concerns related to a suicide prediction model. To do this, we first
53
54
55
56
57
58
59
60

1
2
3 assess participants' general privacy concerns and their privacy knowledge. Then we evaluate participants'
4 opinions and concerns towards various contexts, which we will then be able to compare to concerns towards the
5 use of socioeconomic data in the suicide prediction algorithm. Overall, we aim to determine whether it is ethical
6 to combine publicly available socioeconomic data with health data to improve a suicide prediction algorithm
7 through an assessment of individuals' awareness of this very possibility as well as their comfort with the
8 appropriateness of this use of said data within this particular context.
9

10
11 We approach this through the following research questions:
12

13 RQ1: To what extent do people understand that their socioeconomic data – such as financial, legal, life
14 event, derogatory, and sociodemographic data – is publicly available?
15

16 RQ2: In which contexts do individuals find the use of publicly available information most concerning?
17

18 RQ3: Regarding suicide risk prediction algorithms specifically, which socioeconomic data points are
19 most problematic for inclusion?
20
21
22
23

24 **Methods**

25 *Survey Instrument*

26
27 A survey was developed to measure public opinion about privacy concerns with using various
28 socioeconomic data points across various contexts. The survey consisted of five sections that allowed us to collect
29 information about demographics of respondents, general privacy concerns, awareness of what types of
30 information are publicly available, privacy concerns associated with specific contexts, and concerns over the use
31 of thirty publicly available data attributes in the context of a suicide risk prediction algorithm. To inquire about
32 personal data use across a broad range of contexts, we presented respondents with ten vignettes that described
33 scenarios situated in medical, private business, and social media contexts (see Appendix A). Following each
34 vignette, we asked participants to rate their level of concern over the context, and we asked what factor
35 contributed most to their level of concern. Specific to suicide prediction, we presented respondents with various
36 data attributes¹ that could potentially be used in the context of a suicide risk algorithm and asked participants to
37 rate how concerned they would be if each attribute was used for this purpose (see Appendix B). The survey was
38 tested for clarity and consistency, and the research protocol received Institutional Review Board approval.
39
40
41
42
43
44
45
46
47

48 *Data Collection*

49 The survey was deployed on Qualtrics from July 17-24, 2020. We contracted Qualtrics to recruit
50 approximately respondents between the ages of 26-99 who live in the United States. Respondents with a response
51
52
53

54 ¹ Data attributes were based on a list of over 400 "Socioeconomic Health Attributes" marketed by LexisNexis to improve predictive modeling:
55 <https://risk.lexisnexis.com/products/socioeconomic-health-attributes>
56
57
58
59
60

1
2
3 time of fewer than 2.35 minutes, which was one-half of the median soft launch time, were excluded from the data.
4
5 Respondents who left comments that were unrelated to what we were asking and exhibited a lack of thoughtful
6
7 consideration were also removed from the data. In total, we had 420 respondents in our data set.
8
9

10 *Data Analysis*

11 Descriptive statistics were first used to analyze all questions from the survey. All descriptive statistical
12
13 procedures were done in Microsoft Excel (Version 2006) and R (Version 3.6.1). We then computed a privacy
14
15 knowledge score by giving a participant 1 point for every question in the privacy knowledge section that they
16
17 answered with “Publicly available” and ½ point for every question they answered “Could be determined based on
18
19 other publicly available information.” Their points were summed to obtain their privacy knowledge score. We
20
21 computed a privacy concern score for each participant by averaging their responses to the nine questions in the
22
23 general privacy section. A Pearson correlation coefficient was computed in IBM SPSS Statistics for Windows,
24
25 version 24 to assess the relationship between privacy concern score and privacy knowledge score.

26 To further assess privacy concern, we created two new data sets from the original, one containing
27
28 participants with a high privacy concern score (greater than 3) and one containing participants with a low privacy
29
30 concern score (less than or equal to 3) to assess how each group responded to all of the survey questions. Chi-
31
32 square tests were conducted to determine whether there was a relationship between all of the demographic factors
33
34 and general privacy concern rating (high or low) and between concern over each vignette and privacy concern
35
36 rating. All chi-square analysis was done in SPSS. A t-test was performed in R to determine if there was a
37
38 significant difference between the privacy knowledge score means of the high general privacy concern group and
39
40 low general privacy concern group. To evaluate the relationship between privacy concern and knowledge over
41
42 specific data attributes, we asked about data attributes from related categories in both the concern and knowledge
43
44 section of the survey. For each concern question, we split the data into two groups, some concern (participants
45
46 who answered extremely, moderately, or somewhat concerned) and little-to-no concern (participants who
47
48 answered slightly or not at all concerned) and analyzed how much knowledge each group had regarding that
49
50 specific data attribute. We isolated all participants an additional time based on their concern toward the suicide
51
52 prediction model vignette. Participants who selected extremely, moderately, or somewhat concerned in response
53
54 to this vignette were placed in a “some concern” group, and participants who selected slightly or not at all
55
56 concerned were placed in a “little-to-no concern” group. We then analyzed the groups' responses to their concern
57
58 over the use of various data attributes in a suicide prediction model.
59
60

52 **Results**

54 *Demographics*

1
2
3 Over half (251/420, 59.8%) of the participants were aged 26-45, while the other 169 participants (40.2%)
4 ranged from 46-66 or more years old. The mean age was 44.97 years (SD = 14.92). Exactly half of the
5 participants were female (210/420, 50.0%). Participants could select multiple ethnicities: 327 participants were
6 white, 46 were African American/Black, 17 were Asian/Pacific Islander, 25 were Hispanic/Latinx, six were
7 Middle Eastern, and seven were Native American/Indigenous. Most participants were married or in a domestic
8 partnership (260/420, 61.9%), and 100 participants (23.8%) were single and never married. Education levels
9 varied; 330 participants (78.6%) had at least some level of college experience, whereas 88 (21.0%) had either less
10 than a high school degree, a high school degree or equivalent, or had attended trade school. Most participants
11 (274/420, 65.2%) were employed, and 72 (17.1%) were retired. There was a broad range of household income;
12 69.8% of participants (293/420) had an income of less than \$100,000, whereas 28.1% of participants (118/420)
13 had an income of greater than \$100,000. Household sizes tended to be small (286/420, 68.1%), although 129
14 participants (30.7%) lived with 3-5 members and 5 participants (1.2%) lived with six or more members.
15
16
17
18
19
20
21
22
23

24 *Knowledge of Publicly Available Information*

25 We assessed the knowledge of participants regarding the publicness of various socioeconomic data
26 attributes by testing their awareness of 15 data elements that are publicly available. Summary results are provided
27 in Table 1. Overall, 4179 of the total 6300 answer responses (66.3%) to the privacy knowledge questions were
28 correct (participant selected "Could be determined based on other publicly available information" or "Publicly
29 available), and 2121 responses (33.7%) were incorrect (participant selected "Not publicly available" or "I don't
30 know"). The publicness of several data attributes was fairly common knowledge: whether you own or rent at your
31 current address (77.4% correctness), whether or not you are registered to vote (76.4% correctness), the amount of
32 time you lived at your previous address (75.2% correctness), the last recorded sale price of your current address
33 (74.3% correctness), and the amount of time since you last moved (72.4% correctness). The most missed question
34 asked about the publicness of your total number of relatives and associates that own a boat or airplane (46.0%
35 correctness). Questions asking about knowledge of the publicness of derogatory record data all had between 64-
36 70% correctness: time since your most recent arrest (64.5% correctness), the total number of misdemeanor
37 convictions (68.3% correctness), the total count of household members with felony convictions (68.6%
38 correctness), whether you have been housed in a correctional facility (69.0% correctness), and total bankruptcy
39 filings (69.3% correctness). Questions with approximately half of the correct responses were: amount of time
40 since your last car accident (53.3% correctness), the total number of relatives and associates who have attended
41 college (56.9% correctness), number of members in your household with licenses for concealed weapons (59.8%
42 correctness), and your estimated household income range (63.6% correctness).
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 [Table 1 here]
4
5

6 *Vignette Privacy Concerns*

7
8 Respondents were presented with ten vignettes to gauge their privacy concerns across different types of
9 information gathered across different contents. Overview results are provided in Table 2. The vignette that
10 generated the most concern was about a public health worker collecting GPS data to assess adherence to stay-at-
11 home orders during the COVID-19 pandemic (319/420, 76% answered that they were somewhat, moderately, or
12 extremely concerned about this vignette). The vignette that generated the least concern was about a restaurant
13 owner conducting surveys to improve the restaurant's quality of service (184/420, 43.8% answered that they were
14 not at all concerned about this vignette).
15
16
17
18

19 The vignette discussing the use of socioeconomic data in a suicide prediction algorithm showed
20 considerable concern, with two-thirds of participants answering that they were somewhat, moderately, or
21 extremely concerned about this use of data, while 32.1% answered that they were not at all concerned or slightly
22 concerned.
23
24
25

26
27 [Table 2 here]
28
29

30 *Factors Contributing to Concern*

31
32 Overall, the factors that contributed most to concern over the vignettes were the *purpose of collecting the*
33 *data* (784/3780 total responses to factor questions, 20.7%) and the *potential future use of data* (737/3780 total
34 responses to factor questions, 19.5%). The *purpose of data collection* was the greatest contributing factor for the
35 vignettes about fitness data being used to develop a weight loss product, purchase transactions being used to stock
36 products, search history being used for targeted advertising, Twitter 'Following' lists being used to identify
37 accounts of people getting information about the Black Lives Matter movement, and email tracking being used to
38 understand a company's target audience. *Future use* was the greatest contributing factor for the vignettes about
39 tracking COVID-19 stay-at-home orders and the use of genomic data to identify a cancer-causing mutation.
40
41
42
43

44 For the vignette discussing the use of socioeconomic data in a suicide prediction algorithm, a technical
45 error prevented users from seeing the options of *potential future use of the data* and *none*. Based on the choices
46 available to respondents for this vignette, the *type of data being collected* generated the most concern (139/420,
47 33.1%) for this vignette.
48
49
50

51 *Concern Over Data Attributes in Suicide Prediction Model*

52
53
54
55
56
57
58
59
60

1
2
3
4 Specific to the development of a suicide prediction algorithm, we presented respondents with various data
5 attributes that would be used for that purpose (see Appendix B). These results are summarized in Table 3.

6 Attributes with the greatest expressed concern (indicating some, moderate, or extreme concern) include those
7 about annual income (78.6%), ownership of assets (71.2%) or value of real estate (71.4%), court appearances
8 (69.0%), arrest records (68.3%), and felony records (67.1%), and whether one holds a license for concealed
9 weapons (64.8%). Concerns were also evident in data attributes about one's entire household, with many
10 exceeding the concern expressed for the individual data attribute.
11
12

13
14 Attributes with the least amount of expressed concern include possessing a hunting or fishing license
15 (50.0%), whether one attended college (51.7%), or the number of times in a car accident (54.5%).
16
17
18

19 [Table 3 here]
20
21

22 *General Concerns Translate to Specific Concerns*

23
24 Based on participants' responses to the general privacy questions, we computed a privacy concern score
25 and assigned participants to a high or low general privacy concern group. Two hundred fifteen participants
26 (51.2%) were assigned to the high general privacy concern group and 205 (48.8%) were assigned to the low
27 general privacy concern group. The mean privacy concern score for all participants was 3.06 out of 5 (SD = 0.60).
28 The mean privacy concern score for the high concern group was 3.53 (SD = 0.33), and the mean privacy concern
29 score for the low concern group was 2.56 (SD = 0.38).
30
31
32

33 Overall, and as expected, those with low general privacy concerns tended to have lower concerns with the
34 vignettes, while those with greater general privacy concerns found the vignettes more concerning. For the vignette
35 discussing the use of socioeconomic data in a suicide prediction algorithm showed considerable concern, 80.0%
36 of respondents with high overall privacy concerns found this particular vignette concerning. And of the
37 respondents with low overall privacy concerns, 55.1% found this vignette concerning.
38
39
40

41 We then isolated the responses of participants in the high privacy concern group and the low privacy
42 concern group to evaluate how each group responded to the appropriateness of specific socioeconomic data
43 attributes being used in a suicide prediction model. The overall privacy concern groupings tended to be indicative
44 of participants' responses to the specific data attributes. Of the 6450 total responses to questions regarding concern
45 over the use of various data attributes from the high privacy concern group, there were 3865 moderately or
46 extremely concerned responses (59.9%). Of the 6150 total responses to questions regarding concern over the use
47 of various data attributes from the low privacy concern group, there were only 1767 moderately or extremely
48 concerned responses (28.7%) and 1842 not at all concerned responses (30.0%).
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 Among the high overall privacy concern group, considerable concern was expressed for all data attributes,
4 with the lowest-rated factor being “Whether you have a hunting or fishing license” with only 64.2% expressing
5 some, moderate, or extreme concern. For those with low overall privacy concerns, some attributes still presented
6 considerable concern: 70.2% of this group expressed some, moderate, or extreme concern about “Your estimated
7 annual income,” and 70.7% expressed similar concern over “The total estimated annual income for your entire
8 household.” Data attributes referencing household members also tended to rate higher than other attributes for this
9 group.
10
11
12
13
14
15

16 Discussion

17 Overall Findings

18
19 While the benefits of big data are manifold, it is necessary to consider the ethical questions and privacy
20 concerns that arise when integrating data pertaining to various contexts of people’s lives into algorithmic and
21 related computational models. Previous studies have approached these questions from a variety of contexts,
22 including personal fitness data (Bietz *et al.*, 2016; Zimmer *et al.*, 2020) and social media monitoring (Reuter *et*
23 *al.*, 2019). This study differs in that we aimed to capture the variation in privacy concerns across several contexts,
24 spanning from medicine to business to social media. We were also specifically interested in assessing the ethical
25 concerns related to a suicide prediction model and determining people’s awareness of the publicness of select data
26 attributes, as well as which of these data attributes generated the most concern in such a context.
27
28
29
30
31

32 Addressing RQ1, we found that overall, two-thirds of our respondents correctly determined that the data
33 elements in the survey were publicly available (either directly available or through some sort of inference). While
34 this can be viewed positively, the fact remains that one-third did not have complete awareness of the extent of the
35 publicness of various socioeconomic data points.
36
37

38 Addressing RQ2, we found considerable concern across the various contexts represented in our vignettes.
39 With the exception of vignette 7 (a restaurant using customer satisfaction surveys to improve quality), the
40 majority of respondents expressed some level of concern about the data use proposed within the hypothetical
41 vignettes. The highest levels of concern centered on GPS tracking for social distancing compliance, and marketers
42 monitoring email and search engine activities. General concerns over the collection and use of personal data
43 during the COVID-19 pandemic might be a contributing factor to concerns over vignette 1. Overall, concern was
44 greatest (68.0% expressing some, moderate, or extreme concern) in vignettes that focused on the use of personal
45 information within the medical context (vignettes 1, 2, and 10).
46
47
48
49
50

51 Our results also show that the most common factors contributing concerns across the various vignettes
52 were the *purpose of data collection* and the *potential future use of data*, a finding supported by existing research
53
54
55
56
57
58
59
60

1
2
3 showing consumers are most concerned about how companies are and might be using their personal information
4 in the future (Hoffman *et al.*, 1999; Phelps *et al.*, 2018) [11].
5

6
7 By splitting the data into a high privacy concern group and a low privacy concern group, we were able to
8 identify trends within and across these groups. While we expected people with generally high levels of overall
9 privacy concerns to therefore express concerns with our vignettes, we were more curious as to whether
10 individuals who typically have low privacy concerns might suddenly express concern for a particular scenario. As
11 with the general findings, even those with low privacy concerns expressed considerable worry about GPS-
12 tracking during the COVID-19 pandemic, as well as having researchers monitor Black Lives Matter activity on
13 Twitter. Here, our low privacy concern respondents expressed similar worries from the high concern group about
14 how such data might be used for other purposes.
15
16
17
18
19
20

21 *Data Concerns with Suicide Risk Prediction Modelling*

22 RQ3 reflects on our specific interest in measuring individuals' comfort with incorporating socioeconomic
23 data within suicide risk prediction models. While this particular vignette ranked in the middle of overall concern,
24 various data elements stood out as particularly problematic among our respondents. Our results show a clear
25 concern from all participants in data attributes related to income, crime and court records, and assets. This is
26 consistent with other research [14,17], indicating that most consumers were unwilling to share information about
27 household income and other financial information. Data about one's household – beyond just the individual –
28 were also particularly concerns for our respondents, suggesting that even if one might be comfortable with their
29 own being used for risk modeling, data about other household members is more problematic. This held true even
30 for respondents with generally low privacy concerns, suggesting these data elements are particularly troublesome
31 when used within this context.
32
33
34
35
36
37

38 Connected to RQ1, a concerning finding is that many attributes that a majority of respondents failed to
39 recognize were publicly-available were also flagged as particularly concerning in the detailed assessment of data
40 used within suicide risk prediction algorithms. For example, 54.0% of respondents didn't recognize that the
41 "count of relatives and associates that own a boat or airplane" was publicly-available, yet 75.5% found it
42 somewhat, moderately, or extremely concerning that the data element "Whether or not anyone in your household
43 owns assets (such as a watercraft, an aircraft, or real estate property)" might be used in a suicide risk prediction
44 algorithm. Similarly, over 40% of respondents did not realize "number of members in your household with
45 licenses for concealed weapons" was publicly-available, while 66.4% found using such data concerning. This
46 suggests many respondents have concerns over the use of certain data elements while underestimating the general
47 availability of the data.
48
49
50
51
52
53
54
55
56
57
58
59
60

Study Limitations

We recognize that participants recruited through Qualtrics are likely digitally savvy individuals and of a high enough socioeconomic status to own a device on which to take the survey. We acknowledge that these characteristics likely had some impact on our results. To help mitigate the effects of these characteristics, we requested that Qualtrics provide us with a specific distribution of individuals across age, income, and gender. Nonetheless, these characteristics undoubtedly had an influence on how participants responded, in particular, to the vignettes.

We also recognize that had we framed the suicide prediction model vignette in slightly different terms, it could have elicited a different response from participants. For example, had we put a greater emphasis on the benefits and societal good of creating such an algorithm and had we clarified that all personal information would be de-identified, perhaps participants would have been less concerned over the use of data in this way. Future work could focus more specifically on participants' concerns over a suicide risk prediction algorithm and include vignettes all with the main purpose of creating a suicide prediction model but altering more minor factors about the vignettes, such as the type of data used, how it was obtained, and whether the algorithm would be used by someone other than clinicians.

Conclusion

In this study, we measured public opinion regarding the use of data in various contexts. In particular, we were interested in assessing opinion over the use of publicly available socioeconomic data in a suicide risk prediction algorithm. To aid in our analysis of these contexts, we also measured public knowledge of select data attributes and concern over the incorporation of these attributes into the suicide prediction model.

Combining socioeconomic data with existing medical records gives researchers the opportunity to improve suicide prediction models. It is clear that the overall goal of this initiative, minimizing suicide attempts, is good and beneficial to society. However, informed by the lens of contextual integrity, the incorporation of socioeconomic data within suicide risk prediction models threatens to violate existing norms of what information is appropriate within the medical context. We found that over two-thirds of participants have at least some concern level toward using socioeconomic data in the suicide prediction algorithm. In comparison to the response to the nine other vignettes, this suicide prediction model vignette fell approximately in the middle in terms of the level of concern. This indicates that while this case is less concerning than some popular uses of data today, such as tracking of search history or email tracking, it is undoubtedly more concerning than researchers accessing genomic data from an ancestry website or fitness data from a wearable device.

We also found that the publicness of some data attributes was well known, such as voter registration records and address records, whereas the publicness of other types of information was less well known, such as

1
2
3 asset records of relatives and accident records. We highlighted certain data attributes that were particularly
4 sensitive to individuals who exhibited both high and low privacy concern, such as data related to income, assets,
5 and criminal records. Taken together, medical patients may have a lack of awareness that their doctors have
6 access to their socioeconomic data and data about their household members which has been aggregated by a third
7 party, and some of those data elements are particularly problematic.
8
9
10

11 Ultimately, we were able to determine that the appropriateness of incorporating personal data within
12 various computational applications is contextually dependent, with the appropriateness often determined by the
13 type of use and concern over future uses of data. We found that the use of certain data attributes is more
14 concerning than others, and that individuals often lack full knowledge of the availability of public data, especially
15 certain sensitive socioeconomic data attributes about our lives and our broader households. Specifically, we
16 determined that participants were most concerned about the use of income records, asset data, and criminal
17 records in suicide risk prediction models, with asset data also being among the data elements participants were
18 least aware were publicly available. Therefore, researchers hoping to rely on such data need to take steps to fully
19 consider the broader ethical and privacy implications of relying on such data, despite their possible predictive
20 value.
21
22
23
24
25
26

27 In the broadest sense, we have shown how confronting the ethical and privacy implications of
28 incorporating publicly available socioeconomic data into algorithmic models presents a unique challenge that
29 requires more than simply relying on the public availability of such data. Researchers – and the general public –
30 are better off when we rely on robust conceptual frameworks such as contextual integrity and engage in social
31 science-based research to better understand the knowledge and expectations of the general public. Algorithmic
32 models like those to help predict suicide risk can be of great public benefit, but only if pursued in an ethically
33 informed manner.
34
35
36
37
38
39

40 **Acknowledgments**

41 This material is based upon work supported by the National Science Foundation REU site grant #IIS-
42 1950826 “Data Science Across the Disciplines.” We also thank Dr. Jordan Smoller and his colleagues at Harvard
43 Medical School and in the Psychiatric & Neurodevelopmental Genetics Unit (PNGU) at Massachusetts General
44 Hospital for their feedback and support.
45
46
47
48

49 **References**

50 Barak-Corren, Y., Castro, V.M., Javitt, S., Hoffnagle, A.G., Dai, Y., Perlis, R.H., Nock, M.K., *et al.* (2017),
51 “Predicting Suicidal Behavior From Longitudinal Electronic Health Records”, *The American Journal of*
52 *Psychiatry*, Vol. 174 No. 2, pp. 154–162.
53
54
55
56
57
58
59
60

- 1
2
3 Bietz, M.J., Bloss, C.S., Calvert, S., Godino, J.G., Gregory, J., Claffey, M.P., Sheehan, J., *et al.* (2016),
4 “Opportunities and challenges in the use of personal health data for health research”, *Journal of the*
5 *American Medical Informatics Association : JAMIA*, Vol. 23 No. e1, pp. e42–e48.
6
7
8 boyd, danah and Crawford, K. (2012), “Critical Questions for Big Data”, *Information, Communication & Society*,
9 Vol. 15 No. 5, pp. 662–679.
10
11 Crain, M. (2018), “The limits of transparency: Data brokers and commodification”, *New Media & Society*, SAGE
12 Publications, Vol. 20 No. 1, pp. 88–104.
13
14 Fiesler, C., Beard, N. and Keegan, B.C. (2020), “No Robots, Spiders, or Scrapers: Legal and Ethical Regulation
15 of Data Collection Methods in Social Media Terms of Service”, *Proceedings of the International AAAI*
16 *Conference on Web and Social Media*, Vol. 14, pp. 187–196.
17
18 Fiesler, C. and Proferes, N. (2018), “‘Participant’ Perceptions of Twitter Research Ethics”, *Social Media +*
19 *Society*, Vol. 4 No. 1, p. 2056305118763366.
20
21
22 Fiscella, K., Franks, P., Gold, M.R. and Clancy, C.M. (2000), “Inequality in Quality: Addressing Socioeconomic,
23 Racial, and Ethnic Disparities in Health Care”, *JAMA*, Vol. 283 No. 19, p. 2579.
24
25 Hoffman, D.L., Novak, T.P. and Peralta, M. (1999), “Building consumer trust online”, *Communications of the*
26 *ACM*, Vol. 42 No. 4, pp. 80–85.
27
28 Hoofnagle, C. (2004), “Big Brother’s Little Helpers: How ChoicePoint and Other Commercial Data Brokers
29 Collect and Package Your Data for Law Enforcement”, *North Carolina Journal of International Law*,
30 Vol. 29 No. 4, p. 595.
31
32
33 Marmot, M., Friel, S., Bell, R., Houweling, T.A. and Taylor, S. (2008), “Closing the gap in a generation: health
34 equity through action on the social determinants of health”, *The Lancet*, Vol. 372 No. 9650, pp. 1661–
35 1669.
36
37
38 Martin, K. and Nissenbaum, H. (2017), “Privacy Interests in Public Records: An Empirical Investigation”,
39 *Harvard Journal of Law & Technology*, Vol. 31, p. 111.
40
41 Metcalf, J. and Crawford, K. (2016), “Where are human subjects in Big Data research? The emerging ethics
42 divide”, *Big Data & Society*, Vol. 3 No. 1, p. 2053951716650211.
43
44 Metcalf, J., Keller, E.F. and boyd, danah. (2016), *Perspectives on Big Data, Ethics, and Society*, Council for Big
45 Data, Ethics, and Society, available at: [https://bdes.datasociety.net/council-output/perspectives-on-big-](https://bdes.datasociety.net/council-output/perspectives-on-big-data-ethics-and-society/)
46 [data-ethics-and-society/](https://bdes.datasociety.net/council-output/perspectives-on-big-data-ethics-and-society/) (accessed 19 April 2019).
47
48
49 Metzger, M.J. (2007), “Communication Privacy Management in Electronic Commerce”, *Journal of Computer-*
50 *Mediated Communication*, Vol. 12 No. 2, pp. 335–361.
51
52 Nissenbaum, H. (2010), *Privacy in Context: Technology, Policy, and the Integrity of Social Life*, Stanford Law
53 Books, Stanford, Calif.
54
55
56
57
58
59
60

- 1
2
3 O'Connor, R.C. and Portzky, G. (2018), "Looking to the Future: A Synthesis of New Developments and
4 Challenges in Suicide Research and Prevention", *Frontiers in Psychology*, Vol. 9, available
5 at:<https://doi.org/10.3389/fpsyg.2018.02139>.
6
7
8 Phelps, J., Nowak, G. and Ferrell, E. (2018), "Privacy Concerns and Consumer Willingness to Provide Personal
9 Information":, *Journal of Public Policy & Marketing*, SAGE PublicationsSage CA: Los Angeles, CA,
10 available at:<https://doi.org/10.1509/jppm.19.1.27.16941>.
11
12
13 Reuter, K., Zhu, Y., Angyan, P., Le, N., Merchant, A.A. and Zimmer, M. (2019), "Public Concern About
14 Monitoring Twitter Users and Their Conversations to Recruit for Clinical Trials: Survey Study", *Journal*
15 *of Medical Internet Research*, Vol. 21 No. 10, p. e15455.
16
17
18 Ribeiro, J., Gutierrez, P., Joiner, T., Kessler, R., Petukhova, M., Sampson, N., Stein, M., *et al.* (2017), "Health
19 care contact and suicide risk documentation prior to suicide death: Results from the Army Study to
20 Assess Risk and Resilience in Servicemembers", *Journal of Consulting and Clinical Psychology*, Vol.
21 85 No. 4, pp. 403–408.
22
23
24 Shilton, K. (2015), "Emerging Ethics Norms in Social Media Research", presented at the Workshop on Beyond
25 IRBs: Ethical Review Processes for Big Data Research, available at:
26 <https://bigdata.fpf.org/papers/emerging-ethics-norms-in-social-media-research/> (accessed 28 December
27 2016).
28
29
30 "Suicide statistics". (2019), *American Foundation for Suicide Prevention*, 15 November, available at:
31 <https://afsp.org/suicide-statistics/> (accessed 4 October 2020).
32
33
34 Tene, O. and Polonetsky, J. (2013), "Big Data for All: Privacy and User Control in the Age of Analytics",
35 *Northwestern Journal of Technology and Intellectual Property*, Vol. 11 No. 5, pp. 239–273.
36
37 Turecki, G. and Brent, D.A. (2016), "Suicide and suicidal behaviour", *Lancet (London, England)*, Vol. 387 No.
38 10024, pp. 1227–1239.
39
40 Vitak, J., Shilton, K. and Ashktorab, Z. (2016), "Beyond the Belmont Principles: Ethical Challenges, Practices,
41 and Beliefs in the Online Data Research Community", *Proceedings of the 19th ACM Conference on*
42 *Computer-Supported Cooperative Work & Social Computing*, ACM, New York, NY, USA, pp. 941–
43 953.
44
45
46 Zimmer, M. (2010), "'But the data is already public': On the ethics of research in Facebook", *Ethics and*
47 *Information Technology*, Vol. 12 No. 4, pp. 313–325.
48
49 Zimmer, M. (2016), "OkCupid Study Reveals the Perils of Big-Data Science", *Wired*, 14 May, available at:
50 <https://www.wired.com/2016/05/okcupid-study-reveals-perils-big-data-science/> (accessed 28 May 2016).
51
52
53
54
55
56
57
58
59
60

1
2
3 Zimmer, M., Kumar, P., Vitak, J., Liao, Y. and Chamberlain Kritikos, K. (2020), “‘There’s Nothing Really They
4 Can Do with This Information’: Unpacking How Users Manage Privacy Boundaries for Personal Fitness
5 Information”, *Information, Communication & Society*, Vol. 23 No. 7, pp. 1020–1037.
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Table I: Knowledge of Publicly Available Information

Data Attribute	I don't know	Not publicly available	Could be determined based on other information	Publicly available
1. The amount of time since your last car accident	83 19.8%	113 26.9%	133 31.7%	91 21.7%
2. Whether you own or rent at your current address	42 10.0%	53 12.6%	157 37.4%	168 40.0%
3. Total count of your relatives and associates that own a boat or airplane	110 26.2%	117 27.9%	125 29.8%	68 16.2%
4. Last recorded sale price of your current address	43 10.2%	65 15.5%	130 31.0%	182 43.3%
5. Time since your most recent arrest	60 14.3%	89 21.2%	101 24.0%	170 40.5%
6. Your total number of misdemeanor convictions	60 14.3%	73 17.4%	126 30.0%	161 38.3%
7. Total count of your household members with felony convictions	73 17.4%	56 14.0%	138 32.9%	150 35.7%
8. Whether or not you've been housed in a correctional facility	69 16.4%	61 14.5%	135 32.1%	155 36.9%
9. Your total number of bankruptcy filings	56 13.3%	73 17.4%	143 34.0%	148 35.2%
10. Total number of relatives and associates who have attended college	83 19.8%	98 23.3%	131 31.2%	108 25.7%
11. Your estimated household income range	49 11.7%	104 24.8%	177 42.1%	90 21.4%
12. The number of members in your household with licenses for concealed weapons	78 18.6%	91 21.7%	143 34.0%	108 25.7%
13. Amount of time since you last moved	40 9.5%	76 18.1%	171 40.7%	133 31.7%
14. Whether or not you are registered to vote	44 10.5%	55 13.1%	141 33.6%	180 42.9%
15. The amount of time you lived at your previous address	38 9.0%	66 15.7%	149 35.5%	167 39.8%

Table II: Vignette Privacy Concerns

Vignette	Not at all concerned	Slightly concerned	Somewhat concerned	Moderately concerned	Extremely concerned
1. GPS data for tracking stay-at-home order adherence	37 8.8%	64 15.2%	72 17.1%	81 19.3%	166 39.5%
2. Genomic data for identification of cancer-causing mutation	74 17.6%	93 22.1%	80 19.0%	72 17.1%	101 24.0%
3. Fitness data for a weight loss product	72 17.1%	76 18.1%	90 21.4%	74 17.6%	108 25.7%
4. Transaction records for stocking purposes	57 13.6%	91 21.7%	80 19.0%	81 19.3%	111 26.4%
5. Search history for targeted advertising	44 10.5%	76 18.1%	75 17.9%	97 23.1%	128 30.5%
6. Following lists for tracking users involved in Black Lives Matter movement	67 16.0%	60 14.3%	77 18.3%	84 20.0%	132 31.4%
7. Customer satisfaction surveys to improve restaurant quality	184 43.8%	54 12.9%	62 14.8%	57 13.6%	63 15.0%
8. Email tracking for understanding target audience	33 7.9%	82 19.5%	92 21.9%	78 18.6%	135 32.1%
9. Comments to improve a social media ad	110 26.2%	78 18.6%	77 18.3%	71 16.9%	84 20.0%
10. Socioeconomic data for a suicide prediction model	62 14.8%	73 17.4%	94 22.4%	95 22.6%	96 22.9%

Table III: Concern Over Data Attributes in Suicide Prediction Model

Data Attribute	Not at all concerned	Slightly concerned	Somewhat concerned	Moderately concerned	Extremely concerned
16. Number of times you have been in a car accident	97 23.1%	94 22.4%	78 18.6%	74 17.6%	77 18.3%
17. Distance (in miles) between you and your nearest relative	78 18.6%	62 14.8%	101 24.0%	88 21.0%	91 21.7%
18. Whether or not you are registered to vote	108 25.7%	54 12.9%	81 19.3%	91 21.7%	86 20.5%
19. The number of phone numbers associated with you	55 13.1%	70 16.7%	84 20.0%	89 21.2%	122 29.0%
20. Whether or not you attended college	124 29.5%	79 18.8%	79 18.8%	70 16.7%	68 16.2%
21. Total number of household member who attended college	117 27.9%	65 15.5%	75 17.9%	89 21.2%	74 17.6%
22. Your estimated annual income	37 8.8%	53 12.6%	93 22.1%	109 26.0%	128 30.5%
23. The total estimated annual income for your entire household	35 8.3%	55 13.1%	102 24.3%	94 22.4%	134 31.9%
24. The original mortgage dollar amount at your current address	74 17.6%	78 18.6%	87 20.7%	86 20.5%	95 22.6%
25. The estimated market value of your previous address	98 23.3%	49 11.7%	93 22.1%	101 24.0%	79 18.8%
26. The difference in neighborhood median household income between your address and your most recent address	81 19.3%	78 18.6%	84 20.0%	97 23.1%	80 19.0%
27. The number of multi-family properties in your neighborhood	128 30.5%	50 11.9%	88 21.0%	87 20.7%	67 16.0%
28. Your current address' neighborhood crime index, based on law enforcement data	110 26.2%	63 15.0%	88 21.0%	85 20.2%	74 17.6%
29. Your previous address' neighborhood crime index, based on law enforcement data	112 26.7%	52 12.4%	82 19.5%	99 23.6%	75 17.9%
30. Whether or not you own assets (such as a watercraft, aircraft, or real estate property)	45 10.7%	76 18.1%	104 24.8%	98 23.3%	97 23.1%
31. Whether or not anyone in your household owns assets (such as a watercraft, an aircraft, or real estate property)	41 9.8%	62 14.8%	110 26.2%	107 25.5%	100 23.8%
32. The total value for all real estate properties you own	55 13.1%	65 15.5%	86 20.5%	103 24.5%	111 26.4%
33. The total value for all real estate properties everyone in your household owns	49 11.7%	59 14.0%	92 21.9%	103 24.5%	117 27.9%
34. Total number of real estate	88	45	87	90	110

1						
2						
3						
4	properties sold within last 5 years	21.0%	10.7%	20.7%	21.4%	26.2%
5	35. The total number of court records	81	49	91	83	116
6	(including felony, misdemeanor,	19.3%	11.7%	21.7%	19.8%	27.6%
7	lien, judgment, bankruptcy, or					
8	eviction) listed in your name					
9	36. The total number of court records	75	54	95	85	111
10	(including felony, misdemeanor,	17.9%	12.9%	22.6%	20.2%	26.4%
11	lien, judgment, bankruptcy, or					
12	eviction) for your entire					
13	household					
14	37. Your total number of arrests	90	43	72	102	113
15		21.4%	10.2%	17.1%	24.3%	26.9%
16	38. Total number of arrests for your	85	50	73	95	117
17	entire household	20.2%	11.9%	17.4%	22.6%	27.9%
18	39. Your total number of felony	88	50	73	94	115
19	convictions	21.0%	11.9%	17.4%	22.4%	27.4%
20	40. Total number of felony	92	52	71	99	106
21	convictions for your entire	21.9%	12.4%	16.9%	23.6%	25.2%
22	household					
23	41. Whether you have a hunting or	140	70	65	79	66
24	fishing license	33.3%	16.7%	15.5%	18.8%	15.7%
25	42. Whether anyone in your	137	51	83	79	70
26	household has a hunting or fishing	32.6%	12.1%	19.8%	18.8%	16.7%
27	license					
28	43. Whether you have a license for	88	60	78	91	103
29	concealed weapons	21.0%	14.3%	18.6%	21.7%	24.5%
30	44. Whether anyone in your	87	54	95	95	89
31	household has a license for	20.7%	12.9%	22.6%	22.6%	21.2%
32	concealed weapons					
33	45. The number of times you have	89	65	87	80	99
34	changed addresses in the last 5	21.2%	15.5%	20.7%	19.0%	23.6%
35	years					
36						
37						
38						
39						
40						
41						
42						
43						
44						
45						
46						
47						
48						
49						
50						
51						
52						
53						
54						
55						
56						
57						
58						
59						
60						

Appendix A: Vignettes

Survey respondents were presented with ten vignettes that described scenarios situated in medical, private business, and social media contexts and were asked to provide their level of concern for their data being used in this way:

1. A public health worker is collecting location data from your phone's GPS system. The data will be used to track your adherence to stay-at-home orders during the COVID-19 pandemic.
2. A medical researcher is collecting genomic data from an ancestry website that collects saliva samples. The data will be used to confirm the identity of a cancer-causing mutation.
3. A research team for a fitness company is collecting data about your activity levels from your wearable fitness tracker. The data will be used to develop and market a new weight loss product.
4. A researcher for a private company is collecting your purchase transaction records from their stores. The data will be used to analyze customer purchase habits and then stock items accordingly.
5. A marketer is collecting your search history from a popular search engine. The data will be used to advertise their products to you based on your interests.
6. A university researcher is collecting the 'Following' list from Twitter accounts that have used the hashtag #BlackLivesMatter. The data will be used to identify the accounts of people who are getting information about the Black Lives Matter movement.
7. A manager at a local restaurant is collecting your responses to a customer satisfaction survey. The data will be used to improve the quality of service.
8. A marketer is collecting information about when and where you opened an email from an email tracking service. The data will be used to gain a better understanding of their target audience.
9. A marketer is collecting comments from one of their social media advertisements. The data will be used to assess the reaction to the ad and improve the ad accordingly.
10. A medical researcher wants to collect socioeconomic data from public databases and combine it with medical records of people in your community to improve an algorithm that identifies patients with suicide risk.

Respondents were also asked what factor contributed most to their level of concern:

- Who is obtaining your data
- The type of data being collected
- Where the data is being obtained from
- The purpose of collecting the data
- The potential future use of the data
- None of these

Appendix B: Data Attributes for Suicide Risk Prediction

Survey respondents were presented various data attributes would be used in the context of a suicide risk algorithm and asked participants to rate how concerned they would be if each attribute was used for that purpose:

1. Number of times you have been in a car accident
2. Distance (in miles) between you and your nearest relative
3. Whether or not you are registered to vote
4. The number of phone numbers associated with you
5. Whether or not you attended college
6. Total number of household members who attended college
7. Your estimated annual income
8. The total estimated annual income for your entire household
9. The original mortgage dollar amount at your current address
10. The estimated market value of your previous address
11. The difference in neighborhood median household income between your address and your most recent address
12. The number of multi-family properties in your neighborhood
13. Your current address' neighborhood crime index, based on law enforcement data
14. Your previous address' neighborhood crime index, based on law enforcement data
15. Whether or not you own assets (such as a watercraft, an aircraft, or real estate property)
16. Whether or not anyone in your household owns assets (such as a watercraft, an aircraft, or real estate property)
17. The total value for all real estate properties you own
18. The total value for all real estate properties everyone in your household owns
19. Total number of real estate properties sold within last 5 years
20. The total number of court records (including felony, misdemeanor, lien, judgment, bankruptcy, or eviction) listed in your name
21. The total number of court records (including felony, misdemeanor, lien, judgment, bankruptcy, or eviction) for your entire household
22. Your total number of arrests
23. Total number of arrests for your entire household
24. Your total number of felony convictions
25. Total number of felony convictions for your entire household
26. Whether you have a hunting or fishing license
27. Whether anyone in your household has a hunting or fishing license
28. Whether you have a license for concealed weapons
29. Whether anyone in your household has a license for concealed weapons
30. The number of times you have changed addresses in the last 5 years