

By Dr Libby Bishop and Simon Parker

# Big data - Ethical best practice

These data often differ from traditional research data (e.g., surveys) in that they have not been generated specifically for research.

- The data collection was not subject to any formal ethical review processes.
- Protections applied when data are collected (e.g., informed consent) and processed (e.g., de-identification), will not have been implemented.
- Using the data for research may substantially differ from the original purpose for which it was collected (e.g., data to improve direct health care used later for research), and this was not anticipated when data were generated.
- Data are less often held as discrete collections and may be linked with diverse sources, indeed a benefit of big data lies in the capacity to link many data sources.

## Ethical data resources

### UK Cabinet Office Data Science Ethical Framework

This is a general guide for all data scientists. The target audience is researchers in government, but the guide is useful for any data scientist and is relevant across disciplines and genres of data. There are detailed questions and excellent examples.



**Six key principles: at a glance view**

- 1 Start with clear user need and public benefit**  
Data science offers huge opportunities to create evidence for policymaking, and make quicker and more accurate operational decisions. Being clear about the public benefit will help you justify the sensitivity of the data (principle 2) and the method that you want to use (principle 3).
- 2 Use data and tools which have the minimum intrusion necessary**  
You should always use the minimum data necessary to achieve the public benefit. Sometimes you will need to use sensitive personal data. There are steps that you can take to safeguard people's privacy e.g. de-identifying or aggregating data to higher levels, querying against datasets or using synthetic data.
- 3 Create robust data science models**  
Good machine learning models can analyse far larger amounts of data far more quickly and accurately than traditional methods. Think through the quality and representativeness of the data, flag if algorithms are using protected characteristics (e.g. ethnicity) to make decisions, and think through unintended consequences. Complex decisions may well need the wider knowledge of policy or operational experts.
- 4 Be alert to public perceptions**  
The Data Protection Act requires you to have an understanding of how people would reasonably expect their personal data to be used. You need to be aware of shifting public perceptions. Social media data, commercial data and data scraped from the web allow us to understand more about the world, but come with different terms and conditions and levels of consent.
- 5 Be as open and accountable as possible**  
Being open allows us to talk about the public benefit of data science. Be as open as you can about the tools, data and algorithms (unless doing so would jeopardise the aim, e.g. fraud). Provide explanations in plain English and give people recourse to decisions which they think are incorrectly made. Make sure your project has oversight and accountability built in throughout.
- 6 Keep data secure**  
We know that the public are justifiably concerned about their data being lost or stolen. Government has a statutory duty to protect the public's data and as such it is vital that appropriate security measures are in place.

More detail in annex below

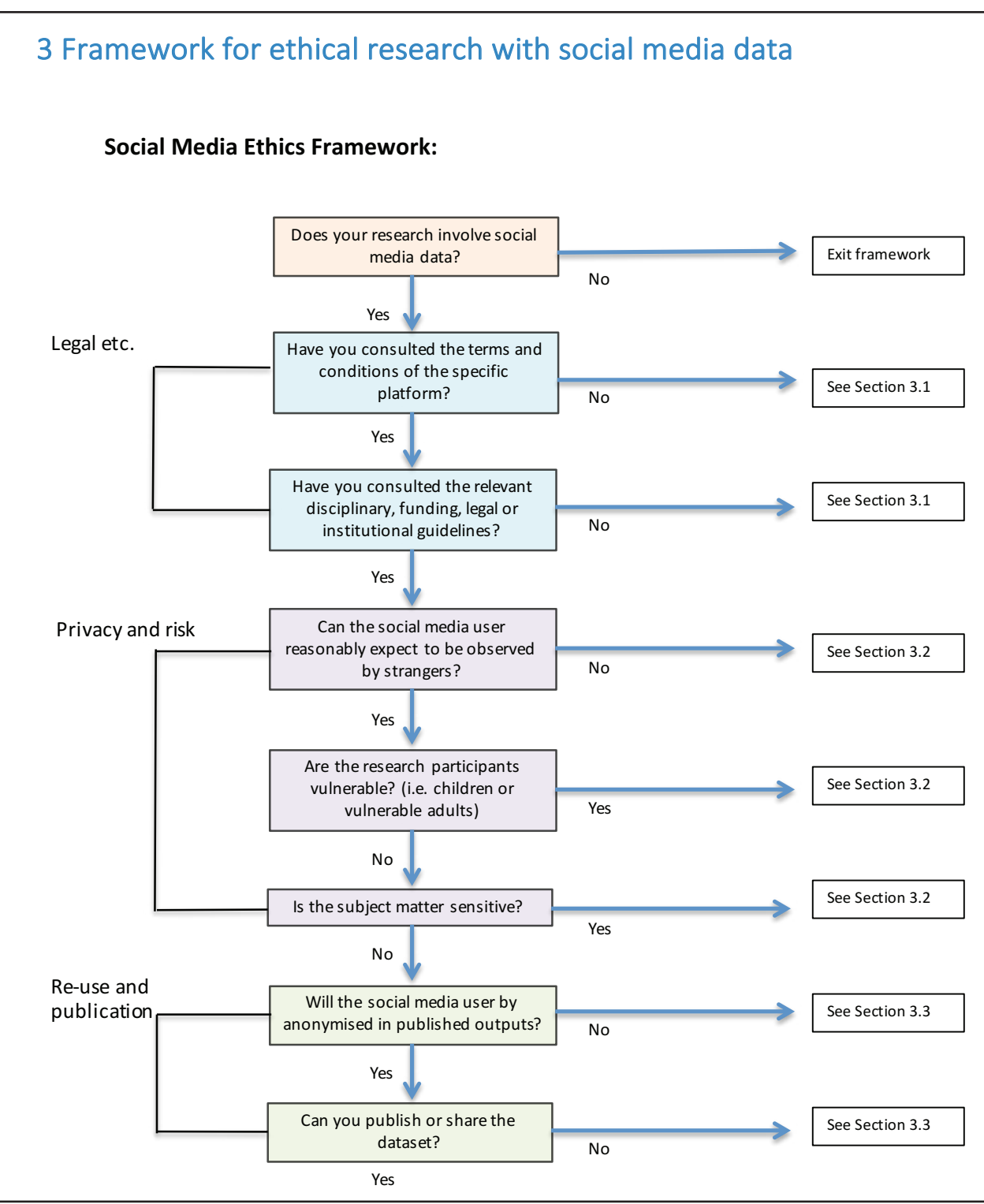
4

### The OECD report

Is the data personal or sensitive (under Data Protection laws)?  
Might the research discriminate against groups, not just individuals?  
In what settings is the information gathered, and what uses are expected in those settings?  
What are data subjects' reasonable expectations concerning the research project's re-contextualisation of their information?

### Social Media Ethical Framework

The Social Data Science Lab at the University of Cardiff brings together social, computer, political, health, statistical and mathematical scientists to study study diverse dimensions of New and Emerging Forms of Data in social and policy contexts. Their extensive and practical research has informed this guide for doing ethical research with social media.



### UK Data Service



Luke Sloan, Jeffrey Morgan, Pete Burnap, Matthew Williams. Who Tweets? Deriving the Demographic Characteristics of Age, Occupation and Social Class from Twitter User Meta-Data; PLOS One. March 2, 2015 <https://doi.org/10.1371/journal.pone.0115545>.

Williams, M., Burnap, P., & Sloan, L. (2017). Towards an ethical framework for publishing Twitter data in social research: taking into account users' views, online context and algorithmic estimation. Sociology (forthcoming).

## Sharing complex data

### Publishing social media

Because social media platforms vary greatly in their terms of use, as well as expectations of users (e.g., about privacy), each must be considered separately when publishing data. The Social Data Science Lab has produced an easy-to-use flowchart for making publishing decisions about aggregated and textual data from Twitter.

