

Available online at www.sciencedirect.com

ScienceDirect

www.compseconline.com/publications/prodclaw.htm

Computer Law &
Security Review

Functional anonymisation: Personal data and the data environment



Mark Elliot ^{a,*}, Kieron O'Hara ^b, Charles Raab ^c, Christine M. O'Keefe ^d, Elaine Mackey ^a, Chris Dibben ^c, Heather Gowans ^e, Kingsley Purdam ^a, Karen McCullagh ^f

- ^a University of Manchester, Manchester, UK
- ^b University of Southampton, Southampton, UK
- ^c University of Edinburgh, Edinburgh, UK
- ^d CSIRO, Canberra, Australia
- e University of Oxford, Oxford, UK
- ^f University of East Anglia, Norwich, Norfolk, UK

ABSTRACT

Keywords:
Anonymisation
Deidentification
Deanonymisation
Statistical disclosure control
Data environment
ADF
DDF
Functional anonymisation
Release-and-forget
Obscurity

Anonymisation of personal data has a long history stemming from the expansion of the types of data products routinely provided by National Statistical Institutes. Variants on anonymisation have received serious criticism reinforced by much-publicised apparent failures. We argue that both the operators of such schemes and their critics have become confused by being overly focused on the properties of the data itself. We claim that, far from being able to determine whether data is anonymous (and therefore non-personal) by looking at the data alone, any anonymisation technique worthy of the name must take account of not only the data but also its environment.

This paper proposes an alternative formulation called *functional anonymisation* that focuses on the relationship between the data and the environment within which the data exists (the *data environment*). We provide a formulation for describing the relationship between the data and its environment that links the legal notion of personal data with the statistical notion of disclosure control. Anonymisation, properly conceived and effectively conducted, can be a critical part of the toolkit of the privacy-respecting data controller and the wider remit of providing accurate and usable data.

© 2018 Mark Elliot, Kieron O'Hara, Charles Raab, Christine M. O'Keefe, Elaine Mackey, Chris Dibben, Heather Gowans, Kingsley Purdam, Karen McCullagh

1. Introduction

Superficially, the notion of anonymisation¹ is straightforward: if information contains the identity of a person, then other facts about them can be revealed by the dissemination

of that information, and this may breach that person's privacy. For example, sentence (1) discloses information about Jane.

- (1) Jane is a 39-year-old female, suffering from diabetes, who presented herself for treatment on 24th May.
- * Corresponding author. CCSR and Social Statistics, Humanities, University of Manchester, Bridgeford Street, Manchester, M13 9PL, UK. E-mail address: mark.elliot@manchester.ac.uk (M. Elliot).
- ¹ In some jurisdictions (for example the US, Canada and Australia) the term 'de-identification' is used to mean what anonymisation means in the EU context. In this paper we will use the term anonymisation throughout. https://doi.org/10.1016/j.clsr.2018.02.001

0267-3649/© 2018 Mark Elliot, Kieron O'Hara, Charles Raab, Christine M. O'Keefe, Elaine Mackey, Chris Dibben, Heather Gowans, Kingsley Purdam, Karen McCullagh

If we can isolate and remove (or replace) that part of the information that contains the person's identity, then that person will not be identifiable from the information, and his or her privacy will no longer be at risk in this context. In our example, we could replace sentence (1) by (2).

(2) A 39-year-old female, suffering from diabetes, presented herself for treatment on 24th May.

Clearly (1) is more specific and has more content than (2), but (2) still retains much of the information that is in (1). The decrease in information content in moving from (1) to (2) may be compensated for by the gain in privacy protection. There may also be practical benefits, in that people may be more willing to provide accurate or sensitive information if they trust that their privacy will be protected (Oswald, 2014). In many fields of public policy such as health, privacy protection is consonant with the public interest in the use of high-quality sources of data, but it can also be a barrier to research.

However, although the basic idea of anonymisation seems straightforward, the procedure is easier said than done (or, rather, done effectively). For example, sentence (2) might easily reveal the identity of the referent, if one knew a little extra information: for example, that Jane was the only woman of that approximate age who presented herself for treatment that day; someone who knew only that about Jane would thereby learn that Jane had diabetes from (2).

Indeed, it can be formally shown that anonymisation can always, in theory, be reversed, as long as there is some informational content remaining in the data (see for example Dwork, 2006). An adversary² attempting such a reversal could have access to an unpredictably wide range of information: for example, some of the information that we wish to protect by anonymisation might have been published on social media by Jane herself.

This does seem to lead us a worrying conclusion, that the only way to be certain of countering the threat of reidentification is to turn the information into noise (e.g. turning all of the values in a database to randomly generated ones). Does this mean that anonymisation is doomed to failure, and thus, legally or ethically, that the anonymiser has no justifiable practical basis for anonymisation? Are we condemned never to redeem any of the value of medical data, which is inherently associated with individuals at the micro-level, because in theory – an opportunity might emerge for an adversary to re-identify the individuals in question? Or is there a tradeoff to be made - as argued for example by Cavoukian and El Emam (2011) and Rubinstein and Hartzog (2016) - between the social (or commercial) value of sharing data, and some risk of identifying people, even if that trade-off has consequences for personal privacy?

It is difficult to answer these questions without making the concept of anonymisation more concrete. We argue in this paper

that (i) there are various interpretations of 'anonymisation', and also of the related notion of the risk of re-identification; (ii) what is deemed to be an acceptable level of risk will affect understandings of anonymisation and (iii) that anonymisation itself is a complex process requiring attention to far more than the data. This line of reasoning leads us to posit that, contrary to a series of influential commentaries, anonymisation, properly conceived and effectively conducted, can be a critical part of the toolkit of the privacy-respecting data controller and the wider remit of providing accurate and usable data.

The question of identifiability, which underlies anonymisation, is prominent in data protection legislation. In the US, 'personally identifiable information' (PII) has a narrower scope, referring to information maintained by a federal agency that can be used to trace an individual's identity or that is linkable to an individual (see McCallister et al., 2010). The definition is supported with examples of such identifiers: names, addresses (including email addresses), and other known identifiers such as the Social Security number. The European Union's (EU) data protection regime incorporates a category of 'personal data', defined as data from which the subject of the data is identifiable, either on its own or in tandem with auxiliary pieces of data.3 This creates a legal, as well as an ethical, driver for anonymisation. If the removal or perturbation of information can transform personal data into non-personal data, or PII into non-PII,4 then the information itself is outside the scope of data protection or privacy regulation, thereby reducing the constraints on the use of that information.

Anonymisation can therefore be seen through the lens of data protection law. If we look at the EU General Data Protection Regulation (GDPR) (EU, 2016/679), important and sometimes onerous restrictions are imposed on personal data, defined (article 4) as:

any information relating to an identified or identifiable natural person ('data subject'); an identifiable natural person is one who can be identified, directly or indirectly, in particular by reference to an identifier such as a name, an identification number, location data, an online identifier or to one or more factors specific to the physical, physiological, genetic, mental, economic, cultural or social identity of that natural person.

Note that identification can be direct or indirect. GDPR Recital 26 includes the explanation:

To ascertain whether means are reasonably likely to be used to identify the natural person, account should be taken of all objective factors, such as the costs of and the amount of time required for identification, taking into consideration the available technology at the time of the processing and technological developments.

² We use the term "adversary" throughout the paper to refer to an agent who attempts to re-identify an individual population unit within a de-identified dataset and the term "attack" to refer to the attempted re-identification. Synonymous terms that are found elsewhere in the literature are "intruder" (e.g. Elliot and Dale, 1999), and "snooper" (e.g. Duncan et al., 2011).

³ Directive 95/46/EC, Art 2(a); GDPR 2016/679, Art 4(1).

⁴ In this paper, except where flagged otherwise, we use the term 'personal data' to mean data from which people are identifiable, and therefore risky in a privacy sense, covering both EU personal data and PII.

Download English Version:

https://daneshyari.com/en/article/6890507

Download Persian Version:

https://daneshyari.com/article/6890507

<u>Daneshyari.com</u>