



## City Research Online

### City, University of London Institutional Repository

---

**Citation:** Hagger-Johnson, G. E., Harron, K., Goldstein, H., Parslow, R., Dattani, N., Borja, M. C., Wijlaars, L. & Gilbert, R. (2014). Making a hash of data: what risks to privacy does the NHS's care.data scheme pose?. *British Medical Journal (BMJ)*, 348(mar25 7), g2264. doi: 10.1136/bmj.g2264

This is the published version of the paper.

This version of the publication may differ from the final published version.

---

**Permanent repository link:** <https://openaccess.city.ac.uk/id/eprint/6956/>

**Link to published version:** <https://doi.org/10.1136/bmj.g2264>

**Copyright:** City Research Online aims to make research outputs of City, University of London available to a wider audience. Copyright and Moral Rights remain with the author(s) and/or copyright holders. URLs from City Research Online may be freely distributed and linked to.

**Reuse:** Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

---

City Research Online:

<http://openaccess.city.ac.uk/>

[publications@city.ac.uk](mailto:publications@city.ac.uk)

---

## LETTERS

## THE NHS'S CARE.DATA SCHEME

## Making a hash of data: what risks to privacy does the NHS's care.data scheme pose?

Gareth E Hagger-Johnson *senior research associate*<sup>1</sup>, Katie Harron *research associate*<sup>2</sup>, Harvey Goldstein *professor of statistics*<sup>2</sup>, Roger Parslow *senior lecturer in epidemiology*<sup>3</sup>, Nirupa Dattani *senior research fellow*<sup>4</sup>, Mario Cortina Borja *senior lecturer*<sup>2</sup>, Linda Wijlaars *research associate*<sup>2</sup>, Ruth Gilbert *professor of clinical epidemiology*<sup>2</sup>

<sup>1</sup>Institute of Child Health/Department of Epidemiology and Public Health, University College London, London, UK; <sup>2</sup>Institute of Child Health, University College London, London, UK; <sup>3</sup>Division of Epidemiology and Biostatistics, Leeds Institute of Genetics, Health and Therapeutics, University of Leeds, Leeds, UK; <sup>4</sup>Centre for Maternal and Child Health Research, City University London, London, UK

Care.data proposes to link individual level hospital episode statistics (HES) and general practice data at the Health and Social Care Information Centre. As is currently the case for HES, linked data will be pseudoanonymised before being released to researchers.<sup>1</sup> A proposed alternative is for identifiers (such as NHS number, date of birth) to be pseudoanonymised at source,<sup>2</sup> using an encrypted hash, before linkage is performed.<sup>3,4</sup>

Pseudoanonymisation at source will increase data linkage errors, where two records belonging to the same patient fail to link (missed match) or two records are incorrectly assigned to the same patient (false match). Duplicate records and "confusions" (two patients sharing a record) often occur in clinical settings (for example, owing to changes of name or address, typographical errors).

Data linkage errors have clinical implications but are also relevant to commissioning and research. False matches lead to overestimation of prevalence (if cases are counted twice). Missed matches lead to underestimation of prevalence (if cases are missed) and loss of statistical power. When healthier subgroups of the population are more likely to link correctly than others, biased estimates of relative risk can occur. Linkage errors lower the quality of information available and can lead to flawed decision making.

Records that can be linked are restricted to those with complete identifiers required by the linkage algorithm, but not all of these will be correctly linked. For example, an NHS number might be present and valid,<sup>3</sup> yet incorrect. Pseudoanonymisation will prevent techniques that overcome identifier errors, such as partial matching on date of birth,<sup>1</sup> and will feedback to providers to prevent it. And if we want to plan for better integration of services across health and social care,<sup>5</sup> we should make best use of patient identifiers, not scramble them and ignore any errors.

Competing interests: GEH-J has an honorary contract with the Health and Social Care Information Centre (HSCIC) as part of a project funded by the Economic and Social Research Council (ESRC) to study data linkage errors. The views stated are his own.

Full response at: [www.bmj.com/content/348/bmj.g1547/rr/689516](http://www.bmj.com/content/348/bmj.g1547/rr/689516).

- 1 HSCIC. Replacement of the HES Patient ID (HESID). Leeds: Health and Social Care Information Centre, 2009.
- 2 Hoeksma J. The NHS's care.data scheme: what are the risks to privacy? *BMJ* 2014;348:g1547. (17 February.)
- 3 Hipisley-Cox J. Validity and completeness of the NHS number in primary and secondary care electronic data in England 1991-2013. University of Nottingham, 2013.
- 4 EMIS National User Group. EMIS NUG proposals for realising the benefits of the GP record. 2014. [www.emisnug.org.uk/article/emis-nug-proposals-realising-benefits-gp-record](http://www.emisnug.org.uk/article/emis-nug-proposals-realising-benefits-gp-record).
- 5 Secretary of State. Health and Social Care Act 2012. Stationery Office, 2010.

Cite this as: *BMJ* 2014;348:g2264

© BMJ Publishing Group Ltd 2014