



# A preprocessing optimization applied to the cell suppression problem in statistical disclosure control

Martin Serpell<sup>a,\*</sup>, Jim Smith<sup>a</sup>, Alistair Clark<sup>a</sup>, Andrea Staggemeier<sup>b</sup>

<sup>a</sup> University of the West of England, Frenchay Campus, Coldharbour Lane, Bristol BS16 1QY, United Kingdom

<sup>b</sup> Head of Center for Statistical and Analytical Intelligence, Office for National Statistics, Newport, United Kingdom

## ARTICLE INFO

### Article history:

Received 28 July 2010

Received in revised form 24 September 2012

Accepted 8 February 2013

Available online 13 March 2013

### Keywords:

Information sharing

Statistical table

Cell suppression problem

Statistical disclosure control

Preprocessing optimization

Mathematical programming

## ABSTRACT

As organizations start to publish the data that they collect, either internally or externally, in the form of statistical tables they need to consider the protection of the confidential information held in those tables. The algorithms used to protect the confidential information in these statistical tables are computationally expensive. However a simple preprocessing optimization applied prior to protection can save time, improve the resultant protection and on occasions enable the use of exact methods where otherwise heuristic methods would have been necessary. The theory behind this preprocessing optimization, how it can be applied and its effectiveness are described in this paper.

© 2013 Elsevier Inc. All rights reserved.

## 1. Introduction

How do we share information in the digital economy? Benefits can be obtained by pooling private data and then analyzing it, however this risks exposing private data. How to protect privacy in such circumstances is an active area of research [21,22]. This is a problem that has been addressed by National Statistics Agencies for some time but is now a problem that needs to be addressed by a wider variety of organizations. Much information that is made public is done so in the form of statistical tables as publishing the source data would break confidentiality. Each cell in the statistical table will typically contain an aggregate of some information, for example it could be the total cost of a given drug dispensed in a given area of a country. The statistical table will also have row and column totals that are referred to as marginals. Like National Statistics Agencies, many organizations are obliged to maintain the confidentiality of the information they hold. If confidentiality is compromised then organizations can lose the co-operation and goodwill of their data contributors and may even be subjected to legal action.

Confidentiality can be compromised if, for example, there is a cell in the published statistical table that is not an aggregate of data from many contributors but from just one or two. To illustrate this point in Figs. 1–3 the cell data is represented in the format  $X_{(y)}$  where  $X$  is the sum of the data provided by  $y$  contributors. In Figs. 1–3 any cell with less than or equal to two contributors has been suppressed to ensure the confidentiality of those contributors. If the information in the cell can be traced back to source then confidentiality has been compromised. Typically in this situation the cell in question is left blank in the published statistical table and this cell is called a primary suppressed cell.

\* Corresponding author.

E-mail address: [Martin2.Serpell@uwe.ac.uk](mailto:Martin2.Serpell@uwe.ac.uk) (M. Serpell).

	1	2	3	4	5	6	Total
A	9 <sub>(1)</sub>	51 <sub>(5)</sub>	41 <sub>(4)</sub>	47 <sub>(5)</sub>	3 <sub>(1)</sub>	48 <sub>(5)</sub>	199 <sub>(21)</sub>
B	8 <sub>(2)</sub>	1 <sub>(1)</sub>	54 <sub>(7)</sub>	44 <sub>(5)</sub>	45 <sub>(2)</sub>	12 <sub>(2)</sub>	164 <sub>(19)</sub>
C	8 <sub>(11)</sub>	70 <sub>(8)</sub>	6 <sub>(2)</sub>	76 <sub>(8)</sub>	64 <sub>(7)</sub>	21 <sub>(2)</sub>	245 <sub>(38)</sub>
D	33 <sub>(6)</sub>	46 <sub>(7)</sub>	45 <sub>(6)</sub>	27 <sub>(6)</sub>	37 <sub>(6)</sub>	60 <sub>(8)</sub>	248 <sub>(39)</sub>
E	87 <sub>(8)</sub>	51 <sub>(6)</sub>	18 <sub>(6)</sub>	35 <sub>(5)</sub>	49 <sub>(7)</sub>	72 <sub>(4)</sub>	312 <sub>(36)</sub>
F	48 <sub>(7)</sub>	59 <sub>(6)</sub>	39 <sub>(5)</sub>	65 <sub>(9)</sub>	35 <sub>(6)</sub>	58 <sub>(9)</sub>	304 <sub>(42)</sub>
Total	193 <sub>(35)</sub>	278 <sub>(33)</sub>	203 <sub>(30)</sub>	294 <sub>(36)</sub>	233 <sub>(29)</sub>	271 <sub>(30)</sub>	1472 <sub>(195)</sub>

**Fig. 1.** An example  $6 \times 6$  statistical table. The eight primary cells have been highlighted with shading. The number of contributors to each cell is given in brackets.

	1	2	3	4	5	6	Total
A	9 <sub>(1)</sub>	51 <sub>(5)</sub>	41 <sub>(4)</sub>	47 <sub>(5)</sub>	3 <sub>(1)</sub>	48 <sub>(5)</sub>	199 <sub>(21)</sub>
B	8 <sub>(2)</sub>	1 <sub>(1)</sub>	54 <sub>(7)</sub>	44 <sub>(5)</sub>	45 <sub>(2)</sub>	12 <sub>(2)</sub>	164 <sub>(19)</sub>
C	8 <sub>(11)</sub>	70 <sub>(8)</sub>	6 <sub>(2)</sub>	76 <sub>(8)</sub>	64 <sub>(7)</sub>	21 <sub>(2)</sub>	245 <sub>(38)</sub>
D	33 <sub>(6)</sub>	46 <sub>(7)</sub>	45 <sub>(6)</sub>	27 <sub>(6)</sub>	37 <sub>(6)</sub>	60 <sub>(8)</sub>	248 <sub>(39)</sub>
E	87 <sub>(8)</sub>	51 <sub>(6)</sub>	18 <sub>(6)</sub>	35 <sub>(5)</sub>	49 <sub>(7)</sub>	72 <sub>(4)</sub>	312 <sub>(36)</sub>
F	48 <sub>(7)</sub>	59 <sub>(6)</sub>	39 <sub>(5)</sub>	65 <sub>(9)</sub>	35 <sub>(6)</sub>	58 <sub>(9)</sub>	304 <sub>(42)</sub>
Total	193 <sub>(35)</sub>	278 <sub>(33)</sub>	203 <sub>(30)</sub>	294 <sub>(36)</sub>	233 <sub>(29)</sub>	271 <sub>(30)</sub>	1472 <sub>(195)</sub>

**Fig. 2.** An example  $6 \times 6$  statistical table with the sequence in which suppressed cells are exposed using an unpicking algorithm. The eight primary cells have been highlighted with shading. The number of contributors to each cell is given in brackets.

	1	2	3	4	Total
A	100 <sub>(1)</sub>	100 <sub>(1)</sub>	100 <sub>(1)</sub>	100 <sub>(4)</sub>	400 <sub>(7)</sub>
B	100 <sub>(4)</sub>	100 <sub>(1)</sub>	100 <sub>(1)</sub>	100 <sub>(4)</sub>	400 <sub>(10)</sub>
C	100 <sub>(1)</sub>	100 <sub>(4)</sub>	100 <sub>(4)</sub>	100 <sub>(1)</sub>	400 <sub>(10)</sub>
D	100 <sub>(1)</sub>	100 <sub>(4)</sub>	100 <sub>(4)</sub>	100 <sub>(1)</sub>	400 <sub>(10)</sub>
Total	400 <sub>(7)</sub>	400 <sub>(10)</sub>	400 <sub>(10)</sub>	400 <sub>(10)</sub>	1600 <sub>(37)</sub>

**Fig. 3.** An example  $4 \times 4$  statistical table, provided by Lawrence Cox, that cannot be unpicked. The nine primary cells have been highlighted with shading. The number of contributors to each cell is given in brackets.

It is the duty of the National Statistics Agency to determine which cells should be protected. For each primary suppressed cell the National Statistics Agency will also determine a range within which the value of the primary suppressed cells must not be able to be calculated. This range is bounded by the lower protection level, *lpl*, and the upper protection level, *upl*. The values assigned to the *lpl* and *upl* are not published. This alone however does not guarantee confidentiality as the cell value may possibly be calculated by subtracting the other cell values in the same row/column from the row/column total. Therefore to guarantee the confidentiality of the information being published in a statistical table it is necessary to suppress other cells in the statistical table. These other suppressed cells are known as secondary suppressed cells. The combination of primary and secondary suppressed cells are known as a suppression pattern. Finding the lowest cost suppression pattern is known as the cell suppression problem in statistical disclosure control, see [23,14] for more details.

The cell suppression problem is a member of the class of NP-hard problems when solving for optimality. In fact, the problem of finding a secondary suppression pattern is easy to be achieved, for example if all cells are suppressed this is a feasible pattern but clearly not optimal. It is when solving the cell suppression problem optimally that as the size of the table to be protected grows the number of possible solutions that need to be evaluated grows much quicker. For a table with  $n$  cells there are  $2^n - 1$  possible suppression patterns. Because of the very large number of constraints that define the cell suppression

Download English Version:

<https://daneshyari.com/en/article/393008>

Download Persian Version:

<https://daneshyari.com/article/393008>

[Daneshyari.com](https://daneshyari.com)