

Information Preserving Regression-based Tools for Statistical Disclosure Control

Øyvind Langsrud

Statistics and Computing, 2019, Volume 29, Issue ?, pp ???-???

This is an author's accepted manuscript version of the article according to the policy of Springer (author's own personal, self-maintained website immediately on acceptance or institutional website/repository 12 months after first publication). This is a post-peer-review, pre-copyedit version of an article published in Statistics and Computing. The final authenticated version is available online at: <http://dx.doi.org/10.1007/s11222-018-9848-9>

An open view-only version is available at <https://rdcu.be/bfeWQ>

Also see the accompanying R package, RegSDC, at <https://cran.r-project.org/package=RegSDC>

Abstract This paper presents a unified framework for regression-based statistical disclosure control for microdata. A basic method, known as information preserving statistical obfuscation (IPSO), produces synthetic data that preserve variances, covariances and fitted values. The data are then generated conditionally according to the multivariate normal distribution. Generalizations of the IPSO method are described in the literature and these methods aim to generate data more similar to the original data. This paper describes these methods in a concise and interpretable way, which is close to efficient implementation. Decomposing the residual data into orthogonal scores and corresponding loadings is an essential part of the framework. Both QR decomposition (Gram Schmidt orthogonalization) and singular value decomposition (principal components) may be used. Within this framework, new and generalized methods are presented. In particular, a method is described by means of which the correlations to the original principal component scores can be controlled exactly. It is shown that a suggested method of random orthogonal matrix masking (ROMM) can be implemented without generating an orthogonal matrix. Generalized methodology for hierarchical categories is presented within the context of microaggregation. Some information can then be preserved at the lowest level and more information at higher levels. The presented methodology is also applicable to tabular data. One possibility is to replace the content of primary and secondary suppressed cells with generated values. It is proposed replacing suppressed

cell frequencies with decimal numbers and it is argued that this can be a useful method.

Keywords microdata · anonymization · synthetic data · microaggregation · hybrid microdata · cell suppression · official statistics

1 Introduction

Microdata are data sets in which each record contains several variables concerning a person or an organization. Usually a microdata data set cannot be made publicly available for reasons of confidentiality. Methods for statistical disclosure control of microdata aim to create protected data that can be released (Hundepool et al. 2012). Two main categories are non-perturbative methods (information reduction, coarsening data) and perturbative methods (changing data). The latter category is often divided into perturbative masking methods and synthetic data generators. In addition, hybrid data generators combine original and synthetic data. Synthetic data are data randomly drawn from a statistical model, often under constraints such that certain statistics or internal relationships of the original data set are preserved.

A broad class of masking methods falls into the category of matrix masking (Duncan and Pearson 1991), which consists of replacing the original data matrix, \mathbf{Y} , with $\mathbf{M}_1 \mathbf{Y} \mathbf{M}_2 + \mathbf{M}_3$. Then, \mathbf{M}_1 is a record-transforming mask, \mathbf{M}_2 is a variable-transforming mask, and \mathbf{M}_3 is a displacing mask. A specific masking method is noise addition where \mathbf{M}_1 and \mathbf{M}_2 are identity matrices and \mathbf{M}_3 contains randomly generated data. As follows below, when the M-matrices are allowed to be stochastically generated, it becomes problematic to dis-

Ø. Langsrud
Statistics Norway, P.O. Box 8131 Dep., 0033 Oslo, Norway
Tel.: +47-21094425
E-mail: Oyvind.Langsrud@ssb.no

tinguish between masking methods and synthetic data generators.

Given a parametric model, making statistical inference from synthetic data the same way as from original data requires that sufficient statistics for the joint distribution of all variables in the data are preserved. This can be achieved by drawing data jointly from the conditional distribution. In general, this is difficult, and data may instead be drawn with fixed values of uncertain parameters. The distributional properties are then violated and variances in particular will be underestimated.

The problems just described are commonly handled by using multiple imputation with fully conditional specifications (Loong and Rubin 2017; Drechsler 2011; Reiter and Raghunathan 2007). Multiple imputation makes it possible to correct for the additional variance from imputation. The use of fully conditional specifications is a way of overcoming the problems of drawing from a multivariate distribution. Then, the assumed joint distribution is specified indirectly. The United States Census Bureau releases synthetic data products prepared by using multiple imputation (Benedetto et al. 2013; Jarmin et al. 2014). The conditional specifications of continuous variables involve normality and linear regression models. These specifications are consistent with the properties of a joint multivariate normal distribution. Sometimes the variables are transformed so that they have an approximately normal distribution.

When all the variables to be synthesized are continuous and assuming multivariate normality and a linear regression model, it is possible to generate synthetic data in a way that exactly preserves the distributional properties. Then there is no need for multiple imputation. Such a basic regression-based synthetic data generator is described by Burridge (2003) as information preserving statistical obfuscation (IPSO). The generated data then constitute a matrix of independent samples from the multivariate normal distribution conditioned on parameter estimates. An equivalent data generation method was described by Langsrud (2005) within the area of Monte Carlo testing. The simulation-based tests were referred to as rotation tests since the generated data can be interpreted as randomly rotated data. This means that the data can be generated as matrix masking with \mathbf{M}_1 as a random rotation matrix and with \mathbf{M}_2 being identity and \mathbf{M}_3 being zero. Strictly speaking, \mathbf{M}_1 is then a uniformly (according to the Haar measure) distributed random orthogonal matrix restricted to preserving fitted values. Given that \mathbf{M}_1 is a restricted orthogonal matrix, this ensures that the sample covariance matrix is preserved. The specific

uniform distribution ensures that original and generated data are independent. Such simulation methodology can be traced back to Wedderburn (1975) in the only intercept case where fitted values are the sample means. In this case Ting et al. (2008) described generalized methodology called random orthogonal matrix masking (ROMM). Their generalization allows \mathbf{M}_1 to be generated in other ways so that the similarity to the original data can be controlled.

Authors of popular software tools for microdata protection emphasize the importance of low information loss or high data utility provided that the disclosure risk is acceptable (Hundepool et al. 2014; Templ et al. 2015). Therefore, it is beneficial to use information preserving methods that generate data more similar to the original data than the IPSO method. In addition to ROMM, such methods are described in Muralidhar and Sarathy (2008) and in Domingo-Ferrer and Gonzalez-Nicolas (2010). In these papers the distributional assumptions are discussed in detail. The IPSO method is made to preserve the sufficient statistics under multivariate normality. Even if the original data are not multivariate normal, the results of common statistical analyses relying on (multivariate) normality will be the same regardless of whether the original or the IPSO generated data are used. The results of analyses not relying on normality will, however, be different. Generating data more similar to the original data also means that the results of any analysis will tend to be more similar to the results based on the original data. The empirical distribution of the original data will also be much better approximated.

The present paper describes several information preserving methods under a common matrix decomposition framework. The aim is twofold. First, to describe existing regression-based methods in a concise and interpretable way, which is close to efficient implementation. Second, to develop new techniques and generalized methodology within the framework. All in all, an important class of tools within the area of statistical disclosure control is described. Disclosure risk and information loss are not investigated in this paper, but both are important when the methodology is applied in practise. We believe that the methodology is well suited to controlling these properties, since the flexibility enables targeted changes.

A requirement that fitted values must be preserved means that we can only change the residual data. In Section 2 the residual data are decomposed into scores and loading using an optional decomposition method (QR or SVD/PCA). Generating data with preserved information can be done by generating new scores and Section 2 formulates an algorithm for the IPSO method.

Section 3 considers the addition of synthetic values to preliminary residuals from arbitrarily predefined data. According to Muralidhar and Sarathy (2008) in particular, it is possible at the variable level to select the degree of similarity to the original data. We will also suggest another and related method that makes it possible to control exactly the correlation between generated and original scores. This method can be viewed as a modified and extended variant of the principal component approach in Calvino (2017).

Section 4 describes an approach in which the IPSO algorithm is modified by using scores from arbitrary residual data. In Section 5 we look into the topic of ROMM and discuss how other methods can be viewed within this context. We show that simulations according to Section 4 can be an efficient and equivalent alternative to doing computations via an orthogonal matrix.

Section 6 considers the generation of microdata within the context of microaggregation. This means that the IPSO method is applied to several clusters/categories as described by Domingo-Ferrer and Gonzalez-Nicolas (2010). New generalized methodology for hierarchical categories is presented. Some information can then be preserved at the lowest level and more information at higher levels.

The methodology is also directly applicable to tabular data obtained by crossing several categorical variables. Statistical disclosure control for tabular data is an important field and several methods exist (Hundepool et al. 2012; Salazar-Gonzalez 2008). In Section 7 we assume that a table suppression method has been applied. It is then possible to replace the content of suppressed (primary and secondary) cells by generated/synthetic values. In particular, it is suggested replacing suppressed cell frequencies with decimal numbers and it is argued that this can be a useful method.

2 Information preserving statistical obfuscation (IPSO)

Assume a multivariate multiple regression model defined by

$$\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{E} \quad (1)$$

where \mathbf{Y} is a $n \times k$ matrix containing n observations of k confidential variables and where the $n \times m$ matrix \mathbf{X} contains corresponding non-confidential variables including a constant term (column of ones). The standard assumptions are that the rows of \mathbf{E} are independent multivariate normal with zero means and a common covariance matrix.

We now want to generate, \mathbf{Y}^* , a synthetic variant of \mathbf{Y} , so that the fitted values and the sample covari-

ance matrix are preserved. Equivalently, we impose the restrictions that $\mathbf{Y}^{*T}\mathbf{Y}^* = \mathbf{Y}^T\mathbf{Y}$ and $\mathbf{X}^T\mathbf{Y}^* = \mathbf{X}^T\mathbf{Y}$. The synthetic data can be drawn from the multivariate regression model conditioned on the restrictions. Burrige (2003) described such simulations as information preserving statistical obfuscation. Langsrud (2005) considered the same problem within the area of Monte Carlo testing and the simulation-based tests were referred to as rotation tests. He described an algorithm based on QR decomposition which is very similar to the algorithm in Burrige (2003). Here we formulate a similar algorithm using a general decomposition:

$$\mathbf{Y} = \hat{\mathbf{Y}} + \hat{\mathbf{E}} = \hat{\mathbf{Y}} + \mathbf{T}\mathbf{W} \quad (2)$$

$$\mathbf{Y}_s = \hat{\mathbf{Y}}_s + \hat{\mathbf{E}}_s = \hat{\mathbf{Y}}_s + \mathbf{T}^*\mathbf{W}_s \quad (3)$$

$$\mathbf{Y}^* = \hat{\mathbf{Y}} + \hat{\mathbf{E}}^* = \hat{\mathbf{Y}} + \mathbf{T}^*\mathbf{W} \quad (4)$$

The data is split into fitted values and residuals. Some method is used to decompose the residuals as a matrix, \mathbf{T} , of orthogonal scores and a matrix, \mathbf{W} , of loadings (2). A simulated data matrix, \mathbf{Y}_s , whose elements are independent standard normal deviates, is decomposed similarly (3) to obtain the matrix \mathbf{T}^* of synthetic scores orthogonal to \mathbf{X} . The synthetic data is obtained by replacing \mathbf{T} by \mathbf{T}^* (4).

The fitted values are computed according to the regression model (1). We allow collinear columns in \mathbf{X} . The fitted values are uniquely defined independently of how this problem is handled. Two obvious candidates for the decomposition method are QR decomposition (Gram Schmidt orthogonalization) and singular value decomposition (SVD) (Strang 1988). In the present paper we refer to a generalized QR decomposition as described in Appendix 1 ($\mathbf{T} = \mathbf{Q}$ and $\mathbf{W} = \mathbf{R}$). This decomposition is unique and dependent columns are allowed. When using SVD as dealt with in Appendix 2, we are essentially performing principal components analysis (PCA) (Jolliffe 2002), but the scores are scaled differently ($\mathbf{T} = \mathbf{U}$ and $\mathbf{W} = \mathbf{\Lambda}\mathbf{V}^T$).

The properties of the synthetic data are independent of the choice of a decomposition method, but this choice is important in some of the modifications below. Note that Burrige (2003) formulated the algorithm using the Choleski decomposition (see Appendix 1). Then, \mathbf{W} is calculated from the covariance matrix estimate and $\hat{\mathbf{E}}^*$ is found as $\hat{\mathbf{E}}_s\mathbf{W}_s^{-1}\mathbf{W}$. In the only intercept case, Mateo-Sanz et al. (2004) formulated another algorithm where \mathbf{W} is calculated by Choleski decomposition and \mathbf{T}^* is generated by an orthogonalization algorithm.

Note that in the special case of a single confidential variable ($k = 1$), then \mathbf{T} is simply the vector of residuals scaled to unit length and \mathbf{W} is the scale factor (scalar). Thus, \mathbf{Y}^* is obtained by replacing the original residuals by residuals obtained from simulated data, scaled

so that the sum of squared residuals is preserved. This yields a simplified procedure compared to the one recently published in Klein and Datta (2018).

A logical question is: Why not simply release the regression coefficients and the covariance matrix? Since \mathbf{X} is released there is no more information in the IPSO-generated data than in these estimates. A good reason for releasing synthetic data is that this is a user-friendly product. Most users would prefer a data set which can be analyzed by ordinary tools. Sophisticated users who would prefer estimates instead can easily calculate them. Anyway, this question is relevant for IPSO, but not for the other methods described below.

3 Arbitrary residual data with a synthetic addition

Instead of drawing synthetic data directly from the model we may start with an arbitrary data set \mathbf{Y}_g and add a synthetic matrix to preserve the required information. However, since we want to preserve the fitted values we will only use the residuals from the arbitrary data set. Then, we generate synthetic data according to

$$\mathbf{Y}_g = \hat{\mathbf{Y}}_g + \hat{\mathbf{E}}_g \quad (5)$$

$$\mathbf{Y}^* = \hat{\mathbf{Y}} + \hat{\mathbf{E}}_g + \mathbf{T}^* \mathbf{C} \quad (6)$$

Here we make use of a synthetic \mathbf{T}^* similar to the one generated in (3). In this case we will generate \mathbf{T}^* in a manner that it is orthogonal to both \mathbf{X} and \mathbf{Y}_g . In practice the QR decomposition of a composite matrix (27) may be used. The matrix \mathbf{C} is calculated from an equation (see below), but a solution may not exist.

With reference to the decomposition of \mathbf{Y} (2) two important special cases of this approach can be expressed as

$$\mathbf{Y}^* = \hat{\mathbf{Y}} + \hat{\mathbf{E}} \mathbf{D} + \mathbf{T}^* \mathbf{C} \quad (7)$$

$$\mathbf{Y}^* = \hat{\mathbf{Y}} + \mathbf{T} \mathbf{D} \mathbf{W} + \mathbf{T}^* \sqrt{\mathbf{I} - \mathbf{D}^2} \mathbf{W} \quad (8)$$

where \mathbf{D} is a diagonal matrix whose i 'th diagonal element ($d_i \in [0, 1]$) controls how much information from the i 'th column of $\hat{\mathbf{E}}$ (7) or \mathbf{T} (8) is re-used. The first method (7) was introduced by Muralidhar and Sarathy (2008) and is described in a different way in their paper. Equation (8) is written without \mathbf{C} since, in that case, the expression for \mathbf{C} is known. The square root of a diagonal matrix means that we take the square root of each diagonal element. Hence, we calculate new unit score vectors as $d_i T_i + \sqrt{1 - d_i^2} T_i^*$ where T_i denotes the i 'th column of \mathbf{T} . With a constant term included in \mathbf{X} , it is easy to show that d_i is a correlation coefficient in both equations. That is, in (7), for each variable, we

control exactly the correlation between the original and the generated residuals. And in (8) we control exactly the correlation between the original score vector and the score vector of the generated data. Note that if all diagonal elements of \mathbf{D} are set to be equal, (7) and (8) become equivalent.

Using (8) in the cases where all d_i 's are either zero or one, we are keeping some score vectors and simulating others. If in addition the decomposition method used in (2) is SVD, the method is very similar to the one described in Calvino (2017), where some PCA score vectors are swapped (randomly permuted). Using swapping instead of the above simulation leads to an approximate instead of an exact preservation of the covariance matrix. Our approach is also more general than the one in Calvino (2017) since we allow regression variables beyond the constant term. Thus, within our framework, PCA means PCA of residuals. In any case, we preserve regression fits and preserving some PCA scores is one way of generating data that is even more similar to the original data.

Working with PCA score vectors can be very useful when there are many highly correlated variables, especially when a few components account for most of the variability in the data. One may want to rotate the components for better interpretability, however. Understanding the components is of particular interest when some variables are considered more sensitive than others. Then, the correlations (d_i 's) can be specified accordingly. Possible rotation of components fits into the general framework of this paper. In fact, going from PCA to QR is a rotation. If the aim is to have good control of the relationship to a few important y-variables, the QR approach would be easier. One may choose the column ordering of the y-variables to correspond to the importance. Then d_1 is the correlation between the original and the generated residuals for the first y-variable, similar to (7). Also note that, when using QR, setting some of the first d_i 's equal to one is equivalent to reclassifying those y-variables as x-variables. In general, the method based on controlling components (8) is more reliable than (7) because it avoids the problem of a possible nonexistent \mathbf{C} .

To discuss the calculation of \mathbf{C} in general we will introduce a scalar parameter α and we re-write equation (6) as

$$\mathbf{Y}^* = \hat{\mathbf{Y}} + \alpha \hat{\mathbf{E}}_g + \mathbf{T}^* \mathbf{C} \quad (9)$$

The matrix \mathbf{C} can now be found from the equation

$$\mathbf{C}^T \mathbf{C} = \hat{\mathbf{E}}^T \hat{\mathbf{E}} - \alpha^2 \hat{\mathbf{E}}_g^T \hat{\mathbf{E}}_g \quad (10)$$

and the Cholesky decomposition may be used to compute \mathbf{C} when possible. An alternative is to find \mathbf{C} via

the eigen decomposition (see Appendix 2). But a solution may not exist since negative eigenvalues of the computed right side can occur (not positive definite).

However, we will present here a way of computing the largest α that makes the equation solvable. By manipulating the characteristic eigenvalue equation it turns out that the limit is found as the square root of the smallest eigenvalue of $\hat{\mathbf{E}}^T \hat{\mathbf{E}} (\hat{\mathbf{E}}_g^T \hat{\mathbf{E}}_g)^{-1}$. A modification is needed when $(\hat{\mathbf{E}}_g^T \hat{\mathbf{E}}_g)$ is not invertible. One possibility is to find the limit as the square root of the inverse of the largest eigenvalue of $(\hat{\mathbf{E}}^T \hat{\mathbf{E}})^{-1} \hat{\mathbf{E}}_g^T \hat{\mathbf{E}}_g$. The calculation of the limit can be useful in combination with the method (7) originally described by Muralidhar and Sarathy (2008). In practice, one may choose the limit if $\alpha = 1$ cannot be chosen. Using $\alpha < 1$ means that the initial residual data are downscaled. In (7) this means that all the required correlations are multiplied by α and that the similarity to the original data is not as high as expected.

4 Using scores from arbitrary residual data

To generate \mathbf{Y}^* in a way that preserves the information without considering its distributional properties, \mathbf{Y}_s (3) can be replaced by anything as long as the residual data have the desired dimension. Necessary scores can then be computed. Using an arbitrary data set \mathbf{Y}_g , synthetic data can be generated by

$$\mathbf{Y}_g = \hat{\mathbf{Y}}_g + \hat{\mathbf{E}}_g = \hat{\mathbf{Y}}_g + \mathbf{T}_g \mathbf{W}_g \quad (11)$$

$$\mathbf{Y}^* = \hat{\mathbf{Y}} + \mathbf{T}_g \mathbf{W} \quad (12)$$

with \mathbf{W} as above (2).

A possible application is when \mathbf{Y}_g is a preliminary data set in some sense. In particular, the methodology can be used in combination with another method for anonymization of microdata. If the result of such a method (\mathbf{Y}_g) preserves the information approximately, the method here can be used as a correction to achieve exactness (\mathbf{Y}^*). In such cases one would expect \mathbf{Y}^* to be close to \mathbf{Y}_g . This is the case when QR is used as the decomposition method in (2), but not necessarily so when SVD is used. The problem is that SVD is unstable when the difference between two singular values is small. Therefore, QR decomposition is preferable.

Another possible application is to combine original and randomly generated data and one may generate \mathbf{Y}_g by

$$\mathbf{Y}_g = \mathbf{Y} \mathbf{D} + \mathbf{Y}_s (\mathbf{I} - \mathbf{D}) \quad (13)$$

using a diagonal matrix, \mathbf{D} , as earlier. To make the parameters in \mathbf{D} interpretable in this case, the variables in the randomly generated data, \mathbf{Y}_s , should be scaled to

have variances similar to those of the original variables. When all diagonal elements are close to one this means that \mathbf{Y}_g is close to \mathbf{Y} . Again, it is important to use QR decomposition to avoid instability.

Since the aim is to generate scores, \mathbf{Y}_g may be created by using original scores instead of original data. This is the case for the ROMM method described below in (17).

5 Random orthogonal matrix masking (ROMM)

Data generated according to ROMM (Ting et al. 2008) can be expressed as

$$\mathbf{Y}^* = \mathbf{M}^* \mathbf{Y} \quad (14)$$

where \mathbf{M}^* is a randomly generated orthogonal matrix. As described by Langsrud (2005) the basic IPSO method of Section 2 can equivalently be formulated as such a ROMM method with

$$\mathbf{M}^* = \mathbf{U}_X \mathbf{U}_X^T + \mathbf{U}_E \mathbf{P}^* \mathbf{U}_E^T \quad (15)$$

where the columns of \mathbf{U}_X form an orthogonal basis for the column space of \mathbf{X} and \mathbf{U}_E form an orthonormal basis for the complement so that $[\mathbf{U}_X \ \mathbf{U}_E]$ is an orthogonal matrix. The matrix \mathbf{P}^* is drawn as a uniformly distributed orthogonal matrix. Ting et al. (2008) stated that all methods that preserve means and the covariance matrix are special cases of ROMM. This also holds beyond the only intercept case. Regardless of the method, as long as the generated data preserve information as above, we can always write $\mathbf{Y}^* = \hat{\mathbf{Y}} + \mathbf{T}^* \mathbf{W} = \mathbf{M}^* \mathbf{Y}$ where

$$\mathbf{M}^* = [\mathbf{U}_X \ \mathbf{T}^* \ \mathbf{U}_A^*] [\mathbf{U}_X \ \mathbf{T} \ \mathbf{U}_A]^T \quad (16)$$

Here, \mathbf{U}_A and \mathbf{U}_A^* are additional orthogonal columns constructed so that the composite matrices are square.

A ROMM approach to generalizing the IPSO method can be expressed by (15) and by letting \mathbf{P}^* be the result of orthogonalizing $\mathbf{I} + \lambda \mathbf{H}$ where \mathbf{H} is a matrix filled with standard normal deviates. Since we are referring to a regression model, this is an extended version of the method (only intercept) originally proposed in Ting et al. (2008). However, in that paper, such an extension was indicated in the appendix. The QR decomposition will be used to orthogonalize $\mathbf{I} + \lambda \mathbf{H}$, and SVD will not work. The identity matrix is obtained when $\lambda = 0$ and the original data are then unchanged. An ordinary uniformly distributed orthogonal matrix is obtained when $\lambda \rightarrow \infty$ and thus we have the IPSO method. Since QR decomposition is sequential, the diagonal elements of \mathbf{P}^* will have a decreasing tendency to approach one (sequentially phenomenon). This

is probably not an intention of the method. This means that the choice of the basis, \mathbf{U}_E , matters. To avoid this phenomenon, a possibility is to generate the basis randomly. To discuss this further we look at the special case without a constant term or any other x-variables. We also assume that \mathbf{T} has full rank (or is temporary extended to have full rank). Now it is possible to use \mathbf{P}^* directly as \mathbf{M}^* so that $\mathbf{T}^* = \mathbf{P}^*\mathbf{T}$. Thus, the above sequentially phenomenon means that the similarity to the original data is highest for the first observation. Another possibility is to make use of (15) without \mathbf{U}_X and with $\mathbf{U}_E = \mathbf{T}$. Then it follows that $\mathbf{T}^* = \mathbf{T}\mathbf{P}^*$ and now the sequentially phenomenon means that the similarity to the original data is highest for the first score vector. Since \mathbf{T} is orthogonal it turns out that we can perform the QR needed to create \mathbf{P}^* from $\mathbf{I} + \lambda\mathbf{H}$ after left multiplication with \mathbf{T} . In addition we have that $\mathbf{T}\mathbf{H}$ is distributed exactly like \mathbf{H} . That is, we can find \mathbf{T}^* by orthogonalizing $\mathbf{T} + \lambda\mathbf{H}$. If we only need the first columns of \mathbf{T}^* (others were included temporarily), we can simplify the method, without changing the result, by only generating as many columns of \mathbf{H} as needed.

The discussion when we have x-variables is similar, but a little trickier. We need to “work” in the correct subspace. One possibility is to generate \mathbf{H} so that it is orthogonal to \mathbf{X} . This can be done by left-multiplying a preliminarily generated matrix by \mathbf{U}_E . An easier and equivalent approach is to generate \mathbf{H} straightforwardly and instead enforce orthogonality to \mathbf{X} afterwards. Thus we obtain a method that is a special case of the one in Section 4 with

$$\mathbf{Y}_g = \mathbf{T} + \lambda\mathbf{H} \quad (17)$$

where \mathbf{H} is a matrix filled with standard normal deviates. This method is more efficient than generating the orthogonal matrix, \mathbf{M}^* . Especially in cases with a large number of observations, ROMM will be very time and memory consuming. Using (11), (12) and (17) now solves this problem.

Note that Ting et al. (2008) have described another ROMM approach that makes use of special block diagonal matrices. We will not discuss this topic in detail. But when IPSO is performed separately within clusters, as described in the section below, this can be interpreted according to (14) where \mathbf{M}^* is a block diagonal matrix where the blocks are independent random orthogonal matrices.

6 Generalized microaggregation

A particular method for statistical disclosure control is microaggregation (Domingo-Ferrer and Mateo-Sanz 2002; Templ et al. 2015). Then, a method is first used to

group the records into clusters. The aggregates within these clusters are released instead of the individual record values. Equivalently, the microdata records can be replaced by averages within the clusters.

The clustered data can also be the starting point for the IPSO method in Section 2 and one may proceed in two ways:

MHa: Run the IPSO procedure separately within each cluster.

MHb: Add dummy variables to \mathbf{X} corresponding to the clustering before running IPSO.

The first method is described by Domingo-Ferrer and Gonzalez-Nicolas (2010) as a hybrid method that combines microaggregation with generation of synthetic data. The notation, MHa, is chosen since this method is known as MicroHybrid. The alternative method introduced here, MHb, is also useful in combination with microaggregation. In both methods sums within clusters, the overall fitted values and the overall sample covariance matrix are preserved. In MHa the fits and the covariance matrix estimates within each cluster are also preserved. An extended variant of MHb is to cross all the original x-variables with the clusters. Then all the fits will be preserved within clusters, but not the covariance matrix estimates.

To discuss this more generally we will present an algorithm in which some regression variables are crossed with clusters and some are not. For this purpose we assume two data matrices of x-variables, \mathbf{X}_A and \mathbf{X}_B . The matrix \mathbf{X}_A consists of all the variables to be crossed with clusters (including the intercept). As before, we let \mathbf{X} denote the full regression model matrix. We will now extend and partition this matrix as

$$\mathbf{X}_{ext} = [\mathbf{X} \ \mathbf{X}_3] = [\mathbf{X}_1 \ \mathbf{X}_2 \ \mathbf{X}_3] \quad (18)$$

- \mathbf{X}_1 is the regression variable matrix obtained by crossing \mathbf{X}_A with clusters.
- \mathbf{X}_2 contains \mathbf{X}_B adjusted for (made orthogonal to) \mathbf{X}_1 .
- \mathbf{X}_3 contains the regression variable matrix obtained by crossing \mathbf{X}_B with clusters followed by adjustment for \mathbf{X}_1 and \mathbf{X}_2 .

Fitted values obtained by regressing \mathbf{Y} onto \mathbf{X} and the corresponding residuals are two orthogonal parts of \mathbf{Y} . According to (18) we also divide each of these parts into two orthogonal parts.

$$\mathbf{Y} = \hat{\mathbf{Y}} + \hat{\mathbf{E}} = (\hat{\mathbf{Y}}_1 + \hat{\mathbf{Y}}_2) + (\hat{\mathbf{E}}_3 + \hat{\mathbf{E}}_4) \quad (19)$$

Here $\hat{\mathbf{Y}}_1$, $\hat{\mathbf{Y}}_2$ and $\hat{\mathbf{E}}_3$ are the regression fits obtained by using \mathbf{X}_1 , \mathbf{X}_2 and \mathbf{X}_3 , respectively.

Performing computations via \mathbf{X} will be very inefficient when we have several variables, observations and

clusters. We also want to (partly) preserve the covariance matrix within clusters. Therefore some of the computations will be performed within clusters and in the following we will call this the local level. As opposed to *globally*, *locally* means that we perform the computations by looping through all the clusters. In particular, we can compute the fitted values \hat{Y}_1 locally and then we can use the original matrix X_A and will not need X_1 . Above, X_2 was made orthogonal to X_1 . In this case it does not matter whether we do the orthogonalization globally or locally. In the algorithm below we perform this computation locally using X_A instead of X_1 . As mentioned above, the regression fits obtained by using X_2 globally are contained in \hat{Y}_2 . If we instead use X_2 to compute regression fits locally, the result is $\hat{Y}_{2ext} = \hat{Y}_2 + \hat{E}_3$. This means that we have sketched how to calculate the four parts in (19). In order to compute synthetic data this will be combined with computations on simulated data. We can use the following algorithm for all the computations:

1. Globally simulate a data matrix Y_s with the same dimension as Y .
2. Locally calculate \hat{Y}_1 , \hat{Y}_{s1} and \hat{X}_B by regressing Y , Y_s and X_B onto X_A .
3. Locally calculate $X_2 = X_B - \hat{X}_B$.
4. Locally calculate \hat{Y}_{2ext} and \hat{Y}_{s2ext} by regressing Y and Y_s onto X_2 .
5. Locally calculate $\hat{E}_4 = Y - \hat{Y}_1 - \hat{Y}_{2ext}$ and $\hat{E}_{s4} = Y_s - \hat{Y}_{s1} - \hat{Y}_{s2ext}$.
6. Locally decompose $\hat{E}_4 = T_4 W_4$ and $\hat{E}_{s4} = T_{s4} W_{s4}$ and calculate $\hat{E}_4^* = T_{s4} W_4$.
7. Globally calculate \hat{Y}_2 and \hat{Y}_{s2} by regressing Y and Y_s onto X_2 .
8. Globally calculate $\hat{E}_3 = \hat{Y}_{2ext} - \hat{Y}_2$ and $\hat{E}_{s3} = \hat{Y}_{s2ext} - \hat{Y}_{s2}$.
9. Globally decompose $\hat{E}_3 = T_3 W_3$ and $\hat{E}_{s3} = T_{s3} W_{s3}$ and calculate $\hat{E}_3^* = T_{s3} W_3$.
10. Globally calculate $Y^* = \hat{Y}_1 + \hat{Y}_2 + \hat{E}_3^* + \hat{E}_4^*$.

In addition to this main algorithm we will discuss five possible modifications.

- a) Drop item 6 and replace item 9 by: Locally decompose $\hat{E}_3 + \hat{E}_4 = TW$ and let $\hat{E}_4^* = T_{s4} W$ and set $\hat{E}_3^* = 0$.
- b) Drop item 6 and instead, after item 9: Find locally $\hat{E}_4^* = T_{s4} C$ where the matrix C satisfies $C^T C = \hat{E}_4^T \hat{E}_4 + \hat{E}_3^T \hat{E}_3 - \hat{E}_3^{*T} \hat{E}_3^*$. This is similar to (6).
- c) This is a generalization with a) and b) as special cases. Drop item 6 and modify item 9 by $\hat{E}_3^* = \alpha T_{s3} W_3$ and find \hat{E}_4^* as in b). This is similar to (9).
- d) If we have groups of clusters, we may want some of the variables in X_B to be crossed with these groups.

At the beginning of the algorithm we then replace X_B with a model matrix that contains such crossing.

- e) If the clusters are partitioned into smaller pieces, we may want some of the variables in X_A (maybe only the constant term) to be crossed with these pieces. Before item 2 in the algorithm we will then locally replace X_A with a model matrix that contains crossing with pieces that are present in the actual cluster.

The main algorithm preserves the fitted values (or regression parameters). However, the covariance matrix is partitioned into two terms, $\hat{E}_3^T \hat{E}_3 + \hat{E}_4^T \hat{E}_4$. The last term is preserved at the local level, but the first term is only preserved at the global level. In the special case where X_B is empty, the main algorithm simplifies to MHa and \hat{E}_3 vanishes.

When X_B is not empty, we can use modification a), b) or c) to preserve the covariance matrix at the local level. In a) we have a side effect regarding the parameter estimates corresponding to X_B . They are preserved globally as required, but within each cluster these parameter estimates have become identical to the global estimate. This may or may not be a required property. Modification b) avoids this property. With many variables in X_B and/or few observations in each cluster, modification a) or b) may not be possible. An extra problem in b) is that the matrix C may not be found. In cases where a) is possible but not b) one may choose c). Within each cluster it is possible, using the method in Section 3, to compute the largest α that makes it possible to find C . The final α can be chosen as the smallest of all these values.

Using modification d) we are only changing the input. This method can be a compromise if one is uncertain whether a variable should belong to X_A or X_B .

Modification e) is especially useful if we want to preserve means in groups that are too small to be used as clusters. In practice modification e) may also be used in the special case of a single cluster. Then X_B is not needed and the method simplifies to MHb. The covariance matrix is only preserved at the top level.

Modifications d) and e) can be combined and they can also be combined with a), b) or c). All in all, the main algorithm with the possibility of the modifications is a flexible approach for generating hybrid microdata within the context of microaggregation.

We will now exemplify the methodology by discussing scenarios where the general methodology may be used in different ways, including the five modifications. We assume business data where Y consists of sensitive continuous economic variables such as employment expenses and taxes. We have several categorical

non-sensitive variables such as region, industry group and employment size classes. We will refer to the latter categorical variable as K_{ESC} . The data can be subjected to microaggregation with clusters created from the non-sensitive variables. We assume that this is done in such a way that K_{ESC} is unique within each cluster. Synthetic data can be generated by the IPSO method in two ways, MHa and MHb, as described at the beginning of this section. In both cases, \mathbf{X}_B is empty. When MHa is used, \mathbf{X}_A contains the intercept. The other variant, MHb, means that input to the method is a single cluster and that \mathbf{X}_A contains dummy variables corresponding to the microaggregation clustering. The choice between MHa and MHb is a question of whether one wants to reproduce the cluster specific covariances or not. If the variability in some clusters is very low, the synthetic data will be close to the original data in those clusters when MHa is used. When MHa is used, a very high correlation between some sensitive variables in some clusters will also be reproduced. To avoid this property, method MHb may be chosen instead.

Now assume that the exact number of employees, N_{ESC} , is a non-sensitive variable. Then, one objective might be to preserve the correlations (or equivalently regression fits) between N_{ESC} and the sensitive variables. This could be done by including N_{ESC} in \mathbf{X}_A even though K_{ESC} is already a variable used within the microaggregation clustering. This way of doing it is still a possibility within MHa since IPSO, in general, can include x-variables (beyond the intercept). This means that the correlations to N_{ESC} will be preserved within each cluster. If one company has a value of N_{ESC} very different from the others in a cluster, the sensitive values of this company will influence the correlations a lot and another method may be preferred. One possibility is to include N_{ESC} in \mathbf{X}_B instead. The correlations between N_{ESC} and the sensitive variables will be preserved, but not within each cluster.

A possibility between including N_{ESC} in \mathbf{X}_A and including N_{ESC} in \mathbf{X}_B is to preserve the correlations within each employment size class (each K_{ESC} category). Since the employment size classes are groups of clusters, this mean that we use modification d). A related method is to use the employment size classes as the clusters within the method and include N_{ESC} in \mathbf{X}_A . Again the correlations are preserved within the employment size classes. This time covariances of the sensitive variables are preserved within the employment size classes, but not within the original microaggregation clustering. We can still preserve the means within the microaggregation clusters by using modification e). This means that the original microaggregation clusters are cluster pieces within the description of e).

Modifications a), b) and c) are relevant when \mathbf{X}_B is in use. This was the case in the example above when N_{ESC} was included directly in \mathbf{X}_B and when modification d) was in use. In these cases the standard method does not preserve covariances exactly within each cluster (as required), but this can be achieved by means of a), b) or c). In the case of a) this means that the relations, in terms of regression coefficients, between N_{ESC} and the sensitive variables are made identical in each cluster. The data may seem strange, because identical regression coefficients will never occur in real data. Therefore, the method is probably not preferable and one would prefer b). Since relations are preserved at the global level, there is some information about this in each cluster. When b) is used one can imagine that someone may misunderstand this and may use the data to analyse how relationships differ between clusters. Method a) will prevent someone publishing misleading results about differences. Although, method b) is the preferable choice. But method b) may not be possible and c) is chosen instead.

At the end of this section, we note that an alternative approach is to set all the residuals to zero. If we consider only intercepts in \mathbf{X}_A and \mathbf{X}_B we obtain ordinary microaggregation. The original values are replaced by means within clusters. Beyond the only intercept cases, we obtain a form of generalized microaggregation which, in some cases, can be a useful alternative to ordinary microaggregation. Based on the above exemplification one could imagine including N_{ESC} in \mathbf{X}_A or \mathbf{X}_B . One reason could be as follows. Suppose many users always divide all numbers by N_{ESC} to obtain per employee data. From their perspective it would be better to divide all data by N_{ESC} before microaggregation so that means per employee are preserved within clusters. Producing two data sets is not an alternative, but including N_{ESC} in the method may be a solution. Means per employee will be closer to the truth at the same time as ordinary means are preserved.

7 Suppressed tabular data

It is possible to combine the IPSO methodology in Section 2 with statistical disclosure control for tabular data. As described in Section 6, \mathbf{X} may be composed of dummy variables corresponding to a categorical variable. Sums within the categories will then be preserved. Instead of just a single categorical variable we may have several variables, which can be used to tabulate the data in various ways. When playing around with which dummy variables are to be included in \mathbf{X} , one is playing around with which sums are to be preserved. Thus, we can include in \mathbf{X} only those variables that correspond

Table 1 The example frequency table where cells to be suppressed are marked with *.

	col1	col2	col3	col4	Total
row1	3*	11*	32	30	76
row2	1*	9*	13*	8	31
row3	12	22	2*	2*	38
row4	18*	19	16	3*	56
Total	34	61	63	43	201

to cells that are found to be safe by means of a cell suppression method.

Instead of using microdata as a starting point we can use the sums obtained by crossing all the categorical variables (cover table). This often means that most of the input data are preserved and only some new data are generated (the suppressed cells). In the following, we apply the methodology to cell frequencies as a y -variable. Even if all the original values are counts, we will generate synthetic values that are not whole numbers. As will be discussed, this can be considered as a nice property.

An example of a frequency table is given in Table 1. We assume that values below 4 cannot be published. This means that five cells in the table are primarily suppressed. To prevent the possibility of these values being calculated from other cell values, four additional cells are secondarily suppressed. These additional cells are found using cell suppression methodology. Here we have 16 inner cell frequencies (totals excluded) and 16 publishable cell frequencies (suppressed cells excluded).

Generally, we let \mathbf{Y} be $n \times k$ consisting of all the n inner multivariate cell elements in a cover table. Furthermore we let \mathbf{Z} be $m \times k$ consisting all the m multivariate elements of the publishable cells.

We focus in particular on the univariate special case ($k = 1$) where the values are frequencies. In Table 1 we have $n = m = 16$.

Obviously \mathbf{Z} can be calculated from \mathbf{Y} and this can be done via a $n \times m$ dummy matrix \mathbf{X} :

$$\mathbf{Z} = \mathbf{X}^T \mathbf{Y} \quad (20)$$

We have one column in \mathbf{X} for each publishable cell. Each publishable cell is either an inner cell or a sum of several inner cells. In the first case this means that the corresponding column of \mathbf{X} has only one element that is one (others are zero). By regressing \mathbf{Y} onto \mathbf{X} we can calculate fitted values by

$$\hat{\mathbf{Y}} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^\dagger \mathbf{X}^T \mathbf{Y} = \mathbf{X}^\dagger^T \mathbf{X}^T \mathbf{Y} = \mathbf{X}^\dagger^T \mathbf{Z} \quad (21)$$

Since the columns of \mathbf{X} are collinear we use the Moore-Penrose generalized inverse here (see Appendix 2). It is clear that the fitted values can be calculated directly

Table 2 The suppressed values of the example frequency table and corresponding generated values.

(row,col)	\mathbf{Y}	\mathbf{Y}^*	$\hat{\mathbf{Y}}$	$\hat{\mathbf{Y}}_{\text{mod}4}$	$\hat{\mathbf{Y}}_{\text{mod}10}$
(2,1)	1	11.1562	4.5217	0.8696	4.0870
(3,3)	2	9.3549	6.6957	1.8261	2.7826
(3,4)	2	-5.3549	-2.6957	2.1739	1.2174
(1,1)	3	0.1986	4.1739	2.9565	0.6957
(4,4)	3	10.3549	7.6957	2.8261	3.7826
(2,2)	9	6.1986	10.1739	8.9565	6.6957
(1,2)	11	13.8014	9.8261	11.0435	13.3043
(2,3)	13	5.6451	8.3043	13.1739	12.2174
(4,1)	18	10.6451	13.3043	18.1739	17.2174

from the published cell values. We also have that \mathbf{Z} is preserved in the sense that $\mathbf{Z} = \mathbf{X}^T \hat{\mathbf{Y}}$. Many elements of \mathbf{Y} are also elements of \mathbf{Z} and they are reproduced exactly in $\hat{\mathbf{Y}}$. Elements of \mathbf{Y} that are suppressed are replaced in $\hat{\mathbf{Y}}$ by estimates computed from the publishable cells. Estimates of possible suppressed cells not included in \mathbf{Y} can be calculated from $\hat{\mathbf{Y}}$. To be more specific, we extend (20) to all elements and not only the publishable ones and let $\mathbf{Z}_{\text{All}} = \mathbf{X}_{\text{All}}^T \mathbf{Y}$. Now we can compute corresponding estimates by

$$\hat{\mathbf{Z}}_{\text{All}} = \mathbf{X}_{\text{All}}^T \hat{\mathbf{Y}} \quad (22)$$

Even though the starting point is different from ordinary regression, synthetic data can be generated by the same method. The IPSO method in Section 2 preserves $\mathbf{X}^T \mathbf{Y}$ and $\mathbf{Y}^T \mathbf{Y}$. We can generate synthetic residuals, $\hat{\mathbf{E}}^*$, to obtain $\mathbf{Y}^* = \hat{\mathbf{Y}} + \hat{\mathbf{E}}^*$. Then the synthetic variant of \mathbf{Z}_{All} is

$$\mathbf{Z}_{\text{All}}^* = \mathbf{X}_{\text{All}}^T \mathbf{Y}^* \quad (23)$$

It is worth noting that

- The method can be simplified by working with reduced versions of \mathbf{Y} and \mathbf{X} . Rows corresponding to publishable cells can be removed and then redundant columns of \mathbf{X} can also be removed.
- If preserving the covariance matrix is not a requirement, the fitted values ($\hat{\mathbf{Y}}$ and $\hat{\mathbf{Z}}_{\text{All}}$) can be an alternative to synthetic values.

Now we will look at the example in Table 1. All the suppressed cells are inner cells and therefore, when considering generated values, we can refer to $\hat{\mathbf{Y}}$ and \mathbf{Y}^* . These values are shown in Table 2. The underlying $\hat{\mathbf{E}}^*$ was generated in the ordinary way.

When variance is not important, $\hat{\mathbf{Y}}$, may be used. However, these values are quite far from the true values. One may wish for values that are closer to the truth and still safe. One possibility is to do the calculations after a modulo operation. With 10 as the divisor, this means that we only consider the last digit of the suppressed values as unsafe. In our example, this means that we

Table 3 The example frequency table with decimal numbers in place of suppressed values.

	col1	col2	col3	col4	Total
row1	0.4449	13.5551	32	30	76
row2	4.2102	6.4449	12.3449	8	31
row3	12	22	2.6551	1.3449	38
row4	17.3449	19	16	3.6551	56
Total	34	61	63	43	201

replace the largest unsafe values, 11, 13 and 18, by 1, 3 and 8. After the calculations, 10 is added back to these cells. The fitted values obtained in this way are presented in the last column of Table 2. Since the starting point was that all values below 4 were unsafe, one may try 4 as the divisor. However, in this case, this results in model fits that are very close to the truth. In fact, the true value is the closest integer in all cases.

Note that if one combines the modulo operation with ordinary generation of synthetic values, the variance of the original values will not be preserved. A solution is to scale the residuals by a constant to achieve the required variance. Even if correct variance is not important, synthetic residuals may still be added. One reason is to increase the differences from the truth. Another reason is to ensure that none of the generated values are whole numbers. Then, a linear combination of synthetic values using integer coefficients cannot be a whole number unless it is possible to rewrite the combination as a combination of publishable values (whole numbers). Have in mind that the columns \mathbf{X}_{All} are linearly dependent so that linear combinations can be written in several ways.

Table 3 presents the results from a method that may be used in practice. The modulo 10 method is used and the ordinary residuals are downscaled by a factor of 10. In general, if we replace the original suppressed values by values generated in a manner that includes synthetic residuals, we have the following characteristics.

1. The values of all the publishable cells are unchanged. In particular, the values add up correctly to all totals and subtotals. In practice, however, we need to take into account numerical precision error.
2. The values of all the suppressed cells, including possible subtotals, will not be whole numbers. Obtaining a value as close to a whole number as the level of the numerical precision error is very unlikely.
3. If a suppressed cell value is a whole number, this means that the suppression algorithm has failed. This can be verified by re-running the random generation part.

4. Assuming a successful suppression algorithm, any sum of values resulting in a whole number is a true and publishable value.
5. Negative values may be generated.
6. After cell suppression and generation of decimal numbers, only the inner cells may be stored. Various totals can be calculated when needed.
7. The methodology can be used when several variables are involved (not only cell frequency). The ordinary method (no modulo and ordinary residuals) may then be preferred.

As mentioned above, one reason to generate decimal numbers may be to generate values to be stored. They may never be shown and only used technically. After computation of any requested totals, the whole numbers may be shown and other numbers hidden. In practice, a numerical precision limit is needed to define whole numbers. Another reason to generate decimal numbers is to give more information to users. Without using the modulo method there is no extra information in fitted values than is already hidden in the published cells. Methods that assume positive integers and produce lower and upper bounds will provide more precise information about what is hidden in the data. However, presenting intervals is more complicated and precise information may not be the goal. Appropriate information may be achieved by combining the modulo method and scaling of residuals. A third reason to use the approach presented here may be to control the result of the suppression method. Suppressing linked tables can be especially challenging (de Wolf and Giessing 2009), and not all algorithms guarantee protection. When working with linked tables, we need to include all cells from all tables in \mathbf{Z}_{All} . A fourth reason to calculate decimal numbers is to permit calculation of arbitrary (user defined) sums of cells afterwards. A sum not included in the suppression method may be publishable even if suppressed cells are involved.

The fitted values contain only information available from publishable cells. When variance is preserved, a piece of information is provided about the suppressed cells and this may be a reason to scale the residuals. One problematic special case to be aware of is when the rank of \mathbf{X} is $n - 1$. Normally, the generated score vectors of unit length underlying $\hat{\mathbf{E}}^*$ are randomly orientated in the subspace orthogonal to the column space of \mathbf{X} . But in this special case, the dimension of this subspace is one. Thus, given preserved variance, $\pm\hat{\mathbf{E}}$ are the only two possible instances of $\hat{\mathbf{E}}^*$. This is typically a problem in small example data sets and not in applications. In the example here, the subspace dimension is two and this problem is avoided.

To summarise, without using the modulo method and using residuals scaled in a manner unknown to the user, replacing ordinary suppression (missing value) with decimal numbers releases no extra information. Decimal numbers may, however, be a user-friendly alternative. Switching to modulo should be done with care as the method has not been thoroughly studied.

8 Concluding remarks

This paper has focused on information preserving regression-based methods, which belong to the area of synthetic and hybrid data. Such an approach to protecting microdata is sometimes preferred, but in many applications other methods are used. In the area of statistical disclosure control, insight into the described methodology is useful regardless. Innovative tools are often created by combining aspects of several methods. In the present paper, combined methods have been described in Sections 6 and 7. In both cases, another method is utilized before data are generated by the synthetic data methodology. In the former case, clusters according to microaggregation are made and in the latter case table suppression is performed.

Since a requirement for all the methods treated in this paper is to preserve fitted values, these methods are mostly about how to generate residuals. In some cases, an alternative is to set these residuals at zero. Then we are replacing all the original data with deterministic imputations (in the sense opposed to random imputations). We discussed this approach in Sections 6 and 7.

Regression-based methods have been criticized for not being able to deal with data with outliers in a satisfactory way (Templ and Meindl 2008). This criticism can be met by combining methods. As with other techniques (Templ and Meindl 2008), one can start by dividing the data into outlying and non-outlying observations. Thereafter, one proceeds by generating data within each group. One may use methodology described in this paper in either group, or only in the non-outlying observations group. In the case of both groups, we can say that such methodology is already covered by the general description in Section 6.

The contribution of this paper has mainly been to describe and develop methods in a concise and interpretable way. Several questions arise concerning the performance of the methods and which method to prefer in specific situations. A wide-ranging comparison of the methods is beyond the scope of the present paper, but studies and performance comparisons can be found in the underlying references. The characteristics of the

methods are different and the preferred choice would depend on the situation.

However, if we limit the discussion to non-combined methods (Sections 2-5), we can make some reflections. The IPSO method is the basic method to be used when the aim is only to preserve variances, covariances and fitted values. Technically, performing the computations via QR decomposition is faster than using SVD. Methodology that controls the relationship with the original data by means of a single parameter is easy to understand. The IPSO method and the original data are thus the two extremes. One such method is obtained with all diagonal elements equal in (7) or (8). These two equations are then equivalent. This method is elegant and exact control is obtained. Furthermore, this method is a special case of the one described in Muralidhar and Sarathy (2008). An easy and efficient implementation is obtained by (8) and QR decomposition. It is also easy adjustable, so that the relationship with some x-variables can be controlled better. However, the methodology in Section 3 makes use of new scores generated under certain orthogonality restrictions. In cases with many variables and few observations this may be impossible. The single parameter ROMM method may then be an alternative. This method is efficiently implemented by (11), (12) and (17) which utilize the theory in Sections 4 and 5.

Appendix 1: Generalized QR decomposition

The QR decomposition of a $n \times m$ matrix \mathbf{A} with rank r can be written as

$$\mathbf{A} = \mathbf{Q}\mathbf{R} \quad (24)$$

where \mathbf{Q} is a $n \times r$ matrix whose columns form an orthonormal basis for the column space of \mathbf{A} . This decomposition can be viewed as the matrix formulation of the Gram-Schmidt orthogonalization process. The Cholesky decomposition of $\mathbf{A}^T\mathbf{A}$ can be read from the QR decomposition of \mathbf{A} as $\mathbf{R}^T\mathbf{R}$.

In this paper, in order to allow linearly dependent columns of \mathbf{A} ($r < m$), we refer to a generalized variant of QR decomposition. In such cases a usual decomposition (Chan 1987) is

$$\mathbf{A}\tilde{\mathbf{P}} = \mathbf{Q}\tilde{\mathbf{R}} \quad (25)$$

where $\tilde{\mathbf{P}}$ is a permutation matrix that reorders the columns (pivoting) in order to make a decomposition so that $\tilde{\mathbf{R}}$ is upper triangular.

To make the decomposition unique, we require the diagonal entries of $\tilde{\mathbf{R}}$ to be positive. Furthermore, we require $\tilde{\mathbf{P}}$ to keep the order of the columns as close to the original order as possible (minimal pivoting). We

now have $\mathbf{A} = \mathbf{Q}\tilde{\mathbf{R}}\tilde{\mathbf{P}}^T$ and in generalized QR decomposition (24) we use

$$\mathbf{R} = \tilde{\mathbf{R}}\tilde{\mathbf{P}}^T \quad (26)$$

The QR decomposition of a composite matrix can be written as

$$[\mathbf{A}_1 \ \mathbf{A}_2] = [\mathbf{Q}_1 \ \mathbf{Q}_2] \begin{bmatrix} \mathbf{R}_1^T & \mathbf{R}_2^T \end{bmatrix}^T \quad (27)$$

Now \mathbf{Q}_1 can be computed by QR decomposition of \mathbf{A}_1 . The matrix \mathbf{Q}_2 can be computed by QR decomposition of $\mathbf{A}_2 - \mathbf{Q}_1\mathbf{Q}_1^T\mathbf{A}_2$, which is the residual part after regressing \mathbf{A}_2 onto \mathbf{A}_1 .

Appendix 2: The singular value decomposition

The singular value decomposition (SVD) of a $n \times m$ matrix \mathbf{A} with rank r can be written as

$$\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T \quad (28)$$

where $\mathbf{\Lambda}$ is a $r \times r$ diagonal matrix of strictly positive singular values in descending order. This is the rank-revealing version of the decomposition (Demmel et al. 1999). Other variants of SVD allow some singular values to be zero, but these can be omitted. The columns of \mathbf{U} form an orthonormal basis for the column space of \mathbf{A} and the columns of \mathbf{V} form an orthonormal basis for the row space.

The singular values are the square root of the eigenvalues of $\mathbf{A}^T\mathbf{A}$ and $\mathbf{A}\mathbf{A}^T$. The eigen decompositions of these two symmetric matrices can be read directly from the SVD of \mathbf{A} as $\mathbf{V}\mathbf{\Lambda}^2\mathbf{V}^T$ and $\mathbf{U}\mathbf{\Lambda}^2\mathbf{U}^T$. It is also worth mentioning that an alternative to the ordinary Cholesky decomposition, $\mathbf{A}^T\mathbf{A} = \mathbf{R}^T\mathbf{R}$, is to let $\mathbf{\Lambda}\mathbf{V}^T$ play the role of \mathbf{R} .

To make the SVD unique we can require all column sums of \mathbf{V} to be positive. In cases with equal singular values, the decomposition is not unique regardless.

There is a close relationship between SVD and PCA. In PCA, the variables are usually centered to zero means, and in many cases standardized to equal variances prior to decomposition. If \mathbf{A} is such a centered/standardized matrix then $\mathbf{U}\mathbf{\Lambda}$ is the matrix of PCA scores and \mathbf{V} is the matrix of PCA loadings.

The Moore-Penrose generalized inverse of \mathbf{A} can be written as

$$\mathbf{A}^\dagger = \mathbf{V}\mathbf{\Lambda}^{-1}\mathbf{U}^T \quad (29)$$

We have

$$\mathbf{A}^\dagger = (\mathbf{A}^T\mathbf{A})^\dagger\mathbf{A}^T = \mathbf{A}^T(\mathbf{A}\mathbf{A}^T)^\dagger \quad (30)$$

When \mathbf{A} is invertible, $\mathbf{A}^\dagger = \mathbf{A}^{-1}$. When $\mathbf{A}^T\mathbf{A}$ or $\mathbf{A}\mathbf{A}^T$ is invertible this means, respectively, that $\mathbf{A}^\dagger = (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T$ or $\mathbf{A}^\dagger = \mathbf{A}^T(\mathbf{A}\mathbf{A}^T)^{-1}$.

References

- Benedetto, G., Stinson, M.H., Abowd, J.M.: The Creation and Use of the SIPP Synthetic Beta. Tech. rep., United States Census Bureau (2013)
- Burridge, J.: Information preserving statistical obfuscation. *Stat. Comput.* **13**(4), 321–327 (2003). doi:10.1023/A:1025658621216
- Calvino, A.: A Simple Method for Limiting Disclosure in Continuous Microdata Based on Principal Component Analysis. *J. Off. Stat.* **33**(1), 15–41 (2017). doi:10.1515/JOS-2017-0002
- Chan, T.F.: Rank revealing QR factorizations. *Linear Alg. Appl.* **88-9**, 67–82 (1987). doi:10.1016/0024-3795(87)90103-0
- de Wolf, P.P., Giessing, S.: Adjusting the tau-ARGUS modular approach to deal with linked tables. *Data Knowl. Eng.* **68**(11), 1160–1174 (2009). doi:10.1016/j.datak.2009.06.005
- Demmel, J., Gu, M., Eisenstat, S., Slapnicar, I., Veselic, K., Drmac, Z.: Computing the singular value decomposition with high relative accuracy. *Linear Alg. Appl.* **299**(1-3), 21–80 (1999). doi:10.1016/S0024-3795(99)00134-2
- Domingo-Ferrer, J., Gonzalez-Nicolas, U.: Hybrid microdata using microaggregation. *Inf. Sci.* **180**(15), 2834–2844 (2010). doi:10.1016/j.ins.2010.04.005
- Domingo-Ferrer, J., Mateo-Sanz, J.M.: Practical data-oriented microaggregation for statistical disclosure control. *IEEE Transactions on Knowledge and Data Engineering* **14**(1), 189–201 (2002)
- Drechsler, J.: *Synthetic Datasets for Statistical Disclosure Control*. Springer, New York (2011)
- Duncan, G.T., Pearson, R.W.: Enhancing Access to Microdata While Protecting Confidentiality: Prospects for the Future. *Stat. Sci.* **6**(3), 219–239 (1991)
- Hundepool, A., de Wolf, P.P., Bakker, J., Reedyjk, A., Francini, L., Polettini, S., Capobianchi, A., Domingo, J.: *mu-ARGUS user's manual*, version 5.1. Tech. rep., Statistics Netherlands (2014)
- Hundepool, A., Domingo-Ferrer, J., Francini, L., Giessing, S., Nordholt, E.S., Spicer, K., de Wolf, P.P.: *Statistical Disclosure Control*. John Wiley & Sons, Ltd (2012). doi:10.1002/9781118348239.ch1
- Jarmin, R.S., Louis, T.A., Miranda, J.: Expanding the role of synthetic data at the U.S. Census Bureau. *Statistical Journal of the IAOS* **30**(1-3), 117–121 (2014)
- Jolliffe, I.: *Principal Component Analysis*, Second Edition. Springer, New York (2002)
- Klein, M.D., Datta, G.S.: Statistical disclosure control via sufficiency under the multiple linear regression model. *Journal of Statistical Theory and Practice* **12**(1), 100–110 (2018)
- Langsrud, Ø.: Rotation tests. *Stat. Comput.* **15**(1), 53–60 (2005). doi:10.1007/s11222-005-4789-5
- Loong, B., Rubin, D.B.: Multiply-Imputed Synthetic Data: Advice to the Imputer. *J. Off. Stat.* **33**(4), 1005–1019 (2017). doi:10.1515/JOS-2017-0047
- Mateo-Sanz, J., Martinez-Balleste, A., Domingo-Ferrer, J.: Fast generation of accurate synthetic microdata. In: Domingo-Ferrer, J and Torra, V (ed.) *Privacy in Statistical Databases*, Proceedings, vol. 3050, pp. 298–306 (2004). Conference on Privacy in Statistical DataBases (PSD 2004), Barcelona, Spain, June 09-11, 2004.
- Muralidhar, K., Sarathy, R.: Generating Sufficiency-based Non-synthetic Perturbed Data. *Trans. Data Priv.* **1**(1), 17–33 (2008)

- Reiter, J.P., Raghunathan, T.E.: The multiple adaptations of multiple imputation. *J. Am. Stat. Assoc.* **102**(480), 1462–1471 (2007). doi:10.1198/016214507000000932
- Salazar-Gonzalez, J.J.: Statistical confidentiality: Optimization techniques to protect tables. *Comput. Oper. Res.* **35**(5), 1638–1651 (2008). doi:10.1016/j.cor.2006.09.007
- Strang, G.: *Linear Algebra and Its Applications*, 3rd edition. Harcourt Brace Jovanovich, San Diego (1988)
- Templ, M., Kowarik, A., Meindl, B.: Statistical Disclosure Control for Micro-Data Using the R Package *sdcmicro*. *J. Stat. Softw.* **67**(4) (2015)
- Templ, M., Meindl, B.: Robustification of microdata masking methods and the comparison with existing methods. In: Domingo-Ferrer, J., Saygin, Y. (eds.) *Privacy in Statistical Databases, Proceedings*, pp. 113–126. Springer Berlin Heidelberg, Berlin, Heidelberg (2008). UNESCO Chair in Data Privacy International Conference (PSD 2008), Istanbul, Turkey, September 24–26, 2008.
- Ting, D., Fienberg, S.E., Trottini, M.: Random orthogonal matrix masking methodology for microdata release. *International Journal of Information and Computer Security* **2**(1), 86–105 (2008). doi:10.1504/IJICS.2008.016823
- Wedderburn, R.W.M.: *Random Rotations and Multivariate Normal Simulation*. Research Report, Rothamsted Experimental Station (1975)