

RESEARCH REPORT SERIES

(Disclosure Avoidance #2014-01)

**On Invariant Post Randomization for Statistical
Disclosure Control**

Tapan K. Nayak and Samson A. Adeshiyan

Center for Disclosure Avoidance Research
U.S. Census Bureau
Washington DC 20233

Report Issued: September 11, 2014

Disclaimer: This report is released to inform interested parties of ongoing research and to encourage discussion of work in progress. The views expressed are those of the author and not necessarily those of the U.S. Census Bureau.

On Invariant Post Randomization for Statistical Disclosure Control

Tapan K. Nayak* and Samson A. Adeshiyan^{†‡}

Abstract

In this paper, we investigate certain operational and inferential aspects of invariant PRAM (post randomization method) as a tool for disclosure limitation of categorical data. Invariant PRAMs preserve unbiasedness of certain estimators, but inflate their variances and distort other attributes. We introduce the concept of strongly invariant PRAM, which does not affect data utility or the properties of any statistical method. However, the procedure seems feasible in limited situations. We review methods for constructing invariant PRAM matrices and prove that a conditional approach, which can preserve the original data on any subset of variables, is an invariant PRAM. For multinomial sampling, we derive expressions for variance inflation due to invariant PRAMing and variances of certain estimators of the cell probabilities and also their tight upper bounds. We discuss estimation of these quantities and thereby assessing statistical efficiency loss due to invariant PRAMing. We find a connection between invariant PRAM and creating partially synthetic data using a nonparametric

*Center for Disclosure Avoidance Research, U.S. Census Bureau, Washington, DC 20233 and Department of Statistics, George Washington University, Washington, DC 20052.

[†]U.S. Energy Information Administration, Washington, DC 20585.

[‡]The views expressed in this article are those of the authors and not necessarily those of the U.S. Census Bureau. The analysis and conclusions contained in this paper are those of the authors and do not represent the official position of the U.S. Energy Information Administration (EIA) or the U.S. Department of Energy (DOE).

approach, and compare estimation variance under the two approaches. Finally, we discuss some aspects of invariant PRAM in a general survey context.

Key words and Phrases: Categorical data; randomized response; sampling design; synthetic data; unbiased estimation; variance inflation.

1. Introduction

The Post Randomization Method (PRAM), introduced by Kooiman et al. (1997) and further discussed by Gouweleeuw et al. (1998), De Wolf et al. (1998), Willenborg and De Waal (2001), Ronning (2005), Van den Hout and Elamir (2006), Van den Hout and Kooiman (2006), Cruyff et al. (2008) and others, is an important technique for categorical data perturbation for confidentiality protection. PRAM stochastically transforms each record in a data set using pre-selected probabilities. This deliberate misclassification of the original responses introduces uncertainty about the true category of any respondent. On the other hand, since the misclassification probabilities are known, valid statistical inferences can be derived from PRAMed data. Let X be a categorical variable with categories c_1, \dots, c_k . The basic ideas of PRAM are to (i) select a transition probability matrix $P = ((p_{ij}))$, where $\sum_i p_{ij} = 1$ for $j = 1, \dots, k$, and then (ii) randomly change any original category c_j to c_i with probability p_{ij} ($i, j = 1, \dots, k$). The randomization step is performed for each record in the data set, independently of all other records. We shall denote the transformed variable by Z , in which case, $p_{ij} = P(Z = c_i | X = c_j)$. In practice, data agencies should release PRAMed data along with the transition probability matrix P , also called the PRAM matrix, to the public, so that users can derive valid statistical inferences from the released data.

As noted by Gouweleeuw et al. (1998) and Van den Hout and Van der Heijden (2002), mathematically, PRAM is similar to randomized response (RR) surveys (e.g., Warner, 1965; Chaudhuri and Mukerjee, 1988; Nayak, 1994; Nayak and Adeshiyani, 2009). Thus, many ideas and mathematical results developed for RR surveys can be applied to PRAM. Both methods deal with the dual goals of protecting privacy and preserving statistical information. One difference is that in RR surveys, each respondent randomizes his response at data gathering stage, whereas in PRAM, randomization is carried out by the surveyor after the data are collected. Another divergence, which has not been noted well earlier, is that in RR surveys the transition probability

matrix P is prespecified and hence fixed, but in PRAM it may depend on the data and hence be random. This causes certain differences in the variance inflation due to the two procedures, which we discuss in Section 4.

Let $\pi_i = P[X = c_i], i = 1, \dots, k$, and $\pi = (\pi_1, \dots, \pi_k)'$. Let n denote the sample size and T_i denote the frequency of category c_i in the original sample. Most papers on PRAM assume multinomial sampling, i.e., the original data are collected by random sampling from an infinite population or by simple random sampling with replacement (SRSWR) if the population is finite. We shall assume multinomial sampling, except in Section 6, where we consider general probability sampling. Under multinomial sampling, $\mathbf{T} = (T_1, \dots, T_k)' \sim Mult(n, \pi)$, but note that T_1, \dots, T_k cannot be ascertained from PRAMed data. Let, S_i denote the frequency of category c_i after PRAMing, $\lambda_i = P(Z = c_i), i = 1, \dots, k$ and $\lambda = (\lambda_1, \dots, \lambda_k)'$. Then, $\mathbf{S} = (S_1, \dots, S_k)' \sim Mult(n, \lambda)$, where

$$\lambda = P\pi. \quad (1.1)$$

From (1.1), it follows that for nonsingular P , any unbiased estimator $\tilde{\lambda}$ of λ yields an unbiased estimator of π , given by $\tilde{\pi} = P^{-1}\tilde{\lambda}$. The MLE (and UMVUE) of λ is $\hat{\lambda} = \mathbf{S}/n$, which yields the following estimator of π :

$$\hat{\pi} = P^{-1}\hat{\lambda} = P^{-1}(\mathbf{S}/n). \quad (1.2)$$

It can be seen that $\hat{\pi}$ is an unbiased estimator of π and

$$Var(\hat{\pi}) = \frac{(D_\pi - \pi\pi')}{n} + \frac{[P^{-1}D_\lambda(P^{-1})' - D_\pi]}{n}, \quad (1.3)$$

where D_π is a diagonal matrix with diagonal elements being π_1, \dots, π_k and D_λ is defined similarly (see Chaudhuri and Mukerjee, 1988, p. 43). The first term on the right side of (1.3) is the variance under no randomization and the last term is the additional variance due to PRAMing.

Without PRAMing, i.e., based on the original data, the MLE (and the UMVUE) of π is $\hat{\pi}_0 = \mathbf{T}/n$. Note that $E[\mathbf{S}|\mathbf{T}] = P\mathbf{T}$ and hence $E[\hat{\pi}|\mathbf{T}] = \hat{\pi}_0$. Thus, the estimator $\hat{\pi}$ in (1.2) can

be regarded as an unbiasedly recovered version of $\hat{\pi}_0$. However, if $E[\mathbf{S}|\mathbf{T}] = \mathbf{T}$, i.e.,

$$P\mathbf{T} = \mathbf{T} \quad \text{or equivalently} \quad P\hat{\pi}_0 = \hat{\pi}_0, \quad (1.4)$$

then \mathbf{S}/n is also an unbiased estimator of π , which can be calculated without using P or its inverse. Also, \mathbf{S}/n is always a probability vector, while $\hat{\pi}$ in (1.2) may not be so. Observing these points, Gouweleeuw et al. (1998) defined a PRAM to be an invariant PRAM if P satisfies (1.4). The main goal of this paper is to examine certain practical and inferential sides of invariant PRAM and present some results to advance it further as a disclosure control method.

PRAM can be applied to more than one categorical variable, independently or jointly on the cross-classification of all variables (see, Gouweleeuw et al., 1998). Conceptually, any PRAM can be regarded as being applied to the compound variable created by cross-classifying all variables, and that view is both convenient and appropriate for estimating the joint probabilities. When there are several variables and they are PRAMed independently and invariantly, the marginal probabilities can be estimated from PRAMed data without any adjustment for PRAMing, but not the joint probabilities. This is because, independent invariant PRAMing of individual variables does not imply invariant PRAMing of the compound variable, unless the variables are independently distributed. So, for estimating the joint probabilities, one will need to obtain the transition probability matrix for the combined variable and use its inverse. Unbiased estimation of all cell probabilities without knowing the PRAM matrix requires an invariant PRAM of the combined variable.

In the next section, we introduce and discuss a stronger version of invariant PRAM that has no effect on statistical inferences. Choosing a suitable P satisfying (1.4) is central to devising an invariant PRAM. This requires methods for solving (1.4) for P and assessing the efficacy of any given P for confidentiality protection and also its effect on variance of estimators. In Section 3, we review two existing methods for solving equations (1.4) and (2.1) and formalize an approach that permits PRAMing some of the variables while keeping others unchanged. This

is useful when there are several variables and the total number of cells is large. In Section 4, we derive the variance covariance matrix $V(\hat{\pi}_*)$ of $\hat{\pi}_* = \mathbf{S}/n$ under multinomial sampling, and discuss how a data agency can estimate it. However, a user cannot estimate $V(\hat{\pi}_*)$, unless P is released along with PRAMed data, which is undesirable from confidentiality perspective. We derive a tight upper bound for $V(\hat{\pi}_*)$, which can be estimated from PRAMed data. In Section 5, we relate one type of synthetic data to invariant PRAM. Specifically, we show that creating nonparametric synthetic data is equivalent to applying a specific invariant PRAM, which also induces the maximum variance inflation among all invariant PRAMs. In Section 6, we discuss invariant PRAM for data coming from a probability sample. We give an expression for variance inflation that can be used by data agencies. However, it is difficult for a user to estimate the variance inflation or variances of estimators. We state some concluding remarks in Section 7.

2. Strongly Invariant PRAM

We note that (1.4) implies a specific invariance property, viz., unbiasedness of the observed relative frequency vector (\mathbf{S}/n) for estimating π , but not necessarily of other estimators. Even under (1.4), the distributions of \mathbf{T} and \mathbf{S} are different (unless $P = I$), and consequently (i) the covariance matrices of (\mathbf{T}/n) and (\mathbf{S}/n) are not the same and (ii) unbiased estimation of nonlinear functions of π (e.g., $\sum \pi_i^2$) may not be possible if P is not given. Thus, invariant PRAMs preserve only certain properties of some estimators.

For stronger implications, we may wish to choose P so that statistical properties of all inferential methods remain invariant under PRAMing, i.e., the validity of any inferential procedure would not depend on whether it is applied to the original data or a PRAMed version of it. In such cases, for making inferences about π , data users would not need to know or use the matrix P and can use any procedure that is valid for the original data. The question of whether it is possible to find such P or not was left open for future research in Gouweleeuw et al. (1998). Note

that invariance of all inferential methods essentially requires the distributions of any statistic calculated from original and PRAMed data, respectively, to be the same. Under multinomial sampling, this amounts to requiring the frequency counts from original and PRAMed data, i.e., (T_1, \dots, T_k) and (S_1, \dots, S_k) respectively, to have the same distribution, viz., $Mult(n, \pi)$. This holds if and only if $\lambda = \pi$ or equivalently P satisfies the condition

$$P\pi = \pi. \tag{2.1}$$

We define a PRAM to be a *strongly invariant* PRAM if P satisfies (2.1). One appealing property of a strongly invariant PRAM is that it does not affect either data utility or data analysis; there is no loss of statistical information and PRAMed data can be analyzed without any adjustment for post-randomization. However, such a strong outcome is usually unattainable in practice because typically π is unknown and so (2.1) cannot be solved for P (except for $P = I$). One situation where the true π is known is when the data come from a census. Likewise, when the sample size is large, π can be fairly well estimated from the original data. Thus, a strongly invariant PRAM is applicable when releasing a random subset of census data (or of a large data set) for public use. For example, the Public Use Microdata Sample (PUMS) files released by the U.S. Census Bureau are subsets of census or large survey data. Thus, strongly invariant PRAMing may serve as another tool for confidentiality protection in creating PUMS files.

3. Constructing Invariant PRAMs

To construct a strongly invariant PRAM, one needs to solve (2.1) for P (for given π) and choose a solution. Similarly, constructing an invariant PRAM requires solving (1.4). However, (1.4) and (2.1) are identical, as mathematical equations, and so we shall examine only (2.1). Note that (2.1) has many solutions for P , with the identity matrix being an obvious solution, and if P and P_* satisfy (2.1), then $aP + (1 - a)P_*$ also satisfies (2.1) for all $0 \leq a \leq 1$. Thus, the

solution space of (2.1) is a non-empty convex set. The solution set is also closed under matrix multiplication, viz. if P and P_* are strongly invariant transition probability matrices, then so is the product PP_* .

In many practical applications, X is a compound variable with some structurally empty cells, i.e., cells with zero probability (see Gouweleeuw et al., 1998). We note some practical points about such situations. For notational simplicity, suppose $\pi = (\pi'_+, \pi'_0)'$, where all components of π_+ are positive and π_0 is a vector of zeros. Then, partitioning P analogously, (2.1) reduces to

$$\begin{pmatrix} P_{11} & P_{12} \\ P_{21} & P_{22} \end{pmatrix} \begin{pmatrix} \pi_+ \\ 0 \end{pmatrix} = \begin{pmatrix} \pi_+ \\ 0 \end{pmatrix}. \quad (3.1)$$

Clearly, (3.1) can hold only if $P_{21} = 0$, and P_{12} and P_{22} have no bearing on the validity of (3.1). In particular, one may take $P_{12} = 0$ and $P_{22} = I$. Essentially, (2.1) reduces to $P_{11}\pi_+ = \pi_+$, and thus one needs to consider only the structurally non-zero cells (which also remain so under any strongly invariant PRAM). Similarly, for invariant PRAM, i.e., for (1.4), only the cells with positive frequencies are relevant. Any cell that has zero frequency for the original data will also have zero frequency after invariant PRAMing. In contrast, a cell with original count of zero and a positive probability may get a positive count after strongly invariant PRAMing. Thus, strongly invariant PRAMs may better protect confidentiality than invariant PRAMs. One concern in micro data perturbation is satisfying edit constraints (see De Wolf et al., 1998; Shlomo and De Waal, 2008). The preceding observations show that invariant (or strongly invariant) PRAMs of the cross classification of all variables automatically satisfy all natural edit constraints. For notational simplicity, we shall continue to use (2.1) considering only structurally non-empty cells, instead of switching to the equation $P_{11}\pi_+ = \pi_+$.

We shall now briefly review, partly for later reference, two methods for solving (2.1), presented in Gouweleeuw et al. (1998). To discuss one class of solutions, assume for notational simplicity that $0 < \pi_k \leq \pi_i$ for $i = 1, \dots, k$. Then, letting $p_{ii} = 1 - \theta(\frac{\pi_k}{\pi_i})$ and $p_{ij} = \theta[\frac{\pi_k}{(k-1)\pi_j}]$

for $i \neq j$, it can be seen that $P = ((p_{ij}))$ satisfies (2.1) for all $0 \leq \theta \leq 1$. These $\{p_{ij}\}$, with any $0 \leq \theta \leq 1$, constitute a strongly invariant PRAM matrix. Here, $\{p_{ij}\}$ can be interpreted nicely. Note that $P(Z \neq c_j | X = c_j) = \theta(\frac{\pi_k}{\pi_j}), j = 1, \dots, k$, which implies that θ is the probability of changing a true response of c_k to another category. Any other response c_j ($j \neq k$), changes to another category with probability $\theta(\frac{\pi_k}{\pi_j})$. Moreover, it can be seen that for $i \neq j, P(Z = c_i | X \neq Z) = 1/(k - 1)$, which means that when a true response is changed, the perturbed value is selected at random from one of the remaining $(k - 1)$ categories.

Another method for solving (2.1) uses a two step process. Let $R = ((r_{ij}))$ be any transition probability matrix such that all components of $R\pi$ are positive. Let $q_{ij} = [r_{ji}\pi_i]/[\sum_l r_{jl}\pi_l]$, $Q = ((q_{ij}))$ and $P = QR$. Then, it can be verified that P satisfies (2.1). This also shows that any PRAMed data set (with a known transition probability matrix R) can be PRAMed again, with a suitable transition probability matrix Q , to make the final outcome a strongly invariant PRAM. In this context, we note that R need not be a square matrix, i.e., X and Z need not have the same number of cells. For any transition probability matrix R of order $m \times k$ with all components of $R\pi$ being positive, the final matrix P is a strongly invariant PRAM matrix. However, P would be singular, if $m < k$ or $rank(R) < k$.

When there are several variables, applying strongly invariant PRAM to each variable separately (and independently) does not amount to using a strongly invariant PRAM on the compound variable, unless all variables are independently distributed. So, the two methods discussed above need to be applied to the cross-classification of all variables, which is unwieldy when the total number of cells is large. Gouweleeuw et al. (1998) and De Wolf et al. (1998) discuss several practical issues, such as how to decide which variables are to be PRAMed and how to choose P . They suggest to group the variables such that any two variables falling in two different groups are nearly independent. Then, consider all variables in each group as a compound variable and partition its categories into clusters in such a way that any category can be replaced only by

categories within its cluster. This approach leads to a block diagonal structure of P for each compound variable. Finally, choose the block matrices suitably to achieve confidentiality protection goals. Shlomo and De Waal (2008) applied these recommendations to an Israel Census sample data set. However, as the compound variables (for the groups) are PRAMed separately, the overall process may not be an invariant PRAM of the cross-classification of all variables.

Next, we describe another mechanism for constructing strongly invariant PRAMs of several variables. The procedure uses a conditional approach and has been used by Shlomo and De Waal (2008). The following discussion gives the procedure a formal footing and helps to bring out a connection between invariant PRAM and partially synthetic data, presented in Section 5.

Theorem 3.1. *Consider two categorical variables X and Y . Suppose X is kept unchanged and within each category of X , Y is strongly invariantly PRAMed into Z . Then, the stochastic transformation from (X, Y) to (X, Z) is a strongly invariant PRAM of (X, Y) .*

Proof. Suppose X has k categories, denoted by c_1, \dots, c_k , and Y and Z have m categories labeled d_1, \dots, d_m . Since Y is strongly invariantly PRAMed within each category of X , we have $P(Z = d_j | X = c_i) = P(Y = d_j | X = c_i)$ for all $i = 1, \dots, k$ and $j = 1, \dots, m$. Then considering the joint probabilities, we get

$$\begin{aligned} P(X = c_i, Z = d_j) &= P(X = c_i)P(Z = d_j | X = c_i) \\ &= P(X = c_i)P(Y = d_j | X = c_i) \\ &= P(X = c_i, Y = d_j) \end{aligned}$$

for all $i = 1, \dots, k$ and $j = 1, \dots, m$, which completes the proof. □

It can be seen similarly that a parallel version of Theorem 3.1 holds for invariant PRAM, i.e., invariant PRAMing of Y within each X category is an invariant PRAM of (X, Y) . Thus, the following mechanism can be used for invariant PRAMing. First, divide all variables into two

groups and let X and Y denote the corresponding compound variables; each variable in the data set must be a part of either X or Y . Then, for each category of X , apply an invariant PRAM to Y . One advantage of the conditional approach is that it fully preserves the original data on X , which may be regarded as a control variable. This is useful for preserving information on some of the variables that are deemed important to users. For example, we may PRAM zipcodes within each county to preserve county level counts. We should also note that the conditional approach requires choosing and using a transition probability matrix for each X category, which may be burdensome, but that also offers the flexibility to perturb Y in varying degree across the categories of X . This is helpful for mitigating different disclosure risks for units in different X categories. When the total number of cells is large and many cells have small counts, invariant PRAMing may not be convenient or adequate for privacy protection. In such cases, invariant PRAMing after cell collapsing could be a practical and effective approach.

We believe, the procedure proposed above, viz., dividing all variables into two compound variables X and Y and then invariantly PRAMing Y within each category of X , would suffice in most practical situations. However, we may note that repeated applications this process also result in an invariant PRAM. This follows from the fact that if P and P_* are invariant transition probability matrices (i.e., satisfy (1.4)), then PP_* is also so. In this context, one may put any of the variables in the control set in each step. For example, for a data set with three categorical variables X_1, X_2 and X_3 , we may first invariantly PRAM X_2 and X_3 within X_1 categories, creating (X_1, X'_2, X'_3) and then invariantly PRAM X_1 and X'_2 within X_3 to form (X'_1, X''_2, X'_3) and so on. Such a perturbation process is equivalent to an invariant PRAM.

4. Estimation from Invariantly PRAMed Data

As noted in Section 2, under multinomial sampling, all analytical methods for original data are equally valid for strongly invariantly PRAMed data. Thus, strongly invariant PRAMs do

not raise any new inferential issues. Also, a strongly invariant PRAM cannot be constructed in many practical situations, as π is unknown. So, here we shall discuss estimation based on invariantly PRAM data, as defined by (1.4). Let $P = [P_1 : \dots : P_k]$ be an invariant PRAM matrix and rewrite (1.4) as

$$\sum_{i=1}^k T_i P_i = \mathbf{T}. \quad (4.1)$$

Note that P depends on \mathbf{T} and hence is a random matrix. Let F_{ij} denote the number of units with original category c_i and perturbed category c_j , and let $\mathbf{F}_i = (F_{i1}, \dots, F_{ik})'$. Then, the frequency vector (\mathbf{S}) from PRAMed data can be expressed as $\mathbf{S} = \sum_i^k \mathbf{F}_i$, where given \mathbf{T} and P , $\mathbf{F}_i \sim Mult(T_i, P_i)$, $i = 1, \dots, k$, and they are independent. We shall find this representation helpful for deriving formulas for variance inflation and overall estimation variance.

First, we note a *subtle point* that it makes sense to talk about unconditional distributions of \mathbf{F}_i and \mathbf{S} only when a probability model for P is given. Note that while P depends on \mathbf{T} , it need not be a function of \mathbf{T} , as (1.4) has many solutions. So, for discussing distributional properties of relevant statistics, we shall assume that an algorithm for choosing P from \mathbf{T} , via (1.4), is given, i.e., there is a probability model for P that is fully determined by \mathbf{T} . The algorithm for choosing P may be probabilistic or deterministic, and in the latter case the distribution of P given \mathbf{T} is degenerate. However, for our discussion, we mostly do not need to know the actual probability model for P ; it suffices to know that it exists and is determined by \mathbf{T} .

In the above setting, we investigate statistical properties of $\hat{\pi}_* = \mathbf{S}/n$ as an estimator of π . Using $\mathbf{S} = \sum_{i=1}^k \mathbf{F}_i$ and above mentioned conditional distributions of \mathbf{F}_i , $i = 1, \dots, k$, we get

$$E(\hat{\pi}_* | \mathbf{T}, P) = \frac{1}{n} \sum_{i=1}^k E[\mathbf{F}_i | \mathbf{T}, P] = \frac{1}{n} \sum_{i=1}^k T_i P_i = \frac{1}{n} \mathbf{T} = \hat{\pi}_0, \quad (4.2)$$

by (4.1), and

$$V(\hat{\pi}_* | \mathbf{T}, P) = \frac{1}{n^2} \sum_{i=1}^k T_i [D_{P_i} - P_i P_i'] = \frac{1}{n} [D_{\hat{\pi}_0} - \sum_{i=1}^k (\frac{T_i}{n}) P_i P_i'], \quad (4.3)$$

which represents the added variation due to PRAMing. We believe, data agencies should examine $V(\hat{\pi}_*|\mathbf{T}, P)$ when selecting a suitable P . Note that for any given P , the data agency can calculate (4.3), based on the original data, before applying PRAM.

Since $\mathbf{T} \sim Mult(n, \pi)$ and $\hat{\pi}_0 = \mathbf{T}/n$, we get $E(\hat{\pi}_0) = \pi$ and $V(\hat{\pi}_0) = [D_\pi - \pi\pi']/n$. From these and (4.2) and (4.3) it follows that $E(\hat{\pi}_*) = \pi$ (for any distribution of P) and

$$\begin{aligned} V(\hat{\pi}_*) &= V[E(\hat{\pi}_*|\mathbf{T}, P)] + E[V(\hat{\pi}_*|\mathbf{T}, P)] \\ &= V(\hat{\pi}_0) + \frac{1}{n}[D_\pi - E\{\sum_{i=1}^k (\frac{T_i}{n})P_iP_i'\}]. \end{aligned} \quad (4.4)$$

Thus, the relative frequency vector $\hat{\pi}_*$ from invariantly PRAMed data is an unbiased estimator of π and the last term of (4.4) is the variance inflation matrix, which is different from the last term of (1.3). This difference arises from the fact that (1.3) holds when P is fixed and prespecified, but in invariant PARM, P is data dependent and random. We can also express (4.4) as

$$V(\hat{\pi}_*) = \frac{1}{n}[2D_\pi - \pi\pi'] - \frac{1}{n}E\{\sum_{i=1}^k (\frac{T_i}{n})P_iP_i'\}. \quad (4.5)$$

The expectation in (4.5) may be difficult to derive, depending on the distribution of P . The data agency knows P and the values of T_1, \dots, T_k and hence can estimate (4.5) by

$$\hat{V}_1(\hat{\pi}_*) = \frac{1}{n}[2D_{\hat{\pi}_0} - \hat{\pi}_0\hat{\pi}_0'] - \sum_{i=1}^k (\frac{T_i}{n^2})P_iP_i'.$$

If P is reported along with the perturbed data, a data user may estimate $V(\hat{\pi}_*)$ (without knowing T_1, \dots, T_k) by

$$\hat{V}_2(\hat{\pi}_*) = \frac{1}{n}[2D_{\hat{\pi}_*} - \hat{\pi}_*\hat{\pi}_*'] - \sum_{i=1}^k (\frac{S_i}{n^2})P_iP_i'.$$

However, the reporting of P is problematic from disclosure perspective. As Gouweleeuw et al. (1998) noted, \mathbf{T} is an eigenvector of P corresponding to the eigenvalue 1 and one might be able to deduce the original frequencies from P without any error. If P is not published, estimating $V(\hat{\pi}_*)$ by data users is challenging. However, a user can apply the following result to obtain an upper bound for the variance of $\hat{\pi}_*$.

Theorem 4.1. (a) An upper bound of $V(\hat{\pi}_*)$ is

$$V_{max}(\hat{\pi}_*) = (2 - \frac{1}{n})[\frac{D_\pi - \pi\pi'}{n}] \quad (4.6)$$

in the sense that $[V_{max}(\hat{\pi}_*) - V(\hat{\pi}_*)]$ is non-negative definite for any invariant PRAM.

(b) The upper bound in (4.6) is tight, i.e., there exists an invariant PRAM for which $V(\hat{\pi}_*) = V_{max}(\hat{\pi}_*)$.

Proof. For any given vector a of dimension k , consider

$$a'V(\hat{\pi}_*)a = \frac{1}{n}a'[2D_\pi - \pi\pi']a - \frac{1}{n^2}E\{\sum_{i=1}^k T_i(a'P_i)^2\}. \quad (4.7)$$

Suppose, without loss of generality, that only the first m ($m \leq k$) components of \mathbf{T} are nonzero. Then, using (4.1) and the Cauchy-Schwarz inequality, $(x'y)^2 \leq (x'Ax)(y'A^{-1}y)$, with $x' = (T_1, \dots, T_m)$, $y' = (a'P_1, \dots, a'P_m)$ and $A = \text{diag}(1/T_1, \dots, 1/T_m)$, we get

$$(a'\mathbf{T})^2 = [\sum_{i=1}^m T_i(a'P_i)]^2 \leq n[\sum_{i=1}^m T_i(a'P_i)^2] = n[\sum_{i=1}^k T_i(a'P_i)^2],$$

which implies that

$$\begin{aligned} \frac{1}{n^2}E\{\sum_{i=1}^k T_i(a'P_i)^2\} &\geq \frac{1}{n^3}E(a'\mathbf{T})^2 \\ &= \frac{1}{n}a'[E(\hat{\pi}_0\hat{\pi}'_0)]a \\ &= \frac{1}{n}a'[(\frac{D_\pi - \pi\pi'}{n}) + \pi\pi']a. \end{aligned} \quad (4.8)$$

Next, using (4.8) in (4.7), we get

$$a'V(\hat{\pi}_*)a \leq (2 - \frac{1}{n})a'[\frac{D_\pi - \pi\pi'}{n}]a. \quad (4.9)$$

Part (a) of the theorem now follows readily as (4.9) holds for all a .

For part (b), we show that $V(\hat{\pi}_*)$ equals $V_{max}(\hat{\pi}_*)$ if the transition probability matrix is

$$P_* = \frac{1}{n} \begin{pmatrix} T_1 & T_1 & \cdots & T_1 \\ T_2 & T_2 & \cdots & T_2 \\ \vdots & \vdots & \cdots & \vdots \\ T_k & T_k & \cdots & T_k \end{pmatrix}.$$

Note that P_* satisfies (1.4) and hence it is an invariant PRAM matrix. As each column of P_* is $\hat{\pi}_0$, the sum in (4.4) equals $\hat{\pi}_0 \hat{\pi}'_0$ and hence $V(\hat{\pi}_*)$ reduces to

$$\begin{aligned} V(\hat{\pi}_*) &= V(\hat{\pi}_0) + \frac{1}{n}[D_\pi - E(\hat{\pi}_0 \hat{\pi}'_0)] \\ &= V(\hat{\pi}_0) + \frac{1}{n}[D_\pi - \{V(\hat{\pi}_0) + \pi \pi'\}] \\ &= (2 - \frac{1}{n})V(\hat{\pi}_0), \end{aligned}$$

which is the same as $V_{max}(\hat{\pi}_*)$ in (4.6). □

Estimating linear combinations of π , say $a'\pi$, is a common problem. One special case is estimating the total probability of a set of categories, which can be expressed as $a'\pi$, where $a_i = 1$ if category c_i is included in the set and $a_i = 0$ otherwise. A natural estimator of $a'\pi$ is $a'\hat{\pi}_*$ and its variance is $a'V(\hat{\pi}_*)a$, for which (4.9) gives an upper bound. Data users can estimate the upper bounds in (4.6) and (4.9) by replacing π by $\hat{\pi}_*$. An obvious lower bound for $V(a'\hat{\pi}_*)$ is $V(a'\hat{\pi}_0) = a'[\{D_\pi - \pi \pi'\}/n]a$. Recall that $\hat{\pi}_0$ and $\hat{\pi}_*$ are estimators of π based on original and invariantly PRAMed data, respectively, and from (4.6), $V_{max}(\hat{\pi}_*) = (2 - 1/n)V(\hat{\pi}_0) \approx 2V(\hat{\pi}_0)$; a similar comparison holds for $V(a'\hat{\pi}_*)$ and $V(a'\hat{\pi}_0)$. As the upper bounds in (4.6) and (4.9) are tight, the variance of an estimator under an invariant PRAM may be almost twice the variance of a corresponding estimator based on the original data.

5. Invariant PRAM and Synthetic Data

One significant approach to disclosure avoidance is to release fully or partially synthetic data; see e.g., Rubin (1993), Raghunathan et al. (2003), Little (1993), Reiter (2002, 2003), An and Little (2007), Hawala (2008) and Slavkovic and Lee (2010). The basic idea in creating partially synthetic data is to replace the observed values of some selected variables by simulated data from pertinent predictive distributions. A synthetic data set consisting of original X and synthetic Y values, where both X and Y may be vectors, is created by replacing the observed Y of each unit by a simulated value from a predictive distribution $f(y|x)$, where x is the unit's observed value of X . Usually, $f(y|x)$ is a posterior predictive distribution, but one may also use $f(y|x, \hat{\theta})$, where $\hat{\theta}$ is some estimate (such as MLE) of the parameter(s) of an *assumed* model; see Hawala (2008) and Reiter and Kinney (2012). In the following, we present a link between invariant PRAM and nonparametric partially synthetic data and compare the variance of estimators under the two approaches.

Suppose a data set has only two categorical variables X and Y (both can be compounded variables) with categories labeled as in Theorem 3.1. Let T_{ij} denote the frequency of $(X = c_i, Y = d_j)$, $i = 1, \dots, k, j = 1, \dots, m$, in the original data and let $T_{i.}$ denote the marginal frequency of $(X = c_i)$. Let $p_{j|i} = T_{ij}/T_{i.}$. Here, $p_{j|i}$ is a natural and nonparametric estimate of $f(y|x = c_i)$. So, a nonparametric approach to creating partially synthetic data set would be as follows. For any unit with $X = c_i$, replace its observed Y category by one of d_1, \dots, d_m , with respective probabilities $p_{1|i}, \dots, p_{m|i}$. However, this is equivalent to ordering the cells as $c_1d_1, \dots, c_1d_m, c_2d_1, \dots, c_2d_m, \dots, c_kd_1, \dots, c_kd_m$ and then applying a PRAM with a block diagonal PRAM matrix, whose i th ($i = 1, \dots, k$) block, say P^i , is an $m \times m$ matrix with each column being $(p_{1|i}, \dots, p_{m|i})'$. It is easy to see that P^i is an invariant PRAM matrix for PRAMing Y within the category c_i of X . The whole process is an invariant PRAM, in view of Theorem 3.1. In summary, for categorical variables, a nonparametric synthetic data set can be viewed as an

invariantly PRAMed data set.

Our results of Section 4 can be extended to the case where Y is invariantly PRAMed within each X category. Here, the PRAMing process divides all cells into k groups, one for each category of X , and then perturbs the frequencies of the cells in each group; it does not alter the marginal frequencies of X categories. To see how the method affects the frequencies in each group, let us consider the cells within $X = c_i$. Let $\vec{T}_i = (T_{i1}, \dots, T_{im})'$ and let $P^i = (P_1^i : \dots : P_m^i)$ be an invariant PRAM matrix for this group, i.e.,

$$P^i \vec{T}_i = \vec{T}_i \quad \text{or} \quad \sum_{j=1}^m T_{ij} P_j^i = \vec{T}_i. \quad (5.1)$$

Let S_{ij} denote the frequency of the (i, j) th cell based on PRAMed data and let $\vec{S}_i = (S_{i1}, \dots, S_{im})'$. Here, \vec{S}_i is a perturbed version of \vec{T}_i and as in Section 4, we have $\vec{S}_i = F_1^i + \dots + F_m^i$, where $F_j^i \sim \text{mult}(T_{ij}, P_j^i)$, $j = 1, \dots, m$, are independent. Then, for any P^i satisfying (5.1), it follows that $E(\vec{S}_i | \vec{T}_i, P^i) = \vec{T}_i$ and

$$V(\vec{S}_i | \vec{T}_i, P^i) = D_{\vec{T}_i} - \sum_{j=1}^m T_{ij} P_j^i (P_j^i)'. \quad (5.2)$$

If $P_j^i = (1/T_{i.})(T_{i1}, \dots, T_{im})'$, $j = 1, \dots, m$, which yields nonparametric partial synthetic data, the sum in (5.2) reduces to $(1/T_{i.})\vec{T}_i\vec{T}_i'$. Moreover, it can be seen, as in the proof of Theorem 4.1, that $[\sum_{j=1}^m T_{ij} P_j^i (P_j^i)' - (1/T_{i.})\vec{T}_i\vec{T}_i']$ is non-negative definite for all P^i satisfying (5.1), which leads to the following:

Theorem 5.1. *The induced variation of the cell frequencies due to replacing observed Y by nonparametric synthetic values is at least as large as the variation induced by any invariant PRAM of Y within the categories of X .*

The preceding discussions show that generating synthetic Y in a nonparametric way is equivalent to invariantly PRAMing Y within X categories, with automatically using the most variance inflationary PRAM matrices. In practice, this extreme invariant PRAMing may be unduly perturbative. We believe, data agencies would prefer to invariantly PRAM Y within X categories,

with more flexible and judicious choice of the PRAM matrices. The general approach also permits to apply different amount of perturbation to data in different categories of X . Operationally, one needs to choose and use a solution P^i of (5.1), for each $i = 1, \dots, k$. Here, the first method that we reviewed in Section 3, due to Gouweleeuw et al. (1998), may be useful as it requires just one number, viz. choosing a value of θ , for each X category. Recall that θ represents the probability of changing any value of Y that has the lowest frequency. Usually, the units falling in cells with small counts have high disclosure risk and require greater protection. So, the value of θ for category c_i of X should decrease as $\min_j \{T_{ij}\}$ increases. In particular, if all $T_{ij}, j = 1, \dots, m$ are fairly large, θ should be close to 0 and $P^i \approx I$.

The preceding discussion should not be viewed as a sweeping criticism of disclosure limitation via synthetic data. Clearly, our observations and Theorem 5.1 are about a specific synthetic approach and they do not apply to synthetic data based on models and priors, which may serve well in some applications. However, the nonparametric approach has a particular appeal as models and priors may lead to much biases when they are not correct. Obviously, PRAM is only for categorical variables whereas synthetic data can be created for both categorical and quantitative variables. Also, synthetic data based on multiple imputation involves additional inferential and disclosure control issues. For a more comprehensive comparison of PRAM with synthetic data one should examine both estimation error and disclosure control, which we leave for future research.

6. Estimation Under Probability Sampling

Practical surveys rarely use simple random sampling with replacement (i.e, multinomial sampling). Most surveys use a combination of stratified, cluster, systematic and multi-stage sampling. Also, nonresponse is a common phenomenon. Consequently, one uses a weighted combination of the data for estimating the population probability vector π . Thus, it is important

to examine properties of PRAM when applied to data coming from a general finite population survey. Let N denote the population size and consider a general survey design $p(s)$, where s is a subset of the population and $p(s)$ is the probability of selecting s . As in Section 1, let X be a categorical variable with k categories c_1, \dots, c_k and let π denote the corresponding probability vector. Let U_i ($i = 1, \dots, N$) be a k -dimensional indicator vector representing the X category of unit i . Specifically, if unit i falls in category c_j , then the j th component of U_i is 1 and the remaining components are 0. A similar indicator vector V_i will represent the category of unit i after PRAMing. Suppose

$$\hat{\pi}_w = \sum_{i \in s} w_{si} U_i \quad (6.1)$$

is an estimator of π based on the original data, where the sampled units are weighted by $\{w_{si}\}$.

Suppose the data are PRAMed using a fixed non-singular PRAM matrix $P = [P_1 : \dots : P_k]$.

Then, one can use

$$\hat{\pi}_{w*} = \sum_{i \in s} w_{si} P^{-1} V_i = P^{-1} \left(\sum_{i \in s} w_{si} V_i \right)$$

to estimate π based on PRAMed data. Note that the conditional distribution of V_i given U_i is $mult(1, PU_i)$, which shows that $E[\hat{\pi}_{w*}|s] = \hat{\pi}_w$ and

$$\begin{aligned} V[\hat{\pi}_{w*}|s] &= P^{-1} \left[\sum_{i \in s} w_{si}^2 \{D_{PU_i} - PU_i U_i' P'\} \right] (P^{-1})' \\ &= P^{-1} \left[\sum_{i \in s} w_{si}^2 D_{PU_i} \right] (P^{-1})' - \sum_{i \in s} w_{si}^2 U_i U_i'. \end{aligned} \quad (6.2)$$

From these, we see that $\hat{\pi}_{w*}$ unbiasedly recovers $\hat{\pi}_w$ and

$$V(\hat{\pi}_{w*}) = V_p(\hat{\pi}_w) + E_p\{V[\hat{\pi}_{w*}|s]\}, \quad (6.3)$$

where V_p and E_p are with respect to the sampling design. The last term of (6.3) is variance inflation due to PRAMing, which the data agency can unbiasedly estimate by (6.2). The data agency should examine the size of the matrix in (6.2), measured for example by its trace or determinant, for judging the suitability of P from data utility perspective.

Next consider using

$$\tilde{\pi}_w = \sum_{i \in s} w_{si} V_i$$

for estimating π , i.e., applying the estimator (6.1) for the original data to PRAMed data. Here,

$$E[\tilde{\pi}_w | s] = \sum_{i \in s} w_{si} P U_i = P \left[\sum_{i \in s} w_{si} U_i \right] = P \hat{\pi}_w,$$

which coincides with $\hat{\pi}_w$ if and only if

$$P \hat{\pi}_w = \hat{\pi}_w. \tag{6.4}$$

Based on this observation, Kooiman et al. (1997) and Gouweleeuw et al. (1998) defined invariant PRAM by (6.4) when the survey weights are unequal. We note some properties of these invariant PRAMs in the following.

We shall assume that $\hat{\pi}_w$ is a proper probability vector, i.e., all of its components are non-negative and add to 1. Then, any of the methods discussed in Section 3 can be used to find solutions of (6.4). For using the conditional approach, we would use relevant conditional probabilities implied by $\hat{\pi}_w$. We can still keep data on some variables (X) unchanged and invariantly PRAM the remaining variables (Y) within the categories of X . This is quite important when the survey weights depend on some of the variables. PRAMing of those variables would introduce randomness in the survey weights and hence the above treatment of $E[\tilde{\pi}_w | s]$ would not be correct. To define invariant PRAM via (6.4) it is necessary to assume that all variables determining the survey weights are kept unchanged.

As before, choosing a P satisfying (6.4) depends on the data, via $\hat{\pi}_w$, and the algorithm for calculating P . Thus, the distribution of $\{V_i\}$ also depends on the procedure for selecting one of many solutions of (6.4). Conditionally given s and $P = [P_1 : \dots : P_k]$, $\{V_i\}$ are independent $mult(1, P U_i)$, which implies that $E[\tilde{\pi}_w | s, P] = \hat{\pi}_w$ and

$$V[\tilde{\pi}_w | s, P] = \sum_{i \in s} w_{si}^2 \{D_{P U_i} - P U_i U_i' P'\}$$

$$= \sum_{i \in s} w_{si}^2 D_{PU_i} - P \left[\sum_{i \in s} w_{si}^2 U_i U_i' \right] P'. \quad (6.5)$$

Consequently, if $\hat{\pi}_w$ is an unbiased estimator of π , then so is $\tilde{\pi}_w$ and

$$V(\tilde{\pi}_w) = V(\hat{\pi}_w) + E\{V[\tilde{\pi}_w|s, P]\},$$

where the expectation is with respect to both the distribution of P given s and the sampling design, which a user cannot evaluate if the mechanism for generating P is not released. However, the data agency can calculate (6.5) and use it to assess the level of variance inflation.

7 Discussion

Of the few methods that are available for perturbing categorical microdata for disclosure control, invariant PRAM is appealing because it allows unbiased estimation of all cell probabilities, without any adjustment for data perturbation. We believe, our results and observations contribute to its further development as a useful tool. We noted that preserving all properties of all estimators require a stronger version of invariant PRAM, which is applicable in limited situations. We gave further reasons for using the conditional approach, i.e., invariantly PRAMing some variables within the categories of the remaining variables, for devising an invariant PRAM. It results in a proper invariant PRAM and allows the data agency to selectively retain the original data on some variables, which is particularly important when some of the variables determine the survey weights. We examined the effect of invariant PRAM on the variance of estimators, which previously had not been appreciated well. Our results are useful to data agencies for assessing variance inflation and selecting suitable PRAM matrices. We found a connection between invariant PRAM and nonparametric synthetic data and observed that the former is more flexible and gives greater control to data agencies to balance disclosure control and data utility.

In this paper, we focused mainly on constructing invariant PRAM matrices and effect of invariant PRAMing on the variance of estimators of cell probabilities. We did not go into

measuring disclosure risk. Some general approaches and related issues have been discussed by Duncan and Lambert (1986, 1989), Greenberg and Zayatz (1992), Lambert (1993), Reiter (2005), Skinner and Elliott (2002), Shlomo and Skinner (2010, 2012) and others. Invariant PRAM (or more generally PRAM) is a significant disclosure control technique, but we believe that additional research and directives are needed for it to become a common method. In particular, further guidance on how to choose the PRAM matrix will be helpful. In principle, the choice should achieve disclosure control goals at minimum loss of information, one consequence of which is variance inflation. De Wolf et al. (1998), Gouweleeuw et al. (1998), Skinner and Elliott (2002), Shlomo (2010), Shlomo and De Waal (2008), Shlomo and Skinner (2010, 2012), Van den Hout and Elamir (2006), Willenborg (2000), Yancey et al. (2002) and others have made significant contribution to disclosure risk assessment and the choice of a PRAM matrix. Additional empirical and theoretical research on effects of invariant PRAM on both information loss and disclosure risk will help expand its the scope of applications for disclosure control.

References

- [1] An, D., and Little, R.J.A. (2007). Multiple imputation: an alternative to top coding for statistical disclosure control. *J. Royal Statist. Soc., Ser. A*, 170, 923-940.
- [2] Chaudhuri, A. and Mukerjee, R. (1988). *Randomized Response: Theory and Techniques*. New York: Marcel Dekker.
- [3] Cruyff, M.J.L.F., Van den Hout, A. and Van der Heijden, P.G.M. (2008). The analysis of randomized response sum score variables. *J. Royal Statist. Soc., Ser. B*, 70, 21-30.
- [4] De Wolf, P.-P., Gouweleeuw, J.M., Kooiman, P., Willenborg, L. (1998). Reflections on PRAM. In: *Statistical Data Protection*. Office for Official Publications of the European Communities, Luxembourg, pp. 337-349

- [5] Duncan, G.T. and Lambert, D. (1986). Disclosure-limited data dissemination. *J. Amer. Statist. Assoc.*, 81, 10-28.
- [6] Duncan, G.T. and Lambert, D. (1989). The risk of disclosure for microdata. *J. Business & Econ. Statist*, 7, 207-217.
- [7] Gouweleeuw, J.M., Kooiman, P., Willenborg, L.C.R.J. and De Wolf, P.-P. (1998). Post randomisation for statistical disclosure control: Theory and implementation. *J. Official Statist.*, 14, 463–478.
- [8] Greenberg, B. and Zayatz, L. (1992). Strategies for measuring risk in public use microdata files. *Statist. Neerland.*, 46, 33-48.
- [9] Hawala, S. (2008). Partially synthetic data via expert knowledge, modeling, and matching. Proceedings of the Survey Methods Section, SSC Annual Meeting.
- [10] Kooiman, P., Willenborg, L., and Gouweleeuw, J. (1997). A method for disclosure limitation of microdata. Research paper 9705, Statistics Netherlands, Voorburg.
- [11] Lambert, D. (1993). Measure of disclosure risk and harm. *J. Official Statist.*, 9, 313-331.
- [12] Little, R.J.A. (1993). Statistical analysis of masked data. *J. Official Statist.*, 9, 407-426.
- [13] Nayak, T.K. (1994). On randomized response surveys for estimating a proportion. *Commun. Statist.- Theory Meth.*, 23, 3303-3321.
- [14] Nayak, T.K. and Adeshiyan, S.A. (2009). A unified framework for analysis and comparison of randomized response surveys of binary characteristics. *J. Statist. Plann. Inference*, 139, 2757-2766.
- [15] Raghunathan, T.E., Reiter, J.P., and Rubin, D.B. (2003). Multiple imputation for statistical disclosure limitation. *J. Official Statist.*, 19, 1-16.

- [16] Reiter, J.P. (2002). Satisfying disclosure restrictions with synthetic datasets. *J. Official Statist.*, 18, 531-544.
- [17] Reiter, J.P. (2003). Inference for partially synthetic, public use microdatasets. *Survey Methodology*, 29, 181-188.
- [18] Reiter, J.P. (2005). Estimating identification risk in microdata. *J. Amer. Statist. Assoc.*, 100, 1101-1113.
- [19] Reiter, J. P. and Kinney, S. K. (2012). Inferentially valid, partially synthetic data: Generating from posterior predictive distributions not necessary. *J. Official Statist.*, 28, 583-590.
- [20] Ronning, G. (2005). Randomized response and the binary probit model. *Economics Letters*, 86, 221-228.
- [21] Rubin, D.B. (1993). Discussion: Statistical disclosure limitation. *J. Official Statist.*, 9, 462-468.
- [22] Skinner, C.J. and Elliot, M.J. (2002). A measure of disclosure risk for microdata. *J. R. Statist. Soc., Ser. B*, 64, 855-867.
- [23] Shlomo, N. (2010). Releasing microdata: Disclosure risk estimation, data masking and assessing utility. *J. Privacy & Confidentiality*, 2, 73-91.
- [24] Shlomo, N. and De Waal, T. (2008). Protection of micro-data subject to edit constraints against statistical disclosure. *J. Official Statist.*, 24, 229-253
- [25] Shlomo, N. and Skinner, C. (2010). Assessing the protection provided by misclassification-based disclosure limitation methods for survey microdata. *Ann. Appl. Statist.*, 4, 1291-1310.

- [26] Shlomo, N. and Skinner, C.J. (2012). Privacy protection from sampling and perturbation in survey microdata. *J. Privacy & Confidentiality*, 4, 155-169
- [27] Slavkovic, A.B. and Lee, J. (2010). Synthetic two-way contingency tables that preserve conditional frequencies. *Statistical Methodology*, 7, 225-239
- [28] Van den Hout, A. and Van der Heijden, P.G.M. (2002). Randomized response, statistical disclosure control and misclassification: A review. *Internat. Statist. Rev.*, 70, 269-288.
- [29] Van den Hout, A. and Elamir, E.A.H. (2006). Statistical disclosure control using post randomisation: Variants and measures for disclosure risk. *J. Official Statist.*, 22, 711731.
- [30] Van den Hout, A. and Kooiman, P.(2006). Estimating the linear regression model with categorical covariates subject to randomized response. *Computational Statistics & Data Analysis*, 50, 3311-3323.
- [31] Warner, S.L. (1965). Randomized response: A survey technique for eliminating evasive answer bias. *J. Amer. Statist. Assoc.*, 60, 63-69.
- [32] Willenborg, L. (2000). Optimality models for PRAM. In *Proceedings in Computational Statistics*, Eds. J.G. Bethlehem and P.G.M. van der Heijden. Heidelberg: Physica-Verlag, pp. 505-510.
- [33] Willenborg, L.C.R.J. and De Waal, T. (2001). *Elements of Statistical Disclosure Control*. New York: Springer.
- [34] Yancey, W., Winkler, W., and Creecy, R. (2002). Disclosure risk assessment in perturbative micro-data protection. In J. Domingo-Ferrer, ed., *Inference Control in Statistical Databases*, vol. 2316 of LNCS. New York: Springer. 135-151.