# Matrix Masking Methods Which Preserve Moments

Ruben N. Mera

Census Bureau

September 30, 1997

## Abstract

The challenge of the problem of data disclosure avoidance lies in finding a compromise between protecting respondents' identities while providing the analyst with useful information. Statistics to be preserved most commonly include, means, variances and covariances by domains of the classification variables; regression coefficients; independence and measures of association; and parameters of loglinear and logistic regression models. In this paper, a linear transformation of the original data is presented, which preserves means, covariances, and eventually higher order moments.

## 1 Introduction

Several methods of disclosure limitation have been proposed in the literature. *Cell supression* was investigated by Greenberg and Zayatz [10], Cox [3], [4], and Carvalho et al. [1]. It involves the elimination of records at high risk of identification. When data is being summarized in cross-classification contingency tables, cells with few counts are deleted but the marginal totals of contingency tables are preserved. *Swapping* observations was proposed originally by Dalenius and Reiss [5], [6]. A variant of this method was implemented by the Census Bureau for the 1990 Census [12]. *Rounding* and *truncation* methods are discussed in [11]. *Data smoothing* and *imputation* [15], [9] consist in releasing a sample from a distribution with similar characteristics as the original data. Kim [13] proposed a method of *addition of random noise*, by which the transformed data mantains many of the statistical characteristics of the original data. *Inclusion of simulated data* has also been implemented in the past. Entire simulated rows are introduced to the data. Cell supression,

random noise, and inclusion of simulated data, are all particular cases of *matrix masking*, a technique that was introduced by Duncan and Pearson [8].

In the sequel, $\mathbf{X}$ will denote an $n \times p$ real valued matrix of the original data. Matrix masking involves the release of mapped data $\mathbf{Y} = \mathbf{AXB} + \mathbf{C}$ for some matrices $\mathbf{A}$, $\mathbf{B}$, and $\mathbf{C}$. In this paper we seek a transformation of the type $\mathbf{Y} = \mathbf{T}\mathbf{X}$, where $\mathbf{T}$ is a suitable $m \times n$ matrix, so that the transformation $\mathbf{T}$ preserves means, covariance matrix, and eventually, depending on $m$, $n$ and $p$, higher order moments. In other words,

$$\frac{1}{m}\sum_{i=1}^{m} y_{ij} = \frac{1}{n}\sum_{i=1}^{n} x_{ij} \quad and \qquad (1)$$

$$\frac{1}{m}\sum_{i=1}^{m} y_{ij}y_{ik} = \frac{1}{n}\sum_{i=1}^{n} x_{ij}x_{ik} \qquad (2)$$

$$j,k = 1,\ldots,p$$

or equivalently,

$$\frac{1}{m}\mathbf{Y}'\mathbf{1}_m = \frac{1}{n}\mathbf{X}'\mathbf{1}_n \quad and \qquad (3)$$

$$\frac{1}{m}\mathbf{Y}'\mathbf{Y} = \frac{1}{n}\mathbf{X}'\mathbf{X} \qquad (4)$$

where $\mathbf{1}_n = (1,1,\ldots,1)'$; the subscript shall be omitted if it is obvious from the context. Given a vector $\mathbf{u} = (u_1, u_2, \ldots, u_n)'$, $\|\mathbf{u}\|$ will denote its usual $\Re^n$ norm, $\|\mathbf{u}\| = (\sum u_i^2)^{1/2}$. Methods of random noise are presented in section 3. While somewhat restrictive in that the transformation $\mathbf{T}$ is given by an $n \times n$ matrix, it offers a simple, intuitive, and easy to implement methodology. In section 4 the general case is proposed.

## 2 Matrix Representations

Some basic results of linear algebra are reviewed here. For a more in depth reading, see for instance [7], [18], [20], [19], [17]. A square matrix $\mathbf{P}$ is said

---

445

to be *orthogonal* if $\mathbf{P}'\mathbf{P} = \mathbf{I}$. Viewed as a linear transformation, an orthogonal matrix produces a rotation of the coordinate axes, preserving therefore norms and angles. If $\mathbf{P}$ is orthogonal, so is $\mathbf{P}'$, and we have $\mathbf{P}^{-1} = \mathbf{P}'$. The spectral decomposition theorem asserts that any symmetric matrix $\mathbf{A}$ can be expressed as $\mathbf{A} = \mathbf{P}\mathbf{\Lambda}\mathbf{P}'$, where $\mathbf{\Lambda}$ is a diagonal matrix of eigenvalues of $\mathbf{A}$ and $\mathbf{P}$ is an orthogonal matrix whose columns are the standardized eigenvectors of $\mathbf{A}$. An $m \times n$ matrix $\mathbf{L}$, $m \leq n$ ($m \geq n$), is said to be column orthonormal (row orthonormal) if $\mathbf{A}'\mathbf{A} = \mathbf{I}$ ($\mathbf{A}\mathbf{A}' = \mathbf{I}$). Any $m \times n$ matrix $\mathbf{X}$, $m \geq n$, can be represented as $\mathbf{X} = \mathbf{L}\mathbf{\Lambda}^{1/2}\mathbf{P}'$, where $\mathbf{\Lambda}$ is an $n \times n$ diagonal matrix of the eigenvalues of $\mathbf{X}'\mathbf{X}$, $\mathbf{P}$ is an $n \times n$ orthogonal matrix of its standardized eigenvectors, and $\mathbf{L}$ is an $m \times n$ column orthonormal matrix whose columns consist of the $n$ normalized eigenvectos of $\mathbf{X}\mathbf{X}'$ associated with the $n$ largest eigenvalues of this matrix. This is called the singular value decomposition theorem. For any matrix $\mathbf{A}$ we have rank$[\mathbf{A}]$ = rank$[\mathbf{A}']$= rank$[\mathbf{A}\,\mathbf{A}']$ = rank$[\mathbf{A}'\,\mathbf{A}]$. For conformable matrices the nonzero eigenvalues of $\mathbf{A}\,\mathbf{B}$ are the same as those of $\mathbf{B}\,\mathbf{A}$. In particular, if $\mathbf{A}$ is $m \times n$ column orthonormal, then the $m$ nonzero eigenvalues of $\mathbf{A}\,\mathbf{A}'$ are all equal to 1. Given $\Omega$, a vector subspace of $\Re^n$, every $n \times 1$ vector $\mathbf{y}$ can be expressed uniquely in the form, $\mathbf{y} = \mathbf{u} + \mathbf{v}$, where $\mathbf{u} \in \Omega$ and $\mathbf{v} \in \Omega^{\perp}$, the last being the orthogonal complement of $\Omega$. There is a unique matrix $\mathbf{P}_{\Omega}$, called the projection onto $\Omega$, such that $\mathbf{u} = \mathbf{P}_{\Omega}\,\mathbf{y}$; this matrix is further symmetric and idempotent. The projection of $\Re^n$ onto $\Omega^{\perp}$ is given by $\mathbf{I} - \mathbf{P}_{\Omega}$. If $\Omega = \mathcal{R}[\mathbf{X}]$ (the range of $\mathbf{X}$), then $\mathbf{P}_{\Omega} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'$, where $(\mathbf{X}'\mathbf{X})^{-}$ is any generalized inverse of $\mathbf{X}'\mathbf{X}$.

# 3 Random Noise

The method of random noise can be viewed as a particular case of matrix masking in which $m$, the number of records of the transformed data, equals $n$, the number of records of the original data. Let $\mathbf{E} = \{\varepsilon_{ij}\}_{n \times p}$ be matrix of random noise. Set

$$\mathbf{Y} = \mathbf{X} + \mathbf{E}$$
$$y_{ij} = x_{ij} + \varepsilon_{ij}.$$

Which in general can be expressed as

$$\mathbf{Y} = \mathbf{T}_D\,\mathbf{X} = (\mathbf{I} - 2\mathbf{D})\,\mathbf{X}.$$

In this section we seek sufficient conditions for $\mathbf{T}_D$ to satisfy (3) and (4).

## 3.1 A general solution

If $\mathbf{T}$ is orthogonal, (4) is satisfied, because

$$\mathbf{Y}'\mathbf{Y} = (\mathbf{T}\mathbf{X})'(\mathbf{T}\mathbf{X}) = \mathbf{X}'\mathbf{T}'\mathbf{T}\mathbf{X} = \mathbf{X}'\mathbf{I}\mathbf{X} = \mathbf{X}'\mathbf{X}.$$

If $\mathbf{D}$ is symmetric and idempotent, thus an orthogonal projection, $\mathbf{T}_D = \mathbf{I} - 2\,\mathbf{D}$ is orthogonal, since

$$(\mathbf{I} - 2\mathbf{D})'(\mathbf{I} - 2\mathbf{D}) = \mathbf{I} - 4\mathbf{D} + 4\mathbf{D}'\mathbf{D} = \mathbf{I}.$$

It follows that the transformation $\mathbf{Y} = \mathbf{T}_D\mathbf{X} = (\mathbf{I} - 2\mathbf{D})\mathbf{X}$ satisfies (4). Assume further that $\mathbf{D}$ is an orthogonal projection into a subspace of $\mathbf{1}_n^{\perp}$. We have

$$\mathbf{Y}'\mathbf{1} = \mathbf{X}'(\mathbf{I} - 2\,\mathbf{D})'\mathbf{1} = \mathbf{X}'\mathbf{1},$$

hince (3) also holds.

## 3.2 An explicit solution

A relatively simple technique of construction of such a random noise matrix is as follows. Let $\varepsilon_{i1} = \varepsilon_{i2} = \cdots = \varepsilon_{ip} = \varepsilon_i$, where $\varepsilon = (\varepsilon_1, \ldots, \varepsilon_n)'$ is an arbitrary zero mean vector, so as (1) holds. Set $y_{ij} = x_{ij} - \lambda_j \varepsilon_i$ for some constants $\lambda_j$ to be determined. The $\lambda$s that satisfy (2) are solutions of the system of $p(p+1)/2$ equations

$$\lambda_j \lambda_k \sum_{i=1}^{n} \varepsilon_i^2 - \lambda_j \sum_{i=1}^{n} x_{ik}\,\varepsilon_i - \lambda_k \sum_{i=1}^{n} x_{ij}\,\varepsilon_i = 0$$

By straighforward calculations such solutions are

$$\lambda_j = \frac{2\sum_{i=1}^{n} x_{ij}\,\varepsilon_j}{\sum_{i=1}^{n} \varepsilon_i^2} = 2\,\frac{\langle \mathbf{x}_j, \varepsilon \rangle}{\|\varepsilon\|_2^2}. \tag{5}$$

This leads to the following

**Theorem 1** *For any $n \times 1$ zero mean vector $\varepsilon$ the transformation $\mathbf{T}_\varepsilon$ given by*

$$\mathbf{T}_\varepsilon\,\mathbf{X} = \left(\mathbf{I} - 2\,\frac{\varepsilon\,\varepsilon'}{\varepsilon'\,\varepsilon}\right)\mathbf{X} = (\mathbf{I} - 2\,\mathbf{D}_\varepsilon)\,\mathbf{X} \tag{6}$$

*preserves means and covariances.*

The proof is straightforward, since equation (6) follows directly from (5).

Remark 1. A trivial solution of the type $\lambda_j = 0$ leaves unchanged the column $\mathbf{X}_j$ of $\mathbf{X}$. To avoid trivial solutions, the vector $\varepsilon$ must be chosen so that $\langle \varepsilon, \mathbf{X}_j \rangle \neq 0$, $j = 1, \ldots, p$. If no column $\mathbf{X}_j$ has a constant values, such choice of $\varepsilon$ is always possible.

Remark 2. Regressions performed in either set will clearly produce identical results provided that no

quadratic term or interactions are included in the model. To mantain unchanged regression outcomes when quadratic terms or interactions are included, the new data must preserve third order moments as well.

$$\sum_{i=1}^{n} y_{ij} y_{ik} y_{il} = \sum_{i=1}^{n} x_{ij} x_{ik} x_{il}$$

The only constraint that the vector $\boldsymbol{\varepsilon}$ must satisfy is $\sum \varepsilon_i = 0$. The choice of $\boldsymbol{\varepsilon}$ has therefore $n-1$ "degrees of freedom". We can impose additional conditions to produce specific results. By a suitable choice of $\boldsymbol{\varepsilon}$, for instance, third order moments may be preserved if $n \geq p^3 + 1$; variables that are positive by design may be transformed into positive values; the change $|y_{ij} - x_{ij}| = |\lambda_j \varepsilon_i|$ may also be controlled. This change depends on the choice of the vector $\boldsymbol{\varepsilon}$. Large values of $|\lambda_j|$ are likely to produce larger differences. Since the values of the $\lambda$s depend on the choice of $\boldsymbol{\varepsilon}$, a good choice of this vector may increase the value of a particular $\lambda_j$, but this is not necessarily possible for all the $\lambda$s altogether (unless the variables $X_j$ are highly correlated), because this method does not provide a different choice of $\boldsymbol{\varepsilon}$ for each column of $\mathbf{X}$. Section 3.4 deals with the relationship between $\boldsymbol{\varepsilon}$ and $\lambda$. Later in section 4.2 a more general transformation is discussed, which allows for separate effects in each of the variables.

## 3.3   Categorical variables

Theorem 1 can readily be extended to the case in which the data set also contains categorical variables. Suppose that the data set contains $q$ discrete variables, $V_1, V_2, \ldots V_q$ in addition to the $p$ numerical variables $X_1, X_2, \ldots X_p$. Let $r$ be the number of cells in the cross-classification table by the $q$ categorical variables. By partitioning the data into $r$ subsets and applying this method to each subgroup, the last results will apply to each of the subdomians. In matricial form, suppose that $m_s$ is the number of observations in cell $s$, $s = 1, \ldots, r$. For any cell $s$ with $m_s \geq 2$ select a random vector $\boldsymbol{\varepsilon}_s = (\varepsilon_{s1}, \ldots, \varepsilon_{s,m_s})'$ as before and apply the transformation $\mathbf{T}_{\boldsymbol{\varepsilon}_s} X^*$ as in Theorem 1, where $X^*$ is the subset of observations in cell $s$ restricted to the numerical variables of the data set. Finally repeat this procedure to all cells with counts greater than 1. Clearly Theorem 1 applies in each transformed cell. To construct this mapping first sort $\mathbf{X}$ by the set of categorical vari-

ables and secondly construct the matrix

$$\mathbf{D}_v = \begin{pmatrix} D_{\varepsilon_1} & 0 & 0 & \ldots & 0 \\ 0 & D_{\varepsilon_2} & 0 & \ldots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \ldots & D_{\varepsilon_r} \end{pmatrix}$$

where $D_{\varepsilon_s}$ are $m_s$ by $m_s$ matrices defined by (6). This leads to

**Theorem 2** *If a data set contains $p$ numerical variables $X_1, X_2, \ldots, X_p$, and $q$ categorical variables $V_1, V_2, \ldots, V_q$, then the transformation*

$$\mathbf{Y} = \mathbf{T}_v \mathbf{X} = (\mathbf{I} - 2\mathbf{D}_v) \mathbf{X}$$

*preserves means and covariance matrices in every domain defined by levels of the categorical variables. In particular, results of regressions and analyses of covariance are preserved.*

## 3.4   Geometric Approach

Consider the $n \times p$ matrix $\mathbf{X}$ as $p$ points of $\Re^n$. For a given variable $X_j$, (2) implies that

$$\sum_{i=1}^{n} (x_{ij} - \lambda_j \varepsilon_i)^2 - \sum_{i=1}^{n} x_{ij}^2 = 0 .$$

It follows that

$$\sum_{i=1}^{n} (\lambda_j \varepsilon_i^2 - 2x_{ij}\varepsilon_i) = \langle \boldsymbol{\varepsilon} , \lambda_j \boldsymbol{\varepsilon} - 2\mathbf{x}_j \rangle = 0 .$$

Thus, as a point of $\Re^n$, $\lambda_j \boldsymbol{\varepsilon}$ must be located on the surface of the $n$-dimensional sphere $\mathcal{S}$ with center at $\mathbf{x}_j$ and radious $\|\mathbf{x}_j\|$. To satisfy (4), $\boldsymbol{\varepsilon}$ must lay in the hyperplane $\alpha$ passing through the origin and perpendicular to the vector $\mathbf{1}$. Hence, the vectors $\lambda_j \boldsymbol{\varepsilon}$ lay in the ellipsoid generated by the intersection of $\alpha$ with $\mathcal{S}$. When $var(\mathbf{x}_j) = 0$, $\mathbf{x}_j$ is colinear to the vector $\mathbf{1}$; the numerator in (5) is equal to zero and the intersection of $\alpha$ with $\mathcal{S}$ is restricted to one single point, the origin. In this case $\lambda_j = 0$ and the only solution of (1) and (2) is $\mathbf{x}_j$ itself. As $var(\mathbf{x}_j)$ increases, $\mathbf{x}_j$ moves away from $\mathbf{1}$ and there are increasing values of $\lambda_j$ that solve these equations. Given $\mathbf{x}_j$, the maximum value of $|\lambda_j|$ in the set of solutions of (1) and (2) is reached when the numerator in (5) is maximum, which in turn will be achieved when $\boldsymbol{\varepsilon}$ is the projection of $\mathbf{x}_j$ onto the plane $\alpha$, that is, when $\varepsilon_{ij} = x_{ij} - \bar{x}_j$, in which case $\lambda_j = 2$. It is clear that the $\varepsilon_i$ values that produce a large change in one variable $X_j$ should also produce large changes in all the other variables that are highly correlated with $X_j$. Variables uncorrelated with $X_j$, by

contrast, should produce only small changes. This suggests a more general transformation of the type

$$y_{ij} = x_{ij} - \lambda_j(x_{ij} - \bar{x}_j),$$

whose solution is $\lambda_1 = \lambda_2 = \cdots = \lambda_p = 2$. Hence,

$$y_{ij} = 2\bar{x}_j - x_{ij}.$$

By straightforward calculations it can be seen that the above transformation satisfies (1) and (2) and produces the maximum change in each of the variables. This approach alone is not of interest, since the original data can be reconstructed from the transformed data. There are variants of this case, however, which allow for large individual changes while preventing disclosure. This in turn requires proper definition of distance between two data sets.

## 3.5   Weighted data

If the records are associated with weights $w_i$, to account for different probabilities of selection, a similar procedure may be used. The weighted mean $\bar{x}_w$ and weighted variance $\sigma_w^2$ are respectively given by

$$\bar{x}_w = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i} \quad \text{and} \quad \sigma_w^2 = \frac{\sum_{i=1}^n w_i(x_i - \bar{x}_w)^2}{\sum_{i=1}^n w_i}.$$

Choose $\varepsilon$ so that $\sum w_i \varepsilon_i = 0$ and $\lambda_j = 2\sum w_i x_{ij}\varepsilon_i / \sum w_i \varepsilon_i^2$. The previous results apply also to the weighted data. The details are left to the reader.

## 3.6   Non-Uniqueness

To provide protection from disclosure, we must insure that, given a transformed data $\mathbf{Y}$, there are at least two different data sets $\mathbf{X}_i$ and two transformations $\mathbf{T}_{\varepsilon_i}$ such that $\mathbf{T}_{\varepsilon_i}\mathbf{X}_i = \mathbf{Y}$. It is easy to see that the inverse of $\mathbf{T}_\varepsilon$ is $\mathbf{T}_{-\varepsilon}$. Hence, for any transformed matrix $\mathbf{Y}$, if $\varepsilon_\nu$, $\nu \in \Gamma$, is the set of all zero mean $n$-dimensional vectors, any of the matrices $\mathbf{X}_{\varepsilon_\nu} = \mathbf{T}_{\varepsilon_\nu}\mathbf{Y}$ could be the one with the original data.

# 4   The general case

In this section we study a general linear transformation $\mathbf{Y} = \mathbf{T}\mathbf{X}$ given by an $m \times n$ matrix $\mathbf{T}$. When $m \neq n$, not only the values of the entries in the original data are changed, but also the number of records is modified, providing an extra protection against disclosure. The transformed records have completely lost their identity. Matching records between these two matrices is meaningless, since each record in the transformed data is obtained as a weighted average of all the records of the original data.

## 4.1   Preliminary Results

The following results will be needed in the sequel.

**Lemma 1** *Given two n-dimensional vectors* $\mathbf{a}$ *and* $\mathbf{b}$, *with* $\|\mathbf{a}\| = \|\mathbf{b}\|$, *there is an orthogonal matrix* $\mathbf{C}$ *such that* $\mathbf{C}\mathbf{a} = \mathbf{b}$.

*Proof.* Let $\mathbf{r}_2, \mathbf{r}_3, \cdots, \mathbf{r}_n$ be $n-1$ $n$-dimensional vectors such that the system $\mathbf{a}, \mathbf{r}_2 \cdots, \mathbf{r}_n$ is orthogonal and let $\mathbf{A}$ be the orthogonal matrix whose columns are these $n$ vectors. Define $\mathbf{B}$ in a similar way, for another orthogonal system $\mathbf{b}, \mathbf{s}_2 \ldots, \mathbf{s}_n$. It follows that

$$\mathbf{A}'\mathbf{a} = \mathbf{B}'\mathbf{b} = \begin{pmatrix} \|\mathbf{a}\|^2 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

The lemma holds with $\mathbf{C} = \mathbf{B}\mathbf{A}'$ $\square$

**Lemma 2** *Given 2 vectors,* $\mathbf{a}(n \times 1)$ *and* $\mathbf{b}(m \times 1)$, *with* $n \leq m$ *and* $\|\mathbf{a}\| = \|\mathbf{b}\|$, *there is a column orthonormal* $m \times n$ *matrix* $\mathbf{H}$ *such that* $\mathbf{H}\mathbf{a} = \mathbf{b}$.

*Proof.* Let $\mathbf{A}$ be an arbitrary $m \times n$ column orthonormal matrix and let $\mathbf{z} = \mathbf{A}\mathbf{a}$. Since

$$\|\mathbf{z}\|^2 = \mathbf{z}'\mathbf{z} = \mathbf{a}'\mathbf{A}'\mathbf{A}\mathbf{a} = \|\mathbf{a}\|^2 = \|\mathbf{b}\|^2,$$

By Lemma 1 there is an orthogonal matrix $\mathbf{C}$ such that $\mathbf{C}\mathbf{z} = \mathbf{b}$. The lemma holds for $\mathbf{H} = \mathbf{C}\mathbf{A}$. $\square$

**Lemma 3** *Let* $p \leq m \leq n$. *For any* $n \times p$ *column orthonormal matrix* $\mathbf{L}$ *there is an* $m \times n$ *row orthonormal matrix* $\mathbf{H}$ *such that* $\mathbf{H}\mathbf{L}$ *is column orthonormal. The matrix* $\mathbf{H}$ *can further be chosen so that* $\mathbf{H}'\mathbf{a} = \mathbf{b}$ *for any two given conformable vectors* $\mathbf{a}$ *and* $\mathbf{b}$, *provided that* $\|\mathbf{a}\| = \|\mathbf{b}\|$.

*Proof.* Since $p \leq m$, it is possible to chose an orthogonal system of $n \times 1$ vectors $v_1, v_2, \cdots, v_n$ such that the last $n - m$ vectors $v_{n-m+1}, \cdots, v_n$ are orthogonal to the $p$ column vectors of $\mathbf{L}$. Denote by $\mathbf{\Lambda}$ the $m \times n$ matrix whose first $m$ columns form the $m$-dimensional identity matrix and the last $n - m$ columns are zero, $\mathbf{\Lambda} = (\mathbf{I}_m \; 0)$, and let $\mathbf{P}$ be the orthogonal matrix whose columns are the $n$ vectors $v_i$. Define $\mathbf{H} = \mathbf{\Lambda}\mathbf{P}'$. It must be proved that $\mathbf{H}\mathbf{L}$ is column orthonormal. For this,

$$(\mathbf{H}\mathbf{L})'(\mathbf{H}\mathbf{L}) = \mathbf{L}'\mathbf{H}'\mathbf{H}\mathbf{L} = (\mathbf{L}'\mathbf{P})(\mathbf{\Lambda}'\mathbf{\Lambda})(\mathbf{P}'\mathbf{L}).$$

Note that

$$\Lambda'\Lambda = \begin{pmatrix} \mathbf{I}_n & 0 \\ 0 & 0 \end{pmatrix},$$

By the selection of the column vectors of $\mathbf{P}$, the last $n - m$ columns of $\mathbf{P}$, being otrhogonal to all $p$ row vectors of $\mathbf{L}'$, we conclude that the last $n - m$ columns of $\mathbf{L}'\mathbf{P}$ are zero. Therefore,

$$(\mathbf{L}'\mathbf{P})(\Lambda'\Lambda)(\mathbf{P}'\mathbf{L}) = (\mathbf{L}'\mathbf{P})(\mathbf{P}'\mathbf{L}) = \mathbf{I}$$

as we wanted to prove. Given now any two vectors $\mathbf{a}$ and $\mathbf{b}$ with $\|\mathbf{a}\| = \|\mathbf{b}\|$, let $\mathbf{H}'\mathbf{a} = \mathbf{z}$.

$$\mathbf{z}'\mathbf{z} = \mathbf{a}'\mathbf{H}\mathbf{H}'\mathbf{a} = \mathbf{a}'\mathbf{a} = \|\mathbf{a}\|^2 = \|\mathbf{b}\|^2$$

and by Lemma 1 there is an orthogonal matrix $\mathbf{C}$ such that $\mathbf{C}\mathbf{z} = \mathbf{b}$. Then choose $\mathbf{H}_1 = \mathbf{H}\mathbf{C}$. $\square$

## 4.2 General Transformation

In this section, a matrix $\mathbf{T}$ satisfying (3) and (4), which is not unique is found. The non-uniqueness of $\mathbf{T}$ is essential to prevent the reconstruction of $\mathbf{X}$ from $\mathbf{Y}$.

**Theorem 3** *Let $p \le n$, $m$. For any $n \times p$ matrix $\mathbf{X}$ of rank $p$ there is an $m \times n$ matrix $\mathbf{T}$ such that the transformation $\mathbf{Y} = \mathbf{T}\mathbf{X}$ satisfies (3) and (4).*

*Proof.* Suppose first that $m \ge n$. Let $\mathbf{H}$ be an $m \times n$ column orthonormal matrix such that $\mathbf{H}(n^{-1/2}\mathbf{1}_n) = m^{-1/2}\mathbf{1}_m$, whose existence is guaranteed by Lemma 2. Let $\mathbf{T} = \sqrt{m/n}\,\mathbf{H}$. Thus, $\mathbf{T}'\mathbf{T} = (m/n)\mathbf{I}$ and (4) readily follows. On the other hand,

$$\mathbf{1}_n = \frac{n}{m}\mathbf{T}'\mathbf{T}\mathbf{1}_n = \frac{n}{m}\mathbf{T}'\mathbf{1}_m\,.$$

Hence,

$$\frac{1}{m}\mathbf{Y}'\mathbf{1}_m = \frac{1}{m}\mathbf{X}'\mathbf{T}'\mathbf{1}_m = \frac{1}{n}\mathbf{X}'\mathbf{1}_n$$

and (3) also holds. Suppose now that $m < n$. By the singular value decomposition theorem, $\mathbf{X}$ can be expressed as $\mathbf{X} = \mathbf{L}\Lambda^{1/2}\mathbf{P}'$, where $\mathbf{P}$ is the orthogonal matrix of eigenvalues of $\mathbf{X}'\mathbf{X}$, $\Lambda$ is the diagonal matrix of its eigenvectors and $\mathbf{L}$ is column orthonormal. Let $\mathbf{H}$ be as in Lemma 3, i.e. $\mathbf{H}\mathbf{L}$ is column orthonormal and $\mathbf{H}'(m^{-1/2}\mathbf{1}_m) = n^{-1/2}\mathbf{1}_n$. Define $\mathbf{T} = \sqrt{n/m}\,\mathbf{H}$. By straightforward calculations, similar to those in the previous case, it can be seen that (3) and (4) are satisfied for this $\mathbf{T}$. $\square$

# References

[1] Carvalho, F. de, Dellaert, N. and Osorio M. de S. (1994) *Statistical Disclosure in two dimensional tables: General tables.* Journal of the American Statisttical Association, **89**, 1547-1557.

[2] Cox, L. and Sande, G. (1978). *Automated statistical disclosure control* American Statistical Association, Proceedings of the Section on Survey Research Methods, 177-182.

[3] Cox, L. (1980) *Supression methodology and Statistical disclosure control.* Journal of the American Statisttical Association, **75**, 191-194.

[4] Cox, L. (1995) *Network models for complementary cell supression.* Journal of the American Statisttical Association, **90**, 177-182.

[5] Dalenius, T. and Reiss, S.P. (1978). *Data-swapping: A techinque for disclosure control.* American Statistical Association, Proceedings of the Section on Survey Research Methods, 191-194.

[6] Dalenius, T. and Reiss, S.P. (1982). *Data-swapping: A techinque for disclosure control.* Journal of Statistical Planning and Inference, 6, 73-85.

[7] Dillon, William R. and Goldstein, Mathew (19994) *Multivariate Analysis*, Wiley.

[8] Duncan, G.T. and Pearson, R.B. (1991) *Enhancing Acces to Microdata While Protecting Confidenciality: Prospects for the Future.* Statistical Science 6, 219-239.

[9] Fineberg, S.E. (1994) *A radical proposal for the provision of micro-datass samples and the preservation of confidentiality.* Techincal Report No. 611, Department

[10] Greenberg, B.V. and Zayatz, L.V. (1992) *Strategies for measuring risk in public use microdata files.* Statistica Neerlandica **46**, 33-48.

[11] Griffin, R. and Thompson, J. (1987). *Confidentiality Techniques for the 1990 Census.* Presented at the Concurrent Session of the Joint Census Advisory Committee of October of 1987.

[12] Griffin, R., Navarro, A., and Flores-Baez, L. (1989). *Disclosure avoidance ofr the 1990 census.* American Statistical Association, Proceedings of the Section on Survey Research Methods, 516-521.

[13] Kim, Jay J. and Winkler, William E. (1995) *Masking Microdata Files* ASA, Proceedings of Section on Survey Research Methods.

[14] Navarro, A., Flores-Baez, L. and Thompson, J. (1988). *Results of data switching simulation.* Presented at the Spring meeting of the American Statistical Association and Population Statistics Census Advisory Committees.

[15] Rubin, D.B. (1987). *Multiple Imputation for Nonresponse in Surveys.* Wiley, New York.

[16] Rubin, D.B. (1983 ). *Discussion, statistical disclosure limitation.* Journal of official Statistics 9, 461-468.

[17] Schott, James R. (1997) *Matrix Analysis for Statistics* Wiley.

[18] Searle, S.R. (1971) *Linear Models* Wiley.

[19] Searle, S.R. (1982) *Matrix Algebra Useful for Statistics* Wiley.

[20] Seber, G.A.F. (1984). *Multivariate Observations* Wiley.