

# A Robust Approach of Regression-Based Statistical Matching for Continuous Data

Sooncheol Sohn<sup>1</sup> · Myoungshic Jhun<sup>2</sup>

<sup>1</sup>Department of Statistics, Korea University; <sup>2</sup>Department of Statistics, Korea University

(Received February 4, 2012; Revised March 20, 2012; Accepted March 30, 2012)

---

## Abstract

Statistical matching is a methodology used to merge microdata from two (or more) files into a single matched file, the variants of which have been extensively studied. Among existing studies, we focused on Moriarity and Scheuren's (2001) method, which is a representative method of statistical matching for continuous data. We examined this method and proposed a revision to it by using a robust approach in the regression step of the procedure. We evaluated the efficiency of our revised method through simulation studies using both simulated and real data, which showed that the proposed method has distinct advantages over existing alternatives.

Keywords: Donor file, recipient file, matched file, common variable, unique variable, statistical matching.

---

## 1. Introduction

Let  $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$  be a random vector with density  $f(\mathbf{x}, \mathbf{y}, \mathbf{z})$ , where  $\mathbf{X} = (X_1, \dots, X_p)$ ,  $\mathbf{Y} = (Y_1, \dots, Y_q)$  and  $\mathbf{Z} = (Z_1, \dots, Z_r)$  are vectors of random variables of dimension  $p, q$  and  $r$ , respectively.  $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$  is assumed to have a nonsingular distribution with mean vector  $\boldsymbol{\mu}$  and covariance matrix  $\Sigma$  as in (1.1).

$$(\boldsymbol{\mu}, \Sigma) = \left( \begin{pmatrix} \boldsymbol{\mu}_X \\ \boldsymbol{\mu}_Y \\ \boldsymbol{\mu}_Z \end{pmatrix}, \begin{pmatrix} \Sigma_{XX} & \Sigma_{XY} & \Sigma_{XZ} \\ \Sigma_{YX} & \Sigma_{YY} & \Sigma_{YZ} \\ \Sigma_{ZX} & \Sigma_{ZY} & \Sigma_{ZZ} \end{pmatrix} \right). \quad (1.1)$$

Suppose that  $A$  and  $B$  are two files consisting of  $n_A$  and  $n_B$ , independent and identically distributed observations generated from  $f(\mathbf{x}, \mathbf{y}, \mathbf{z})$ . Assume that  $\mathbf{X}$  is observed in both files and  $\mathbf{Y}$  appears only in one file ( $A$ ), while  $\mathbf{Z}$  is observed exclusively in the second file ( $B$ ) as in (1.2).  $\mathbf{X}$  are thus common variables, whereas  $\mathbf{Y}$  and  $\mathbf{Z}$  are unique variables.

$$\begin{aligned} \text{File } A : (\mathbf{x}_a^A, \mathbf{y}_a^A) &= (x_{a1}^A, \dots, x_{ap}^A, y_{a1}^A, \dots, y_{aq}^A), \quad a = 1, \dots, n_A, \\ \text{File } B : (\mathbf{x}_b^B, \mathbf{z}_b^B) &= (x_{b1}^B, \dots, x_{bp}^B, z_{b1}^B, \dots, z_{br}^B), \quad b = 1, \dots, n_B. \end{aligned} \quad (1.2)$$

---

This research was supported by Basic Science Research Program through the National Research Foundation of Korea(NRF) funded by the Ministry of Education, Science and Technology (2010-0009204).

<sup>2</sup>Corresponding author: Department of Statistics, Korea University, Seoul 136-701, Korea.

E-mail: [jhun@korea.ac.kr](mailto:jhun@korea.ac.kr)

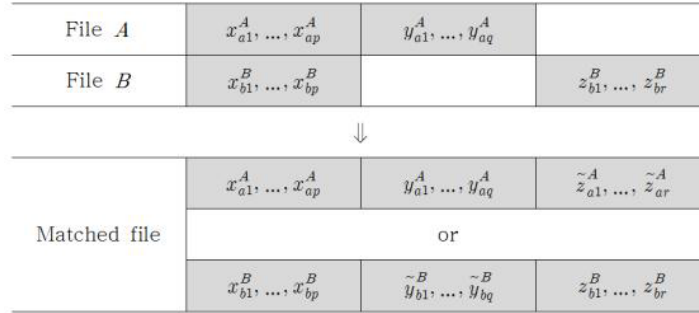


Figure 1.1. Illustration of statistical matching

Record linkage, or exact matching, is a methodology designed to merge the same entities from two (or more) different files. On the other hand, statistical matching, as a methodology designed to allocate microdata from similar entities, is used to merge microdata from two (or more) files into a single matched file (see Figure 1.1), where  $\tilde{z}_{ak}^A, k = 1, \dots, r$  and  $\tilde{y}_{bj}^B, j = 1, \dots, q$  denote the matched unique variables from file *A* and *B*, respectively.

Statistical matching and its variants have been discussed by Okner (1972, 1974), Sims (1972, 1974), Paass (1982), and Rodgers (1984). Of particular note, Kadane (1978) and Rubin (1986) described regression-based procedures to produce a matched file. When using Kadane’s (1978) method, Moriarity and Scheuren (2001) observed that the correlation coefficients between unique variables are not preserved during the matching procedure. Also, Moriarity and Scheuren (2003) indicated that Rubin (1986) did not consider the preservation of the correlation matrix structure, and asserted that the way in which his method estimated the secondary predicted values was redundant. In other words, it was proposed that Moriarity and Scheuren’s (2001) method could improve on the existing regression-based procedures of Kadane (1978) and Rubin (1986).

In this paper, we focused on Moriarity and Scheuren’s method as being paradigmatic of statistical matching for continuous data, and proposed a revision by using a robust approach in the regression step of the procedure. Furthermore, through simulation studies, using both simulated and real data, we showed that our proposed method represents an improvement on that of Moriarity and Scheuren. The rest of this paper is organized as follows: Section 2 briefly examines Moriarity and Scheuren’s method, highlights its problems, and proposes an improved statistical matching technique; Section 3 and Section 4 present the results of the simulation and case studies; and Section 5 concludes.

## 2. Statistical Matching Methods

In a statistical matching framework, we naturally assume the conditional independence of  $(\mathbf{Y}, \mathbf{Z})$ , given  $\mathbf{X}$ . This was assumed, either explicitly or implicitly, in all previous statistical matching applications. This assumption is usually referred to as the conditional independence assumption(CIA), and  $f(\mathbf{x}, \mathbf{y}, \mathbf{z}) = f_{\mathbf{X}}(\mathbf{x})f_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{x})f_{\mathbf{Z}|\mathbf{X}}(\mathbf{z}|\mathbf{x})$  when the CIA holds. Moreover, given the multivariate normality of  $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$ ,  $\Sigma_{\mathbf{YZ}} = \Sigma_{\mathbf{YX}}(\Sigma_{\mathbf{XX}})^{-1}\Sigma_{\mathbf{XZ}}$  and the conditional expectations of  $\mathbf{Z}$  and  $\mathbf{Y}$  in files *A* and *B* will be simplified as follows:

$$E(\mathbf{Z}|\mathbf{X}, \mathbf{Y}) = E(\mathbf{Z}|\mathbf{X}) = \boldsymbol{\mu}_{\mathbf{Z}} + \Sigma_{\mathbf{ZX}}\Sigma_{\mathbf{XX}}^{-1}(\mathbf{X} - \boldsymbol{\mu}_{\mathbf{X}}), \tag{2.1}$$

$$E(\mathbf{Y}|\mathbf{X}, \mathbf{Z}) = E(\mathbf{Y}|\mathbf{X}) = \boldsymbol{\mu}_{\mathbf{Y}} + \Sigma_{\mathbf{YX}}\Sigma_{\mathbf{XX}}^{-1}(\mathbf{X} - \boldsymbol{\mu}_{\mathbf{X}}). \tag{2.2}$$

The previous argument can easily be extended to multivariate linear regression models

$$\mathbf{y}_a^A = (\mathbf{X}_a^A)^T \boldsymbol{\beta}_Y + \boldsymbol{\epsilon}_a^A, \quad a = 1, \dots, n_A, \tag{2.3}$$

$$\mathbf{z}_b^B = (\mathbf{X}_b^B)^T \boldsymbol{\beta}_Z + \boldsymbol{\epsilon}_b^B, \quad b = 1, \dots, n_B, \tag{2.4}$$

where  $\mathbf{X}_a^A$  and  $\mathbf{X}_b^B$  are  $(p+1)q \times q$  and  $(p+1)r \times r$  block diagonal matrices representing  $p$  explanatory variables, with the diagonal elements being  $(1, \mathbf{x}_a^A)$  and  $(1, \mathbf{x}_b^B)$  and the off-diagonal elements zeros, respectively.  $\boldsymbol{\beta}_Y$  and  $\boldsymbol{\beta}_Z$  are  $(p+1)q \times 1$  and  $(p+1)r \times 1$  vectors of regression coefficients, and  $\boldsymbol{\epsilon}_a^A$  and  $\boldsymbol{\epsilon}_b^B$  are  $q \times 1$  and  $r \times 1$  error vectors which are distributed as multivariate normal with mean vectors  $\mathbf{0}$  and covariance matrices  $\Sigma_{YY}$  and  $\Sigma_{ZZ}$ , respectively. The conditional expectations (2.1) and (2.2) can then be estimated as (2.5) and (2.6), respectively.

$$\hat{\mathbf{z}}_a^A = (\mathbf{X}_a^A)^T \hat{\boldsymbol{\beta}}_Z, \quad a = 1, \dots, n_A, \tag{2.5}$$

$$\hat{\mathbf{y}}_b^B = (\mathbf{X}_b^B)^T \hat{\boldsymbol{\beta}}_Y, \quad b = 1, \dots, n_B, \tag{2.6}$$

where  $\hat{\boldsymbol{\beta}}_Z$  and  $\hat{\boldsymbol{\beta}}_Y$  are least squares estimators(LSEs) in the files  $B$  and  $A$ , respectively.

### 2.1. Moriarity and Scheuren's (2001) method

Moriarity and Scheuren (2001) defined file  $A$  as the donor file and file  $B$  as the recipient file in (1.2). They assumed that  $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$  have a nonsingular multivariate normal distribution with mean vector and covariance matrix as in (1.1), where  $p = q = r = 1$ . Their method consists of a regression step and a matching step, where the regression step produces estimates of the “missing” values as in (2.1) and (2.2), and the matching step produces the matched file. When the CIA holds, their method based on the models (2.3) and (2.4) can be summarized as follows:

<Moriarity and Scheuren's method>

- Regression step

Stage 1: In models (2.3) and (2.4), obtain the LSEs  $\hat{\boldsymbol{\beta}}_Y$  and  $\hat{\boldsymbol{\beta}}_Z$  for  $\boldsymbol{\beta}_Y$  and  $\boldsymbol{\beta}_Z$ , respectively.

Stage 2: Calculate the predicted values  $\hat{\mathbf{z}}_a^A = (\mathbf{X}_a^A)^T \hat{\boldsymbol{\beta}}_Z$ ,  $a = 1, \dots, n_A$  and  $\hat{\mathbf{y}}_b^B = (\mathbf{X}_b^B)^T \hat{\boldsymbol{\beta}}_Y$ ,  $b = 1, \dots, n_B$  in the donor and recipient files, respectively.

Stage 3: Add the random residuals  $\mathbf{r}_a^A \sim \text{MVN}(\mathbf{0}, \Sigma_{ZZ} - T_1)$  and  $\mathbf{r}_b^B \sim \text{MVN}(\mathbf{0}, \Sigma_{YY} - T_2)$  to  $\hat{\mathbf{z}}_a^A$  and  $\hat{\mathbf{y}}_b^B$  for  $a = 1, \dots, n_A$  and  $b = 1, \dots, n_B$ , respectively. If  $\Sigma_{ZZ} - T_1(\Sigma_{YY} - T_2)$  is not nonsingular, then  $\mathbf{r}_a^A = \mathbf{0}$  ( $\mathbf{r}_b^B = \mathbf{0}$ ), where  $T_1 = \Sigma_{ZZ} \Sigma_{XX}^{-1} \Sigma_{XZ}$  and  $T_2 = \Sigma_{YY} \Sigma_{XX}^{-1} \Sigma_{XY}$ .

- Matching step

Stage 1: Calculate the Mahalanobis distance in  $\{(\mathbf{Y}^A, \hat{\mathbf{Z}}^A + \mathbf{R}^A), (\hat{\mathbf{Y}}^B + \mathbf{R}^B, \mathbf{Z}^B)\}$ , where  $\mathbf{R}^A = (\mathbf{r}_1^A, \dots, \mathbf{r}_{n_A}^A)^T$  and  $\mathbf{R}^B = (\mathbf{r}_1^B, \dots, \mathbf{r}_{n_B}^B)^T$ .

Stage 2: Carry out the constrained match using RELAX-IV software (Bertsekas, 1991; Bertsekas and Tseng, 1994), which is able to conduct a constrained match by solving a “transportation” linear programming problem.

Improving on Kadane (1978), Moriarity and Scheuren added random residuals  $\mathbf{r}_a^A$  and  $\mathbf{r}_b^B$  to the  $\hat{\mathbf{z}}_a^A$  and  $\hat{\mathbf{y}}_b^B$  for Stage 3 of the regression step. When  $\Sigma_{ZZ} - T_1(\Sigma_{YY} - T_2)$  is nonsingular, in

contrast with Kadane (1978), the covariance matrices of  $(\mathbf{X}^A, \mathbf{Y}^A, \hat{\mathbf{Z}}^A + \mathbf{R}^A)$  and  $(\mathbf{X}^B, \hat{\mathbf{Y}}^B + \mathbf{R}^B, \mathbf{Z}^B)$  are nonsingular matrices and the expected values of these are equal to the population covariance matrix  $\Sigma$ . Furthermore, in order to preserve the covariance structure, Moriarity and Scheuren conducted a constrained match excepting common variable  $\mathbf{X}$  in the matching step. Consequently, compared with Kadane (1978) and Rubin (1986), they showed improved results, especially in terms of preservation of the correlation between  $\mathbf{Y}$  and  $\mathbf{Z}$ .

However, Moriarity and Scheuren assumed that  $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$  follows a multivariate normal distribution, which can lead to two problems. First, since the LSE in the regression step is known to be best under the normality assumption of  $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$ , the method may not be appropriate when the distribution of  $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$  has outliers or a heavy tail. Secondly, because the method used random residuals  $\mathbf{r}_a^A$  and  $\mathbf{r}_b^B$  from the normal distribution, if the distribution of  $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$  is not normal, these random residuals could be problematic.

## 2.2. Proposed method

First, in model (2.3), the estimator  $\hat{\beta}_{\mathbf{Y}}^{LDE}$  for  $\beta_{\mathbf{Y}}$  that minimizes over all possible values of  $\beta_{\mathbf{Y}}$

$$\sum_{a=1}^{n_A} \left\| \mathbf{y}_a^A - (\mathbf{X}_a^A)^T \beta_{\mathbf{Y}} \right\| \quad (2.7)$$

is called the least distance estimator(LDE), where  $\|\cdot\|$  denotes the usual Euclidean distance. In model (2.4), we can obtain the estimator  $\hat{\beta}_{\mathbf{Z}}^{LDE}$  for  $\beta_{\mathbf{Z}}$  analogously. The LDE is an estimator considering the relationship among response variables, and the relative efficiency of the LDE with respect to the least absolute estimator increases as the correlation between the response variables increases (Jhun and Choi, 2009). Additionally, Bai *et al.* (1990) found that the LDE is robust when the data contain outliers. Although the LSE is the most widely used estimator in regression modeling, it can be seriously affected when outliers exist and its estimation process fails to take into account the relationship among response variables. Thus, if the LSE is used, valuable information contained in the interdependence structure of the response variables may be overlooked, unlike when the LDE is employed. For these reasons, in statistical matching the LDE can be considered as an alternative to the LSE when the distribution of  $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$  has outliers or a heavy tail.

Next, consider the empirical residuals  $\mathbf{u}_a^A = \mathbf{y}_a^A - (\mathbf{X}_a^A)^T \hat{\beta}_{\mathbf{Y}}^{LDE}$ ,  $a = 1, \dots, n_A$  and  $\mathbf{u}_b^B = \mathbf{z}_b^B - (\mathbf{X}_b^B)^T \hat{\beta}_{\mathbf{Z}}^{LDE}$ ,  $b = 1, \dots, n_B$  in the donor and recipient files, respectively. Then,  $\mathbf{U}^A = (\mathbf{u}_1^A, \dots, \mathbf{u}_{n_A}^A)^T$  and  $\mathbf{U}^B = (\mathbf{u}_1^B, \dots, \mathbf{u}_{n_B}^B)^T$  are distribution-free and have the same mean vectors and covariance matrices as those of  $\mathbf{R}^A$  and  $\mathbf{R}^B$ , respectively. Because of these properties,  $\mathbf{U}^A$  and  $\mathbf{U}^B$  can be valid alternatives to the  $\mathbf{R}^A$  and  $\mathbf{R}^B$  even when the distribution of  $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$  is normal. Therefore, we propose an alternative statistical matching method which revises that of Moriarity and Scheuren as follows:

<Proposed method>

- Regression step

Stage 1: In models (2.3) and (2.4), obtain the LDEs  $\hat{\beta}_{\mathbf{Y}}^{LDE}$  and  $\hat{\beta}_{\mathbf{Z}}^{LDE}$  for  $\beta_{\mathbf{Y}}$  and  $\beta_{\mathbf{Z}}$ , respectively.

Stage 2: Calculate the predicted values  $\hat{\mathbf{z}}_a^A = (\mathbf{X}_a^A)^T \hat{\beta}_{\mathbf{Z}}^{LDE}$ ,  $a = 1, \dots, n_A$  and  $\hat{\mathbf{y}}_b^B = (\mathbf{X}_b^B)^T \hat{\beta}_{\mathbf{Y}}^{LDE}$ ,  $b = 1, \dots, n_B$  in the donor and recipient files, respectively.

Stage 3: Add randomly the empirical residuals  $\mathbf{u}_a^A = \mathbf{y}_a^A - (\mathbf{X}_a^A)^T \hat{\boldsymbol{\beta}}_{\mathbf{Y}}^{LDE}$  and  $\mathbf{u}_b^B = \mathbf{z}_b^B - (\mathbf{X}_b^B)^T \hat{\boldsymbol{\beta}}_{\mathbf{Z}}^{LDE}$  to  $\hat{\mathbf{y}}_b^B$  and  $\hat{\mathbf{z}}_a^A$  for  $a = 1, \dots, n_A$  and  $b = 1, \dots, n_B$ , respectively.

- Matching step

Stage 1: Calculate the Mahalanobis distance in  $\{(\mathbf{Y}^A, \hat{\mathbf{Z}}^A + \mathbf{U}^A), (\hat{\mathbf{Y}}^B + \mathbf{U}^B, \mathbf{Z}^B)\}$ .

Stage 2: For each  $b = 1, \dots, n_B$ , impute the live value  $y_{a^*}$  corresponding to the closest entity  $a^*$  in file  $A$  with respect to the Mahalanobis distance, *i.e.* the unconstrained match.

For matching, either the constrained or the unconstrained match can be used in the matching step, but we employ the unconstrained match in order to obtain the closest entity to the reality. In fact, the proposed method replaces the LSE used in that of Moriarity and Scheuren with the LDE in Stage 1, and also swaps random residuals for empirical residuals in Stage 3 of the regression step. That is to say, we propose a robust revision of Moriarity and Scheuren’s method. We thus expect that the proposed method will create a matched file which is closer to the reality than is produced by that of Moriarity and Scheuren, especially when the distribution of  $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$  is not normal.

### 3. Simulation Study

In order to compare the proposed method with that of Moriarity and Scheuren in performance, we carried out a simulation study. We considered two cases: one in which  $\mathbf{X}, \mathbf{Y}$  and  $\mathbf{Z}$  are all univariate (Case 1), and another in which  $\mathbf{X}, \mathbf{Y}$  and  $\mathbf{Z}$  are all multivariate (Case 2). For both cases, we set up the mean vector and covariance matrix of  $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$  as

$$(\boldsymbol{\mu}_1, \Sigma_1) = \left( \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \mathbf{1} & \Sigma_{\mathbf{XY}} & \Sigma_{\mathbf{XZ}} \\ \Sigma_{\mathbf{XY}} & \mathbf{1} & \Sigma_{\mathbf{YZ}} \\ \Sigma_{\mathbf{XZ}} & \Sigma_{\mathbf{YZ}} & \mathbf{1} \end{pmatrix} \right), \tag{3.1}$$

where  $\Sigma_{\mathbf{YZ}} = \Sigma_{\mathbf{XY}} \Sigma_{\mathbf{XZ}}$  and the elements of  $(\Sigma_{\mathbf{XY}}, \Sigma_{\mathbf{XZ}})$  are taken as random draws from a Uniform  $(0, 1)$  distribution, but  $\Sigma_1$  should be nonsingular.

#### 3.1. Case 1: $\mathbf{X}, \mathbf{Y}$ and $\mathbf{Z}$ are all univariate

We considered four types of distributions for  $(X_1, Y_1, Z_1)$ , where  $p = q = r = 1$ :

- i) a multivariate normal distribution  $MVN(\boldsymbol{\mu}_1, \Sigma_1)$ , where the mean vector  $\boldsymbol{\mu}_1$  and covariance matrix  $\Sigma_1$  are the same as in (3.1).
- ii) a contaminated multivariate normal distribution  $0.9MVN(\boldsymbol{\mu}_1, \Sigma_1) + 0.1MVN(\boldsymbol{\mu}_1, \Sigma_2)$ , where the mean vector  $\boldsymbol{\mu}_1$  and covariance matrix  $\Sigma_1$  are as in i) and  $\Sigma_2$  is the same as  $\Sigma_1$  except the variance of  $Y_1$  is 100
- iii) a multivariate  $t$  distribution with 3 degrees of freedom with covariance matrix  $\Sigma_1$
- iv) a multivariate Cauchy distribution

We generated  $n = 2000$  observations of  $(X_1, Y_1, Z_1)$  from each of the four types of distributions and randomly separated them into two files of the same size,  $n_A = n_B = 1000$ . We then deleted the variable  $Y_1$  in one file and  $Z_1$  in another, in order to create the recipient and donor files, respectively. We performed statistical matching by using both our proposed method and that of

**Table 3.1.** The means and the standard deviations(SD) of evaluation factors for Case 1

Dist.	Matching methods	$Y_1^D$		$\text{Corr}(X_1, \tilde{Y}_1)$		$\text{Corr}(\tilde{Y}_1, Z_1)$	
		Mean	SD	Mean	SD	Mean	SD
MVN	(Real value)	0	0	0.440	0.272	0.225	0.202
	M & S	1125	24	0.437	0.281	0.220	0.206
	Proposed	1127	22	0.436	0.279	0.221	0.214
CMN	(Real value)	0	0	0.135	0.201	0.077	0.129
	M & S	2473	38	0.117	0.191	0.171	0.191
	Proposed	2364	41	0.139	0.211	0.076	0.103
Mt(3)	(Real value)	0	0	0.487	0.329	0.213	0.279
	M & S	1656	31	0.363	0.292	0.315	0.326
	Proposed	1599	30	0.455	0.304	0.206	0.223
MC	(Real value)	0	0	0.446	0.681	0.207	0.567
	M & S	14325	20639	0.102	0.459	0.082	0.504
	Proposed	10470	5820	0.304	0.439	0.164	0.234

M & S: Moriarity and Scheuren's (2001) method

Proposed: the proposed method

MVN: multivariate normal

CMN: contaminated multivariate normal

Mt(3): multivariate  $t$  with 3 degrees of freedom

MC: multivariate Cauchy

Moriarity and Scheuren. For the assessment of sampling variation, this procedure was repeated 100 times independently.

Usually, the efficiency of a matched file is determined by two factors. First, how close the values of the matched unique variable are to the real values, and secondly, how similar the correlation coefficients between the matched unique variable and the other variables are to the real correlation coefficients. To assess the former, we calculated the sum of absolute differences between the real  $y_{b1}$  ( $b = 1, \dots, n_B$ ) in the original recipient file and the matched  $\tilde{y}_{b1}$  ( $b = 1, \dots, n_B$ ), denoting this sum as  $Y_1^D$ . For the latter, we compared the  $\text{Corr}(X_1, \tilde{Y}_1)$  and  $\text{Corr}(\tilde{Y}_1, Z_1)$  in the matched file with the real  $\text{Corr}(X_1, Y_1)$  and  $\text{Corr}(Y_1, Z_1)$ , respectively.

Table 3.1 shows the means and the standard deviations of  $Y_1^D$ ,  $\text{Corr}(X_1, \tilde{Y}_1)$ , and  $\text{Corr}(\tilde{Y}_1, Z_1)$  based on 100 independent repetitions. This shows that when the distribution of  $(X_1, Y_1, Z_1)$  is multivariate normal, for  $Y_1^D$ , the results from Moriarity and Scheuren's (2001) method and those from the proposed method are not significantly different. The correlation coefficients suggest that both methods well preserved the real values, as expected. However, when the distribution of  $(X_1, Y_1, Z_1)$  is not multivariate normal, the means of  $Y_1^D$  from the proposed method are smaller than those from Moriarity and Scheuren's method, and the correlation coefficients from the proposed method are closer to the real values in the proposed method. Generally, since the variables  $\mathbf{Y}$  and  $\mathbf{Z}$  are never jointly observed, researchers tend to be more interested in the preserving of the correlation of  $(\mathbf{Y}, \mathbf{Z})$  than that of  $(\mathbf{X}, \mathbf{Y})$ . In this respect, the proposed method has greater potential utility.

### 3.2. Case 2: $\mathbf{X}, \mathbf{Y}$ and $\mathbf{Z}$ are all multivariate

We extended the simulation study to a multivariate situation. We considered four types of distributions for  $(X_1, X_2, Y_1, Y_2, Z_1, Z_2)$ , where  $p = q = r = 2$ . The features of the distributions are similar to those in Case 1, where the variances of  $Y_j$ ,  $j = 1, 2$  equal 100 in  $\Sigma_2$  for a contaminated multivariate normal distribution. As in Case 1, we generated  $n = 2000$  observations of  $(X_1, X_2, Y_1, Y_2, Z_1, Z_2)$

**Table 3.2.** The means and the standard deviations (SD) of evaluation factors for Case 2

Dist.	Matching methods	$Y_1^D$	$Y_2^D$	Correlation coefficient of				Correlation coefficient of			
		Mean (SD)	Mean (SD)	$(X_1, \tilde{Y}_1)$ Mean (SD)	$(X_1, \tilde{Y}_2)$ Mean (SD)	$(X_2, \tilde{Y}_1)$ Mean (SD)	$(X_2, \tilde{Y}_2)$ Mean (SD)	$(\tilde{Y}_1, Z_1)$ Mean (SD)	$(\tilde{Y}_1, Z_2)$ Mean (SD)	$(\tilde{Y}_2, Z_1)$ Mean (SD)	$(\tilde{Y}_2, Z_2)$ Mean (SD)
MVN	(Real value)	0 (0)	0 (0)	0.422 (0.230)	0.422 (0.227)	0.418 (0.227)	0.421 (0.228)	0.239 (0.143)	0.244 (0.157)	0.243 (0.159)	0.240 (0.145)
	M & S	1049 (119)	1045 (103)	0.234 (0.176)	0.235 (0.176)	0.104 (0.115)	0.110 (0.126)	0.098 (0.105)	0.092 (0.089)	0.105 (0.121)	0.098 (0.099)
	Proposed	870 (178)	886 (175)	0.405 (0.226)	0.401 (0.223)	0.403 (0.223)	0.408 (0.221)	0.234 (0.140)	0.238 (0.151)	0.237 (0.154)	0.227 (0.140)
CMN	(Real value)	0 (0)	0 (0)	0.175 (0.101)	0.171 (0.102)	0.174 (0.102)	0.174 (0.107)	0.104 (0.074)	0.110 (0.076)	0.101 (0.078)	0.103 (0.081)
	M & S	1961 (212)	1949 (206)	0.076 (0.059)	0.075 (0.068)	0.032 (0.045)	0.037 (0.053)	0.038 (0.050)	0.038 (0.047)	0.035 (0.056)	0.040 (0.048)
	Proposed	1738 (223)	1754 (231)	0.152 (0.090)	0.162 (0.088)	0.159 (0.092)	0.159 (0.093)	0.102 (0.070)	0.100 (0.072)	0.100 (0.069)	0.104 (0.077)
Mt(3)	(Real value)	0 (0)	0 (0)	0.412 (0.257)	0.402 (0.258)	0.405 (0.250)	0.422 (0.242)	0.270 (0.222)	0.254 (0.200)	0.244 (0.201)	0.264 (0.223)
	M & S	1506 (189)	1527 (204)	0.244 (0.176)	0.252 (0.171)	0.108 (0.128)	0.122 (0.127)	0.115 (0.119)	0.117 (0.127)	0.120 (0.118)	0.121 (0.132)
	Proposed	1188 (266)	1284 (243)	0.349 (0.208)	0.347 (0.204)	0.362 (0.197)	0.360 (0.21)	0.227 (0.168)	0.221 (0.156)	0.212 (0.157)	0.231 (0.165)
MC	(Real value)	0 (0)	0 (0)	0.453 (0.564)	0.445 (0.568)	0.477 (0.537)	0.324 (0.581)	0.220 (0.693)	0.265 (0.645)	0.162 (0.657)	0.197 (0.637)
	M & S	33633 (101940)	50517 (20101)	0.135 (0.322)	0.151 (0.346)	0.104 (0.295)	0.070 (0.317)	0.037 (0.310)	0.006 (0.332)	0.106 (0.373)	0.106 (0.360)
	Proposed	9840 (11717)	9832 (8379)	0.176 (0.286)	0.281 (0.328)	0.243 (0.301)	0.194 (0.336)	0.140 (0.294)	0.125 (0.276)	0.109 (0.360)	0.156 (0.331)

M & S: Moriarity and Scheuren's (2001) method  
 Proposed: the proposed method  
 MVN: multivariate normal  
 CMN: contaminated multivariate normal  
 Mt(3): multivariate  $t$  with 3 degrees of freedom  
 MC: multivariate Cauchy

from each of the four types of distributions and randomly separated them into two files of the same size,  $n_A = n_B = 1000$ . We then deleted the variable  $(Y_1, Y_2)$  in one file and  $(Z_1, Z_2)$  in another, in order to create the recipient and donor files, respectively. We performed statistical matching by using the proposed method and that of Moriarity and Scheuren. For the reduction of sampling variation, this procedure was repeated 100 times independently, as in Case 1.

Table 3.2 shows the means and the standard deviations of the evaluation factors from 100 independent repetitions. Overall, the results are similar to those in Table 3.1. However, when the distribution of  $(X_1, X_2, Y_1, Y_2, Z_1, Z_2)$  is multivariate normal, the means of  $Y_1^D$  and  $Y_2^D$  based on the proposed method are lower than those based on Moriarity and Scheuren's method. Besides this, the means of  $\text{Corr}(X_i, \tilde{Y}_j)$ ,  $i, j = 1, 2$  and  $\text{Corr}(\tilde{Y}_j, Z_k)$ ,  $j, k = 1, 2$  based on the proposed method are closer to the real values than those based on Moriarity and Scheuren's method. Since the LDE takes into account the relationship among response variables, such results seem to hold even when the distribution of  $(X_1, X_2, Y_1, Y_2, Z_1, Z_2)$  is multivariate normal.

#### 4. Real Example

In order to demonstrate the usefulness of the proposed method in a real situation, the method was applied to the Boston Housing Data found in Harrison and Rubinfeld (1978). The data set

File <i>A</i> (donor file)	DIS ( $X_1$ ), LSTAT ( $X_2$ )	NOX ( $Y_1$ ), MEDV ( $Y_2$ )	
File <i>B</i> (recipient file)	DIS ( $X_1$ ), LSTAT ( $X_2$ )		RM ( $Z_1$ ), CRIM ( $Z_2$ )

Figure 4.1. The donor and recipient files using Boston Housing Data

Table 4.1. The means and the standard deviations (SD) of evaluation factors for the real example

Matching methods	$Y_1^D$ Mean (SD)	$Y_2^D$ Mean (SD)	Correlation coefficient of				Correlation coefficient of			
			$(X_1, Y_1)$ Mean (SD)	$(X_1, Y_2)$ Mean (SD)	$(X_2, Y_1)$ Mean (SD)	$(X_2, Y_2)$ Mean (SD)	$(Y_1, Z_1)$ Mean (SD)	$(Y_1, Z_2)$ Mean (SD)	$(Y_2, Z_1)$ Mean (SD)	$(Y_2, Z_2)$ Mean (SD)
(Real value)	0 (0)	0 (0)	-0.771 (0.011)	0.249 (0.044)	0.593 (0.029)	-0.738 (0.016)	-0.300 (0.031)	0.437 (0.044)	0.691 (0.038)	0.240 (0.145)
M & S	27 (1)	2345 (82)	-0.332 (0.055)	0.126 (0.054)	0.221 (0.058)	-0.101 (0.062)	-0.112 (0.060)	0.189 (0.061)	0.057 (0.058)	0.098 (0.099)
Proposed	19 (1)	1657 (84)	-0.662 (0.046)	0.225 (0.059)	0.539 (0.068)	-0.581 (0.049)	-0.268 (0.063)	0.286 (0.076)	0.383 (0.067)	0.227 (0.140)

M & S: Moriarity and Scheuren's (2001) method

Proposed: the proposed method

involves housing conditions of 506 observations, with 14 variables. We used 6 of these variables: DIS (weighted distances to five Boston employment centers), LSTAT (percentage of the population that is lower status), NOX (nitric oxides concentration), MEDV (median value of owner-occupied homes in \$1000s), RM (average number of rooms per dwelling), and CRIM (per capita crime rate by town). We randomly separated the data file into two equally sized files,  $n_A = n_B = 253$  to make the donor and recipient files as in Figure 4.1. Then, we performed statistical matching by using our proposed method and that of Moriarity and Scheuren. As in the previous simulation studies, this procedure was repeated 100 times independently.

Table 4.1 shows the means and the standard deviations for the evaluation factors, based on these 100 independent repetitions. Considering that the standard deviations of NOX ( $Y_1$ ) and MEDV ( $Y_2$ ) were 0.12 and 9.20 in the original complete data set, respectively, the means of  $Y_1^D$  and  $Y_2^D$  based on the proposed method were significantly lower than those based on Moriarity and Scheuren's method. Moreover, the means of  $\text{Corr}(X_i, \tilde{Y}_j)$ ,  $i, j = 1, 2$  and  $\text{Corr}(\tilde{Y}_j, Z_k)$ ,  $j, k = 1, 2$  based on the proposed method were much closer to the real values than those based on Moriarity and Scheuren's method.

## 5. Conclusion

In this research, we examined and proposed revisions to Moriarity and Scheuren's (2001) method, which is a representative method of statistical matching for continuous data. Their method improved on Kadane's (1978) and Rubin's (1986) methods by using the random residuals in the regression step and conducting a constrained match excepting the common variables in the matching step. However, since it assumed that  $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$  follows multivariate normal distribution, it may not be appropriate when the distribution of  $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$  has outliers or a heavy tail. This study's proposed method employs the LDE, which is an alternative to the LSE and robust to outliers, and takes into account the relationship among response variables, thus forming an alternative to the LSE. The



employment of empirical residuals as an alternative to the use of random residuals can also make our proposed method distribution-free. Through simulation studies using both simulated and real data, we found that the proposed method has distinct advantages over existing alternatives.

## References

- Bai, Z. D., Chen, X. R., Miao, B. Q. and Rao, C. R. (1990). Asymptotic theory of least distance estimate in multivariate linear model, *Statistics*, **21**, 503–519.
- Bertsekas, D. P. (1991). *Linear Network Optimization: Algorithms and Codes*, Massachusetts: MIT Press, Cambridge.
- Bertsekas, D. P. and Tseng, P. (1994). RELAX-IV: A faster version of the RELAX code for solving minimum cost flow problems. Available on the Internet at <http://web.mit.edu/dimitrib/www/home.html>.
- Harrison, D. and Rubinfeld, D. L. (1978). Hedonic prices and the demand for clean air, *Journal of Environmental Economics and Management*, **5**, 81–102.
- Jhun, M. and Choi, I. (2009). Bootstrapping least distance estimator in the multivariate regression model, *Computational Statistics & Data Analysis*, **53**, 4221–4227.
- Kadane, J. B. (1978). Some statistical problems in merging data files. 1978 Compendium of Tax Research, U.S. Department of the Treasury, 159–171. (Reprinted in *Journal of Official Statistics*, **17**, 423–433).
- Moriarity, C. and Scheuren, F. (2001). Statistical matching: a paradigm for assessing the uncertainty in the procedure, *Journal of Official Statistics*, **17**, 407–422.
- Moriarity, C. and Scheuren, F. (2003). A note on Rubin's statistical matching using file concatenation with adjusted weights and multiple imputations, *Journal of Business & Economic Statistics*, **21**, 65–73.
- Okner, B. (1972). Constructing a new data base from existing micro-data sets: The 1966 merge file, *Annals of Economic and Social Measurement*, **1**, 325–362.
- Okner, B. (1974). Data matching and merging: An overview, *Annals of Economic and Social Measurement*, **3**, 347–352.
- Paass, G. (1982). Statistical match with additional information. Internal Report IPES.82.0204, Gesellschaft für Mathematik und Datenverarbeitung, Bonn, W. Ger.
- Rodgers, W. L. (1984). An evaluation of statistical matching, *Journal of Business & Economic Statistics*, **2**, 91–102.
- Rubin, D. B. (1986). Statistical matching using file concatenation with adjusted weights and multiple imputations, *Journal of Business and Economic Statistics*, **4**, 87–94.
- Sims, C. A. (1972). Comments and rejoinder, *Annals of Economic and Social Measurement*, **1**, 343–345; 355–357.
- Sims, C. A. (1974). Comment, *Annals of Economic and Social Measurement*, **3**, 395–397.