

6-2004

Regression Analysis with Linked Data

P. Lahiri

Iowa State University

Michael D. Larsen

Iowa State University

Follow this and additional works at: http://lib.dr.iastate.edu/stat_las_preprints

 Part of the [Statistics and Probability Commons](#)

Recommended Citation

Lahiri, P. and Larsen, Michael D., "Regression Analysis with Linked Data" (2004). *Statistics Preprints*. 56.
http://lib.dr.iastate.edu/stat_las_preprints/56

This Article is brought to you for free and open access by the Statistics at Iowa State University Digital Repository. It has been accepted for inclusion in Statistics Preprints by an authorized administrator of Iowa State University Digital Repository. For more information, please contact digirep@iastate.edu.

Regression Analysis with Linked Data

Abstract

Record linkage, or exact matching, can be used to join together two files that contain information on the same individuals, but lack unique personal identification codes. The possibility of errors in linkage causes problems for estimating the relationships between variables on the two files. The effect is analogous to the impact of measurement error. A model of a linear regression relationship between variables in linked files is proposed. Assuming the probabilities that pairs of records are links are known, an unbiased estimator of the regression coefficients is derived. Methods for estimating the linkage probabilities by using mixture models are discussed. A consistent estimator of the covariance matrix of the proposed estimator is proposed. A bootstrap estimator is used to reflect the impact of the uncertainty in record linkage model parameters on the estimators of the regression parameters. A simulation study compares the performance of the proposed estimator and alternatives.

Keywords

Fellegi-Sunter, file matching, latent class, measurement error, mixture model, propagation of error, record linkage

Disciplines

Statistics and Probability

Comments

This preprint was published as P. Lahiri and Michael D. Larsen, "Regression Analysis with Linked Data", *Journal of the American Statistical Association* (2005): 222-230, doi: [10.1198/016214504000001277](https://doi.org/10.1198/016214504000001277)

Regression Analysis with Linked Data,

Iowa State University, Department of Statistics,
Preprint #04-9

P. Lahiri and Michael D. Larsen*

June 18, 2004

Abstract

Record linkage, or exact matching, can be used to join together two files that contain information on the same individuals, but lack unique personal identification codes. The possibility of errors in linkage causes problems for estimating the relationships between variables on the two files. The effect is analogous to the impact of measurement error. A model of a linear regression relationship between variables in linked files is proposed. Assuming the probabilities that pairs of records are links are known, an unbiased estimator of the regression coefficients is derived. Methods for estimating the linkage probabilities by using mixture models are discussed. A consistent estimator of the covariance matrix of the proposed estimator is proposed. A bootstrap estimator is used to reflect the impact of the uncertainty in record linkage model parameters on the estimators of the regression parameters. A simulation study compares the performance of the proposed estimator and alternatives.

Key words: Fellegi-Sunter; File matching; Latent class; Measurement error; Mixture models; Propagation of Error; Record Linkage.

1 Introduction

A goal of record linkage is to join together two files that contain information on the same individuals, but lack unique personal identification codes. Computerized record linkage (CRL) methods

*P. Lahiri is Professor, Joint Program in Survey Methodology, University of Maryland, College Park MD 20742, U.S.A. (E-mail: plahiri@survey.umd.edu) and Michael D. Larsen is Assistant Professor, Department of Statistics, Iowa State University, Ames IA 50011, U.S.A (E-mail: larsen@iastate.edu). The work of the first author was supported in part by National Science Foundation Grant SES-9978145 and a grant from the Gallup Organization. The work of the second author was supported in part by the U. S. Bureau of the Census through Census Contract No. 50-YABC-7-66021 under Census Task Order #46-YABC-8-00004. The authors thank the editor, an anonymous associate editor and two referees for constructive comments that led to a substantial improvement of an earlier version of the paper.

are used in many federal statistical systems (Alvey and Jamerson 1997) and often in medical studies (Newcombe 1988), in which the data bases are very large and processing time and accuracy are concerns. Sophisticated software has been developed for large applications by organizations including Statistics Canada (CANLINK software), the U.S. Census Bureau (Winkler 1994, 1995, and Jaro 1989, 1995), and the Oxford Medical Record Linkage Study (Gill 1997). Since CRL utilizes already existing databases, it enables new statistical analyses without the substantial time and resources needed to collect new data.

Fellegi and Sunter (1969), formalizing ideas of Newcombe *et al.*(1959), proposed a model for record linkage. In the Fellegi-Sunter (1969) model, the two files being compared are called File A and File B . The set of pairs of records $A \times B = \{(a, b), a \in A, b \in B\}$ is composed of two disjoint subsets: the set of true links, M , and the set of true nonlinks, U . Most CRL software attaches weights, similar in nature to weights described in Fellegi and Sunter (1969), reflecting the likelihood that a pair of records, one from each of the two files, corresponds to the same subject.

Mixture models are useful when the population being studied is composed of two or more subpopulations that are not clearly identified (McLachlan and Peel 2000). In the case of record linkage, before clerical review has been completed and in the absence of unique identifying information, the status of pairs as true links and true nonlinks is unknown, but real. Before clerical review is undertaken, mixture models can be applied to measurements of the similarity among pairs of records in order to estimate probabilities used in calculating record linkage weights. In some applications (Larsen and Rubin 2001, Winkler 1988, 1994, 1995, Jaro 1989, 1995), the mixture classes correspond very closely to the sets of true links and true nonlinks.

If mismatch errors are introduced by CRL, statistical analyses based on linked data can be adversely affected. Neter *et al.* (1965) studied the effect of mismatch errors in finite population sampling. They observed that relatively small mismatch error could lead to a substantial bias in estimating the relationship between response errors and true values. Scheuren and Winkler (1993), henceforth referred to as SW (1993), investigated the effect of mismatch errors on the bias of ordinary least squares estimators of regression coefficients in a standard regression model and proposed a method of adjusting for the bias. Scheuren and Winkler (1997) advanced the work further with an iterative procedure that modified the regression and matching results for apparent outliers. See also Scheuren and Winkler (1991).

In this paper we consider an alternative to the bias correction method of SW (1993). For known linkage probabilities, SW obtained their estimator of regression coefficient by adjusting the bias of the ordinary least square estimator for the regression model with mismatch errors, whereas our proposed method provides an unbiased estimator directly for a transformed regression model. In Section 2, we describe the record linkage problem and model. In section 3, we consider the use of mixture models for estimating relevant linkage probabilities and three implementation issues. In Section 4, we review the SW method and then propose a new method of estimating regression coefficients in the presence of mismatch errors. In Section 5, we propose a variance estimator for our regression estimator. We also discuss a bootstrap addition to the variance estimator to account for uncertainty in mixture model parameters. Simulation results are presented in Section 6. Our method improves on a naive, a robust, and the SW (1993) method. Technical proofs are deferred to the Appendix.

2 Record Linkage

In record linkage, the records in two files are compared to one another using available information, which typically does not include unique, error-free personal codes. Individuals can be compared on surname, first name, age or date of birth, and other variables. Some of these matching variables carry a lot of information for identifying individuals, whereas others (e.g., race or sex) contain very little. Some comparisons, however, are useful for discriminating between certain people, such as individuals living in the same household. Information can be missing or recorded with typographical or spelling errors.

The comparisons made on available fields of information result in measurements of agreement between the records in the two files. The outcome of agreement versus disagreement or of the level of correspondence measured in some manner (e.g., see Winkler 1990) on comparison k is stored in comparison variable γ_k . In the simplified case of dichotomous agreement/disagreement outcomes, let $\gamma_k = 1$ if the pair agrees on comparison k and 0 otherwise. The set of K comparisons creates a comparison vector $\gamma = (\gamma_k, k = 1, \dots, K)$ for each pair of records. For the special case of three available fields (i.e., $K = 3$), the possible comparison vectors are (0,0,0), i.e. all disagreements, (0,0,1), (0,1,0), (0,1,1), (1,0,0), (1,0,1), (1,1,0), and (1,1,1), i.e. all agreements.

The pattern of 0's and 1's in a comparison vector γ contains information about whether or not the pair is a true link or true nonlink. True links tend to have more agreements than do true nonlinks. If Γ is the space of all comparison vectors γ and the probabilities of seeing a vector γ among true links and true nonlinks is known, then Fellegi and Sunter's (1969) decision rule for designating pairs as links and nonlinks is based on the ratio

$$R = P(\gamma \in \Gamma|M)/P(\gamma \in \Gamma|U).$$

Intuitively, if R is large, pairs should be designated as links. On the other hand, if R is small, pairs should be called nonlinks. Some values of R , however, are moderate and do not clearly suggest link or nonlink. In practice, pairs with moderate values of R can be sent to clerks for review or can be subjected to further comparisons. Fellegi and Sunter (1969) showed that, at pre-specified error levels for false links and false nonlinks, optimal cutoffs can be determined. The cutoffs are optimal in that they minimize the set of pairs that are sent to clerical review for deciding link status at prespecified error levels. The decision rule can be characterized as the following:

If $R \geq upper$, then designate the pair as a link.

If $upper > R > lower$, then postpone the decision pending clerical review.

If $R < lower$, then designate the pair as a nonlink.

Three issues arise in practice. First, not all possible pairs of records are compared. Instead, pairs are compared within blocks of records that are similar in terms of basic characteristics, such as geography or first letter of last name. Forming blocks, or "blocking" as it is called, greatly reduces the number of pairs compared. If individuals rarely move between blocks, then few true links are lost by implicitly treating all pairs excluded by blocking as nonlinks.

Second, probabilities of comparison vectors by link status, $P(\gamma|M)$ and $P(\gamma|U)$, are not known; they must be estimated under a model using certain assumptions. Given prespecified error rates and estimates of these probabilities, the Fellegi-Sunter (1969) method determines corresponding values of *upper* and *lower*. The performance of the procedure in terms of actual versus specified error rates is sensitive to estimates of probabilities and choice of *upper* and *lower* (Belin 1993, Belin and Rubin 1995). Section 3 describes the use of mixture models for estimating these probabilities.

Third, for a record in File A there might be several candidate links within a particular block in File B. It is assumed in this work that only one of the records in File B is a true link for the record in File A. Given estimated probabilities, in practice, single links for individual records are chosen according to some procedure. Many applications, such as those at the U.S. Census Bureau (e.g., Jaro 1989) use a one-to-one, linear-sum assignment procedure (Burkard and Derigs 1980) to choose individual links. The one-to-one assignment procedure can effectively eliminate many candidate links that have some degree of similarity, but actually are nonlinks. On the other hand, forcing one-to-one matching could remove the true link if one member of the record pair has a better matching record in the other file. The possibility of false-matches and false-nonmatches has serious implications in many record linkage applications, such as counterterrorism (Gomatam and Larsen 2004).

3 Mixture Models

Let G be the number of subpopulations. In our application, we have $G = 2$ subpopulations - one consists of links and the other consisting of nonlinks. The comparison vector γ is assumed to follow a finite mixture model with probability mass function given by:

$$P(\gamma) = \sum_{g=1}^G \pi_g P(\gamma | \text{class } g),$$

where π_g is the probability that a pair of records belongs to the mixture class g and $P(\gamma | \text{class } g)$ is the probability mass function of the comparison vector in class g .

The model in each mixture class makes simplifying assumptions about the relationship between fields of comparison. For example, a common assumption suggested by Fellegi and Sunter (1969) is to assume the fields of comparison are independent within a given class. If there are K comparison fields that are conditionally independent, then in class g the probability of observing a comparison vector is

$$P(\gamma | \text{class } g) = \prod_{k=1}^K P(\gamma_k | \text{class } g),$$

where $P(\gamma_k | \text{class } g)$ is the probability of outcome γ_k on comparison k in class g . Other modeling assumptions are possible and, in some cases, better correspond to the observed data (Larsen and Rubin 2001, Armstrong and Mayda 1993, Thibaudeau 1993). A few authors in other contexts have used mixture models applied to discrete data with modeling assumptions other than conditional independence (see Becker and Yang 1998, references therein, and references in Larsen and Rubin 2001).

The parameters of the mixture model can be estimated using the Expectation-Maximization (EM; Dempster, Laird, and Rubin 1977) and Expectation-Conditional maximization (ECM; Meng and Rubin 1993) algorithms. Several authors, including Larsen and Rubin (2001) and references therein, have implemented these algorithms for the purposes of record linkage. The estimated probability that the pair that produced the comparison vector γ belongs to class g is, by Bayes' Theorem,

$$P(\text{class } g|\gamma) = \pi_g P(\gamma|\text{class } g) / \sum_{h=1}^G \pi_h P(\gamma|\text{class } h).$$

As described in Larsen and Rubin (2001), the estimated probabilities can be used to partition the record pairs into designated links and nonlinks and to estimate error rates. Larsen and Rubin (2001) study model selection and use partial clerical review to improve estimation.

Certain logical inequality constraints can be incorporated into the estimation algorithms. For example, the estimated size of the mixture class should be less than the size of the smaller of the two files, file A or B . Additionally, the probability of agreeing on a comparison should be higher among links than among nonlinks. These inequality constraints are implemented as part of the estimation in the simulation. Winkler (1993) has considered more extensive restrictions on parameters and estimated probabilities.

4 Estimation of Regression Coefficients

After File A and File B have been linked together, it might be of interest to analyze the relationship between a response variable (y) that was originally in File A and a set of covariates (x) that were originally in File B . Neter, Maynes, and Ramanathan (1965), Scheuren and Winkler (1993, 1997), Larsen (1999, 2001), and Lahiri and Larsen (2000), henceforth referred to as LL, discussed this problem. Table 1 illustrates the structure of the files. In the illustration, the matching variables v_k and w_k , $k = 1, \dots, K$, are used for comparison of record pairs. These two files are linked using a computerized record linkage technique. Thus, the true data pairs (x_i, y_i) are not observable. Instead, the record linkage procedure produces pairs (x_i, z_i) in which z_i may or may not correspond to y_i .

Consider the following regression model for $y = (y_1, \dots, y_n)'$:

$$y_i = x_i' \beta + \epsilon_i, i = 1, \dots, n, \quad (1)$$

where $x_i = (x_{i1}, \dots, x_{ip})'$ is a column vector of p known covariates, $\beta = (\beta_1, \dots, \beta_p)'$ is a column vector of p unknown regression coefficients, $E(\epsilon_i) = 0$, $\text{Var}(\epsilon_i) = \sigma^2$, and $\text{Cov}(\epsilon_i, \epsilon_j) = 0$ for $i \neq j$, $i, j = 1, \dots, n$. SW (1993) considered the following model for $z = (z_1, \dots, z_n)'$ given y :

$$z_i = \begin{cases} y_i & \text{with probability } q_{ii}, \\ y_j & \text{with probability } q_{ij} \text{ for } j \neq i, j = 1, \dots, n, \end{cases} \quad (2)$$

where $\sum_{j=1}^n q_{ij} = 1$, $i = 1, \dots, n$. Define $q_i = (q_{i1}, \dots, q_{in})'$, $i = 1, \dots, n$, and $Q = (q_1, \dots, q_n)'$. The naive least-squares estimator of β , which ignores mismatch errors, is given by

$$\hat{\beta}_N = (X'X)^{-1} X'z,$$

where $X = (x_1, \dots, x_n)'$ is a $n \times p$ matrix. An alternative to this naive estimator would be to use a robust estimator such as an estimator that minimizes the sum of absolute deviations (Press *et al.* 1992; page 698). Robust estimators that decrease the influence of outliers should decrease the impact of erroneously paired predictor and response values.

First note that under the model described by (1) and (2),

$$E(z_i) = w_i' \beta,$$

where $w_i = q_i' X = \sum_{j=1}^n q_{ij} x_j$, $i = 1, \dots, n$, is a $p \times 1$ column vector (see Theorem A.1). Thus, bias of $\hat{\beta}_N$ is given by

$$\text{Bias}(\hat{\beta}_N) = E[\hat{\beta}_N - \beta] = [(X'X)^{-1} X'W - I] \beta = [(X'X)^{-1} X'QX - I] \beta,$$

where I is an identity matrix of dimension p and $W = (w_1, \dots, w_n)'$. Thus, in general $\hat{\beta}_N$ is not an unbiased of β and the magnitude of the bias depends on W and β . Evidently, in the special case of no mismatch error, i.e. when $q_{ii} = 1$, $i = 1, \dots, n$, $\text{Bias}(\hat{\beta}_N) = 0$.

In order to reduce the bias of $\hat{\beta}_N$, SW (1993) first investigated the bias of $\hat{\beta}_N$ conditional on the values of y . It can be seen from the calculations of SW that

$$\text{Bias}(\hat{\beta}_N | y) = E[(\hat{\beta}_N - \beta) | y] = (X'X)^{-1} X' B, \quad (3)$$

where $B = (B_1, \dots, B_n)'$ and $B_i = (q_{ii} - 1)y_i + \sum_{j \neq i} q_{ij} y_j = q_i' y - y_i$. Thus B_i is the difference between a weighted average of responses from all observations and the actual response y_i . If an estimator of B , say \hat{B} , is available, then the SW estimator is given by

$$\hat{\beta}_{SW} = \hat{\beta}_N - (X'X)^{-1} X' \hat{B}. \quad (4)$$

Let q_{ij_1} and q_{ij_2} denote the highest and the second highest elements of the vector q_i respectively and let z_{j_1} and z_{j_2} denote the corresponding elements in the vector z . Then the estimator of B_i given by SW (1993) is $\hat{B}_i^{\text{TR}} = (q_{ij_1} - 1)z_{j_1} + q_{ij_2}z_{j_2}$, where TR denotes ‘‘truncation.’’ Note that $\hat{\beta}_{SW}$ is not an unbiased estimator of β in general. However, if the probability is high that the best candidate link is the true link, then the truncation might produce a very small bias. We can suggest a further modification that would use as many candidate links as necessary ($2 \leq s \leq n$) in the calculation of B_i so that $\sum_{r=1}^s q_{ij_r} > p$, where $p = 0.90$ or similar cumulative probability.

Since $E(z) = W\beta$, we can propose an exactly unbiased estimator of β as follows:

$$\hat{\beta}_U = (W'W)^{-1} W' z.$$

As with the SW estimator, it is possible to ‘‘truncate’’ our estimator to utilize only the best candidate links for a record. Instead of w_i , such a procedure could use $w_i^{\text{TR}} = q_{ij_1} x_{j_1} + q_{ij_2} x_{j_2}$.

Note that Q , or equivalently W , is a function of the parameters of the mixture distribution defined in Section 3, i.e. $\hat{Q} = Q(\psi)$ or $\hat{W} = W(\psi)$, where $\psi = \{(P(\gamma|M), \gamma \in \Gamma), (P(\gamma|U), \gamma \in \Gamma), \pi_M\}$. Thus, we can write $\hat{\beta}_{SW} = \hat{\beta}_{SW}(\psi)$ and $\hat{\beta}_U = \hat{\beta}_U(\psi)$. In practice ψ is unknown and a reasonable estimator $\hat{\psi}$ (e.g., maximum likelihood estimation) of ψ is used. In this case, we obtain the SW and our estimator as $\hat{\beta}_{SW}(\hat{\psi})$ and $\hat{\beta}_U(\hat{\psi})$ respectively. Interesting, $\hat{\beta}_U(\hat{\psi})$ is an unbiased estimator of β whenever $\hat{\psi}$ can be assumed to be independent of z . Such a situation is expected in most applications since the distribution of the matching variables (e.g., last name, phone number, etc.) which determines the distribution of $\hat{\psi}$ is usually independent of the response variable y (e.g., income) and hence of z .

5 Variance Estimation

5.1 Estimation of $\text{Var}(\hat{\beta}_{SW})$

SW (1993) suggested estimating the variance of their estimator, i.e. $\hat{\beta}_{SW}(\hat{\psi})$, by modifying the usual least-squares regression variance estimators. The usual variance estimator of $\hat{\beta}_{yx} = (X'X)^{-1}X'y$ is $\hat{\sigma}_0^2(X'X)^{-1}$, where $(n-p)\hat{\sigma}_0^2$ is given by

$$\begin{aligned} (y - X\hat{\beta}_{xy})'(y - X\hat{\beta}_{xy}) &= y'y - \hat{\beta}'_{xy}X'y \\ &= y'y - \hat{\beta}'_{xy}X'X\hat{\beta}_{xy}. \end{aligned} \quad (5)$$

Since $z = y + B$, it is readily seen that $y'y = z'z - 2B'z + B'B$. The SW (1993) estimator of σ^2 is taken to be

$$\hat{\sigma}_{SW}^2 = (z'z - 2B'z + B'B - \hat{\beta}'_{SW}X'X\hat{\beta}_{SW})/(n-p).$$

The corresponding estimator of $\text{Var}(\hat{\beta}_{SW})$ is taken to be $\hat{\sigma}_{SW}^2(X'X)^{-1}$. Clearly, the efficiency of this variance estimator depends on the degree of agreement between the vectors z and y , reflected by the elements of the matrix Q . When it is difficult to match two files, this variance estimator is likely to perform poorly. The estimator of the truncated version of the SW estimator will use \hat{B}_i^{TR} in place of \hat{B}_i and $\hat{\beta}_{SW}^{\text{TR}}$ in place of $\hat{\beta}_{SW}$.

5.2 Estimation of $\text{Var}(\hat{\beta}_U)$

First consider the case when ψ is known. Note that

$$\text{Var}(\hat{\beta}_U) = (W'W)^{-1}W'\Sigma W(W'W)^{-1}, \quad (6)$$

where $\text{Var}(z) = \Sigma = ((\sigma_{ij}))$ and the expressions for σ_{ij} are given in Theorem A.1. We stress that $\Sigma = \Sigma(\beta, \sigma^2, \psi)$, i.e. it depends on both the parameters of the regression model (1), i.e. β and σ^2 , and those of the mixture model, i.e. ψ . Since in this case ψ is known, we simply replace β and σ^2 by their estimators to obtain a variance estimator in the known ψ case. For example, we can use the unbiased estimator $\hat{\beta}_U$ to estimate β .

We shall now consider an alternate estimator of σ^2 . A naive estimator of mean squared error (MSE) based on z is given by

$$\text{MSE} = \frac{1}{n-p}z'(I - X(X'X)^{-1}X')z.$$

Note that MSE is not unbiased for σ^2 under the model described by (1) and (2); it would be unbiased for σ^2 if the real data were (X, z) . In order to obtain an alternate estimator of σ^2 , consider

$$S^2 = z'(I - W(W'W)^{-1}W')z,$$

where I is an n -dimensional identity matrix. According to Theorem A.2,

$$\text{E}(S^2) = (n-p)\sigma^2 + \text{tr}[(I - W(W'W)^{-1}W')H], \quad (7)$$

where $H = ((h_{ij}))$ with $h_{ii} = \beta' A_i \beta$ and $h_{ij} = \beta' A_{ij} \beta$.

Equation (7) motivates us to consider the following estimator of σ^2 :

$$\hat{\sigma}^2 = \max \left(0, \frac{S^2 - \text{tr}[(I - W(W'W)^{-1}W')\hat{H}]}{n - p} \right),$$

where $\hat{H} = ((\hat{h}_{ij}))$ with $\hat{h}_{ii} = \hat{\beta}'_u A_i \hat{\beta}_u$ and $\hat{h}_{ij} = \hat{\beta}'_u A_{ij} \hat{\beta}_u$. It can be shown that $\hat{\sigma}^2$ is consistent for σ^2 under the model described by (1) and (2) and mild regularity conditions (see Theorem A.3).

Now consider the most practical situation when ψ is unknown. In this case, one may naively use $\text{Var}(\hat{\beta}_U)$ with estimated β , σ^2 and ψ as a variance estimator of $\hat{\beta}_U(\hat{\psi})$. However, this variance estimator fails to incorporate the uncertainties due to the estimation of ψ and thus underestimates the true variability of $\hat{\beta}_U(\hat{\psi})$. The same comment applies to the variance estimator of $\hat{\beta}_{SW}(\hat{\psi})$ given in the previous subsection.

Parametric bootstrap methods have been quite effective in providing accurate variance estimators in many complex settings. For example, see Lahiri (2003) for a review of parametric bootstrap methods for complex multi-level models. We now develop a parametric bootstrap method in our context to obtain a reliable variance estimator that captures the additional variability due to the estimation of ψ . Under such a method, we draw B bootstrap samples from the mixture distribution with ψ replaced by $\hat{\psi}$. A single sample is a table of size 2^K of counts. Let $\hat{\psi}^*$ denote the estimator of ψ obtained from the procedure used for $\hat{\psi}$ but based on the bootstrap sample instead of the original sample. Let $\hat{\beta} = \hat{\beta}(\hat{\psi})$ denote any arbitrary estimator of β which depends on $\hat{\psi}$. Then we propose our bootstrap variance estimator as:

$$v_{\text{boot}} = E_{\star}[\text{var}(\hat{\beta}(\hat{\psi}^*))] + V_{\star}[\hat{\beta}(\hat{\psi}^*)], \quad (8)$$

where E_{\star} and V_{\star} denote the expectation and the variance with respect to the bootstrap distribution, $\text{var}[\hat{\beta}(\hat{\psi}^*)]$ is an estimator of $\text{Var}(\hat{\beta})$ with $\hat{\psi}^*$ substituted for ψ . In practice, we propose to use the Monte Carlo method to approximate the bootstrap expectation and variance. Thus,

$$E_{\star}[\text{var}(\hat{\beta}(\hat{\psi}^*))] \approx \frac{1}{B} \sum_{b=1}^B \text{var}[\hat{\beta}(\hat{\psi}^{*b})]$$

and

$$V_{\star}[\hat{\beta}(\hat{\psi}^*)] \approx \frac{1}{B} \sum_{b=1}^B [\hat{\beta}(\hat{\psi}^{*b}) - \hat{\beta}(\hat{\psi})][\hat{\beta}(\hat{\psi}^{*b}) - \hat{\beta}(\hat{\psi})]',$$

where $\hat{\psi}^{*b}$ is the estimator of ψ from the b^{th} bootstrap sample, $b = 1, \dots, B$.

6 A Monte Carlo Simulation

In this section, we use a Monte Carlo simulation to investigate the performances of different estimators of a regression coefficient and the associated variance estimators for a simple linear

regression model in the presence of mismatch errors. Our simulation study includes the naive estimator, $\hat{\beta}_N$, the Scheuren-Winkler estimator, $\hat{\beta}_{SW}(\hat{\psi})$, our proposed estimator, $\hat{\beta}_U(\hat{\psi})$, and a robust estimator mentioned in Section 4. The simulation conditions are first described and then results are presented.

6.1 Simulation Conditions

Four hundred replications are performed under each of two sets of conditions. Table 2 describes the main conditions. In both sets of conditions, the sizes of the files vary between 200 and 10000 records, but are the same for files A and B . The regression slope β varies between 0.20 and 0.80 with the simulated data generated based on a regression model having error variance σ^2 equal to $1-\beta^2$. The X variable is univariate normal with mean zero and variance 1.

In case 1 files A and B have eight to twelve matching fields, whereas in case 2 they have six to ten. Agreements on the fields of information are independent of one another. The probability of agreement among matches varies between 0.55 and 0.95 in case 1 and between 0.55 and 0.85 in case 2. The probability of agreement among nonmatches varies between 0.10 and 0.50 in case 1 and between 0.20 and 0.50 in case 2. The size of blocks affects how many potential links there are between the two files. Blocks are assumed to be linked together correctly, as they would be if they correspond to geographical areas. Pairs from different blocks are nonlinks and not used to estimate probabilities. Block sizes in case 1 range from 10 to 40 records (100 to 1600 potential links per block), whereas in case 2 they range from 20 to 40 records. Thus, Case 2 yields more nonmatches than Case 1, allowing us to understand the effect of nonmatches on different estimation methods.

The files A and B were generated and comparison vectors calculated. The EM algorithm (Dempster, Laird, Rubin 1977) was used to fit a two-class conditional independence mixture model to the comparison vectors to estimate probabilities for the Fellegi-Sunter (1969) algorithm. One product of the EM algorithm in this case are weights that represent the likelihood that a pair of records is a match. Estimated Fellegi-Sunter weights for links and nonlinks overlap more in case 2 than in case 1. The inequality constraints were used in the estimation, but one-to-one assignment was not enforced. It is not entirely clear how to force one-to-one matches and consider probabilities of matching in which two records in one file have a nonzero probability of matching a record in the second file. The use of one-to-one assignment in the analysis of linked files will be studied in future work.

6.2 Simulation Results

We compute four estimates of the slope for each of the four hundred simulation runs and compare with the true slope in terms of the absolute and squared deviations. We then compute average absolute deviation (AAD) and average squared deviation (ASD) for each of the four estimators, the average being taken over the four hundred simulation runs. Our proposed estimator outperforms all the rival estimators in all cases. In order to summarize our results, we define the percent

improvement with respect to AAD of our proposed estimator $\hat{\beta}_U$ over a rival estimator $\hat{\beta}$ as

$$100 \times \frac{AAD_{\hat{\beta}} - AAD_{\hat{\beta}_U(\hat{\psi})}}{AAD_{\hat{\beta}_U(\hat{\psi})}},$$

the percent improvement with respect to ASD is defined similarly.

Table 3 displays the percent relative improvements with respect to both AAD and ASD. The performance of our estimator is impressive. The naive estimator has the worse performance followed by the robust and the Scheuren-Winkler estimators. Since the second set of conditions had less powerful matching information and more difficult settings (e.g., larger blocks), matching was less successful. As expected the performances of all the three rival estimators relative to our proposed estimator get worse in this situation.

The coverage, reported in Table 4, is the percentage of times out of 400 that the following form of a nominal 95% confidence interval,

$$\text{estimate} \pm 2SE,$$

covers the true regression slope. The naive and the robust confidence intervals have the worst coverages. The confidence interval based on the Scheuren-Winkler method improves the coverage with respect to both the naive and the robust methods but it is considerably worse than our proposed method. All the three rival methods are sensitive to the simulation condition. In contrast our method is very stable under different simulation conditions.

For each data set in the simulation, 400 bootstrap comparison vector sets were generated. For each of these, the maximum likelihood estimates of the mixture model parameters were found. Based on the matching probabilities determined by these mixture model parameters, the regression estimates were recomputed. When the bootstrap procedure is used, the coverage of the proposed method improves. The other estimators are hardly affected. The naive and robust estimators, as implemented here, do not use the estimated probabilities determined by the mixture models directly in either estimation or variance estimation. They rely simply on the x- and y-values associated with the best matches. Although the probability estimates change slightly with each bootstrap, the best matches are rarely changed. It would be possible to compute the actual variance of these estimators under the model of (1) and (2). If one were relying on a naive estimator, however, in practice one would not do so. These estimators also are affected by severe bias.

The SW estimator uses the weights, but only in an estimate of bias. The SW estimator of variance used here is the one suggested by SW (1993). The actual variance under the model of (1) and (2) would be different. As such, the variance estimate is largely determined by the naive variance. It is possible that the SW estimator would have better coverage if a bootstrap of the entire data set including the x- and y-values in addition to the comparison vector counts were attempted. Although inclusion of the bootstrap variance estimate has improved the coverage of our estimator, it is still somewhat below the nominal 95% level, and further work is needed to produce additional improvements.

The panels in figure 1 plot the 400 regression estimates using the four methods versus the true simulation values under the first simulation conditions. If all the dots are close to the line

with slope one, then estimators are doing very well. The naive estimates underestimate the true slope most of the time. The robust estimates improve on the naive estimates slightly but still underestimate the true slope. The SW and our estimates appear to be centered in the correct location around the line. Our estimates seem to have less spread about the line. The panels in figure 2 show the decreased performance of all estimates. Our estimates seem to be the least affected.

7 Conclusion

Computerized record linkage can introduce errors into the composite file when errors are made in matching records. The mismatch errors can cause problems for analyses of variables brought together from different source files. In the presence of matching errors, naive estimators of linear regression coefficients are biased toward zero because the errors attenuate the relationship between the predictors and response. In simulations, least median regression was not sufficient to guard against matching errors, whereas the method of Scheuren and Winkler (1993) as applied here made a useful adjustment. Our unbiased method seemed to perform very well across a range of situations. The bootstrap procedure we described is useful for reflecting uncertainty due to matching for our estimation procedure.

Future work will involve comparing our method to Scheuren and Winkler's (1997) iterative method, which we have not implemented, and incorporating iterative clerical review as in Larsen and Rubin (2001). We also plan to investigate alternative bootstrap, jackknife, and multiple imputation options for propagation of error in matching through analyses.

8 References

- Alvey, W., and Jamerson, B. (1997), *Record Linkage Techniques – 1997*, Proceedings of an International Workshop and Exposition. Federal Committee on Statistical Methodology, Office of Management of the Budget.
- Armstrong, J .B., and Mayda, J. E. (1993), "Model-Based Estimation of Record Linkage Error Rates," *Survey Methodology*, 19, 137-147.
- Becker, M.P. and Yang, I. (1998), "Latent Class Marginal Models for Cross-Classifications of Counts," *Sociological Methodology*, 28, 293-325.
- Belin, T. (1993), "Evaluation of Sources of Variation in Record Linkage Through a Factorial Experiment," *Survey Methodology*, 19, 13-29.
- Belin, T. R., and Rubin, D. B. (1995), "A Method for Calibrating False-Match Rates in Record Linkage," *Journal of the American Statistical Association*, 90, 694-707.

- Burkard, R.E., and Derigs, U. (1980), "Assignment and Matching Problems: Solution Methods with FORTRAN-Programs," in *Lecture Notes in Economics and Mathematical Systems, No. 184*, Springer-Verlag: Berlin, Heidelberg, New York, pp. 1-11.
- Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977), "Maximum Likelihood from Incomplete Data Via the EM Algorithm," (with comments), *Journal of the Royal Statistical Society, Ser. B*, 39, 1-37.
- Fellegi, I. P., and Sunter, A. B. (1969), "A Theory for Record Linkage," *Journal of the American Statistical Association*, 64, 1183-1210.
- Gill, L. E. (1997), "OX-LINK: The Oxford Medical Record Linkage System Demonstration of the PC Version," in *Record Linkage Techniques – 1997*, Proceedings of an International Workshop and Exposition. Federal Committee on Statistical Methodology, Office of Management of the Budget, page 491.
- Gomatam, S., and Larsen, M.D. (2004), "Record Linkage and Counterterrorism," *Chance*, 17(1): 25-29.
- Jaro, M. A. (1989), "Advances in Record-Linkage Methodology as Applied to Matching the 1985 Census of Tampa, Florida," *Journal of the American Statistical Association*, 84, 414-420.
- (1995), "Probabilistic Linkage of Large Public Health Data Files," *Statistics in Medicine*, 14, 491-498.
- Lahiri, P. (2003), "On the impact of bootstrap in survey sampling and small-area estimation," *Statistical Science*, 18, 199-210.
- Lahiri, P. and Larsen, M.D. (2000), "Regression analysis with linked data," *Proceedings of the Section on Survey Research Methods, American Statistical Association*, 11-19.
- Lahiri, P., and Rao, J.N.K. (1995), "Robust Estimation of Mean Squared Error of Small Area Estimators," *Journal of American Statistical Association*, 90, 758-766.
- Larsen, M.D. (1999), "Multiple Imputation Analysis of Records Linkage Using Mixture Models," *Proceedings of the Statistical Society of Canada, Survey Methods Section*, 65-71.
- Larsen, M.D. (2001), "Methods For Model-Based Record Linkage And Analysis Of Linked Files," *Proceedings of the Government Statistics Section, American Statistical Association*.
- Larsen, M. D., and Rubin, D. B. (2001), "Iterative automated record linkage using mixture models," *Journal of the American Statistical Association*, 96, 32-41.

- McLachlan, G. J., and Peel, D. (2000), *Finite mixture models*, New York: John Wiley & Sons, Inc.
- Meng, X. L., and Rubin, D. B. (1993), "Maximum Likelihood Estimation Via the ECM Algorithm: A General Framework," *Biometrika*, 80, 267-278.
- Neter, John, Maynes, E. S., and Ramanathan, R. (1965), "The effect of mismatching on the measurement of response errors," *JASA*, 60, 1005-1027.
- Newcombe, H. B. (1988), *Handbook of record linkage: Methods for health and statistical studies, administration, and business*, Oxford University Press: Oxford.
- Newcombe, H.B., Kennedy, J.M., and Axford, S.J. and James, A.P. (1959), "Automatic Linkage of Vital Records," *Science*, 954-959.
- Press, W. H., Teukolsky, S. A., Vetterling, W. T., and Flannery, B. P. (1992), *Numerical Recipes in Fortran*, 2nd edition. Cambridge University Press: Cambridge, pp. 698-700.
- Scheuren, F., and Winkler, W. E. (1991), "Regression analysis of data files that are computer matched," *Proceedings of the Annual Research Conference*, the U.S. Census Bureau, 669-687.
- Scheuren, F., and Winkler, W. E. (1993), "Regression analysis of data files that are computer matched," *Survey Methodology*, 19, 39-58.
- Scheuren, F., and Winkler, W. E. (1997), "Regression analysis of data files that are computer matched – Part II," *Survey Methodology*, 23, 157- 165.
- Sen, A., and Srivastava, M. (1990), *Regression Analysis*. Springer-Verlag: New York, pp. 277-278.
- Thibaudeau, Y. (1993), "The Discrimination Power of Dependency Structures in Record Linkage," *Survey Methodology*, 19, 31-38.
- Winkler, W. E. (1988), "Using the EM Algorithm for Weight Computation in the Fellegi-Sunter Model of Record Linkage," *American Statistical Association Proceedings of Survey Research Methods Section*, pp. 667- 671.
- Winkler, W. E. (1990), "String comparator metrics and enhanced decision rules in the Fellegi-Sunter model of record linkage", *American Statistical Association Proceedings of the Section on Survey Research Methods*, pp. 354-359.
- Winkler, W. E. (1993), "Improved decision rules in the Fellegi-Sunter model of record linkage," in *American Statistical Association Proceedings of Survey Research Methods Section*, pp. 274-279.

Winkler, W. E. (1994), "Advanced Methods for Record Linkage," in *American Statistical Association Proceedings of Survey Research Methods Section*, pp. 467-472.

— (1995), "Matching and Record Linkage," in *Business Survey Methods*, ed. Cox, B. G., Binder, D. A., Chinnappa, B. N., Christianson, A., Colledge, M. J., and Kott, P. S., New York: Wiley Publications, pp. 355-384.

9 Appendix: Theorems and Proofs

Theorem A.1:

Under the model described by (1) and (2), we have for $i, j = 1, \dots, n$ ($i \neq j$)

- (a) $E(z_i) = w'_i \beta$;
- (b) $\text{Var}(z_i) = \sigma^2 + \beta' A_i \beta$, where $A_i = \sum_{j=1}^n q_{ij} d_{ij} d'_{ij}$ and $d_{ij} = x_j - w_i$;
- (c) $\text{Cov}(z_i, z_j) = \beta' A_{ij} \beta$, where $A_{ij} = \sum_{u=1}^n \sum_{v \neq u}^n q_{iu} q_{jv} d_{iu} d'_{jv}$.

Proof:

(a): Part (a) follows by noting that for $i = 1, \dots, n$, $E(z_i) = E[E(z_i|y)]$, $E(z_i|y) = q_{ii} y_i + \sum_{l \neq i} q_{il} y_l = \sum_{l=1}^n q_{il} y_l$, and $E(y_i) = x'_i \beta$.

(b): To prove part (b), we will use the fact that

$$\text{Var}(z_i) = E[\text{Var}(z_i|y)] + \text{Var}[E(z_i|y)]. \quad (9)$$

Using the intermediate step in (a),

$$\text{Var}[E(z_i|y)] = \text{Var}\left(\sum_{l=1}^n q_{il} y_l\right) = \sum_{l=1}^n q_{il}^2 \sigma^2 \quad (10)$$

By the definition of variance,

$$\begin{aligned} \text{Var}(z_i|y) &= \left(y_i - \sum_{l=1}^n q_{il} y_l\right)^2 q_{ii} + \sum_{j \neq i} \left(y_j - \sum_{l=1}^n q_{il} y_l\right)^2 q_{ij} \\ &= \sum_{j=1}^n \left(y_j - \sum_{l=1}^n q_{il} y_l\right)^2 q_{ij} \end{aligned} \quad (11)$$

Now, we compute the expectation for each j :

$$\begin{aligned} E\left(y_j - \sum_{l=1}^n q_{il} y_l\right)^2 &= \text{Var}\left(y_j - \sum_{l=1}^n q_{il} y_l\right) + E\left(y_j - \sum_{l=1}^n q_{il} y_l\right)^2 \\ &= \sigma^2 \left(1 - 2q_{ij} + \sum_{l=1}^n q_{il}^2\right) + (d'_{ij} \beta)^2, \end{aligned} \quad (12)$$

because

$$\begin{aligned}
\text{Var}(y_j - \sum_{l=1}^n q_{il}y_l) &= \text{Var}(y_j(1 - q_{ij}) - \sum_{l \neq j} q_{il}y_l) \\
&= \sigma^2(1 - q_{ij})^2 + \sum_{l \neq j} q_{il}^2 \sigma^2 \\
&= \sigma^2(1 - 2q_{ij} + \sum_{l=1}^n q_{il}^2)
\end{aligned}$$

and

$$\mathbf{E}(y_i - \sum_{l=1}^n q_{il}y_l) = x'_i \beta - \sum_{l=1}^n q_{il} x'_l \beta = d'_{ij} \beta.$$

Part (b) follows by (10) into (9), resultant (9) and (8) into (7), and simplifying.

(c): Turning to part (c), we will use the fact that, for $i \neq j$,

$$\text{Cov}(z_i, z_j) = \mathbf{E}[\text{Cov}(z_i, z_j | y)] + \text{Cov}[\mathbf{E}(z_i | y), \mathbf{E}(z_j | y)].$$

Since $\mathbf{E}(z_i | y) = \sum_{l=1}^n q_{il}y_l$ and $\mathbf{E}(z_j | y) = \sum_{l=1}^n q_{jl}y_l$,

$$\text{Cov}[\mathbf{E}(z_i | y), \mathbf{E}(z_j | y)] = \sum_{l=1}^n q_{il}q_{jl}\sigma^2. \tag{13}$$

By the definition of covariance,

$$\text{Cov}(z_i, z_j | y) = \sum_{u=1}^n \sum_{v=1}^n q_{iu}q_{jv}(y_u - \sum_{k=1}^n q_{ik}y_k)(y_v - \sum_{l=1}^n q_{jl}y_l),$$

where $q_{iu}q_{jv} = \text{Prob}(z_i = y_u, z_j = y_v | y)$. Note that $q_{iu}q_{ju} = q_{iv}q_{jv} = 0$. Now, for $u \neq v$,

$$\begin{aligned}
\mathbf{E}[(y_u - \sum_{k=1}^n q_{ik}y_k)(y_v - \sum_{l=1}^n q_{jl}y_l)] &= \text{Cov}[(y_u - \sum_{k=1}^n q_{ik}y_k), (y_v - \sum_{l=1}^n q_{jl}y_l)] \\
&\quad + \mathbf{E}[(y_u - \sum_{k=1}^n q_{ik}y_k)]\mathbf{E}[(y_v - \sum_{l=1}^n q_{jl}y_l)] \\
&= 0 - \sigma^2(1 - q_{iu})q_{ju} - \sigma^2(1 - q_{jv})q_{iv} \\
&\quad + \sigma^2 \sum_{k \neq u,v}^n q_{ik}q_{jk} + (d'_{iu}\beta)(d'_{jv}\beta)
\end{aligned}$$

since $\mathbf{E}(y_u - \sum_{k=1}^n q_{ik}y_k) = (x'_u - w'_i)\beta = d'_{iu}\beta$ and $\mathbf{E}(y_v - \sum_{l=1}^n q_{jl}y_l) = d'_{jv}\beta$. Thus

$$\begin{aligned}
\mathbf{E}[\text{Cov}(z_i, z_j | y)] &= \sum_{u=1}^n \sum_{v \neq u}^n [q_{iu}q_{jv}[\sigma^2(-(1 - q_{iu})q_{ju} - (1 - q_{jv})q_{iv} \\
&\quad + \sum_{k \neq u,v}^n q_{ik}q_{jk})] + (d'_{iu}\beta)(d'_{jv}\beta)].
\end{aligned}$$

Noting that $\sum_{v,v \neq u}^n q_{iuv} = q_{iu}$, $\sum_{u,u \neq v}^n q_{iuv} = q_{jv}$, $\sum_{u=1}^n \sum_{v,v \neq u}^n q_{iuv} = 1$, and

$$\begin{aligned} \sum_{u=1}^n \sum_{v,v \neq u}^n q_{iuv} \sum_{k \neq u,v}^n q_{ik}q_{jk} &= \sum_{u=1}^n \sum_{v,v \neq u}^n q_{iuv} \left(\sum_{k=1}^n q_{ik}q_{jk} - q_{iu}q_{ju} - q_{iv}q_{jv} \right) \\ &= \sum_{k=1}^n q_{ik}q_{jk} - \sum_{u=1}^n q_{iu}^2 q_{ju} - \sum_{v=1}^n q_{iv}q_{jv}^2, \end{aligned}$$

we arrive at

$$\begin{aligned} \mathbf{E}[\text{Cov}(z_i, z_j | y)] &= \sigma^2 \left[- \sum_{u=1}^n q_{iu}q_{ju} + \sum_{u=1}^n q_{iu}^2 q_{ju} - \sum_{v=1}^n q_{iv}q_{jv} + \sum_{v=1}^n q_{iv}q_{jv}^2 \right. \\ &\quad \left. + \sum_{k=1}^n q_{ik}q_{jk} - \sum_{u=1}^n q_{iu}^2 q_{ju} - \sum_{v=1}^n q_{iv}q_{jv}^2 \right] \\ &\quad + \sum_{u=1}^n \sum_{v \neq u}^n (d'_{iu}\beta)(d'_{jv}\beta). \end{aligned} \quad (14)$$

Adding the two parts, (13) and (14), yields

$$\text{Cov}(z_i, z_j) = \sum_{u=1}^n \sum_{v \neq u}^n q_{iu}q_{jv} (d'_{iu}\beta)(d'_{jv}\beta). \quad (15)$$

Theorem A.2

Under the model described by (1) and (2), we have

$$\mathbf{E}(S^2) = (n-p)\sigma^2 + \text{tr}[I - W(W'W)^{-1}W']H, \quad (16)$$

where $S^2 = z'[I - W(W'W)^{-1}W']z$ and $H = ((h_{ij}))$ with $h_{ii} = \beta'A_i\beta$ and $h_{ij} = \beta'A_{ij}\beta$.

Proof:

First note that $S^2 = z'[I - W(W'W)^{-1}W']z = \text{tr}[(I - W(W'W)^{-1}W')zz']$ and that $\mathbf{E}(zz') = \text{Var}(z) + \mathbf{E}(z)\mathbf{E}(z') = \Sigma + (W\beta)(W\beta)'$, where $\Sigma = ((\sigma_{ij}))$ with $\sigma_{ii} = \text{Var}(z_i)$ and $\sigma_{ij} = \text{Cov}(z_i, z_j)$, $i \neq j$. So

$$\begin{aligned} \mathbf{E}(S^2) &= \text{tr}[(I - W(W'W)^{-1}W')(\Sigma + W\beta\beta'W')] \\ &= \text{tr}[(I - W(W'W)^{-1}W')\Sigma] \\ &= \sigma^2 \text{tr}[(I - W(W'W)^{-1}W')] + \text{tr}[(I - W(W'W)^{-1}W')H] \\ &= (n-p)\sigma^2 + \text{tr}[(I - W(W'W)^{-1}W')H]. \end{aligned} \quad (17)$$

The second term on the first line of (17) vanishes since $\text{tr}[W\beta\beta'W' - W(W'W)^{-1}W'W\beta\beta'W'] = \text{tr}[W\beta\beta'W' - W\beta\beta'W'] = 0$. The conversion to the last line of (17) follows from usual regression algebra (e.g., Sen and Srivastava 1990, page 278).

Theorem A.3

Under the model described by (1) and (2) and the following regularity assumptions

(i) $\sup_{ij} |x_{ij}| \leq c < \infty$

(ii) $W'\Sigma W = O(n)$

the estimator $\hat{\sigma}^2 = \max[0, S^2 - \text{tr}(I - W(W'W)^{-1}W')\hat{H}/(n - p)]$, where $\hat{H} = ((\hat{h}_{ij}))$ with $\hat{h}_{ii} = \hat{\beta}_u' A_i \hat{\beta}_u$ and $\hat{h}_{ij} = \hat{\beta}_u' A_{ij} \hat{\beta}_u$, is consistent for σ^2 as $n \rightarrow \infty$.

Proof:

Let

$$\tilde{\sigma}^2 = u + \frac{1}{n - p} \left(\text{tr}[(I - W(W'W)^{-1}W')H] - \text{tr}[(I - W(W'W)^{-1}W')\hat{H}] \right), \quad (18)$$

where

$$u = \frac{S^2 - \text{tr}[I - W(W'W)^{-1}W']H}{n - p}.$$

Since $E(\hat{\beta}_u) = \beta$ and $\text{Var}(\hat{\beta}_u) = (W'W)^{-1}W'\Sigma W(W'W)^{-1} = O(n^{-1})$, by assumptions (i) and (ii), we have $\hat{\beta}_u \xrightarrow{p} \beta$ as $n \rightarrow \infty$. Thus the second term in (18) tends to 0 as $n \rightarrow \infty$.

We will show that $u \xrightarrow{p} \sigma^2$ as $n \rightarrow \infty$. By Theorem A.2, $E(u) = \sigma^2$. Now,

$$\text{Var}(u) = \frac{\text{Var}(S^2)}{(n - p)^2}.$$

Note that $\text{Var}(S^2) = \text{Var}(\eta'(I - W(W'W)^{-1}W')\eta)$, where $\eta = z - W\beta$. Using part (d) of Lemma C.4 of Lahiri and Rao (1995), we have

$$\text{Var}(\eta'(I - W(W'W)^{-1}W')\eta) = O(n).$$

So, $\text{Var}(S^2) = O(n)$ and $\text{Var}(u) = O(n^{-1})$. Thus, the result is established.

Table 3: The percent relative improvement of the proposed estimator over rival estimators

Method	AAD	ASD
<i>Simulation Case 1</i>		
Naive	84	170
Robust	51	86
Scheuren-Winkler	33	72
<i>Simulation Case 2</i>		
Naive	293	960
Robust	216	590
Scheuren-Winkler	109	327

Table 4: Percent coverage of 95% confidence intervals with and without bootstrap adjustment of standard errors.

	Coverage Before Bootstrap	Coverage After Bootstrap
<i>Simulation Case 1</i>		
Naive	34	34
Robust	50	50
Scheuren-Winkler	59	60
Lahiri-Larsen	83	88
<i>Simulation Case 2</i>		
Naive	4	4
Robust	8	8
Scheuren-Winkler	40	41
Lahiri-Larsen	85	89

Figure 1: Comparison of Four Estimators on Four Hundred Data Sets, first set of simulation conditions. Plots of Naive, Robust, Scheuren-Winkler, and Lahiri-Larsen estimators versus the truth. Diagonal lines have slope 1.

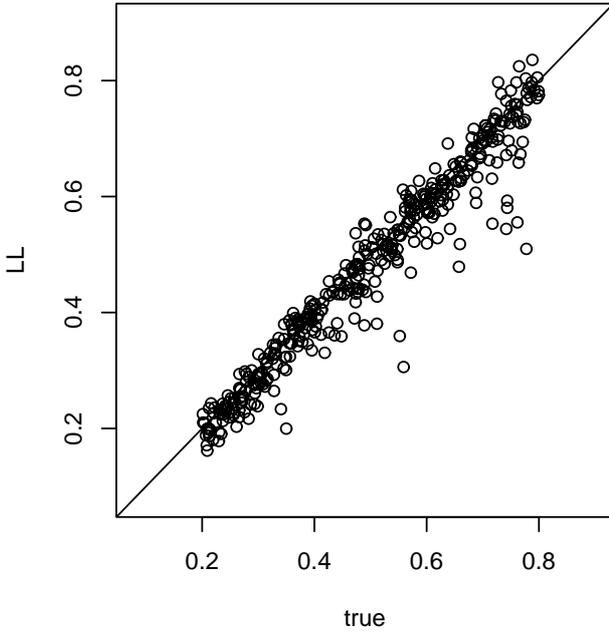
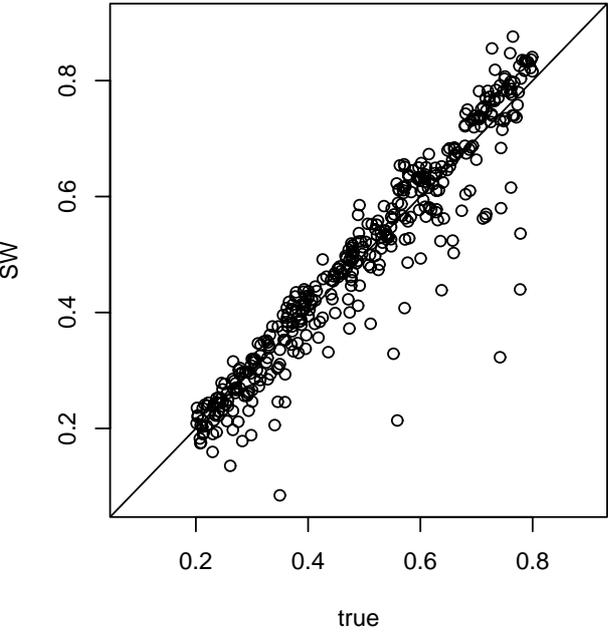
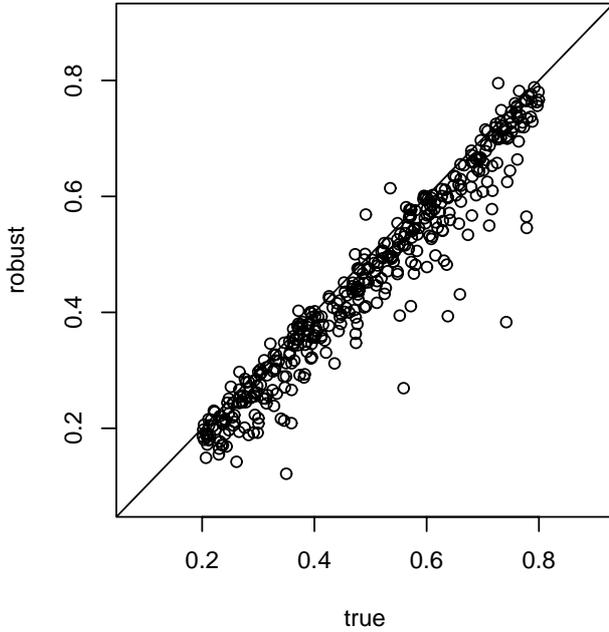
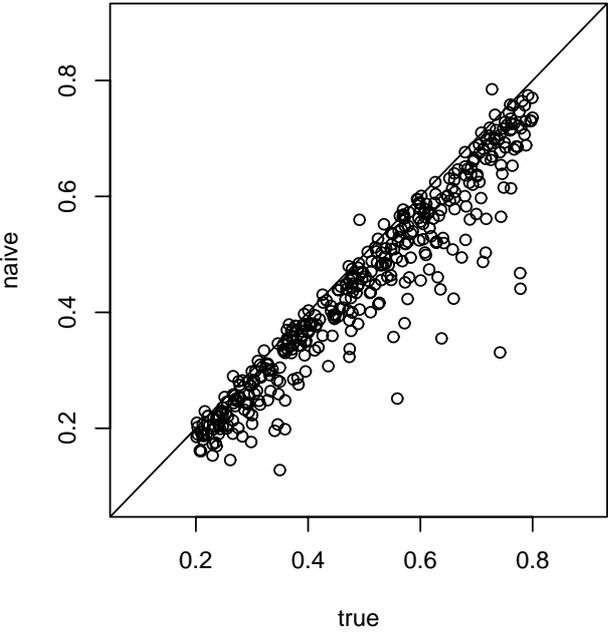


Figure 2: Comparison of Four Estimators on Four Hundred Data Sets, second set of simulation conditions. Plots of Naive, Robust, Scheuren-Winkler, and Lahiri-Larsen estimators versus the truth. Diagonal lines have slope 1.

