

Monte Carlo Estimation of $g(\mu)$ from Normally Distributed Data with Applications

Leonard A. Stefanski
Department of Statistics
North Carolina State University
Raleigh, NC 27695-8203, USA

Steven J. Novick
Schering-Plough Research Institute
Kenilworth, NJ, USA 07033-1300

Viswanath Devanarayan
Eli Lilly & Company
Indianapolis, IN 46285, USA

Institute of Statistics Mimeo Series No. 2569

Abstract We derive Monte Carlo-amenable solutions to the problem of unbiased estimation of a nonlinear function of the mean of a normal distribution. For most nonlinear functions the maximum likelihood estimator is biased. Our method yields a Monte Carlo approximation to the uniformly minimum variance unbiased estimator for a wide class of nonlinear functions. Applications to problems arising in the analysis of data measured with error and the secondary analysis of estimated data are described.

Keywords: Corrected Score, Deconvolution, Jackknife, Measurement Error, Nonlinear, Resampling, SIMEX, Spherical Uniform Distribution, UMVUE.

1 Introduction

We present simulation-based solutions to the problem of estimating $g(\mu)$ given an iid sample $W_1, \dots, W_m \sim N(\mu, \sigma^2)$ with both μ and σ^2 unknown, and describe applications to the analysis of data measured with error. This problem and closely related problems have a long history dating from the 1950s (Kolmogorov, 1950; Kitagawa, 1956; Washio et al., 1956; Neyman and Scott, 1960; Gray, Watkins and Schucany 1973; Woodward, 1977; Watkins, 1979), and solutions have been derived using a variety of approaches. The solution most relevant to our approach is that of Gray, et al. (1973) for which Fortran subroutines were presented by Woodward (1977).

Gray et al. (1973) used the generalized jackknife to derive the uniformly minimum variance unbiased estimator of $g(\mu)$ for polynomial $g(\cdot)$. Then, noting the form of the solution, they generalized the estimator to the case in which $g(\cdot)$ has an infinite power series representation, resulting in the estimator

$$\hat{g}(\mu) = g(\bar{W}) + \sum_{k=1}^{\infty} \frac{(-1)^k \Gamma((m-1)/2) g^{(2k)}(\bar{W})}{k! \Gamma((m-1)/2 + k)} \left(\frac{\hat{\sigma}^2}{4} \right)^k.$$

Like Gray et al. (1973) our estimator is derived from initial consideration of generalized jackknife estimators. However, the path we take differs from that of Gray et al. (1973) and leads to characterizations of the estimator as expectations of relatively simple functions with respect to certain uniform distributions. This in turn suggests a simple Monte Carlo estimation procedure. A significant theoretical and practical advantage of our method relative to that of Gray et al. (1973) is that our method also produces an unbiased estimator of the variance of the Monte Carlo estimator of $g(\mu)$.

Our interest in the problem lies in its relevance to problems arising in the analysis of data measured with error. Specifically, the construction of corrected score functions (Nakamura, 1990; Stefanski, 1989; Novick and Stefanski, 2002) and the identification

of deconvoluting kernels (Carroll and Hall, 1988; Stefanski and Carroll, 1990; Fan, 1991c) when the measurement error variance is unknown but replicate measurements are available. In such applications a common design is $m = 2$ replicates and we examine this case in detail. A very useful and significant advantage of the new method for dealing with measurement error is the automatic accommodation of heteroscedastic measurement errors. The new method is an alternative to the empirical simulation extrapolation method (simex) for replicate measurements studied by Devanarayan and Stefanski (2002). Main results on the fundamental estimation problem are developed in Section 2. Applications to measurement error problems are described in Section 4 and certain technical details are given in Section 5.

2 Main Results

Our assumptions are that the components of $\mathbf{W} = (W_1, \dots, W_m)^T$ are iid $N(\mu, \sigma^2)$, and $\theta = g(\mu)$ is the parameter of interest. The estimators we propose are closely related to, and motivated by, jackknife estimators and this is where we start.

2.1 The Jackknife Connection

Consider the jackknife approach to reducing bias in the maximum likelihood estimator of $g(\mu)$, $\tilde{\theta} = g(\bar{W})$. The leave- b -out version of $\tilde{\theta}$ is

$$\tilde{\theta}_{(-b)} = \frac{1}{\binom{m}{b}} \sum g(\bar{W}_{(i_1, \dots, i_b)}), \quad 1 \leq b < m, \quad (1)$$

where $\bar{W}_{(i_1, \dots, i_b)}$ is the data sample mean with W_{i_1}, \dots, W_{i_b} removed. By definition, the leave-0-out estimator is $\tilde{\theta}$. The first-order, reduced-bias jackknife estimator, $m\tilde{\theta} - (m-1)\tilde{\theta}_{(-1)}$, is the intercept of the line fit to the points $\{(a_0, \tilde{\theta}), (a_1, \tilde{\theta}_{(-1)})\}$, where $a_b = 1/(m-b)$. Higher-order jackknife corrections are obtained by fitting b^{th} -order

polynomials to the pairs $\{(a_0, \tilde{\theta}), (a_1, \tilde{\theta}_{(-1)}), \dots, (a_b, \tilde{\theta}_{(-b)})\}$. The intercept again provides the reduced-bias estimator.

Because $0 < a_0 < a_1 < \dots$, the intercepts of the polynomials used in jackknife estimation lie outside the range of the data used to fit the polynomials. Thus the jackknife estimators are obtained by extrapolating the polynomials to the origin, which can be identified as $a_{-\infty}$ in the sense that $a_{-\infty} = \lim_{b \rightarrow -\infty} a_b = 0$.

Note that in (1), $\overline{W}_{(i_1, \dots, i_b)} = \overline{W} + (\overline{W}_{(i_1, \dots, i_b)} - \overline{W})$ where $(\overline{W}_{(i_1, \dots, i_b)} - \overline{W})$ is an unbiased estimator of 0, and thus is a contrast in $\mathbf{W} = (W_1, \dots, W_m)^T$. After normalizing these contrasts to unit norm, (1) can be rewritten as

$$\tilde{\theta}_{(-b)} = E \left\{ g \left(\overline{W} + \lambda_b^{1/2} m^{-1/2} \mathbf{T}_b^T \mathbf{W} \right) \mid \mathbf{W} \right\} \quad (2)$$

where $\lambda_b = ma_b - 1$ and \mathbf{T}_b is a random $m \times 1$ vector, independent of \mathbf{W} , distributed uniformly on the set of normalized contrasts \mathcal{C}_b consisting of the $\binom{m}{b}$ vectors with b entries equal to $-\gamma m^{-1}$ and $(m - b)$ entries equal to $\gamma \{(m - b)^{-1} - m^{-1}\}$. The normalizing constant is $\gamma = \{m(m - b)/b\}^{1/2}$. The contrasts in \mathcal{C}_b have norm 1, but are not orthogonal. Normalization of the contrasts in \mathcal{C}_b essentially eliminates their dependence on a_b (although \mathcal{C}_b does depend on b both in cardinality of the set and the composition of the contrasts), and thus (2) reveals that $\hat{\theta}_{(-b)}$ is essentially a function of a_b only through $\lambda_b = ma_b - 1$. Note that $\lambda_b = -1$ when $a_b = 0$. The higher-order jackknife can now be described as fitting a b^{th} -order polynomial to the points $\{(\lambda_0, \tilde{\theta}), (\lambda_1, \tilde{\theta}_{(-1)}), \dots, (\lambda_b, \tilde{\theta}_{(-b)})\}$ and extrapolating to $\lambda_{-\infty} = -1$.

Why bother fitting polynomials and extrapolating? Equation (2) explicitly reveals the dependence of $\hat{\theta}_{(-b)}$ on λ_b . Evaluation at $\lambda_b = -1$ results in what we call the *exact-extrapolant, jackknife “estimator,”*

$$\tilde{\tilde{\theta}}_{(-b)} = E \left\{ g \left(\overline{W} + im^{-1/2} \mathbf{T}_b^T \mathbf{W} \right) \mid \mathbf{W} \right\}, \quad (3)$$

where $i = \sqrt{-1}$ (the quotation marks indicate that $\tilde{\theta}_{(-b)}$ is not real-valued and thus is not an estimator in the strict sense). The expectation of $\tilde{\theta}_{(-b)}$ is

$$E\left(\tilde{\theta}_{(-b)}\right) = E\left\{g\left(\bar{W} + im^{-1/2}\mathbf{T}_b^T\mathbf{W}\right)\right\} = E\left[E\left\{g\left(\bar{W} + im^{-1/2}\mathbf{T}_b^T\mathbf{W}\right) \mid \mathbf{T}_b\right\}\right].$$

Consider the inner conditional expectation. Because \mathbf{T}_b is a normalized contrast vector, the argument of the function $g()$

$$\eta = \bar{W} + im^{-1/2}\mathbf{T}_b^T\mathbf{W} = \mu + (\bar{W} - \mu) + im^{-1/2}\mathbf{T}_b^T\mathbf{W} = \mu + \varepsilon_1^* + i\varepsilon_2^*,$$

where $\varepsilon_1^* = \bar{W} - \mu$ and $\varepsilon_2^* = m^{-1/2}\mathbf{T}_b^T\mathbf{W}$ are conditionally independent and identically distributed $N(0, \sigma^2/m)$ given \mathbf{T}_b . Note that the conditional distribution of ε_1^* and ε_2^* does not depend on \mathbf{T}_b . Suppose now that $g()$ is an entire function, so that we can write

$$g(\eta) = g(\mu + \varepsilon_1^* + i\varepsilon_2^*) = \sum_{k=0}^{\infty} \frac{g^{(k)}(\mu) \{\varepsilon_1^* + i\varepsilon_2^*\}^k}{k!}.$$

Assuming that the interchange of summation and expectation are justified, an appeal to Lemma 5.0.1 shows that

$$E\{g(\eta) \mid \mathbf{T}_b\} = \sum_{k=0}^{\infty} \frac{g^{(k)}(\mu) E\{(\varepsilon_1^* + i\varepsilon_2^*)^k\}}{k!} = g(\mu).$$

It follows that the unconditional expectation of $g(\eta)$ equals $g(\mu)$, that is, $E\left(\tilde{\theta}_{(-b)}\right) = g(\mu)$. In other words, the complex-valued, exact-extrapolant, jackknife “estimator” $\tilde{\theta}$ is unbiased for θ .

There are two problems with this “estimator.” The first, that it is complex-valued, is easily solved by dropping the imaginary part, which is an unbiased estimator of 0 as the analysis above shows. The second is lack of efficiency resulting from the particular averaging in (2). Except for special cases $\text{Re}\left\{\tilde{\theta}_{(-b)}\right\}$ is not a function of the data only via $\left\{\bar{W}, \sum_j (W_j - \bar{W})^2\right\}$, and thus need not be fully efficient. Examination of

(3) provides a clue to increasing efficiency. The only properties of the contrasts \mathbf{T}_b used in proving unbiasedness were that they are contrasts. In other words, there is nothing special about the set of contrasts \mathcal{C}_b . One might expect that variation can be reduced by averaging over the largest possible set of normalized contrasts. The theorem in the next section shows this to be the case.

2.2 The First Form on the Estimator

The estimators we describe are defined in terms of uniformly distributed random contrast vectors. Such vectors can be defined, and generated, in terms of standard normal random variables. Let Z_1, \dots, Z_m be independent and identically distributed $N(0, 1)$ independently of W_1, \dots, W_m . Define

$$T_j = \frac{Z_j - \bar{Z}}{\sqrt{\sum_{j=1}^m (Z_j - \bar{Z})^2}}, \quad j = 1, \dots, m, \quad (4)$$

where \bar{Z} is the sample mean of Z_1, \dots, Z_m . Then $\mathbf{T} = (T_1, \dots, T_m)^T$ is a random vector uniformly distributed on \mathcal{C} , the set of $m \times 1$ normalized, contrast vectors.

Theorem 2.2.1 *Assume that W_1, \dots, W_m are independent and identically distributed $N(\mu, \sigma^2)$ with both μ and σ^2 unknown. Suppose that $g(\cdot)$ is an entire function over the complex plane that also satisfies certain integrability conditions that will be made clear in the proof. Then the estimator,*

$$\hat{\theta} = E \left\{ g \left(\bar{W} + im^{-1/2} \mathbf{T}^T \mathbf{W} \right) \mid \mathbf{W} \right\}, \quad (5)$$

is real-valued, unbiased for $\theta = g(\mu)$, and depends on \mathbf{W} only through the statistics $\left\{ \bar{W}, \sum_j (W_j - \bar{W})^2 \right\}$. Consequently, if $\hat{\theta}$ has finite variance it is the uniformly minimum variance unbiased estimator of θ .

Proof. Note that the random vector \mathbf{T} is independent of \mathbf{W} . The proof uses the facts that, almost surely,

$$\mathbf{T}^T \mathbf{1} = \sum T_j = 0, \quad \text{and} \quad \mathbf{T}^T \mathbf{T} = \sum T_j^2 = 1. \quad (6)$$

Upon taking expectations via conditioning,

$$E(\hat{\theta}) = E\left\{g\left(\bar{W} + im^{-1/2}\mathbf{T}^T\mathbf{W}\right)\right\} = E\left[E\left\{g\left(\bar{W} + im^{-1/2}\mathbf{T}^T\mathbf{W}\right) \mid \mathbf{T}\right\}\right].$$

Consider the inner conditional expectation. Invoking the constraints in (6) shows that the argument of the function $g()$

$$\eta = \bar{W} + im^{-1/2}\mathbf{T}^T\mathbf{W} = \mu + (\bar{W} - \mu) + im^{-1/2}\mathbf{T}^T\mathbf{W} = \mu + \varepsilon_1^* + i\varepsilon_2^*,$$

where $\varepsilon_1^* = \bar{W} - \mu$ and $\varepsilon_2^* = m^{-1/2}\mathbf{T}^T\mathbf{W}$ are conditionally independent and identically distributed $N(0, \sigma^2/m)$ given \mathbf{T} .

By assumption $g(\eta) = g(\mu + \varepsilon_1^* + i\varepsilon_2^*) = \sum_{k=0}^{\infty} g^{(k)}(\mu) \{\varepsilon_1^* + i\varepsilon_2^*\}^k / k!$. Assuming that the interchange of summation and expectation are justified, an appeal to Lemma 5.0.1 shows that

$$E\{g(\eta) \mid \mathbf{T}\} = \sum_{k=0}^{\infty} g^{(k)}(\mu) E\left\{(\varepsilon_1^* + i\varepsilon_2^*)^k\right\} / k! = g(\mu). \quad (7)$$

It follows that $E\{g(\eta)\} = g(\mu)$ and thus $E(\hat{\theta}) = g(\mu)$.

It remains to show that $\hat{\theta}$ is the uniformly minimum variance unbiased estimator.

Again by assumption on $g()$

$$g(\eta) = g\left(\bar{W} + im^{-1/2}\mathbf{T}^T\mathbf{W}\right) = \sum_{k=0}^{\infty} \frac{g^{(k)}(\bar{W}) \left\{im^{-1/2}\mathbf{T}^T\mathbf{W}\right\}^k}{k!},$$

so that

$$\hat{\theta} = E\{g(\eta) \mid \mathbf{W}\} = \sum_{k=0}^{\infty} \frac{g^{(k)}(\bar{W}) i^k m^{-k/2} E\left\{(\mathbf{T}^T\mathbf{W})^k \mid \mathbf{W}\right\}}{k!}. \quad (8)$$

Appealing to Lemma 5.0.2 with $p = 1$ and $\mathbf{P}_x = m^{-1}\mathbf{1}\mathbf{1}^T$ shows that for odd k , $E\{(\mathbf{T}^T \mathbf{W})^k | \mathbf{W}\} = 0$, from which it follows that $\hat{\theta}$ is real-valued; and for even k $E\{(\mathbf{T}^T \mathbf{W})^k | \mathbf{W}\} = Q_{m,k}\left(\sum_j (W_j - \bar{W})^2\right)$, for some real-valued function $Q_{m,k}()$ that does not depend on \mathbf{W} , from which it follows that $\hat{\theta}$ is a function of \bar{W} and $\sum_j (W_j - \bar{W})^2$. Completeness of the $N(\mu, \sigma^2)$ family justifies the efficiency assertion. ♠

Sufficient conditions for the interchanges of summation and expectation in (7) and (8) can be readily derived (Stefanski, 1989).

2.2.1 Monte Carlo Estimation

From equation (8) it follows that

$$\hat{\theta} = E\{g(\eta) | \mathbf{W}\} = \sum_{k=0}^{\infty} \frac{g^{(2k)}(\bar{W})(-1)^k E\{(\mathbf{T}^T \mathbf{W})^{2k} | \mathbf{W}\}}{m^k (2k)!}. \quad (9)$$

For polynomials the sum in (9) is finite, otherwise it is generally infinite; and even for simple functions, closed form expressions, when they can be determined, generally involve special functions. For example, when $g(\mu) = \exp(\mu)$ the uniformly minimum variance unbiased estimator estimator (Neyman and Scott, 1960) is

$$\exp(\bar{W}) \Gamma\left(\frac{m-1}{2}\right) J_{(m-1)/2-1}\left(\hat{\sigma}\sqrt{\frac{m-1}{m}}\right) \left(\frac{\hat{\sigma}\sqrt{\frac{m-1}{m}}}{2}\right)^{1-(m-1)/2},$$

where $J_\nu(z)$ is the Bessel function, $J_\nu(z) = \sum_{k=0}^{\infty} \frac{(-1)^k}{k! \Gamma(\nu+k+1)} \left(\frac{z}{2}\right)^{\nu+2k}$.

The technical details of calculating $\hat{\theta}$ can be avoided by replacing the conditional expectation in (5) by a Monte Carlo mean. Define

$$\hat{\theta}_B = \frac{1}{B} \sum_{b=1}^B \text{Re}\left\{g\left(\bar{W} + im^{-1/2}\mathbf{T}_b^T \mathbf{W}\right)\right\} \quad (10)$$

where $\text{Re}()$ “real part” and the vectors $\mathbf{T}_b = (T_{b,1}, \dots, T_{b,m})^T$ are independent and identically distributed Monte Carlo replicates of \mathbf{T} . These are easily constructed

from pseudo-random, standard normal variables using the relationships in (4). As indicated in the proof of the theorem the imaginary part of $g(\overline{W} + im^{-1/2}\mathbf{T}_b^T\mathbf{W})$ has expectation 0, so that taking the real part in (10) does not affect bias and yields in a real-valued estimator.

The estimator simplifies in the special case $m = 2$. For then $T_2 = -T_1$, $\Pr(T_1 = 1/\sqrt{2}) = 1/2 = \Pr(T_1 = -1/\sqrt{2})$, and $T_1^2 = 1/2$ almost surely. In this case

$$\begin{aligned}
\hat{\theta}_B &= \frac{1}{B} \sum_{b=1}^B \operatorname{Re} \left\{ g \left(\overline{W} + i2^{-1/2} (T_{b,1}W_1 + T_{b,2}W_2) \right) \right\} \\
&= \frac{1}{B} \sum_{b=1}^B \operatorname{Re} \left\{ \sum_{k=0}^{\infty} \frac{g^{(2k)}(\overline{W})}{(2k)!} \left\{ im^{-1/2} T_{b,1} (W_1 - W_2) \right\}^{2k} \right\} \\
&= \frac{1}{B} \sum_{b=1}^B \operatorname{Re} \left\{ \sum_{k=0}^{\infty} \frac{(-1)^k g^{(2k)}(\overline{W})}{4^k (2k)!} (W_1 - W_2)^{2k} \right\} \\
&= \operatorname{Re} \left\{ \sum_{k=0}^{\infty} \frac{(-1)^k g^{(2k)}(\overline{W})}{4^k (2k)!} (W_1 - W_2)^{2k} \right\} \\
&= \operatorname{Re} \left\{ g \left(\overline{W} + (i/2) (W_1 - W_2) \right) \right\}
\end{aligned}$$

independent of B . In other words, when $m = 2$ the uniformly minimum variance unbiased estimator can be obtained exactly in terms of $\operatorname{Re} \{g(\cdot)\}$.

For example, the uniformly minimum variance unbiased estimator of $g(\mu) = \exp(\mu)$ is $\hat{\theta} = \exp(\overline{W}) \cos((W_1 - W_2)/2)$, and the uniformly minimum variance unbiased estimator of $g(\mu) = \mu^r$ is

$$\begin{aligned}
\hat{\theta}_r &= \operatorname{Re} \left[\left\{ \overline{W} + (i/2) (W_1 - W_2) \right\}^r \right] \\
&= \operatorname{Re} \left\{ \sum_{k=0}^r \binom{r}{k} \overline{W}^{r-k} (i^k/2^k) (W_1 - W_2)^k \right\} \\
&= \sum_{k=0}^{\lfloor r/2 \rfloor} \binom{r}{2k} \overline{W}^{r-2k} \frac{(-1)^k}{4^k} (W_1 - W_2)^{2k}.
\end{aligned}$$

For $r = 2, 3$ and 4 , the estimators are $\hat{\theta}_2 = \overline{W}^2 - (1/4)(W_1 - W_2)^2$, $\hat{\theta}_3 = \overline{W}^3 - (3/4)\overline{W}(W_1 - W_2)^2$, and $\hat{\theta}_4 = \overline{W}^4 - (3/2)\overline{W}^2(W_1 - W_2)^2 + (1/16)(W_1 - W_2)^4$, and their unbiasedness can be verified directly.

Note that in the case of $m = 2$, the only possible jackknife estimator (1) is the linear jackknife with $b = 1$. In this case, the set of contrasts (for $b = 1$) described in Section 2.1 is $\mathcal{C}_1 = \{(1/\sqrt{2}, -1/\sqrt{2}), (-1/\sqrt{2}, 1/\sqrt{2})\}$. So for $m = 2$, \mathcal{C}_1 is the largest set of possible contrasts over which to average, and thus the exact extrapolant, jackknife estimators are fully efficient for $m = 2$.

2.2.2 Extensions to Regression Models

Now assume that $\mathbf{W}_{m \times 1} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_m)$ where \mathbf{X} is a fixed $m \times p$ design matrix of rank p with $m \geq 2p$ and both $\boldsymbol{\beta}$ and σ^2 unknown. Suppose that the parameter of interest is a nonlinear function of a linear function of $\boldsymbol{\beta}$, say $\theta = g(\mathbf{a}^T \boldsymbol{\beta})$ for some fixed, non-random \mathbf{a} . Let $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}$. Then the generalization of (5) to the regression case results in the uniformly minimum variance unbiased estimator

$$\hat{\theta} = E \left[g \left(\mathbf{a}^T \hat{\boldsymbol{\beta}} + i \mathbf{a}^T (\mathbf{X}^T \mathbf{X})^{-1/2} \mathbf{T}^T \mathbf{W} \right) \mid \mathbf{W} \right]$$

where the random $m \times p$ matrix \mathbf{T} can be defined (and hence generated) in terms of standard normal random variables. Specifically,

$$\mathbf{T} = (\mathbf{I}_m - \mathbf{P}_\mathbf{X}) \mathbf{Z} \{ \mathbf{Z}^T (\mathbf{I}_m - \mathbf{P}_\mathbf{X}) \mathbf{Z} \}^{-1/2} \quad (11)$$

where the matrix $\mathbf{Z}_{m \times p}$ has iid entries $Z_{i,j} \sim N(0, 1)$, and $\mathbf{P}_\mathbf{X}$ is the projection matrix of \mathbf{X} . It is shown in Lemma 5.0.2 that the matrix $\{ \mathbf{Z}^T (\mathbf{I}_m - \mathbf{P}_\mathbf{X}) \mathbf{Z} \}$ is positive definite almost surely when $m \geq 2p$ so that its negative square root exists almost surely. Note also that \mathbf{T} is such that

$$\mathbf{T}^T \mathbf{X} = \mathbf{0} \quad \text{and} \quad \mathbf{T}^T \mathbf{T} = \mathbf{I}_p. \quad (12)$$

The proof is virtually identical to the proof of Theorem 2.2.1.

2.3 The Second Form on the Estimator

The form of the estimator in Theorem 2.2.1 makes apparent its close relationship to jackknife estimation. However, it hides the dependence on \bar{W} and $\sum_j (W_j - \bar{W})^2$. The following result gives a simpler and more general form of the estimator.

Theorem 2.3.1 *Suppose that $\hat{\mu}$ and $\hat{\sigma}^2$ are independent random variables such that $\hat{\mu} \sim N(\mu, \tau\sigma^2)$ and $(d\hat{\sigma}^2/\sigma^2) \sim \text{Chi-Squared}(d)$, where τ and d are known. Define the estimator,*

$$\hat{\theta} = E \left\{ g \left(\hat{\mu} + i(\tau d)^{1/2} \hat{\sigma} \left(\frac{Z_1}{\sqrt{Z_1^2 + \dots + Z_d^2}} \right) \right) \middle| \hat{\mu}, \hat{\sigma}^2 \right\}, \quad (13)$$

where the components of $\mathbf{Z} = (Z_1, \dots, Z_d)^T$ are independent and identically distributed $N(0, 1)$ independent of $\hat{\mu}$ and $\hat{\sigma}^2$. Then $E(\hat{\theta}) = g(\mu)$.

Proof. Define $\mathbf{T} = \mathbf{Z}/\|\mathbf{Z}\|$. The random vector \mathbf{T} is not a contrast but is uniformly distributed on the d -dimensional unit ball ($\mathbf{T}^T \mathbf{T} = 1$, almost surely). Let the components of $\mathbf{W}_* = (W_{*1}, \dots, W_{*d})^T$ be iid $N(0, \sigma^2)$ random variables independent of $\hat{\mu}$ and \mathbf{Z} . Then $(\mathbf{T}^T, d^{-1} \mathbf{W}_*^T \mathbf{W}_*, \hat{\mu})^T$ and $(\mathbf{T}^T, \hat{\sigma}^2, \hat{\mu})^T$ have the same distribution. Thus it is sufficient to prove that

$$E \left\{ g \left(\hat{\mu} + i\tau^{1/2} \|\mathbf{W}_*\| \mathbf{T}^T \mathbf{e}_1 \right) \middle| \hat{\mu}, \mathbf{W}_*^T \mathbf{W}_* \right\}$$

has expectation $g(\mu)$, where \mathbf{e}_1 is the $d \times 1$ unit vector with a 1 in the first position.

The distribution of $\mathbf{T}^T \mathbf{v}$ is the same for any \mathbf{v} with $\|\mathbf{v}\| = 1$ (Watson, 1983; Lemma 5.0.2), thus $\mathbf{T}^T \mathbf{e}_1$ and $\mathbf{T}^T \mathbf{W}_*/\|\mathbf{W}_*\|$ have the same distribution, and it is sufficient to show that $\hat{\theta}_* = E \left\{ g \left(\hat{\mu} + i\tau^{1/2} \mathbf{T}^T \mathbf{W}_* \right) \middle| \hat{\mu}, \mathbf{W}_*^T \mathbf{W}_* \right\}$ has mean $g(\mu)$.

Upon taking expectations and then conditioning,

$$E(\hat{\theta}_*) = E \left\{ g \left(\hat{\mu} + i\tau^{1/2} \mathbf{T}^T \mathbf{W}_* \right) \right\} = E \left[E \left\{ g \left(\hat{\mu} + i\tau^{1/2} \mathbf{T}^T \mathbf{W}_* \right) \middle| \mathbf{T} \right\} \right].$$

The argument of the function $g(\cdot)$, $\eta = \hat{\mu} + i\tau^{1/2}\mathbf{T}^T\mathbf{W}_* = \mu + \varepsilon_1^* + i\varepsilon_2^*$, where $\varepsilon_1^* = \hat{\mu} - \mu$ and $\varepsilon_2^* = \tau^{1/2}\mathbf{T}^T\mathbf{W}_*$ are conditionally iid $N(0, \tau\sigma^2)$ given \mathbf{T} . The proof concludes in the same fashion as the proof of Theorem 2.2.1. \spadesuit

Note that under sufficient regularity conditions on $g(\cdot)$, as $d \rightarrow \infty$ and $\hat{\sigma} \rightarrow \sigma$ the estimator in (13) converges to $\hat{\theta} = E \left\{ g \left(\hat{\mu} + i\tau^{1/2}\sigma Z_1 \right) \middle| \hat{\mu} \right\}$.

2.3.1 Monte Carlo Estimation

The Monte Carlo version of (13) is simpler than (10). It is

$$\hat{\theta}_B = \frac{1}{B} \sum_{b=1}^B \operatorname{Re} \left\{ g \left(\hat{\mu} + i(\tau d)^{1/2} \hat{\sigma} T_b \right) \right\}, \quad (14)$$

where T_1, \dots, T_B are iid replicates of the random variable $T = Z_1 / \sqrt{Z_1^2 + \dots + Z_d^2}$, and Z_1, \dots, Z_d are independent and identically distributed $N(0, 1)$. When $d = 1$, $T^2 = 1$ almost surely and (14) is equal to $\operatorname{Re} \left\{ g \left(\hat{\mu} + i\tau^{1/2} \hat{\sigma} \right) \right\}$ independent of B .

3 Variance and Variance Estimation

We now derive variance expressions and variance estimators for the second form of the estimator (13) which we write as

$$\hat{\theta} = E \left\{ g \left(\hat{\mu} + i(\tau d)^{1/2} \hat{\sigma} T \right) \middle| \hat{\mu}, \hat{\sigma} \right\}, \quad (15)$$

where $T = Z_1 / \sqrt{Z_1^2 + \dots + Z_d^2}$.

Let Y_1, Y_2 and Y_3 be real random variables. For the complex random variable, $Y_1 + iY_2$, define $\operatorname{Var}_*(Y_1 + iY_2) = E \{(Y_1 + iY_2)^2\} - \{E(Y_1 + iY_2)\}^2$ and $\operatorname{Var}_*(Y_1 + iY_2 | Y_3) = E \{(Y_1 + iY_2)^2 | Y_3\} - \{E(Y_1 + iY_2 | Y_3)\}^2$. It can be verified that the familiar variance decomposition holds for Var_* ,

$$\operatorname{Var}_*(Y_1 + iY_2) = \operatorname{Var}_* \{E(Y_1 + iY_2 | Y_3)\} + E \{\operatorname{Var}_*(Y_1 + iY_2 | Y_3)\}. \quad (16)$$

When $Y_2 = 0$ almost surely Var_* is the usual variance operator.

When $g(\cdot)$ is an entire function so is $g^2(\cdot)$. Thus, subject to integrability conditions, not only does $E \left\{ g \left(\hat{\mu} + i(\tau d)^{1/2} \hat{\sigma} T \right) \right\} = g(\mu)$, but $E \left\{ g^2 \left(\hat{\mu} + i(\tau d)^{1/2} \hat{\sigma} T \right) \right\} = g^2(\mu)$ as well. It follows that

$$\begin{aligned} \text{Var}_* \left\{ g \left(\hat{\mu} + i(\tau d)^{1/2} \hat{\sigma} T \right) \right\} &= E \left\{ g^2 \left(\hat{\mu} + i(\tau d)^{1/2} \hat{\sigma} T \right) \right\} \\ &\quad - \left[E \left\{ g \left(\hat{\mu} + i(\tau d)^{1/2} \hat{\sigma} T \right) \right\} \right]^2 \\ &= g^2(\mu) - g^2(\mu) = 0. \end{aligned} \quad (17)$$

Combining (16), (17) and the fact that $\hat{\theta} = E \left\{ g \left(\hat{\mu} + i(\tau d)^{1/2} \hat{\sigma} T \right) \mid \hat{\mu}, \hat{\sigma} \right\}$ is real-valued leads to the identity

$$\begin{aligned} \text{Var} \left(\hat{\theta} \right) &= -E \left[\text{Var}_* \left\{ g \left(\hat{\mu} + i(\tau d)^{1/2} \hat{\sigma} T \right) \mid \hat{\mu}, \hat{\sigma}^2 \right\} \right] \\ &= -E \left[E \left\{ g^2 \left(\hat{\mu} + i(\tau d)^{1/2} \hat{\sigma} T \right) \mid \hat{\mu}, \hat{\sigma}^2 \right\} - \hat{\theta}^2 \right]. \end{aligned} \quad (18)$$

Because both $\hat{\theta}$ and $E \left\{ g^2 \left(\hat{\mu} + i(\tau d)^{1/2} \hat{\sigma} T \right) \mid \hat{\mu}, \hat{\sigma}^2 \right\}$ are real-valued functions of $\hat{\mu}$ and $\hat{\sigma}^2$, it follows that

$$\widehat{\text{Var}} \left(\hat{\theta} \right) = \hat{\theta}^2 - E \left\{ g^2 \left(\hat{\mu} + i(\tau d)^{1/2} \hat{\sigma} T \right) \mid \hat{\mu}, \hat{\sigma}^2 \right\} \quad (19)$$

is an unbiased estimator of $\text{Var} \left(\hat{\theta} \right)$. A Monte Carlo version of this is obtained as

$$\widetilde{\text{Var}}_B \left(\hat{\theta} \right) = -\text{Re} \left[\frac{1}{B-1} \sum_{b=1}^B \left\{ g \left(\hat{\mu} + i(\tau d)^{1/2} \hat{\sigma} T_b \right) - \bar{g}_\bullet \right\}^2 \right] \quad (20)$$

where $\bar{g}_\bullet = B^{-1} \sum_{b=1}^B g \left(\hat{\mu} + i(\tau d)^{1/2} \hat{\sigma} T_b \right)$. To see that $\widetilde{\text{Var}}_B \left(\hat{\theta} \right)$ is unbiased note that

$$\begin{aligned} E \left\{ \widetilde{\text{Var}}_B \left(\hat{\theta} \right) \mid \hat{\mu}, \hat{\sigma}^2 \right\} &= -\text{Re} \left(E \left[\frac{1}{B-1} \sum_{b=1}^B \left\{ g \left(\hat{\mu} + i(\tau d)^{1/2} \hat{\sigma} T_b \right) - \bar{g}_\bullet \right\}^2 \mid \hat{\mu}, \hat{\sigma}^2 \right] \right) \\ &= -\text{Re} \left[\text{Var}_* \left\{ g \left(\hat{\mu} + i(\tau d)^{1/2} \hat{\sigma} T \right) \mid \hat{\mu}, \hat{\sigma}^2 \right\} \right], \end{aligned}$$

and thus

$$\begin{aligned}
E \left\{ \widetilde{\text{Var}}_B(\hat{\theta}) \right\} &= E \left(-\text{Re} \left[\text{Var}_* \left\{ g \left(\hat{\mu} + i(\tau d)^{1/2} \hat{\sigma} T \right) \mid \hat{\mu}, \hat{\sigma}^2 \right\} \right] \right) \\
&= -E \left[\text{Var}_* \left\{ g \left(\hat{\mu} + i(\tau d)^{1/2} \hat{\sigma} T \right) \mid \hat{\mu}, \hat{\sigma}^2 \right\} \right] \\
&= \text{Var}(\hat{\theta}).
\end{aligned}$$

Note that $\widetilde{\text{Var}}_B(\hat{\theta})$ is an estimator of $\text{Var}(\hat{\theta})$ and not an estimator of $\text{Var}(\hat{\theta}_B)$ where $\hat{\theta}_B$ is defined in (14). To obtain an estimator of the latter note that

$$\begin{aligned}
\text{Var}(\hat{\theta}_B) &= \text{Var} \left\{ E(\hat{\theta}_B \mid \hat{\mu}, \hat{\sigma}^2) \right\} + E \left\{ \text{Var}(\hat{\theta}_B \mid \hat{\mu}, \hat{\sigma}^2) \right\} \\
&= \text{Var}(\hat{\theta}) + E \left\{ \text{Var}(\hat{\theta}_B \mid \hat{\mu}, \hat{\sigma}^2) \right\}.
\end{aligned} \tag{21}$$

The second component in (21) is unbiasedly estimated by $B^{-1}S_G^2$ where

$$S_G^2 = \frac{1}{B-1} \sum_{b=1}^B (G_b - \bar{G})^2$$

with $G_b = \text{Re} \left\{ g \left(\hat{\mu} + i(\tau d)^{1/2} \hat{\sigma} T_b \right) \right\}$ and $\bar{G} = B^{-1} \sum_{b=1}^B G_b$. Thus an unbiased estimator of $\text{Var}(\hat{\theta}_B)$ is

$$\widehat{\text{Var}}_B(\hat{\theta}_B) = \widetilde{\text{Var}}_B(\hat{\theta}) + B^{-1}S_G^2 \tag{22}$$

A simple example helps understand these unusual formulas. Consider the function $g(\mu) = \mu^2$. The facts $E(T) = E(T^3) = 0$, $E(T^2) = 1/d$ and $E(T^4) = 3/(d^2 + 2d)$ are used repeatedly in the following. The estimator defined in (15) is, $\hat{\theta} = E \left\{ \left(\hat{\mu} + i(\tau d)^{1/2} \hat{\sigma} T \right)^2 \mid \hat{\mu}, \hat{\sigma}^2 \right\} = \hat{\mu}^2 - \tau \hat{\sigma}^2$, and the Monte Carlo version (14) of the estimator is $\hat{\theta}_B = \hat{\mu}^2 - \tau d \hat{\sigma}^2 \left(B^{-1} \sum_{b=1}^B T_b^2 \right)$, From (18), (19) and (20) one finds $\text{Var}(\hat{\theta}) = 4\mu^2\tau\sigma^2 + 2\tau^2\sigma^4 + 2\tau^2\sigma^4/d$, $\widetilde{\text{Var}}(\hat{\theta}) = 4\hat{\mu}^2\tau\hat{\sigma}^2 + \tau^2\hat{\sigma}^4 - 3d\tau^2\hat{\sigma}^4/(d+2)$, and $\widetilde{\text{Var}}_B(\hat{\theta}) = 4\hat{\mu}^2\tau d \hat{\sigma}^2 S_T^2 - \tau^2 d^2 \hat{\sigma}^4 S_{T^2}^2$, where S_T^2 and $S_{T^2}^2$ are the sample variances of $\{T_1, \dots, T_B\}$ and $\{T_1^2, \dots, T_B^2\}$. Finally, from (21) and (22) one finds $\text{Var}(\hat{\theta}_B) = \text{Var}(\hat{\theta}) + 2(d-1)\tau^2\sigma^4/(Bd)$ and $\widehat{\text{Var}}_B(\hat{\theta}_B) = \widetilde{\text{Var}}_B(\hat{\theta}) + \tau^2 d^2 \hat{\sigma}^4 S_{T^2}^2/B$.

4 Applications

We now present sample applications of the basic estimation method. A simulation study of the fundamental problem of estimating a nonlinear function of a normal mean is given in Section 4.1. In Sections 4.2 and 4.3 we describe the relevance of the method, and outline its application, to regression measurement error models and deconvolution respectively. Thorough development and study of the latter two applications is beyond the scope of the current paper and will be the topic of forthcoming papers. A key feature of the current method applied to measurement error problems is the automatic handling of heteroscedastic measurement error.

4.1 Estimating $g(\mu)$

We conducted a small simulation to illustrate the performance of the Monte Carlo estimator (14) for different values of B . The function chosen for the simulation is $g(\mu) = \mu \sin(\mu) - \mu(\mu - \pi)$. Other details of the simulation design and the results are reported in Table 1. The bias in the maximum likelihood estimator is evident as is the lack of bias in the Monte Carlo estimators. The variance estimator is from (22). For this particular example there is little advantage of taking $B > 32$.

4.2 Errors-in-Variables Regression

The corrected score method is a generally applicable method for the analysis of data measured with error (Stefanski, 1989; Nakamura, 1990; Carroll et al., 1995). The Monte Carlo estimators in (10) and (14) can be used to extend the method of Monte Carlo corrected score equations (Novick and Stefanski, 2002) to the case of replicate measurements with heteroscedastic errors. We describe the application in the context of regression modelling. The new method is an alternative to, and in a certain sense an exact implementation of, the empirical simulation extrapolation method (simex)

Estimator	MLE	$B = 2$	$B = 8$	$B = 32$	$B = 128$	$B = 512$
Mean	-1.60	0.01	0.01	0.01	0.01	0.01
Std Err	0.02	0.02	0.02	0.02	0.02	0.02
t_{CALC}	-93.04	0.40	0.56	0.65	0.55	0.54
MC Variance	29.65	60.44	54.11	52.52	52.03	51.92
Mean $\widehat{\text{Var}}$		60.41	54.42	51.80	51.74	51.50
Std Err $\widehat{\text{Var}}$		0.79	0.51	0.42	0.43	0.44

Table 1: Simulation design: $g(\mu) = \mu \sin(\mu) - \mu(\mu - \pi)$; $\mu = \pi$; $\theta = g(\pi) = 0$; $\tau = 1$; $\sigma^2 = 1$; $d = 10$; $\widehat{\mu} \sim N(\mu, \tau\sigma^2)$; $\widehat{\sigma}^2 \sim \sigma^2\chi^2(d)/d$; Number of replications = 100,000. Estimator: MLE = $g(\widehat{\mu})$; $B = \bullet$, Monte Carlo estimator (14) with $B = k$. Mean, Std Err and t_{CALC} , Monte Carlo mean, standard error, and t -statistic for testing mean = 0. MC Variance, sample variance of 100,000 replicates; Mean $\widehat{\text{Var}}$, Mean of 100,000 variance estimates calculated using (22); Std Err $\widehat{\text{Var}}$, standard error of Mean $\widehat{\text{Var}}$.

studied by Devanarayan and Stefanski (2002).

Suppose that in the absence of measurement error the regression parameter β is consistently estimated by the m-estimator $\widehat{\beta}$ solving

$$\sum_{r=1}^n \psi(Y_r, \mathbf{Z}_r, X_r, \beta) = \mathbf{0}, \quad (23)$$

where Y_r , \mathbf{Z}_r and X_r are a response variable, a vector of predictors measured without error and a predictor subject to measurement error respectively. Suppose that iid replicate measurements $W_{r,1}, \dots, W_{r,m_r}$ of X_r are available such that $W_{r,j} = X_r + U_{r,j}$ where $U_{r,1}, \dots, U_{r,m_r}$ are independent and identically distributed $N(0, \sigma_r^2)$. Let \overline{W}_r and $\widehat{\sigma}_r^2$ denote the sample mean and variance of the components of $\mathbf{W}_r = (W_{r,1}, \dots, W_{r,m_r})^T$.

Define $T_{b,r} = Z_{b,1} / \sqrt{Z_{b,1}^2 + \dots + Z_{b,m_r-1}^2}$, and the $Z_{b,j}$ are independent and identically distributed $N(0, 1)$. Then provided that the components of $\psi(\cdot)$ are entire functions of their third argument, the Monte Carlo score function

$$\psi_B(Y_r, \mathbf{Z}_r, \mathbf{W}_r, \beta) = \frac{1}{B} \sum_{b=1}^B \text{Re} \left\{ \psi \left(Y_r, \mathbf{Z}_r, \overline{W}_r + i(m_r - 1)^{1/2} \widehat{\sigma}_r T_{b,r}, \beta \right) \right\},$$

has the property that

$$E \{ \boldsymbol{\psi}_B (Y_r, \mathbf{Z}_r, \mathbf{W}_r, \boldsymbol{\beta}) \mid Y_r, \mathbf{Z}_r, X_r \} = \boldsymbol{\psi} (Y_r, \mathbf{Z}_r, X_r, \boldsymbol{\beta}). \quad (24)$$

This is the key property of a corrected score. It follows from (24) that the system of equations $\sum_{r=0}^n \boldsymbol{\psi}_B (Y_r, \mathbf{Z}_r, \mathbf{W}_r, \boldsymbol{\beta}) = \mathbf{0}$, admits a consistent sequence of m-estimators; and thus asymptotic inference is relatively straight forward (Novick and Stefanski, 2002).

For example, for Poisson regression the likelihood score function, $\psi(Y, X, \boldsymbol{\beta}) = (Y - \exp(\boldsymbol{\beta}_1 + \boldsymbol{\beta}_x X))(1, X)^T$, satisfies the smoothness requirements, as do score functions for regression models that are polynomial in \mathbf{Z} and X . An important advance over the methods in Novick and Stefanski (2002) is the accommodation of heteroscedastic measurement errors (σ_r^2) with a finite number (m_r) of replicates.

As an illustration of the method for polynomials we present results from an analysis of a subset of data from the Framingham Heart Study described in Carroll et al., (1995). The variable measured with error is systolic blood pressure (SBP), for which there are two measurements available, and the response variable in our analysis is the average of two measurements of serum cholesterol (SC). Preliminary analyses suggested that the relationship of SC to SBP is nonlinear and that the error in the SBP measurements is not homoscedastic (note that we did not transform SBP as is done by Carroll et al., 1995). The results from several quadratic regressions are presented in Table 2.

We fit three so-called naive quadratic models using as predictor variable: i) the first exam measurement W_1 ; ii) the second exam measurement W_2 ; and iii) the average of the two measurements \bar{W} . These allowed us to construct a crude jackknife corrected estimate, $2\hat{\boldsymbol{\beta}}_{\bar{W}} - (\hat{\boldsymbol{\beta}}_{W_1} + \hat{\boldsymbol{\beta}}_{W_2})/2$, with which to compare the Monte Carlo corrected score estimate. Attenuation in the naive estimates is apparent from the

	W_1	W_2	\bar{W}	JK	MCCS
Intercept	1.80 (1.09)	1.43 (1.09)	1.97 (1.09)	2.33	2.51 (1.18)
Linear	0.31 (0.06)	0.30 (0.06)	0.37 (0.07)	0.44	0.47 (0.08)
Quadratic	-0.0046 (0.0013)	-0.0036 (0.0012)	-0.0056 (0.0015)	-0.0072	-0.0081 (0.0020)

Table 2: Coefficient estimates from the Framingham data quadratic regression model with standard errors in parentheses. W_1 , regression on Exam I measurement; W_2 , regression on Exam II measurement; \bar{W} , regression on average; JK, jackknife bias corrected; MCCS, Monte Carlo corrected score.

trend in the latter three (\bar{W} , JK, MCCS) sets of regression coefficients. The Monte Carlo estimates are consistent with the crude jackknife estimates in light of their approximate (conservative) bias corrections.

4.3 An Application of Estimating $\{\sin(\mu)/\mu\}^{2k}$

We now show that our results can be used for estimating a density function from replicate measurements. The estimator is a generalization of the deconvoluting-kernel density estimator for normal measurement error studied by Carroll and Hall (1988), Fan (1991a,b,c, 1992), Hesse (1999), Stefanski (1989), Stefanski (1990), Stefanski and Carroll (1990), and Wand (1998).

Suppose that the measurements $W_{j,r}$ and true values X_r are distributed as in the previous section, with the additional assumption that the X_r are independent and identically distributed with density f_X . For a suitable constant c the function $Q(\mu) = c\{\sin(\mu)/\mu\}^{2k}$ is a probability density function. It is also an entire function to which Theorem 2.3.1 applies and thus for fixed x and $\lambda > 0$, conditionally on the

true value X_r ,

$$E \left\{ \frac{1}{B} \sum_{b=1}^B \frac{1}{\lambda} \operatorname{Re} \left(Q \left[\frac{x - \{\overline{W}_r + i(1 - 1/m_r)^{1/2} \hat{\sigma}_r T_{b,r}\}}{\lambda} \right] \right) \right\} = \frac{1}{\lambda} Q \left(\frac{x - X_r}{\lambda} \right)$$

It follows that

$$\hat{f}_X(x) = \frac{1}{n} \sum_{r=1}^n \frac{1}{B} \sum_{b=1}^B \frac{1}{\lambda} \operatorname{Re} \left(Q \left[\frac{x - \{\overline{W}_r + i(1 - 1/m_r)^{1/2} \hat{\sigma}_r T_{b,r}\}}{\lambda} \right] \right) \quad (25)$$

is, conditionally on X_1, \dots, X_n , an unbiased estimator of $\frac{1}{n\lambda} \sum_{r=1}^n Q \left(\frac{x - X_r}{\lambda} \right)$, which is a kernel density estimator of f_X constructed with the kernel Q . So \hat{f}_X is an estimator of f_X that has the same bias as a kernel density estimator based on the true data. Its construction and conditional expectation properties are essentially identical to the deconvoluting-kernel density estimators studied in the references cited above. The key advantage is that the new estimator allows for heteroscedastic measurement error. A generalization and a numerical illustration is given in the next section.

4.3.1 Analysis of Processed Data

Consider a situation where for each of $r = 1, \dots, n$ subjects a full-rank linear model $\mathbf{Y}_r = \mathbf{X}_r \boldsymbol{\beta}_r + \boldsymbol{\epsilon}_r$, $\boldsymbol{\epsilon}_r \sim N(\mathbf{0}, \sigma_r^2 \mathbf{I})$, is fit resulting in the n statistics $\hat{\boldsymbol{\beta}}_r$, $\hat{\sigma}_r^2$ and $(\mathbf{X}_r^T \mathbf{X}_r)^{-1}$, where $\hat{\sigma}_r^2$ is the mean squared error on $d_r > 0$ degrees of freedom. Then for a fixed \mathbf{a} , $\mathbf{a}^T \hat{\boldsymbol{\beta}}_r \sim N(\mathbf{a}^T \boldsymbol{\beta}_r, \sigma_r^2 \mathbf{a}^T (\mathbf{X}_r^T \mathbf{X}_r)^{-1} \mathbf{a})$. Note that there are two sources of heteroscedasticity in this case, subject-specific heteroscedasticity (σ_r^2), and design-induced heteroscedasticity (\mathbf{X}_r).

Suppose now that the regression coefficients $\boldsymbol{\beta}_r$ are modelled as independent and identically distributed random vectors and interest lies in estimating the density of $\mathbf{a}^T \boldsymbol{\beta}_r$. Define $\tau_r = \mathbf{a}^T (\mathbf{X}_r^T \mathbf{X}_r)^{-1} \mathbf{a}$. Then appealing to Theorem 2.3.1 to generalize (25) to the case under consideration shows that

$$\hat{f}_B(x) = \frac{1}{n} \sum_{r=1}^n \frac{1}{B} \sum_{b=1}^B \frac{1}{\lambda} \operatorname{Re} \left(Q \left[\frac{x - \{\mathbf{a}^T \hat{\boldsymbol{\beta}}_r + i(d_r \tau_r)^{1/2} \hat{\sigma}_r T_{b,r}\}}{\lambda} \right] \right) \quad (26)$$

is an estimator with the property that conditioned on $\mathbf{a}^T \boldsymbol{\beta}_1, \dots, \mathbf{a}^T \boldsymbol{\beta}_n$

$$E \{ \hat{f}_B(x) \} = \frac{1}{n\lambda} \sum_{r=1}^n Q \left(\frac{x - \mathbf{a}^T \boldsymbol{\beta}_r}{\lambda} \right). \quad (27)$$

The expectation estimated by the r^{th} Monte Carlo mean in (26) is

$$Q^* \left(\frac{x - \mathbf{a}^T \hat{\boldsymbol{\beta}}_r}{\lambda}, \lambda, d_r, \tau_r, \hat{\sigma}_r \right)$$

where

$$Q^*(z, \lambda, d_r, \tau_r, \hat{\sigma}_r) = \frac{\Gamma\left(\frac{d_r}{2}\right)}{2\pi} \int_{-\infty}^{\infty} e^{-itz} \left\{ \frac{\theta(t)}{2} \right\}^{1-d_r/2} J_{d_r/2-1}(\theta(t)) C(t) dt,$$

with $\theta(t) = -it(d_r \tau_r)^{1/2} \hat{\sigma}_r / \lambda$, and $C(t)$ is the characteristic function of $Q(t)$. Thus as $B \rightarrow \infty$ the estimator $\hat{f}_B(x)$ converges to

$$\hat{f}_\infty(x) = \frac{1}{n\lambda} \sum_{r=1}^n Q^* \left(\frac{x - \mathbf{a}^T \hat{\boldsymbol{\beta}}_r}{\lambda}, \lambda, d_r, \tau_r, \hat{\sigma}_r \right).$$

The function $Q^*(z, \lambda, d_r, \tau_r, \hat{\sigma}_r)$ plays the same role in the current case (heteroscedastic measurement error with independent variance estimates) as the deconvoluting kernel defined in Stefanski and Carroll (1990) does in the homoscedastic measurement error, know-variance case. In fact, if $\hat{\sigma}_r$ is set equal to σ_r , then in the limit as $d_r \rightarrow \infty$, $Q^*(z, \lambda, d_r, \tau_r, \hat{\sigma}_r)$ converges to the deconvoluting kernel defined in Stefanski and Carroll (1990).

The estimators defined above assume heteroscedastic residual variation in the subject-specific linear models. For the case in which the assumption $\sigma_1^2 = \dots = \sigma_n^2$ is warranted, the estimators defined above apply upon replacing $\hat{\sigma}_r^2$ and d_r with the pooled estimator of variance, and pooled degrees of freedom for each r .

We illustrate the estimator (26) using a simulated data set of size $n = 300$. Random simple linear regression models were generated by: 1) generating a random sample size m_r uniformly distributed on $\{3, 4, 5, 6, 7\}$; 2) generating random design

points by adding $N(0, .04)$ errors to $\{1, 2, \dots, m_r\}$; 3) generating a random slope from a shifted (to mean 1) and scaled (to variance 1) chi-squared(4) distribution; 4) generating Y_r according to the linear model with zero slope and random intercept, and common $\sigma_r = 2.0$ (i.e., homoscedastic residual error variation). The deconvoluting estimator of the random slopes was calculated as in (26) with $B = 50$ under the (correct) assumption of homoscedastic error variation (using the pooled estimate of variance) and under the (also correct but less precise) assumption of heteroscedastic residual error variances (using separate mean squared errors and degrees of freedom).

We used the kernel proportional to $\{\sin(x)/x\}^6$ standardized to have variance 1. Rates of convergence for deconvolution estimators with normal measurement error and known error variance are notoriously slow (Carroll and Hall, 1988), because variance increases exponentially as the bandwidth decreases. The intent of the example is to illustrate the unbiasedness property (27), and we did not attempt to estimate the bandwidth, but rather fixed it at the relatively large value $\lambda = 1.0$. A forthcoming paper will study convergence rates of (26) and examine the feasibility of deconvolution with normal errors and estimated variances in moderate sample sizes, as Wand (1998) has done for the known variance case.

Plots of estimated densities are displayed in Figure 1. The first panel contains the true density and the kernel estimator constructed using the error-free data, and makes evident the large amount of smoothing induced by the bandwidth. The second panel contains the estimator from the error-free data and the kernel estimator constructed from the error-contaminated data. The third and fourth panels display the two estimates from the second panel and the deconvolution estimates calculated under the assumptions of heteroscedastic and homoscedastic errors respectively.

We also calculated estimated moments, mean, variance, third and fourth central moments, for these simulated data using the true data, the observed data, and the

	True Data	Observed Data	Heteroscedastic	Homoscedastic
Mean	1.05	1.09	1.09	1.09
Variance	1.22	1.66	1.27	1.24
Third	2.30	2.20	2.22	2.18
Fourth	10.06	13.08	9.89	9.41

Table 3: Estimated moments of the slope from the simulated random coefficient model data.

m-estimation method described in Section 4.2 (in this case there is no response Y_r or error-free covariate \mathbf{Z}_r , and ψ is chosen to yield the indicated moments). The estimates are given in Table 3.

4.3.2 Comments on the Examples

A notable feature of the examples presented above which is not apparent in the results is the ease with which the Monte Carlo estimates were obtained. We exploited the complex number capabilities available in GAUSS to use the same code that would be used for standard analyses, evaluating functions at the complex arguments as indicated, and invoking GAUSS’ “real-part” function as necessary. For the regression model example, the estimating equations were passed to an equation solver that provided solutions and sandwich-formula, asymptotic variance matrix estimates.

4.4 Robustness

In applications robustness is often an issue. The estimators we propose depend on normality and the smoothness assumptions on g . Novick and Stefanski (2002) provided evidence of robustness for a similar class of estimators. Whether our estimators possess a similar level of robustness is an open question that we plan to address in forthcoming papers that develop the applications in greater detail.

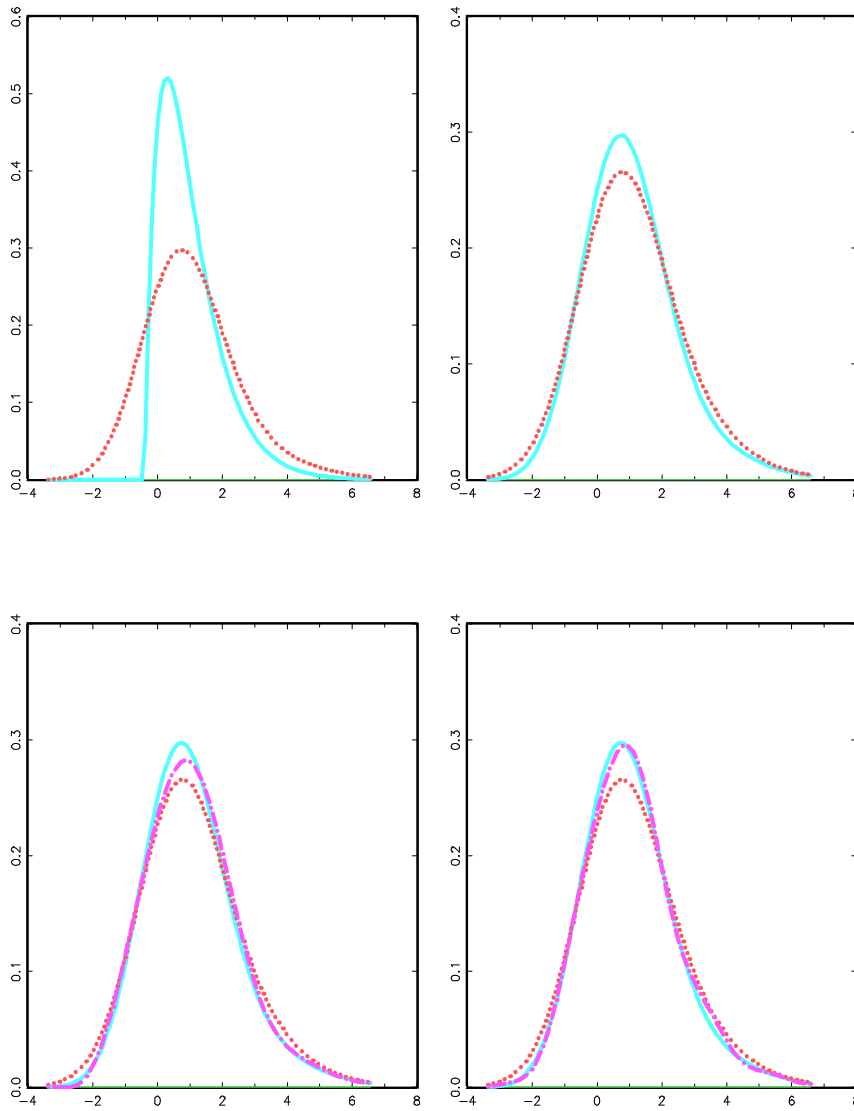


Figure 1: Density Estimates. Panel 1: solid, true density; dotted, true-data estimate. Panel 2: solid, true-data estimate; dotted, observed-data estimate. Panel 3: solid, true-data estimate; dotted, observed-data estimate; dot-dash, deconvolution estimate assuming heteroscedasticity. Panel 4: solid, true-data estimate; dotted, observed-data estimate; dot-dash, deconvolution estimate assuming homoscedasticity. Note that Panels 2-4 are scaled differently than Panel 1.

5 Appendix: Technical Details

Lemma 5.0.1 *If ε_1 and ε_2 are iid $N(0, \tau^2)$, then $E\{(\varepsilon_1 + i\varepsilon_2)^k\} = 0$, for $k = 1, 2, 3, \dots$*

Proof. The moment generating function of $\varepsilon_1 + i\varepsilon_2$ is equal to 1 and thus all of its derivatives are 0. ♠

Lemma 5.0.2 *Let \mathbf{Z} , \mathbf{T} and $\mathbf{P}_\mathbf{X}$ be defined as in (11). Then: (i) $\mathbf{Z}^T (\mathbf{I}_m - \mathbf{P}_\mathbf{X}) \mathbf{Z}$ is positive definite almost surely if and only if $m \geq 2p$; (ii) for fixed non-random vectors \mathbf{u} and \mathbf{s} the distribution of $\mathbf{u}^T \mathbf{T}^T \mathbf{s}$ is symmetric, possesses moments of all orders, and depends on \mathbf{s} only through $\mathbf{s}^T (\mathbf{I}_m - \mathbf{P}_\mathbf{X}) \mathbf{s}$.*

Proof. Because $\text{rank}(\mathbf{Z}^T (\mathbf{I}_m - \mathbf{P}_\mathbf{X}) \mathbf{Z}) \leq \text{rank}(\mathbf{I}_m - \mathbf{P}_\mathbf{X}) = m - p$, $m - p \geq p$ is necessary. Now assume that $m \geq 2p$. Write $(\mathbf{I}_m - \mathbf{P}_\mathbf{X}) = \mathbf{Q} \mathbf{D} \mathbf{Q}^T$, where \mathbf{Q} is orthogonal and \mathbf{D} is diagonal with first $m - p$ entries equal to 1, and last p entries equal to 0. Because $\mathbf{Q}^T \mathbf{Z}$ has the same distribution as \mathbf{Z} , $(\mathbf{Z}^T (\mathbf{I}_m - \mathbf{P}_\mathbf{X}) \mathbf{Z})$ and $\mathbf{Z}^T \mathbf{D} \mathbf{Z}$ are identically distributed. The latter is equal to $\sum_{i=1}^{m-p} \mathbf{r}_i \mathbf{r}_i^T$ where \mathbf{r}_i^T is the i^{th} row of \mathbf{Z} . The matrix $\sum_{i=1}^{m-p} \mathbf{r}_i \mathbf{r}_i^T$ has a Wishart distribution and is positive definite almost surely provided $m - p \geq p$ (Srivastava and Khatri, 1979, p 72).

Symmetry of the distribution of $\mathbf{u}^T \mathbf{T}^T \mathbf{s}$ follows from the fact that

$$\mathbf{u}^T \mathbf{T}^T \mathbf{s} = \mathbf{u}^T \{ \mathbf{Z}^T (\mathbf{I}_m - \mathbf{P}_\mathbf{X}) \mathbf{Z} \}^{-1/2} \mathbf{Z}^T (\mathbf{I}_m - \mathbf{P}_\mathbf{X}) \mathbf{s} = h(\mathbf{Z}), \quad (28)$$

where $h(-\mathbf{z}) = -h(\mathbf{z})$ and $h(\mathbf{Z})$ and $h(-\mathbf{Z})$ have the same distribution. From the second property in (12), $|\mathbf{u}^T \mathbf{T}^T \mathbf{s}| \leq \|\mathbf{u}\| \|\mathbf{s}\|$ from which it follows that all moments of $\mathbf{u}^T \mathbf{T}^T \mathbf{s}$ exist.

It remains to show that the distribution of $\mathbf{u}^T \mathbf{T}^T \mathbf{s}$ depends on \mathbf{s} only through $\mathbf{s}^T (\mathbf{I}_m - \mathbf{P}_\mathbf{X}) \mathbf{s}$. Note that by (28) if \mathbf{s} is in the column space of \mathbf{X} then $\mathbf{u}^T \mathbf{T}^T \mathbf{s} =$

0 almost surely and the claim is trivially true. Now consider the case in which $(\mathbf{I}_m - \mathbf{P}_x) \mathbf{s} \neq \mathbf{0}$. Define

$$\mathbf{P}_\dagger = \frac{(\mathbf{I}_m - \mathbf{P}_x) \mathbf{s} \mathbf{s}^T (\mathbf{I}_m - \mathbf{P}_x)}{\|(\mathbf{I}_m - \mathbf{P}_x) \mathbf{s}\|^2}.$$

From the decomposition $(\mathbf{I}_m - \mathbf{P}_x) \mathbf{Z} = \mathbf{P}_\dagger (\mathbf{I}_m - \mathbf{P}_x) \mathbf{Z} + (\mathbf{I}_m - \mathbf{P}_\dagger) (\mathbf{I}_m - \mathbf{P}_x) \mathbf{Z}$, it follows that $\mathbf{Z}^T (\mathbf{I}_m - \mathbf{P}_x) \mathbf{Z} = \mathbf{Z}_1^T \mathbf{Z}_1 + \mathbf{Z}_2^T \mathbf{Z}_2$ where $\mathbf{Z}_2 = (\mathbf{I}_m - \mathbf{P}_\dagger) (\mathbf{I}_m - \mathbf{P}_x) \mathbf{Z}$, and $\mathbf{Z}_1 = \mathbf{s}^T (\mathbf{I}_m - \mathbf{P}_x) \mathbf{Z} / \|(\mathbf{I}_m - \mathbf{P}_x) \mathbf{s}\|$. Thus

$$\mathbf{s}^T \mathbf{T} = \|(\mathbf{I}_m - \mathbf{P}_x) \mathbf{s}\| \mathbf{Z}_1 (\mathbf{Z}_1^T \mathbf{Z}_1 + \mathbf{Z}_2^T \mathbf{Z}_2)^{-1/2}. \quad (29)$$

We now show that: \mathbf{Z}_1^T has a p -variate standard normal distribution; \mathbf{Z}_1^T and \mathbf{Z}_2^T are independent; and the distribution of $\mathbf{Z}_2^T \mathbf{Z}_2$ does not depend on \mathbf{s} . It follows that the distribution of $\mathbf{s}^T \mathbf{T}$ depends on \mathbf{s} only through the scalar multiplier $\|(\mathbf{I}_m - \mathbf{P}_x) \mathbf{s}\|$ in (29).

Note that $\mathbf{Z}_1^T = \mathbf{Z}^T \mathbf{v} / \|\mathbf{v}\|$ where $\mathbf{v} = (\mathbf{I}_m - \mathbf{P}_x) \mathbf{s}$. Because the rows of \mathbf{Z}^T are iid it follows that the rows of \mathbf{Z}_1^T are also iid. Let \mathbf{c}_1^T denote the first row of \mathbf{Z}^T . Then the first component of \mathbf{Z}_1^T is $\mathbf{c}_1^T \mathbf{v} / \|\mathbf{v}\| \sim N(0, 1)$.

Independence of the rows of \mathbf{Z}^T implies that the i^{th} row of \mathbf{Z}_1^T is independent of the j^{th} row of \mathbf{Z}_2^T for $i \neq j$. Also the rows are identically distributed so it is only necessary to establish that the first rows of \mathbf{Z}_1^T and \mathbf{Z}_2^T are independent. These are $\mathbf{c}_1^T (\mathbf{I}_m - \mathbf{P}_x) \mathbf{s} / \|\mathbf{v}\|$ and $\mathbf{c}_1^T (\mathbf{I}_m - \mathbf{P}_x) (\mathbf{I}_m - \mathbf{P}_\dagger)$ respectively, which are uncorrelated and hence independent.

The matrix $\mathbf{Z}_2^T \mathbf{Z}_2 = \mathbf{Z}^T \mathbf{M} \mathbf{Z}$ where $\mathbf{M} = (\mathbf{I}_m - \mathbf{P}_x) (\mathbf{I}_m - \mathbf{P}_\dagger) (\mathbf{I}_m - \mathbf{P}_x)$. By the spectral composition we can find a diagonal matrix \mathbf{D}_* with 0 – 1 entries such that $\mathbf{M} = \mathbf{Q}_* \mathbf{D}_* \mathbf{Q}_*^T$ where \mathbf{Q}_*^T is orthonormal and thus $\mathbf{Q}_*^T \mathbf{Z}$ and \mathbf{Z} are identically distributed. It follows that $\mathbf{Z}_2^T \mathbf{Z}_2$ and $\mathbf{Z}^T \mathbf{D}_* \mathbf{Z}$ are identically distributed and the distribution of the latter obviously does not depend on \mathbf{s} . ♠

When $p = 1$ and $\mathbf{P}_x = m^{-1}\mathbf{1}\mathbf{1}^T$, the vector $\mathbf{T} = (T_1, \dots, T_m)^T$ is such that $\sum T_j = 0$ and $\sum T_j^2 = 1$. Thus \mathbf{T} lies on the intersection of an $(m - 1)$ -dimensional hyperplane and the m -dimensional unit ball which is an $(m - 1)$ -dimensional ball. In this case \mathbf{T} is equal in distribution to $\mathbf{A}\mathbf{T}_*$ where $\mathbf{A}_{m \times (m-1)}$ is a matrix such that $\mathbf{A}\mathbf{A}^T = \mathbf{I}_m - m^{-1}\mathbf{1}\mathbf{1}^T$ and $\mathbf{A}^T\mathbf{A} = \mathbf{I}_{m-1}$, and \mathbf{T}_* is uniformly distributed on the $(m - 1)$ -dimensional unit ball. So $\mathbf{T}^T\mathbf{s}$ and $\mathbf{T}_*^T\mathbf{A}^T\mathbf{s}$ have the same distribution and the conclusions of Lemma 5.0.2 can be deduced from known properties of the distribution of \mathbf{T}_* . (Watson, 1983, p 45).

REFERENCES

- Carroll, R. J. & Hall, P. (1988). Optimal rates of convergence for deconvolving a density. *Journal of the American Statistical Association*, 83, 1184-1186.
- Carroll, R. J., Ruppert, D. and Stefanski, L. A. (1995). *Measurement Error in Non-linear Models*, London: Chapman Hall.
- Devanarayan, V. & Stefanski, L. A. (2002), Empirical Simulation Extrapolation for Measurement error Models with Replicate Measurements, *Statistics and Probability Letters*, 59, 219-225.
- Fan, J. (1991a). On the optimal rates of convergence for nonparametric deconvolution problems. *The Annals of Statistics*, 19, 1257-1272.
- Fan, J. (1991b). Asymptotic normality for deconvolving kernel density estimators, *Sankhyā, Series A*, 53, 97-110.
- Fan, J. (1991c). Global behavior of deconvolution kernel estimates. *Statistica Sinica*, 1, 541-551.
- Fan, J. (1992). Deconvolution with supersmooth distributions. *Canadian Journal of Statistics*, 20, 23-37.
- Gray H. L., Watkins, T. A., and Schucany W. R., (1973). On the Jackknife Statistic and its Relation to UMVU Estimators in the Normal Case, *Communications in Statistics*, 2(4), 285-320.

- Hesse, C. H. (1999). Data-driven deconvolution. *Journal of Nonparametric Statistics*, 10, 343-373.
- Kitagawa, T. (1956). The operational calculus and the estimations of functions or parameters admitting sufficient statistics, *Bulletin of Mathematical Statistics*, 6, 95-108.
- Kolmogorov, A. N. (1950). Unbiased estimates, *Izvestiya Akademii Nauk SSSR, Ser. Mat.*, 14, 303-326; see also *Am. Math. Soc. Translation* No. 98.
- Nakamura, T. (1990). Corrected score functions for errors-in-variables models: methodology and application to generalized linear models. *Biometrika*, 77, 127-137.
- Neyman, J., and Scott, E. L. (1960). Corrections for bias introduced by a transformation of variables, *Annals of Mathematical Statistics*, 31, 643-655.
- Novick, S. J. and Stefanski, L. A. (2002). Corrected Score Estimation via Complex Simulation Extrapolation, *Journal of the American Statistical Association*, 97, 472-481.
- Srivastava, M. S. and Khatri, C. G. (1979). *An Introduction to Multivariate Statistics*, North Holland, New York.
- Stefanski, L. A. (1989). Unbiased estimation of a nonlinear function of a normal mean with application to measurement error models. *Communications in Statistics, Theory & Methods*, 18, 4335-4358.
- Stefanski, L. A. (1990). Rates of convergence of some estimators in a class of deconvolution problems. *Statistics & Probability Letters*, 9, 229-235.
- Stefanski, L. A. & Carroll, R. J. (1990). Deconvoluting kernel density estimators. *Statistics*, 21, 165-184.
- Wand, M. P. (1998). Finite sample performance of deconvoluting density estimators. *Statistics & Probability Letters*, 37, 131-139.
- Washio, Y., Morimoto, H., and Ikeda, N. (1956). Unbiased estimation based on sufficient statistics, *Bulletin of mathematical Statistics*, 6, 69-93.
- Watkins T. A. (1979). Estimating A Function Of The Mean Of A Multivariate Normal-Distribution, *Communications in Statistics A, Theory*, 8, 245-256.
- Watson, Geoffrey S. (1983). *Statistics on Spheres*, John Wiley & Sons, Inc., New

York.

Woodward, W. A. (1977). Subroutines For Computing Minimum Variance Unbiased Estimators Of Functions Of Parameters In Normal And Gamma Distributions, *Communications in Statistics B, Simulation*, 6, 63-73.