INSTITUTO NACIONAL DE ESTATÍSTICA
STATISTICS PORTUGAL

# REVSTAT
## Statistical Journal

Special issue on
«AISC 2018: Advances in Interdisciplinary
Statistics and Combinatorics»

# REVSTAT

## Statistical Journal

# PREFACE

This special issue features specially invited papers from those who presented at the International Conference on Advances in Interdisciplinary Statistics and Combinatorics (AISC 2018) held at the University of North Carolina – Greensboro, USA, during October 5–7, 2018.

The contributions to this special issue cover several very significant areas of statistics such as Randomized Response Models, Small Area Estimation, Genetics, Statistical Tests, Distribution Theory, and Spatial Statistics.

The guest editors are grateful to the contributors to this issue as well as the editors of REVSTAT for their support during the review process. We also wish to acknowledge the help of the referees who reviewed the papers very promptly and diligently.

Guest Editors:

Sat Gupta
Professor & Head
Department of Mathematics and Statistics
University of North Carolina at Greensboro
Greensboro, NC 27412, USA

Kumer Pial Das
Professor of Mathematics &
AVP for Research, Innovation and Economic Development and Assistant Provost
University of Louisiana at Lafayette
Lafayette, LA 70504, USA

# INDEX

# SIMULTANEOUS INFERENCE OF GENE ISOFORM EXPRESSION FOR RNA SEQUENCING DATA

Author:    BO LI
– Department of Mathematical Sciences,
The Citadel, The Military College of South Carolina,
Charleston, SC 29409, South Carolina, USA
bli@citadel.edu

Abstract:

- In this article, we describe simultaneous inferential methods in detecting differentially expressed gene isoforms based on the Poisson generalized linear models. We derive the joint asymptotic distribution of pivotal quantities. The sample size of RNA sequencing data is often small in practice. Using multiple comparison procedures based on large-sample approximation becomes problematic. The parametric bootstrap method based on pivotal quantities is outlined as a robust alternative. Moreover, we observe the validity of robustness of the bootstrap method when mild overdispersion presents in RNA-sequencing data. We demonstrate the validity of the proposed method in detecting differentially expressed isoforms through Monte Carlo simulation. It shows the proposed method controls the family-wise error rate for large-scale inference. Even though the proposed method can be extended to many experimental designs, we focus on factorial designs in this article.

## 1.    INTRODUCTION

Studies of Gene isoform expression have not only been concentrated on detecting differentially expressed genes with known gene bank ID but also their isoforms due to the development of RNA sequencing technology. RNA sequencing technology, also known as Next Generation Sequencing (NGS), counts how many copies of nucleotide sequence for hundreds to thousands of gene isoforms.

To detect which genes are differentially expressed among hundreds even thousands of genes, researchers often conduct large-scale multiple hypotheses tests simultaneously, see Dudoit *et al.* [3]. One of the major concerns of gene expression analysis is to control the family-wise error rate (FWER). When the multiplicity is overlooked, researchers may claim dozens even hundreds of genes which are differentially expressed but in fact, they are false positives. Concerted efforts have been devoted to controlling FWER for microarray gene expression analysis. Dudoit *et al.* [4] applied Westfall and Young step-down method (Westfall and Young, [13]) based on two-sample Welch's *t*-tests to detect differentially expressed genes in microarray experiments. Alternatively, simultaneous confidence intervals based on the linear models of Kerr *et al.* [7] are constructed, see Hsu *et al.* [6]. Li and Mansouri [8] proposed simultaneous rank tests to search differentially expressed genes when microarray data violate normality assumption and contain a large number of outliers.

Auer and Doerge [1] proposed factorial designs for RNA sequencing experiments. To account for a variety of sources of variations, the resulting observations are fit to the Poisson generalized linear models, see Auer and Doerge [1]. Under this framework, we propose the simultaneous testing procedure to detect differentially expressed gene isoforms such that it controls FWER. Simultaneous test based on large-sample approximation is outlined. The sample size for RNA sequencing study is often small. As it will be shown in Section 4 that the large-sample approximation method does not provide a satisfactory solution in terms of controlling FWER. Monte Carlo simulation of Mansouri and Li [9] shows that percentile-*t* bootstrap method based on pivotal quantities provides a viable method in microarray gene expression analysis. Extension of bootstrap method to RNA sequencing gene expression analysis is hence appealing. In this article, we propose the simultaneous inferential method based on pivotal quantities to detect differentially expressed isoforms using parametric bootstrap. We investigate the performance of the proposed method in controlling the overall error rates through a simulation study.

## 2.    PROBLEM FORMULATION AND PIVOTAL QUANTITIES

### 2.1.    Experimental design and generalized linear model

To account for different sources of variations in observations from treatment, batch, flow cell, and lane, we consider factorial designs for the Next Generation Sequencing. In brief, bar-coded mRNA samples are pooled and assigned to different lanes of a sequencing device in such a way that there are $n$ biological replicates randomly assigned at each combination

of treatment, lane, and flow cell. For details, see Auer and Doerge [1]. Since we can assign an ID to each isoform sequence in RNA sequencing data file, we may use the term "gene" instead of "isoform" in the following.

For gene $l$, $l = 1, ..., g$ we let $Y_{lijkm}$ be the the count of readings from the $i$-th treatment, the $j$-th flow-cell, the $k$-th lane, and the $m$-th biological replicate, $i = 1, ..., a$, $j = 1, ..., b$, $k = 1, ..., c$, and $m = 1, ..., n$. We assume $Y_{lijkm}$'s are independent random observations and the expected value $E(Y_{lijkm}) = \mu_{lijk}$, for $m = 1, ..., n$ follow a per gene Poisson model with log-link (Auer and Doerge, [1]) that

$$(2.1) \qquad \log(\mu_{lijk}) - \log(c_{jk}) = \alpha_l + \tau_{li} + \nu_{lj} + \omega_{lk}$$

where $\alpha_l$ is the overall gene $l$ effect; $\tau_{li}$ is the $i$-th treatment effect on gene $l$ with $\sum_i \tau_{li} = 0$; $\nu_{lj}$ is the $j$-th flow cell effect on gene $l$ with $\sum_j \nu_{lj} = 0$; $\omega_{lk}$ is the $k$-th lane effect on gene $l$ with $\sum_k \omega_{lk} = 0$; $c_{jk}$ is a known constant, namely library size, $j = 1, ..., b$, $k = 1, ..., c$ to normalize the readings from $j$-th flow-cell and $k$-th lane, see Section 6 and Chen *et al.* [2]. We assume that $\alpha_l, \tau_{li}, \nu_{lj}, \omega_{lk}$, for $l = 1, ..., g$, $i = 1, ..., a$, $j = 1, ..., b$, and $k = 1, ..., c$ in (2.1) are fixed effects. Let $N = abcn$ be the total number of readings from each gene.

We let vector

$$\mathbf{Y}_l = \left[ Y_{l1111}, ..., Y_{l111n}, ..., Y_{lijk1}, ..., Y_{lijkn}, ..., Y_{labc1}, ..., Y_{labcn} \right]'$$

be a collection of all readings from gene $l$ and let $\boldsymbol{\mu}_l = E(\mathbf{Y}_l)$, $l = 1, ..., g$. It is useful to write the model in (2.1) in the form of matrix representation that

$$(2.2) \qquad \log(\boldsymbol{\mu}_l / c_{jk}) = X\boldsymbol{\beta}_l$$

where $\boldsymbol{\beta}_l = \left[ \alpha_l, \tau_{l1}, ..., \tau_{l(a-1)}, \nu_{l1}, ..., \nu_{l(b-1)}, \omega_{l1}, ..., \omega_{l(c-1)} \right]'$ and $X$ is the corresponding $N \times (a + b + c - 2)$ design matrix.

Since we use per gene generalized linear model, the model for all genes can be written as

$$(2.3) \qquad \mathbf{1}_g \otimes \log(\boldsymbol{\mu}_l / c_{jk}) = \mathbf{1}_g \otimes X\boldsymbol{\beta}_l \,.$$

## 2.2. Pivotal quantities

For gene $l$, $l = 1, ..., g$ we assume

$$(2.4) \qquad Y_{lijkm} \sim \text{Poisson}(\mu_{lijk}), \qquad \text{for} \quad m = 1, ..., n \,,$$

where

$$(2.5) \qquad \mu_{lijk} = \exp\left[ (\alpha_l + \tau_{li} + \nu_{lj} + \omega_{lk}) + \log(c_{jk}) \right]$$

with $i = 1, ..., a$, $j = 1, ..., b$, $k = 1, ..., c$, and $m = 1, ..., n$.

Let $\widehat{\boldsymbol{\beta}}_{l,N}$ be the maximum likelihood estimation of $\boldsymbol{\beta}_l$, $l = 1, ..., g$. We apply Newton–Raphson method using Fisher Scoring to compute the estimation. We may suppress the notation of the dependence on $N$ and denote the estimation by $\widehat{\boldsymbol{\beta}}_l$.

Now, we define a $q \times (a + b + c - 2)$ comparison matrix $C$ to detect differential gene expression among treatments. In gene expression studies, researchers often interest in $i$) all-pairwise comparisons of gene expression over treatments, or $ii$) comparing gene expression for several treatments versus a control, Hsu *et al.* [6]. We focus on all-pairwise comparisons in this article and analogous results should hold for multiple comparisons to a control. As an example of comparison matrix $C$ for all-pairwise comparisons, see (4.1) in Section 4.

Let $W_l$ be $N \times N$ diagonal weight matrix whose diagonal elements are given by $\mu_{l1111}, ...,$ $\mu_{l111n}, ..., \mu_{lijk1}, ..., \mu_{lijkn}, ..., \mu_{labc1}, ..., \mu_{labcn}$ in order. The vector containing pivotal quantities is given by

$$(2.6) \qquad \boldsymbol{T}(\boldsymbol{\beta}_l) = \Sigma_l^{-1/2}\big[C(\widehat{\boldsymbol{\beta}}_l - \boldsymbol{\beta}_l)\big]$$

where $\Sigma_l$ is a diagonal matrix whose diagonal elements equal to the diagonal elements in $C(X'W_lX)^{-1}C'$, $l = 1, ..., g$.

In relation to the Poisson generalized linear model in (2.3), (2.4), and (2.5), consider gene expression by letting

$$\boldsymbol{T}(\boldsymbol{\beta}) = \big[\boldsymbol{T}(\boldsymbol{\beta}_1)', ..., \boldsymbol{T}(\boldsymbol{\beta}_l)', ..., \boldsymbol{T}(\boldsymbol{\beta}_g)'\big]'.$$

The joint limiting distribution of $\boldsymbol{T}(\boldsymbol{\beta})$ is given by the following Theorem.

**Theorem 2.1.** *Suppose* $\boldsymbol{Y}_1, ..., \boldsymbol{Y}_g$ *are independent vectors, for* $\frac{1}{N}(X'W_lX) \overset{N \to \infty}{\longrightarrow} \mathcal{W}_l$, *which is positive definite, for* $l = 1, ..., g$, *then*

$$(2.7) \qquad \sqrt{N}\boldsymbol{T}(\boldsymbol{\beta}) \overset{D}{\longrightarrow} \mathrm{MVN}(\mathbf{1}_g \otimes \mathbf{0}_q, \Lambda), \qquad as \quad N \to \infty,$$

*where* $\Lambda$ *is a* $gq \times gq$ *block diagonal matrix such that the* $l$*-th* $(q \times q)$ *diagonal block matrix* $\Lambda_l = \lim_{N \to \infty} N\Sigma_l^{-1/2}C(X'W_lX)^{-1}C'\Sigma_l^{-1/2}$, $l = 1, ..., g$.

Proof of Theorem 2.1 immediately follows equation (5.25) and (S.17) of McCulloch *et al.* [10]. Note: since $\Lambda_l$ is unknown in practice, we use a consistent estimator $\widehat{\Lambda}_l = N\widehat{\Sigma}_l^{-1/2}C(X'\widehat{W}_lX)^{-1}C'\widehat{\Sigma}_l^{-1/2}$ where $\widehat{\Sigma}_l$ is a diagonal matrix whose elements equal to the diagonal elements in $C(X'\widehat{W}_lX)^{-1}C'$, and $\widehat{W}_l$ has diagonal elements given by $\exp\{(\widehat{\alpha}_l + \widehat{\tau}_{l1} + \widehat{\nu}_{l1} + \widehat{\omega}_{l1}) + \log(c_{11})\}, ..., \exp\{(\widehat{\alpha}_l + \widehat{\tau}_{li} + \widehat{\nu}_{lj} + \widehat{\omega}_{lk}) + \log(c_{jk})\}, ..., \exp\{(\widehat{\alpha}_l + \widehat{\tau}_{la} + \widehat{\nu}_{lb} + \widehat{\omega}_{lc}) + \log(c_{bc})\}$ in order, $l = 1, ..., g$. In the expression, $\widehat{\alpha}_l$, $\widehat{\tau}_{li}$, $\widehat{\nu}_{lj}$, and $\widehat{\omega}_{lk}$ are maximum likelihood estimation of the parameters, $i = 1, ..., a$; $j = 1, ..., b$; $k = 1, ..., c$. Application of the large-sample approximation method is not trivial since the multivariate normal distribution in Theorem 2.1 has mean and variance with dimension $(gq) \times 1$ and $(gq) \times (gq)$ respectively and the total number of genes $g$, in RNA-sequencing experiments, is typically very large. We propose an Algorithm in Section 4 to reduce the computational burden in RNA-sequencing gene expression analysis.

A challenge besetting RNA-sequencing gene expression analysis may be the overdispersion among counting data, Auer and Doerge [1] and Wang *et al.* [11]. To proceed, we let $\phi_l$ be the dispersion parameter and overdispersion occurs when $\phi_l > 1$, $l = 1, ..., g$.

It is suggested in Auer and Doerge [1] that statistics for detecting differential gene expression should be scaled by the dispersion parameter. Hence, a sequence of pivotal quantities, considering overdispersion, are given by

$$(2.8) \qquad \boldsymbol{T}(\boldsymbol{\beta}_l, \phi_l) = (\phi_l\Sigma_l)^{-1/2}\big[C(\widehat{\boldsymbol{\beta}}_l - \boldsymbol{\beta}_l)\big].$$

The pivotal quantities in (2.6) can be considered as a special case of (2.8) when $\phi_l = 1$. We focus on gene expression analysis for RNA-sequencing data, which presents mild overdispersion such that $\phi_l$ is in a neighborhood of 1, and examine the validity of robustness of the large-sample approximation method through a simulation study in Section 4 in this article.

## 3.  SIMULTANEOUS INFERENCE USING BOOTSTRAP

### 3.1.  Simultaneous inference

In relation to the generalized linear model in (2.3), let the relative gene expression be $\tau_{li} - \tau_{li'}$, $i \neq i' = 1, ..., a$, $l = 1, ..., g$. Detecting all-pairwise differential gene expression can be formulated as testing a sequence of hypotheses that:

$$(3.1) \qquad H_{0l,ii'}: \tau_{li} - \tau_{li'} = 0 \qquad \text{vs.} \qquad H_{1l,ii'}: \tau_{li} - \tau_{li'} \neq 0$$

for $i \neq i' = 1, ..., a$, $l = 1, ..., g$. Hence we conduct $q \times g$ tests simultaneously, where $q$ is the number of rows in comparison matrix $C$ such that $C\boldsymbol{\beta}_l = [\tau_{l1} - \tau_{l2}, ..., \tau_{l(a-1)} - \tau_{la}]'$, see (4.1) for example.

The resulting test statistics are given by

$$(3.2) \qquad \boldsymbol{T}(\widehat{\boldsymbol{\beta}}_l, \widehat{\phi}_l) = (\widehat{\phi}_l \widehat{\Sigma}_l)^{-1/2} C \widehat{\boldsymbol{\beta}}_l$$

for $l = 1, ..., g$ where the plug-in estimation of $\phi_l$ in Auer and Doerge [1] is given by

$$(3.3) \quad \widehat{\phi}_l = \left( \sum_{i,j,k,m} \frac{\left( Y_{lijkm} - \exp\{(\widehat{\alpha}_l + \widehat{\tau}_{li} + \widehat{\nu}_{lj} + \widehat{\omega}_{lk}) + \log(c_{jk})\} \right)^2}{\exp\{(\widehat{\alpha}_l + \widehat{\tau}_{li} + \widehat{\nu}_{lj} + \widehat{\omega}_{lk}) + \log(c_{jk})\}} \right) \Big/ \left( N - (a+b+c-2) \right).$$

For gene $l$, write

$$\boldsymbol{T}(\widehat{\boldsymbol{\beta}}_l, \widehat{\phi}_l) = \left[ T_{12}(\widehat{\boldsymbol{\beta}}_l, \widehat{\phi}_l), ..., T_{ii'}(\widehat{\boldsymbol{\beta}}_l, \widehat{\phi}_l), ..., T_{(a-1)a}(\widehat{\boldsymbol{\beta}}_l, \widehat{\phi}_l) \right]'$$

in association to the hypotheses in (3.1) and the test statistics in (3.2). For all-pairwise comparisons, the total number of comparisons (the total number of elements in $\boldsymbol{T}(\widehat{\boldsymbol{\beta}}_l, \widehat{\phi}_l)$) $q = \binom{a}{2}$.

Simultaneous level-$\alpha$ tests reject hypothesis $H_{0l,ii'}$, $i \neq i' = 1, ..., a$, $l = 1, ..., g$ if:

$$(3.4) \qquad \left| T_{ii'}(\widehat{\boldsymbol{\beta}}_l, \widehat{\phi}_l) \right| > q_\alpha$$

where $q_\alpha$ is the upper $\alpha$-th quantile of the distribution of maximum modulus statistics $\max\limits_{\substack{i \neq i' = 1, ..., a \\ l = 1, ..., g}} \left\{ |T_{ii'}(\widehat{\boldsymbol{\beta}}_l, \widehat{\phi}_l)| \right\}$.

When the magnitude of differential gene expression is of interest, a $(1-\alpha)\,100\%$ simultaneous confidence interval of $\tau_{li} - \tau_{li'}$, $i \neq i' = 1, ..., a$, $l = 1, ..., g$ is given by

$$(3.5) \qquad \boldsymbol{c}'_{ii'} \widehat{\boldsymbol{\beta}}_l \; \pm \; q_\alpha \{ \widehat{\phi}_l \, \boldsymbol{c}'_{ii'} (X' \widehat{W}_l X)^{-1} \boldsymbol{c}_{ii'} \}^{1/2}$$

where $\boldsymbol{c}_{ii'}$ is the row vector of $C$ in association to $\tau_{li} - \tau_{li'}$, $i \neq i' = 1, ..., a$ for all $l = 1, ..., g$.

### 3.2. Bootstrap based on pivotal quantities

It can be shown that the upper $\alpha$-th quantile of the multivariate normal distribution defined in (2.7) is a consistent estimator of $q_\alpha$. RNA sequencing data analysis is often complicated by a large number of unknown parameters but a limited number of observations. Using the large-sample approximation method indicated by Theorem 2.1 can be problematic in the estimation of $q_\alpha$ as it will be shown in Section 4. We propose the parametric bootstrap method based on pivotal quantities to approximate quantiles $q_\alpha$ in detecting differentially expressed genes for RNA sequencing data.

For $r = 1, ..., B$, we define the $q \times 1$ vector of pivotal quantities based on the $r$-th bootstrap sample by

$$(3.6) \qquad \boldsymbol{T}^{(r)}(\widehat{\boldsymbol{\beta}}_l, \widehat{\phi}_l) = \big(\widehat{\phi}_l^{(r)} \widehat{\Sigma}_l^{(r)}\big)^{-1/2} C\big[\widehat{\boldsymbol{\beta}}_l^{(r)} - \widehat{\boldsymbol{\beta}}_l\big], \qquad l = 1, ..., g,$$

where $\widehat{\phi}_l^{(r)}$, $\widehat{\Sigma}_l^{(r)}$, and $\widehat{\boldsymbol{\beta}}_l^{(r)}$ are estimated based on the $r$-th bootstrap data set. Analogously, we write

$$\boldsymbol{T}^{(r)}(\widehat{\boldsymbol{\beta}}_l, \widehat{\phi}_l) = \Big[ T_{12}^{(r)}(\widehat{\boldsymbol{\beta}}_l, \widehat{\phi}_l), ..., T_{ii'}^{(r)}(\widehat{\boldsymbol{\beta}}_l, \widehat{\phi}_l), ..., T_{(a-1)a}^{(r)}(\widehat{\boldsymbol{\beta}}_l, \widehat{\phi}_l) \Big]'.$$

We use the following Algorithm to approximate quantiles $q_\alpha$. For each $r$, $r = 1, ..., B$,

(i) for each $l$, $l = 1, ..., g$ generate random variables $\{Y_{lijkm}\}$ from $\text{Poisson}\big(\exp\big\{(\widehat{\alpha}_l + \widehat{\tau}_{li} + \widehat{\nu}_{lj} + \widehat{\omega}_{lk}) + \log(c_{jk})\big\}\big)$, $i = 1, ..., a$, $j = 1, ..., b$, $k = 1, ..., c$, and $m = 1, ..., n$;

(ii) obtain maximum modulus statistics

$$T_M^{(r)}(\widehat{\boldsymbol{\beta}}_l, \widehat{\phi}_l) = \max_{i \neq i' = 1, ..., a} \big\{ \big| T_{ii'}^{(r)}(\widehat{\boldsymbol{\beta}}_l, \widehat{\phi}_l) \big| \big\}, \qquad l = 1, ..., g,$$

and

$$T_M^{(r)}(\widehat{\boldsymbol{\beta}}, \widehat{\phi}) = \max_{l = 1, ..., g} \big\{ T_M^{(r)}(\widehat{\boldsymbol{\beta}}_l, \widehat{\phi}_l) \big\}.$$

Repeat (i) and (ii) $B$ times, and the upper $\alpha$-th quantile of the sampling distribution of $T_M^{(r)}(\widehat{\boldsymbol{\beta}}, \widehat{\phi})$ is an approximation of $q_\alpha$.

As it will be shown in Section 4, the bootstrap method provides a viable alternative of the large-sample approximation method when the overdispersion parameter is in a neighborhood of $\phi_l = 1$, $l = 1, ..., g$.

## 4.   SIMULATION STUDY

In this section, we investigate the performance of the proposed method in terms of controlling the family-wise error rate (FWER) using Monte Carlo simulation.

We assign the following values to the parameters of the model in (2.1). Let

$$\tau_{li} = 0, \quad \text{for } l = 1, ..., 20, \quad i = 1, 2, 3, 4 \quad (\textit{Complete Null}),$$
$$\tau_{li} = 0, \quad \text{for } l = 1, ..., 15, \quad i = 1, 2, 3, 4 \quad (\textit{Partial Null}).$$

To study the power rates under partial null hypotheses, we let $\tau_{l1} = -0.02$, $\tau_{l2} = 0.01$, $\tau_{l3} = 0.01$, and $\tau_{l4} = 0$, for $l = 16, ..., 20$.

For nuisance parameters, we let $\alpha_l = -3$ and

$$\nu_{lj} = \begin{cases} 0.5, & \text{if } j = 1, \\ -1, & \text{if } j = 2, \\ 0.5, & \text{if } j = 3, \end{cases}$$

for $l = 1, ..., 20$. Let

$$\omega_{lk} = \begin{cases} 0.25, & \text{if } k = 1, \\ -0.5, & \text{if } k = 2, \\ 0.75, & \text{if } k = 3, \\ -1.25, & \text{if } k = 4, \\ 1.5, & \text{if } k = 5, \\ -0.75, & \text{if } k = 6, \end{cases}$$

for $l = 1, ..., 20$.

Assume the library size for each lane and flow cell $c_{jk} = 1,000,000$ for all $j = 1, 2, 3$ and $k = 1, ..., 6$.

We may rewrite the model in (2.1) as $\log(\lambda_{lijk}) = \alpha_l + \tau_{li} + \nu_{lj} + \omega_{lk}$, where the sampling rate $\lambda_{lijk} = E(Y_{lijk}/c_{jk})$ and $c_{jk}$ is a given constant. The observations $Y'_{lijkm}$ are generated from Poisson($\mu_{lijk}$) where $\mu_{lijk} = c_{jk}\lambda_{lijk}$, for $m = 1, 2$. To exam the performance of the proposed method under mild overdispersion, we add Gaussian noise $\epsilon_{lijkm} \sim N(0, (\phi_l - 1)\mu_{lijk})$ ($\phi_l > 1$) to the observations that $Y_{lijkm} = Y'_{lijkm} + [\epsilon_{lijkm}]$, $i = 1, ..., 4$, $j = 1, 2, 3$, $k = 1, ..., 6$, $m = 1, 2$ for gene $l$, $l = 1, ..., 20$ as it is treated in Auer and Doerge [1]. Note that $E(Y'_{lijkm} + \epsilon_{lijkm}) = \mu_{lijk}$ and $\text{Var}(Y'_{lijkm} + \epsilon_{lijkm}) = \phi_l\mu_{lijk}$. We choose $\phi_l = 1.1$, 1.05, 1.01, and 1.001 respectively and let $Y_{lijkm} = Y'_{lijkm}$ for $\phi_l = 1$. In addition, we let the observations equal to zero if it generates "negative" counts, though the chance of generating "negative" counts is rare when the value of $(\phi_l - 1)$ is small.

Hence, the vector of parameters $\boldsymbol{\beta}_l = [\alpha_l, \tau_{l1}, \tau_{l2}, \tau_{l3}, \nu_{l1}, \nu_{l2}, \omega_{l1}, ..., \omega_{l5}]'$, $l = 1, ..., g$. Let $X$ be the corresponding design matrix for all genes. Consider all-pairwise comparisons among treatments. Let $C$ be the $6 \times 11$ comparison matrix given by

$$(4.1) \qquad C = \begin{bmatrix} 0 & 1 & -1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & -1 & 0 & \cdots & 0 \\ 0 & 2 & 1 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & -1 & 0 & \cdots & 0 \\ 0 & 1 & 2 & 1 & 0 & \cdots & 0 \\ 0 & 1 & 1 & 2 & 0 & \cdots & 0 \end{bmatrix}.$$

We run simultaneous tests in (3.4) 1,000 times and compute the empirical overall error rates. Widely used measures of the overall error rates in gene expression analysis are the family-wise error rate (FWER) and the false discovery rate (FDR). Let FWER$_0$ be the probability that at least one true null hypotheses rejected under complete null hypotheses. Let FWER$_1$ be the probability that at least one true null hypotheses rejected under partial null hypotheses. The false discovery rate (FDR) is computed as the average proportion of

wrongly rejected null hypotheses among all rejected hypotheses. FDR is defined as 0 if no rejection were made. To investigate the power of the simultaneous tests, we compute the proportional power rate by obtaining the average proportion of genes found differentially expressed among all misexpressed genes, Dudoit *et al.* [3].

To evaluate the performance of the large-sample approximation method, we use the following Algorithm to generate quantiles based on the multivariate normal distribution defined in Theorem 2.1. In specific, for each $r$, $r = 1, ..., B$,

(**i\***)   generate random variables $\boldsymbol{T}_l^{(r)}$ from $\mathrm{MVN}(\boldsymbol{0}, \widehat{\Lambda}_l)$, for all $l = 1, ..., g$;

(**ii\***)   obtain maximum modulus statistics $T_{M_l}^{(r)} = \max\{|\boldsymbol{T}_l^{(r)}|\}$, $l = 1, ..., g$ and $T_M^{(r)} = \max\{T_{M_l}^{(r)}\}$.

Repeat (i\*) and (ii\*) $B$ times, and the upper $\alpha$-th quantile of the empirical distribution of $T_M^{(r)}$ is an approximation of $q_\alpha$ based on Theorem 2.1.

The performance of the large-sample approximation method and the bootstrap method in the simulation study are summarized in Table 1.

**Table 1**:   Error rates of detecting differentially expressed genes/isoforms — nominal type-1 error rate $\alpha = 0.05$.

| Method | $\phi_l$ | $\mathrm{FWER}_0$ | $\mathrm{FWER}_1$ | FDR | Prop. Power |
|---|---|---|---|---|---|
| No Adjustment | 1.000 | 0.993 | 0.970 | 0.146 | — |
| MVN | $1.000^\dagger$ | 0.072 | 0.059 | 0.003 | 0.889 |
|  | $1.050\,(1.1)^\ddagger$ | 0.084 | 0.061 | 0.003 | 0.861 |
|  | 1.010 | 0.073 | 0.045 | 0.002 | 0.887 |
|  | 1.001 | 0.065 | 0.065 | 0.003 | 0.886 |
| Bootstrap Method | 1.000 | 0.052 | 0.037 | 0.002 | 0.878 |
|  | $1.050\,(1.1)$ | 0.051 | 0.035 | 0.002 | 0.849 |
|  | 1.010 | 0.049 | 0.034 | 0.002 | 0.874 |
|  | 1.001 | 0.050 | 0.045 | 0.002 | 0.872 |

Notes:

i)   Simulation size = 1,000. Bootstrap size $B = 200$.

ii)   $\mathrm{FWER}_0$ denotes the family-wise error rate under complete null hypotheses.

iii)   $\mathrm{FWER}_1$ denotes the family-wise error rate under partial null hypotheses.

iv)   MVN denotes the method of large-sample approximation in Section 3.1.

v)   "Bootstrap Method" means the parametric bootstrap method in Section 3.2.

vi)   † The same value of $\phi_l$ is assigned to all genes.

vii)   ‡ The first 15 genes have $\phi_l = 1.05$ and the last 5 genes have $\phi_l = 1.1$.

viii)   The total computation user time was about 16 hours on a desktop with processor with the following specifications: Intel(R) Core(TM) i5-7600 CPU @ 3.50GHz, 3504 Mhz and Installed physical memory (RAM): 16.0 GB.

It shows that the bootstrap method based on pivotal quantities controls FWER under both complete and partial null hypotheses. This implies the proposed method controls FWER strongly, see Dudoit *et al.* [3]. Without adjustment of multiplicity, it is well known that the overall error rates often exceed the nominal level, particularly in large-scale tests.

Simultaneous tests based on large-sample approximation fail to control FWER in the strong sense in RNA sequencing data analysis. While the overall error rates are controlled at nominal level $\alpha = 0.05$, in average more than 85% of "real" misexpressed genes are detected as differentially expressed genes using the bootstrap method in Section 3.2. Note that it is not useful to address the power rates when the method does not control FWER.

To investigate the performance of the bootstrap method in estimation of quantiles, we generate $1,000$ samples as described above and obtain the $(1-\alpha)$-th quantile of the sampling distribution of pivotal quantities in (2.8). Since the quantiles are generated from a given underlying distribution of maximum modulus distribution empirically, it can be used as a benchmark to evaluate the performance of the proposed method. The results are summarized in Table 2.

**Table 2**:  Quantiles $q_\alpha$ of detecting differentially expressed genes/isoforms — nominal type-1 error rate $\alpha = 0.05$.

| $\phi_l$ | Simulation | MVN | Bootstrap |
|---|---|---|---|
| 1.000 | 3.604 | 3.519 (0.090) | 3.606 (0.094) |
| 1.050 (1.1) | 3.617 | 3.522 (0.086) | 3.611 (0.094) |
| 1.010 | 3.605 | 3.524 (0.090) | 3.609 (0.095) |
| 1.001 | 3.604 | 3.516 (0.084) | 3.604 (0.097) |

Notes:

i)   Simulation size $= 1,000$. Bootstrap size $B = 200$.

ii)  MVN denotes the method of large-sample approximation in Section 3.1 and the Algorithm in Section 4. The quantile is generated from $B = 200$ samples. We repeat the process for $1,000$ times. The mean value of these repeats is included outside of the parentheses and standard deviation is tabulated in the parentheses.

iii) "Bootstrap" means the parametric bootstrap method in Section 3.2. The quantile is generated from $B = 200$ bootstrap samples. We repeat the process for 1,000 times. The mean value of these repeats is included outside of the parentheses and standard deviation is tabulated in the parentheses.

iv)  "Simulation" means: we generate observations from the model in (2.1) with the parameter value assigned in Section 4 and given underlying distributions for $1,000$ times; the upper $\alpha$-th quantile of maximum modulus statistics based on pivotal quantiles in (2.8) is tabulated in the table.

v)   The total computation user time was about 8 hours on a desktop with processor with the following specifications: Intel(R) Core(TM) i5-7600 CPU @ 3.50GHz, 3504 Mhz and Installed physical memory (RAM): 16.0 GB.

It shows from Table 2 that the bootstrap quantiles in Section 3.2 are closer to the simulated quantiles as compared to that generated from MVN. A closer examination sees the quantiles based on normal theory are generally below the simulated quantiles. Therefore, the large-sample approximation method provides a liberal estimation of FWER, as evidenced in Table 1.

## 5.    CONCLUSION AND FUTURE WORK

In this article, we have proposed the parametric bootstrap method based on pivotal quantities in detecting differentially expressed genes for RNA-sequencing data. We have formulated the problem using the Poisson generalized linear models. We have derived the joint limiting distribution of the vector containing pivotal quantities. We have conducted an empirical study to show that the proposed method controls FWER and FDR strongly in detecting differentially expressed genes. The bootstrap method requires a large computation time, parallel computation is recommended particularly for large-scale inference. When data "apparently" violate Poisson distributional assumption, we will investigate the methods involving a large value of overdispersion. To capture the within genes' variation and between genes' variation, we will study the resampling methods, such as moving block bootstrap method in the future work.

## 6.    SOFTWARE

We use the function `glm()` in `R` to obtain maximum likelihood estimation of the parameters in model (2.1). Note that computation of the estimation using `glm()` in `R` may encounter non-convergence. Alternatively, iterative weighted least squares method of Wedderburn [12] may be used in the estimation. Our experience in the simulation study (results not shown) shows that using 20-step iterative weighted least squares method provides satisfactory approximation of the overall Type-I error rates. We use the function `rmvnorm()` of Genz *et al.* [5] in `R` to generate multivariate normal random variables. We use the function `calcNormFactors()` of Chen *et al.* [2] to obtain the library size. Software in the form of `R` code is available on request from the author (bli@citadel.edu).

## REFERENCES

[1]    AUER, P.L. and DOERGE, R.W. (2010). Statistical Design and Analysis of RNA Sequencing Data, *Genetics*, **185**, 405–416.

[2]    CHEN, Y.; LUN, A.T.L. and SMYTH, G.K. (2014). *Differential Expression Analysis of Complex RNA-seq Experiments Using edgeR.* In: "Statistical Analysis of Next Generation Sequencing Data" (S. Datta and D. Nettleton, Eds.), New York, Springer.

[3]    DUDOIT, S.; SHAFFER, J.P. and BOLDRICK, J.C. (2003). Multiple Hypothesis Testing in Microarray Experiments, *Statistical Science*, **18**(1), 71–103.

[4]    DUDOIT, S.; YANG, Y.H.; CALLOW, M.J. and SPEED, T.P. (2002). Statistical methods for identifying differentially expressed genes in replicated cDNA microarray experiments, *Statistica Sinica*, **12**, 111–139.

[5]    GENZ, A.; BRETZ, F.; MIWA, T.; MI, X.; LEISCH, F.; SCHEIPL, F. and HOTHORN, T. (2017). mvtnorm: Multivariate Normal and t Distributions, R package version 1.0-6.
URL: http://CRAN.R-project.org/package=mvtnorm

[6]    HSU, J.C.; CHANG, J.Y. and WANG, T. (2006). Simultaneous confidence intervals for differential gene expressions, *Journal of Statistical Planning and Inference*, **136**(7), 2182–2196.

[7]    KERR, M.K.; MARTIN, M. and CHURCHILL, G.A. (2000). Analysis of Variance for Gene Expression Microarray Data, *Journal of Computational Biology*, **7**(6), 818–837.

[8]    LI, B. and MANSOURI, H.G. (2016). Simultaneous Rank Tests for Detecting Differentially Expressed Genes, *Journal of Statistical Computation and Simulation*, **86**(5), 959–972.

[9]    MANSOURI, H.G. and LI, B. (2019). On simultaneous confidence intervals based on rank-estimates with application to analysis of gene expression data, *Communications in Statistics – Theory and Methods*, **48**(17), 4339–4349.
DOI: 10.1080/03610926.2018.1494287

[10]   MCCULLOCH, C.E.; SEARLE, S.E. and NEUHAUS, J.M. (2008). *Generalized, Linear, and Mixed Models*, New Jersey, Wiley, 2008.

[11]   WANG, J.; HUANG, M.; TORRE, E.; DUECK, H.; SHAFFER, S.; MURRAY, J.; RAJ, A.; LI, M. and ZHANG, N.R. (2018). Gene expression distribution deconvolution in single-cell RNA sequencing, *Proceedings of the National Academy of Sciences*, **115**(28), E6437–E6446.
DOI: 10.1073/pnas.1721085115

[12]   WEDDERBURN, R.W.M. (1974). Quasi-likelihood Functions, Generalized Linear Models, and the Gauss–Newton Method, *Biometrika*, **61**(3), 439–447.

[13]   WESTFALL, P.H. and YOUNG, S.S. (1993). *Resampling-Based Multiple Testing: Examples and Methods for p-Value Adjustment*, New York, Wiley, 1993.

# VARIANCE ESTIMATION USING
# RANDOMIZED RESPONSE TECHNIQUE

Authors: SAT GUPTA
– Department of Mathematics and Statistics, University of North Carolina at
Greensboro, North Carolina, USA
sngupta@uncg.edu

BADR ALORAINI
– Department of Mathematics and Statistics, University of North Carolina at
Greensboro, North Carolina, USA
boalorai@uncg.edu

MUHAMMAD NOUMAN QURESHI
– National College of Business Administration and Economics,
Lahore, Pakistan
nqureshi633@gmail.com

SADIA KHALIL
– Department of Statistics, Lahore College for Women University,
Lahore, Pakistan
sadia_khalil@hotmail.com

Abstract:

• Variance estimation is a well-studied topic in survey sampling but not much work has been done in
this area in the context of Randomized Response Technique (RRT) models. We propose here some
variance estimators for sensitive variables using auxiliary information. We examine the performance
of the proposed estimators through a simulation study and through a numerical example.

## 1.    INTRODUCTION

When conducting surveys, it is sometimes difficult to make a direct observation on the variable of interest. This is more so in the case where the research involves a topic that is a taboo in nature. In surveys on such topics, some of the respondents might give false responses. To offer a solution to this, a Randomized Response Technique (RRT) was developed by Warner [7]. The technique allows respondents to provide a response while maintaining their privacy.

The problem of mean and variance estimation is a topic that has been explored very well by researchers, although less so the problem of variance estimation. This is particularly the case in the context of RRT models. This is the main focus of this study where we examine variance estimation of a sensitive study variable using a highly correlated but non-sensitive auxiliary variable. According to Collins *et al.* [1], the auxiliary variables when combined with the main study variable help to achieve more efficient estimators.

In this paper, three variance estimators have been proposed under RRT using one auxiliary variable and two scrambling variables. In Section 2, some of the variance estimators in literature are reviewed. In Section 3, we propose a new class of variance estimators under RRT and derive their Bias as well as their MSE. We provide a comparison of the proposed estimators in Section 4. A numerical study is conducted in Section 5 based on real data. Some concluding remarks are given in Section 6.

## 2.    ESTIMATORS IN LITERATURE

Let a simple random sample of size $n$ be extracted without replacement from a finite population $U = \{U_1, U_2, ..., U_N\}$. Let $Y$ be a sensitive variable of interest and $X$ be a positively correlated auxiliary variable. Let $(x_i, y_i)$ be the observed $(X, Y)$ values for the $i$-th population unit $U_i$. Let $(\bar{x}, \bar{y})$ and $(\bar{X}, \bar{Y})$ be the sample and population means, and $(s_x^2, s_y^2)$ and $(\sigma_x^2, \sigma_y^2)$ be the sample and population variances respectively. Let

$$s_y^2 = \frac{1}{n-1} \sum_{i=1}^{n} (y_i - \bar{y})^2\,, \qquad s_x^2 = \frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{x})^2\,,$$

$$\sigma_y^2 = \frac{1}{N-1} \sum_{i=1}^{N} (Y_i - \bar{Y})^2\,, \qquad \sigma_x^2 = \frac{1}{N-1} \sum_{i=1}^{N} (X_i - \bar{X})^2\,,$$

$$\bar{X} = \frac{1}{N} \sum_{i=1}^{N} X_i\,, \quad \bar{Y} = \frac{1}{N} \sum_{i=1}^{N} Y_i\,, \quad \bar{x} = \frac{1}{N} \sum_{i=1}^{N} x_i\,, \quad \bar{y} = \frac{1}{N} \sum_{i=1}^{N} y_i\,.$$

An unbiased estimator for the finite population variance is the sample variance given by:

$$t_0 = s_y^2\,.$$

Up to the first degree of approximation, its variance is given by

$$V(t_0) = \theta \sigma_y^4 (\lambda_{40} - 1),$$

where

$$\lambda_{rs} = \frac{\mu_{rs}}{\mu_{20}^{\frac{r}{2}} \mu_{02}^{\frac{s}{2}}}, \qquad \mu_{rs} = \frac{1}{N-1} \sum_{i=1}^{N} (Y_i - \bar{Y})^r (X_i - \bar{X})^s, \qquad \text{and} \qquad \theta = \frac{1}{n}.$$

Also '$r$' and '$s$' are non-negative integers, $\mu_{20}$ and $\mu_{02}$ are the second order moments and $\lambda_{rs}$ is the moment ratio.

Isaki [4] proposed the following ratio estimator of population variance using auxiliary information:

$$t_1 = s_y^2 \left( \frac{\sigma_x^2}{s_x^2} \right).$$

The expressions for Bias and Mean Square Error (MSE) of the estimator, up to the first order of approximation, are given by

$$\mathrm{B}(t_1) = \theta \sigma_y^2 (\lambda_{04} - 1) [1 - f_{04}]$$

and

$$\mathrm{MSE}(t_1) = \theta \sigma_y^4 (\lambda_{40} - 1) + (\lambda_{04} - 1) [1 - 2 f_{04}],$$

where

$$f_{04} = \frac{(\lambda_{22} - 1)}{(\lambda_{04} - 1)}.$$

The regression estimator of population variance was also proposed by Isaki [4] as

$$t_2 = s_y^2 + \alpha(\sigma_x^2 - s_x^2), \qquad \text{where} \quad \alpha = \left( \frac{\sigma_y^2}{\sigma_x^2} \right) f_{04}.$$

The MSE of $t_2$ is given by

$$\mathrm{MSE}(t_2) = \theta \sigma_y^4 (\lambda_{40} - 1)(1 - p^2), \qquad \text{where} \quad p = (\lambda_{22} - 1)/\sqrt{(\lambda_{40} - 1)(\lambda_{04} - 1)}.$$

## 3.  PROPOSED ESTIMATORS

Since $Y$ is sensitive in nature, and hence subject to social desirability bias, we observe only a scrambled version of $Y$ as given by Diana and Perri [2]. This is given by $Z = TY + S$, where $T$ and $S$ are scrambling variables. We also assume that $Y$, $T$ and $S$ are mutually uncorrelated. We also assume $E(S) = 0$ and $E(T) = 1$.

To obtain the Bias and MSE expressions for the proposed estimators, we define the following error terms:

$$s_z^2 = \sigma_z^2 (1 + \delta_z) \qquad \text{and} \qquad \bar{z} = \bar{Z}(1 + e_z),$$

where

$$\delta_z = \frac{s_z^2 - \sigma_z^2}{\sigma_z^2} \qquad \text{and} \qquad e_z = \frac{\bar{z} - \bar{Z}}{\bar{Z}}$$

such that

$$E(\delta_z) = E(e_z) = 0, \quad E(\delta_z^2) = \theta(\lambda_{40} - 1), \quad \text{and} \quad E(e_z^2) = \theta C_z^2; \quad \text{and} \quad E(\delta_z e_z) = \theta \lambda_{30} C_z$$

where

$$C_z^2 = C_y^2 \sigma_T^2 + \left( \frac{\sigma_S^2}{\bar{Y}^2} \right).$$

We now propose several population variance estimators under RRT.

## 3.1.  A basic variance estimator under RRT

Based on the RRT model $Z = TY + S$, we have $\sigma_z^2$ as

$$\begin{aligned}
\sigma_z^2 &= \sigma_{TY+S}^2 = \sigma_{TY}^2 + \sigma_S^2 \\
&= \left( \sigma_T^2 * \sigma_Y^2 + \sigma_T^2 * \left( E[Y] \right)^2 + \left( E[T] \right)^2 * \sigma_Y^2 \right) + \sigma_S^2 \\
&= \left( \sigma_T^2 * \sigma_Y^2 + \sigma_T^2 * (\mu_Y)^2 + \sigma_Y^2 \right) + \sigma_S^2 \\
&= \sigma_T^2 * \sigma_Y^2 + \sigma_T^2 * \mu_Y^2 + \sigma_Y^2 + \sigma_S^2.
\end{aligned}$$

Rearranging, we get

$$\sigma_y^2 = \frac{\sigma_z^2 - \sigma_S^2 - (\sigma_T^2 * \bar{Z}^2)}{\sigma_T^2 + 1}.$$

Estimating $\sigma_z^2$ by its unbiased estimator $s_z^2$, we have our first proposed estimator given by

$$(3.1) \qquad t_0(R) = \frac{s_z^2 - \sigma_S^2 - \sigma_T^2 * \bar{z}^2}{\sigma_T^2 + 1}.$$

Rewriting (3.1), we have

$$t_0(R) = \frac{\sigma_z^2(1 + \delta_z) - \sigma_S^2 - \sigma_T^2 \left[ \bar{Z}(1 + e_z) \right]^2}{\sigma_T^2 + 1}.$$

Subtracting $\sigma_y^2$ on both sides, we obtain

$$(3.2) \qquad \left( t_0(R) - \sigma_y^2 \right) = \frac{\sigma_z^2 \delta_z - 2\sigma_T^2 \bar{Z}^2 e_z - \sigma_T^2 \bar{Z}^2 e_z^2}{\sigma_T^2 + 1}.$$

By taking the expectation on both sides of (3.2), the Bias of $t_0(R)$ is obtained as

$$\text{Bias}\left( t_0(R) \right) = -\theta \left( \frac{\sigma_T^2 \bar{Z}^2}{\sigma_T^2 + 1} \right) C_z^2.$$

By squaring both sides of (3.2) and using the first order approximation, the MSE is obtained as

$$(3.3) \qquad \text{MSE}\left( t_0(R) \right) = \theta \left( \frac{1}{(\sigma_T^2 + 1)^2} \right) \left( \sigma_z^4 (\lambda_{40} - 1) + 4\sigma_T^4 \bar{Z}^4 C_z^2 - 4\sigma_z^2 \sigma_T^2 \bar{Z}^2 \lambda_{30} C_z \right).$$

## 3.2. The ratio estimator under RRT

Isaki [4] proposed the classical ratio estimator $t_1 = s_y^2 \left( \frac{\sigma_x^2}{s_x^2} \right)$. The RRT version of $t_1$ is

$$(3.4) \qquad t_1(R) = \frac{s_z^2 - \sigma_S^2 - \sigma_T^2 * \bar{z}^2}{\sigma_T^2 + 1} * \left( \frac{\sigma_x^2}{s_x^2} \right).$$

To obtain the Bias and MSE, we define the following error terms:

$$s_x^2 = \sigma_x^2 (1 + \delta_x), \qquad \text{where} \quad \delta_x = \frac{s_x^2 - \sigma_x^2}{\sigma_x^2},$$

such that

$$E(\delta_x) = 0, \qquad E(\delta_x^2) = \theta(\lambda_{04} - 1) \qquad \text{and} \qquad E(\delta_x e_z) = \theta \lambda_{12} C_z.$$

Rewriting (3.4), we have

$$t_1(R) = \frac{\sigma_z^2 - \sigma_S^2 - \sigma_T^2 \bar{Z}^2}{\sigma_T^2 + 1} + \frac{2 \sigma_T^2 \bar{Z}^2 e_z \delta_x - \sigma_z^2 \delta_z \delta_x - \sigma_T^2 \bar{Z}^2 e_z^2}{\sigma_T^2 + 1}.$$

Subtracting $\sigma_y^2$ and taking the expectation on both sides, the Bias of $t_1(R)$ is obtained as

$$(3.5) \qquad \text{Bias}(t_1(R)) = \theta \left( \frac{2 \sigma_T^2 \bar{Z}^2 \lambda_{12} C_z - \sigma_z^2 (\lambda_{22} - 1) - \sigma_T^2 \bar{Z}^2 C_z^2}{\sigma_T^2 + 1} \right).$$

For MSE, we have

$$t_1(R) = \frac{\sigma_z^2 + \sigma_z^2 \delta_z - \sigma_S^2 - \sigma_T^2 \bar{Z}^2 - 2 \sigma_T^2 \bar{Z}^2 e_z - \sigma_T^2 \bar{Z}^2 e_z^2}{\sigma_T^2 + 1}$$
$$\frac{- \sigma_z^2 \delta_x - \sigma_z^2 \delta_z \delta_x + \sigma_S^2 \delta_x + \sigma_T^2 \bar{Z}^2 \delta_x + 2 \sigma_T^2 \bar{Z}^2 e_z \delta_x + \sigma_T^2 \bar{Z}^2 e_z^2 \delta_x}{\sigma_T^2 + 1}.$$

Simplifying and ignoring second and higher order terms,

$$t_1(R) = \frac{\sigma_z^2 - \sigma_S^2 W - \sigma_T^2 \bar{Z}^2}{\sigma_T^2 + 1} + \frac{\sigma_z^2 \delta_z - 2 \sigma_T^2 \bar{Z}^2 e_z - \sigma_z^2 \delta_x + \sigma_S^2 \delta_x + \sigma_T^2 \bar{Z}^2 \delta_x}{\sigma_T^2 + 1}.$$

Squaring and taking the expectation on both sides, we have

$$\text{MSE}(t_1(R)) = E \left( \frac{\sigma_z^2 \delta_z}{\sigma_T^2 + 1} - \frac{2 \sigma_T^2 \bar{Z}^2 e_z}{\sigma_T^2 + 1} - \sigma_y^2 \delta_x \right)^2.$$

After some simplifications, the MSE of $t_1(R)$ is obtained as

$$\text{MSE}(t_1(R)) = \theta \frac{1}{(\sigma_T^2 + 1)^2} \left[ \sigma_z^4 (\lambda_{40} - 1) - 2 \sigma_z^2 \sigma_y^2 (\lambda_{22} - 1)(\sigma_T^2 + 1) + \sigma_y^4 (\lambda_{04} - 1)(\sigma_T^2 + 1)^2 \right.$$
$$(3.6)$$
$$\left. + 4 C_z \left( \sigma_T^4 \bar{Z}^4 C_z - \sigma_z^2 \sigma_T^2 \bar{Z}^2 \lambda_{30} + \sigma_T^2 \sigma_y^2 \bar{Z}^2 \lambda_{12}(\sigma_T^2 + 1) \right) \right].$$

### 3.3.  A generalized variance estimator under RRT

We now propose the following class of generalized population variance estimators:

$$(3.7) \quad t_p(R) = \left[ \left( \left( \frac{s_z^2 - \sigma_S^2 - \sigma_T^2 * \bar{z}^2}{\sigma_T^2 + 1} \right) + (\sigma_x^2 - s_x^2) \right) * \left( \frac{(\alpha \sigma_x^2 + \beta)}{\omega(\alpha s_x^2 + \beta) + (1 - \omega)(\alpha \sigma_x^2 + \beta)} \right)^g \right],$$

where $g$, $\alpha$, $\beta$ and $\omega$ are suitably chosen constants. We would choose $g = 1$ for positive correlation between $Y$ and $X$, and $-1$ for negative correlation. $\alpha$ and $\beta$ are known parameters associated with the auxiliary variable and $\omega$ is obtained from optimality consideration.

Using Taylor series approximation, we obtain the bias of the generalized estimator $t_p(R)$ as

$$(3.8) \quad \text{Bias}\big(t_p(R)\big) = \frac{-\theta \sigma_T^2 \bar{Z}^2}{\sigma_T^2 + 1} C_z^2 - (g\omega\psi_i)\,\theta \left( \frac{\sigma_z^2(\lambda_{22} - 1) - 2\,\sigma_T^2 \bar{Z}^2 \lambda_{12} C_z}{\sigma_T^2 + 1} - \sigma_x^2(\lambda_{04} - 1) \right),$$

where $\psi_i = \frac{\alpha \sigma_x^2}{\alpha \sigma_x^2 + \beta}$.

The mean square error is given by

$$(3.9) \quad \begin{aligned} \text{MSE}\big(t_p(R)\big) = \theta & \left[ \left( \frac{\sigma_z^4(\lambda_{40} - 1) + 4\,\sigma_T^4 \bar{Z}^4 C_z^2 - 4\,\sigma_z^2 \sigma_T^2 \bar{Z}^2 \lambda_{30} C_z}{(\sigma_T^2 + 1)^2} \right) \right. \\ & + \left( (\sigma_x^2 + Q\,\sigma_y^2)^2\,(\lambda_{04} - 1) \right) \\ & \left. - 2\left( \frac{\sigma_z^2(\lambda_{22} - 1) - 2\,\sigma_T^2 \bar{Z}^2 \lambda_{12} C_z}{\sigma_T^2 + 1} \right) (\sigma_x^2 + Q\,\sigma_y^2) \right], \end{aligned}$$

where $Q = g\omega\psi_i$.

Differentiate (3.9) w.r.t $Q$:

$$2\,\sigma_y^2(\sigma_x^2 + Q\,\sigma_y^2)\,(\lambda_{04} - 1) = 2\sigma_y^2 \left( \frac{\sigma_z^2(\lambda_{22} - 1) - 2\,\sigma_T^2 \bar{Z}^2 \lambda_{12} C_z}{\sigma_T^2 + 1} \right),$$

$$Q_{\text{opt}} = \frac{1}{\sigma_y^2} \left[ \left( \frac{\sigma_z^2(\lambda_{22} - 1) - 2\,\sigma_T^2 \bar{Z}^2 \lambda_{12} C_z}{\sigma_T^2 + 1} \right) \left( \frac{1}{(\lambda_{04} - 1)} \right) - \sigma_x^2 \right].$$

The MSE at this optimum value is given by

$$(3.10) \quad \begin{aligned} \text{MSE}\big(t_p(R)\big)_{\text{opt}} = \frac{\theta}{(\sigma_T^2 + 1)^2} & \left[ \left( \sigma_z^4(\lambda_{40} - 1) + 4\sigma_T^4 \bar{Z}^4 C_z^2 - 4\sigma_z^2 \sigma_T^2 \bar{Z}^2 \lambda_{30} C_z \right) \right. \\ & \left. - \frac{1}{(\lambda_{04} - 1)} \left( \sigma_z^2(\lambda_{22} - 1) - 2\sigma_T^2 \bar{Z}^2 \lambda_{12} C_z \right)^2 \right]. \end{aligned}$$

## 4.   SIMULATION STUDY

In this section, we use a simulation study to evaluate how efficient the generalized estimator $t_p(R)$ is as compared to both the basic estimator $t_0(R)$ and the ratio estimator $t_1(R)$. We first consider samples of size $N = 1000$ each from three bivariate normal populations determined by the following means and covariance matrices:

$$\text{Population I:} \quad \mu = \begin{bmatrix} 6 \\ 4 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} 4 & 1.6 \\ 1.6 & 1 \end{bmatrix}, \quad \rho_{yx} = 0.80\,;$$

(4.1)  $$\text{Population II:} \quad \mu = \begin{bmatrix} 6 \\ 4 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} 4 & 2.25 \\ 2.25 & 2 \end{bmatrix}, \quad \rho_{yx} = 0.80\,;$$

$$\text{Population III:} \quad \mu = \begin{bmatrix} 6 \\ 4 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} 4 & 1.2 \\ 1.2 & 1 \end{bmatrix}, \quad \rho_{yx} = 0.60\,.$$

These 1000 observations are treated as our finite populations. For the 1000 values generated from these distributions, the means, variances, covariances, and correlations are given by

$$\text{Population I:} \quad \mu_x = 6.029, \quad \mu_y = 4.007, \quad \sigma_x^2 = 3.8862, \quad \sigma_y^2 = 0.9450,$$
$$\sigma_{xy} = 1.5284, \quad \rho_{yx} = 0.7975389\,;$$

$$\text{Population II:} \quad \mu_x = 6.021, \quad \mu_y = 3.9836, \quad \sigma_x^2 = 3.9467, \quad \sigma_y^2 = 1.9998,$$
$$\sigma_{xy} = 2.2382, \quad \rho_{yx} = 0.7967094\,;$$

$$\text{Population III:} \quad \mu_x = 5.962, \quad \mu_y = 3.971, \quad \sigma_x^2 = 4.1149, \quad \sigma_y^2 = 0.9560,$$
$$\sigma_{xy} = 1.2442, \quad \rho_{yx} = 0.5927674\,.$$

For each population, we consider samples of sizes 200 and 500. The scrambling variables $S$ and $T$ are assumed to have normal distributions with $E(S) = 0$ and $E(T) = 1$. We have used different values for $\text{Var}(S)$ and $\text{Var}(T)$.

Before presenting the simulation results, we would like to note that in most studies, researchers have compared estimators only with respect to the *Percent Relative Efficiency* which is defined as

$$\text{PRE} = \frac{\text{MSE}\big(t_0(R)\big)}{\text{MSE}\big(t_i(R)\big)} \times 100, \qquad \text{where} \quad i = 0, 1 \text{ and } p\,.$$

However, for estimators based on RRT methodology, one needs to also consider the *Privacy Protection* offered by the RRT model. With that in mind, Gupta *et al.* [3] introduced a unified measure of estimator quality ($\delta$) given by

$$\delta = \frac{\text{Theoretical MSE}}{\Delta_{DP}}, \qquad \text{where} \quad \Delta_{DP} = E(Z - Y)^2 = \sigma_T^2(\mu_y^2 + \sigma_y^2) + \sigma_s^2$$

is the privacy level for the model $Z = TY + S$, as per Yan *et al.* [8]. A smaller value of ($\delta$) is to be preferred. Khalil *et al.* [6] used this unified measure to compare the performance of various mean estimators under RRT.

**Table 1**: Theoretical (**bold**) and empirical MSEs and PREs of the estimators for Population I with $\sigma_T^2 = 0.5$, $\sigma_y^2 = 1$ and $\rho_{yx} = 0.80$.

| Var(S) | $n$ | Estimator | Mean($\hat{\sigma}_y^2$) | MSE | PRE | $\delta$ |
|--------|-----|-----------|--------------------------|-----|-----|----------|
| 0.2 | 200 | $t_0(R)$ | 1.018416 | **0.4593715** 0.4629093 | **100** 100 | 0.052801 |
| | | $t_1(R)$ | 0.9873038 | **0.4166811** 0.4137594 | **110.2453** 111.8788 | 0.04789438 |
| | | $t_p(R)$ | 0.9708478 | **0.3685766** 0.3689481 | **124.6339** 125.4673 | 0.04236513 |
| | 500 | $t_0(R)$ | 1.021572 | **0.1995375** 0.2100302 | **100** 100 | 0.02293534 |
| | | $t_1(R)$ | 0.9846944 | **0.1430092** 0.146612 | **139.5277** 143.2558 | 0.01643784 |
| | | $t_p(R)$ | 0.9999683 | **0.0946721** 0.0957580 | **210.7669** 219.3343 | 0.01088185 |
| 0.5 | 200 | $t_0(R)$ | 1.034554 | **0.5512713** 0.5593184 | **100** 100 | 0.06125237 |
| | | $t_1(R)$ | 0.9986482 | **0.4943552** 0.5045654 | **111.5131** 110.8515 | 0.05492836 |
| | | $t_p(R)$ | 0.9854447 | **0.4320352** 0.4187965 | **127.5987** 133.5537 | 0.04800391 |
| | 500 | $t_0(R)$ | 1.023019 | **0.2022691** 0.1991505 | **100** 100 | 0.02247434 |
| | | $t_1(R)$ | 0.9816713 | **0.1866725** 0.182478 | **108.3550** 109.1367 | 0.02074139 |
| | | $t_p(R)$ | 1.00232 | **0.1686173** 0.1685935 | **119.9575** 118.1246 | 0.01873526 |
| 1 | 200 | $t_0(R)$ | 1.032376 | **0.6313128** 0.6288249 | **100** 100 | 0.06645398 |
| | | $t_1(R)$ | 0.9967019 | **0.5716984** 0.5582806 | **110.4275** 112.6359 | 0.06017878 |
| | | $t_p(R)$ | 0.9682892 | **0.494106** 0.5058227 | **127.7686** 124.3172 | 0.05201116 |
| | 500 | $t_0(R)$ | 1.040029 | **0.2705931** 0.2652877 | **100** 100 | 0.02848348 |
| | | $t_1(R)$ | 0.9968461 | **0.212635** 0.2254085 | **127.2570** 117.6919 | 0.02238263 |
| | | $t_p(R)$ | 0.9791635 | **0.1965888** 0.204013 | **137.6442** 130.0347 | 0.02069356 |

Tables 1, 2 and 3 show the values of the theoretical MSEs and empirical MSEs. The values from the table confirm that the basic estimator $t_0(R)$ and the ratio estimator $t_1(R)$ are less efficient as compared to the generalized estimator $t_p(R)$. Also, while comparing the generalized estimator $t_p(R)$ with the ratio estimator $t_1(R)$ and basic estimator $t_0(R)$, we note that as the variance of $T$ or variance of $S$ increase, the MSEs increase. This is expected since adding more noise makes the MSE increase. However, if we look at the unified measure ($\delta$), we find that it does not always increase as variance of $T$ or variance of $S$ increase, or at least not to the same extent as does the MSE. For example, for the generalized estimator

$t_p(R)$, theoretical MSE for Population II, with sample size 500, is 0.09227229 for $\sigma_T^2 = 0.2$ but increases to 0.3790013 for $\sigma_T^2 = 1$. In contrast, the ($\delta$) value decreases from 0.023659 to 0.020499. Admittedly, this is not a big drop in ($\delta$) value but at least it is not going up. The important point here is that the 310% increase in MSE (from 0.09227229 to 0.3790013) is more than offset by the significant increase in privacy level in using $\sigma_T^2 = 1$ as compared to $\sigma_T^2 = 0.2$. In another example, for the generalized estimator $t_p(R)$, theoretical MSE for Population III, with sample size 500, is 0.1877209 for $\sigma_T^2 = 0.5$ but increases to 0.3634541 for $\sigma_T^2 = 1$. In contrast, the ($\delta$) value decreases from 0.021453 to 0.021069.

**Table 2**:  Theoretical (**bold**) and empirical MSEs and PREs of the estimators for Population II with $\sigma_s^2 = 0.5$, $\sigma_y^2 = 2$ and $\rho_{yx} = 0.80$.

| Var($T$) | $n$ | Estimator | Mean($\hat{\sigma}_y^2$) | MSE | PRE | $\delta$ |
|---|---|---|---|---|---|---|
| 0.2 | 200 | $t_0(R)$ | 1.961504 | **0.3353948** 0.3330506 | **100** 100 | 0.085998 |
| | | $t_1(R)$ | 1.938223 | **0.3086746** 0.310405 | **108.6564** 107.2955 | 0.079147 |
| | | $t_p(R)$ | 1.97547 | **0.2604031** 0.2696629 | **128.7983** 123.5062 | 0.066770 |
| | 500 | $t_0(R)$ | 1.984015 | **0.1299197** 0.1273284 | **100** 100 | 0.033312 |
| | | $t_1(R)$ | 1.999045 | **0.1057879** 0.1067183 | **122.8114** 119.3126 | 0.027125 |
| | | $t_p(R)$ | 1.985764 | **0.09227229** 0.09218931 | **140.8003** 138.1162 | 0.023659 |
| 0.5 | 200 | $t_0(R)$ | 1.997112 | **0.8036853** 0.7958328 | **100** 100 | 0.084651 |
| | | $t_1(R)$ | 1.988183 | **0.7195406** 0.694571 | **111.6942** 114.5790 | 0.075788 |
| | | $t_p(R)$ | 1.98627 | **0.624445** 0.6421061 | **128.7039** 123.9410 | 0.065772 |
| | 500 | $t_0(R)$ | 1.991561 | **0.2858802** 0.2751116 | **100** 100 | 0.030111 |
| | | $t_1(R)$ | 1.982515 | **0.2471334** 0.232594 | **115.6784** 118.2797 | 0.026030 |
| | | $t_p(R)$ | 1.968053 | **0.1816638** 0.1885275 | **157.3677** 145.9265 | 0.019134 |
| 1 | 200 | $t_0(R)$ | 1.981875 | **1.170947** 1.167372 | **100** 100 | 0.063335 |
| | | $t_1(R)$ | 2.002721 | **1.014171** 0.5582806 | **115.4585** 112.8290 | 0.054855 |
| | | $t_p(R)$ | 1.988997 | **0.955732** 0.969496 | **122.5183** 120.4101 | 0.051694 |
| | 500 | $t_0(R)$ | 1.979819 | **0.5567679** 0.531363 | **100** 100 | 0.030114 |
| | | $t_1(R)$ | 1.998328 | **0.4837988** 0.4790216 | **115.0825** 110.9267 | 0.026168 |
| | | $t_p(R)$ | 1.971607 | **0.3790013** 0.3843118 | **146.9039** 138.2635 | 0.020499 |

**Table 3**:   Theoretical (**bold**) and empirical MSEs and PREs of the estimators
for Population III with $\sigma_s^2 = 0.25$, $\sigma_y^2 = 1$ and $\rho_{yx} = 0.60$.

| Var($T$) | $n$ | Estimator | Mean($\hat{\sigma}_y^2$) | MSE | PRE | $\delta$ |
|---|---|---|---|---|---|---|
| 0.2 | 200 | $t_0(R)$ | 1.021512 | **0.2249759** 0.223441 | **100** 100 | 0.061637 |
| | | $t_1(R)$ | 1.037037 | **0.1962207** 0.1958733 | **114.6545** 114.0742 | 0.053759 |
| | | $t_p(R)$ | 0.979563 | **0.1733191** 0.1752187 | **129.8044** 127.5212 | 0.047484 |
| | 500 | $t_0(R)$ | 0.99568 | **0.09192312** 0.09384772 | **100** 100 | 0.025184 |
| | | $t_1(R)$ | 1.035195 | **0.08558669** 0.08575554 | **107.4035** 109.4363 | 0.023448 |
| | | $t_p(R)$ | 0.995747 | **0.06216159** 0.06279393 | **147.8776** 149.4534 | 0.017030 |
| 0.5 | 200 | $t_0(R)$ | 0.9830188 | **0.6333537** 0.6304459 | **100** 100 | 0.072383 |
| | | $t_1(R)$ | 1.039288 | **0.5491218** 0.5699384 | **115.3393** 110.6164 | 0.062756 |
| | | $t_p(R)$ | 0.971143 | **0.4907475** 0.5044131 | **129.0589** 124.9860 | 0.056085 |
| | 500 | $t_0(R)$ | 0.9941702 | **0.2469968** 0.2442127 | **100** 100 | 0.028228 |
| | | $t_1(R)$ | 0.9846135 | **0.2070803** 0.2115374 | **119.2758** 115.4465 | 0.023666 |
| | | $t_p(R)$ | 0.9992722 | **0.1877209** 0.1827657 | **131.5766** 133.6206 | 0.021453 |
| 1 | 200 | $t_0(R)$ | 0.9571123 | **1.166476** 1.148805 | **100** 100 | 0.067621 |
| | | $t_1(R)$ | 0.9954355 | **1.092394** 1.087534 | **106.7816** 105.6339 | 0.063327 |
| | | $t_p(R)$ | 0.9794743 | **0.9463649** 0.9256485 | **123.2585** 124.1081 | 0.054861 |
| | 500 | $t_0(R)$ | 1.009706 | **0.5152219** 0.4923866 | **100** 100 | 0.029867 |
| | | $t_1(R)$ | 0.9918212 | **0.4304643** 0.458314 | **119.6898** 107.4343 | 0.024954 |
| | | $t_p(R)$ | 0.9856029 | **0.3634541** 0.3569531 | **141.7570** 137.9415 | 0.021069 |

## 5.   APPLICATION

In this section, we use a real data to show the performance of the generalized estimator $t_p(R)$ in comparison to other estimators. For this data which can be obtained from James *et al.* [5], the population size is ($N = 777$). The study variable $Y$ is the reported percent of alumni who donate. The auxiliary variable $X$ is the student to faculty ratio. The scrambling variable $S$ is taken to be a normal random variable with mean equal to zero and variance equal to 0.5. The scrambling variable $T$ is taken to be a normal random variable with mean equal to 1 and variance equal to 0.2, 0.5, and 1.

Population Characteristics are given by

$$N = 777, \quad n = 200, \quad \mu_X = 14.08, \quad \mu_Y = 22.74,$$

$$\sigma_X = 3.95, \quad \sigma_Y = 12.39, \quad \sigma_{XY} = 19.7641, \quad \rho_{yx} = 0.40.$$

From the Table 4, it can be observed that the generalized estimator $t_p(R)$ performs better than the other estimators $t_0(R)$ and $t_1(R)$. Also, we can observe that the unified measure $(\delta)$ does not always increase as variance of $T$ increases, or at least not to the same extent as does the MSE. For example, for the generalized estimator $t_p(R)$, theoretical MSE is 301.0716 for $\sigma_T^2 = 0.2$ but increases to 1196.559 for $\sigma_T^2 = 1$. In contrast, the $(\delta)$ value decreases from 2.23565474 to 1.78234135.

**Table 4:** Theoretical (**bold**) and empirical MSEs and PREs of the estimators.

| $n$ | $\mathbf{Var}(T)$ | Estimator | MSE | PRE | $\delta$ |
|---|---|---|---|---|---|
| 500 | 0.2 | $t_0(R)$ | **519.1796** 490.2126 | **100** 100 | 3.85525016 |
| | | $t_1(R)$ | **435.4705** 437.4432 | **119.2226** 112.0631 | 3.23365501 |
| | | $t_p(R)$ | **301.0716** 297.7625 | **172.4438** 164.6320 | 2.23565474 |
| | 0.5 | $t_0(R)$ | **896.6322** 888.2846 | **100** 100 | 2.66917897 |
| | | $t_1(R)$ | **643.4997** 620.5305 | **139.3368** 143.1492 | 1.91563036 |
| | | $t_p(R)$ | **596.1386** 570.0859 | **150.4066** 155.8159 | 1.77464139 |
| | 1 | $t_0(R)$ | **1805.427** 1876.467 | **100** 100 | 2.68928418 |
| | | $t_1(R)$ | **1618.569** 1650.915 | **111.5446** 113.6622 | 2.41094877 |
| | | $t_p(R)$ | **1196.559** 1105.511 | **150.8849** 169.7375 | 1.78234135 |

## 6.  CONCLUSION

We propose here some variance estimators under RRT. These are the basic estimator $t_0(R)$, ratio estimator $t_1(R)$ and the generalized estimator $t_p(R)$. The simulation study reveals that the generalized estimator $t_p(R)$ is more efficient than the other estimators $t_0(R)$ and $t_1(R)$. We also examine the efficiency of the estimators relative to not just the MSE values, but also with respect to the unified measure of estimators quality $(\delta)$ and observe that while MSE always increases as the noise level increases, the $(\delta)$ value does not necessary follow this pattern. This highlights the significance of respondent under privacy.

## ACKNOWLEDGMENTS

The authors would like to thank the two reviewers for their constructive comments which helped improve the presentation of this paper.

## REFERENCES

[1]   COLLINS, L.M.; SCHAFER, J.L.; *et al.* (2001). A comparison of inclusive and restrictive strategies in modern missing data procedures, *Psychological Methods*, **6**(4), 330–351.

[2]   DIANA, G. and PERRI, P.F. (2011). A class of estimators for quantitative sensitive data, *Statistical Papers*, **52**(3), 633–650.

[3]   GUPTA, S.; SAMRIDHI, M.; SHABBIR, J. and KHALIL, S. (2018). A unified measure of respondent privacy and model efficiency in quantitative RRT models, *Journal of Statistical Theory and Practice*, **12**(3), 506–511.

[4]   ISAKI, C.T. (1983). Variance estimation using auxiliary information, *Journal of the American Statistical Association*, **78**(381), 117–123.

[5]   JAMES, G.; WITTEN, D.; HASTIE, T. and TIBSHIRANI, R. (2013). *An Introduction to Statistical Learning: with Applications in R*, Springer, New York.

[6]   KHALIL, S.; ZHANG, Q. and GUPTA, S. (2019). Mean estimation of sensitive variables under measurement errors using optional RRT models. To appear in *Communications in Statistics – Simulation and Computation*.

[7]   WARNER, S.L. (1965). Randomized response: A survey technique for eliminating evasive answer bias, *Journal of the American Statistical Association*, **60**(309), 63–69.

[8]   YAN, Z.; WANG, J. and LAI, J. (2009). An efficiency and protection degree-based comparison among the quantitative randomized response strategies, *Communications in Statistics – Theory and Methods*, **38**, 404–408.

# ON ARNOLD–VILLASEÑOR CONJECTURES FOR CHARACTERIZING EXPONENTIAL DISTRIBUTION BASED ON SAMPLE OF SIZE THREE

Author:    GEORGE P. YANEV
        – School of Mathematical and Statistical Sciences,
          The University of Texas Rio Grande Valley, Texas, USA
          george.yanev@utrgv.edu
          and
          Institute of Mathematics and Informatics,
          Bugarian Academy of Sciences, Bulgaria

Abstract:

• Arnold and Villaseñor [4] obtain a series of characterizations of the exponential distribution based on random samples of size two. These results were already applied in constructing goodness-of-fit tests in [7]. Extending the techniques from [4], we prove some of Arnold and Villaseñor's conjectures for samples of size three. An example with simulated data is discussed.

Key-Words:

• *exponential distribution; characterizations; order statistics.*

AMS Subject Classification:

• 62G30, 62E10.

## 1.    INTRODUCTION

In general, the problem of characterization of probability distributions is described as follows. Suppose a family of distributions $\mathcal{F}$ possesses a property $\mathcal{A}$. If, conversely, a distribution has property $\mathcal{A}$ only if it is a member of that family, then property $\mathcal{A}$ characterizes the family $\mathcal{F}$. This result is referred to as a characterization of the distributions in $\mathcal{F}$. Primary motivation for characterizations problems is due to statistical applications. If a statistical procedure assumes that property $\mathcal{A}$ holds, then the underlying distribution must be a member of the family $\mathcal{F}$. Naturally, first characterizations results are for the normal family of distributions. The exponential distribution is one of the non-normal distributions, which has received a lot of attention as well. Comprehensive surveys of exponential characterizations can be found in [1], [3], [5], [6], and [8].

More recently, Arnold and Villaseñor [4] obtained a series of characterizations of the exponential distribution based on random samples of size two and conjectured possible generalizations for samples of size three. They provide motivation for their results by pointing out an example of a goodness-of-fit construction. A test for exponentiality based on the characterizations in [4] was recently constructed in [7]. Another possible use of the results in [4] and their generalizations, is in verifying modeling assumptions and in simulations (see also [8]). Extending the techniques from [4], we will prove some of Arnold and Villaseñor's conjectures.

Assume throughout that $X_1, X_2$, and $X_3$ are independent random variables with a common absolutely continuous cumulative distribution function (cdf) $F$, such that $F(0) = 0$ and probability density function (pdf) $f$. Denote $X_{2:2} := \max\{X_1, X_2\}$, $X_{3:3} := \max\{X_1, X_2, X_3\}$, and $\bar{F} = 1 - F$. Consider the relations:

$$(1.1) \qquad \sum_{j=1}^{3} \frac{1}{j} X_j \quad \text{has pdf} \quad \sum_{j=1}^{3} \binom{3}{j}(-1)^{j-1} j\, f(j\,x)\,,$$

$$(1.2) \qquad X_{3:3} \quad \text{has pdf} \quad \sum_{j=1}^{3} \binom{3}{j}(-1)^{j-1} j\, \bar{F}(j\,x)\,,$$

$$(1.3) \qquad \sum_{j=1}^{3} \binom{3}{j}(-1)^{j-1} j\, f(jx) = \sum_{j=1}^{3} \binom{3}{j}(-1)^{j-1} j\, \bar{F}(j\,x)\,,$$

$$(1.4) \qquad X_{2:2} + \frac{1}{3} X_3 \stackrel{d}{=} X_{3:3} \quad \text{and} \quad \sum_{j=1}^{3} \frac{1}{j} X_j \stackrel{d}{=} X_{3:3}\,,$$

where $\stackrel{d}{=}$ denotes equality in distribution. We will prove, under some regularity assumptions on $F$, that each one of these five conditions, on its own, is sufficient for $X_1, X_2$, and $X_3$ to be exponentially distributed.

We organize this paper as follows. Using Laplace transforms, in Section 2 we prove the characterization (1.1). In Section 3, we establish characterization (1.2) utilizing the Taylor series expansion of the cdf $F$. In Section 4, using a recurrent relation, we prove that (1.3) is a sufficient condition for having exponential parent. Section 5 contains characterization results based on (1.4). In Section 6 we provide an example with simulated data. In the concluding section, we discuss possible extensions of the given results.

$$\sum_{j=1}^{3} \frac{1}{j} X_j \text{ has pdf } \sum_{j=1}^{3} \binom{3}{j} (-1)^{j-1} j f(jx)$$

$$\sum_{j=1}^{3} \frac{1}{j} X_j \stackrel{d}{=} X_{3:3}$$

$$X_1, X_2, X_3 \text{ i.i.d. Exp}(\lambda)$$

$$X_{2:2} + \tfrac{1}{3} X_3 \stackrel{d}{=} X_{3:3}$$

$$X_{3:3} \quad \text{has pdf} \quad \sum_{j=1}^{3} \binom{3}{j} (-1)^{j-1} j \bar{F}(jx)$$

## 2.  SUM OF THREE INDEPENDENT VARIABLES

To prove that (1.1) characterizes the exponential distribution, we will convert it into an equation for the Laplace transform $\varphi(t) := E[e^{-tX_1}]$.

**Theorem 2.1.**  *Assume $\varphi(t)$ is finite for all $t$ in a neighbourhood of zero. If for $x > 0$*

$$(2.1) \qquad \sum_{j=1}^{3} \frac{1}{j} X_j \quad has\ pdf \quad \sum_{j=1}^{3} \binom{3}{j} (-1)^{j-1} j f(jx) \,,$$

*then $X_1 \sim \exp(\lambda)$ for some $\lambda > 0$.*

**Proof:**  It follows by (2.1), interchanging the order of summation and integration, that

$$
\begin{aligned}
\varphi(t)\, \varphi\!\left(\frac{t}{2}\right) \varphi\!\left(\frac{t}{3}\right) &= E\left[ e^{-t \sum_{j=1}^{3} \frac{1}{j} X_j} \right] \\
&= \int_0^\infty e^{-tx} \left( \sum_{j=1}^{3} \binom{3}{j} (-1)^{j-1} j f(jx) \right) dx \\
&= \sum_{j=1}^{3} \binom{3}{j} (-1)^{j-1} \int_0^\infty e^{-tx} j f(jx)\, dx \\
&= \sum_{j=1}^{3} \binom{3}{j} (-1)^{j-1} \varphi\!\left(\frac{t}{j}\right).
\end{aligned}
$$

(2.2)

Dividing both sides of (2.2) by $\varphi(t)\,\varphi(t/2)\,\varphi(t/3)$, we obtain

$$(2.3) \qquad 1 \;=\; \alpha(t)\,\alpha\!\left(\frac{t}{2}\right) - 3\,\alpha(t)\,\alpha\!\left(\frac{t}{3}\right) + 3\,\alpha\!\left(\frac{t}{2}\right)\alpha\!\left(\frac{t}{3}\right),$$

where for $t > 0$

$$(2.4) \qquad \alpha(t) \;:=\; \frac{1}{\varphi(t)} \;=\; \sum_{k=0}^{\infty} a_k t^k \,.$$

Note that, the series in (2.4) is convergent in a neighbourhood of zero, by assumption. To prove the theorem, it is sufficient to show that

$$(2.5) \qquad \alpha(t) = 1 + \lambda t, \qquad \lambda > 0\,.$$

We will prove (2.5) by calculating the coefficients of the series in (2.4) to be: $a_0 = 1$, $a_1 = \lambda > 0$, and $a_k = 0$ for $k \geq 2$. It is clear that $a_0 = \varphi^{-1}(0) = 1$. Applying Cauchy formula for multiplication of two power series, we have for any nonzeros $p$ and $q$,

$$(2.6) \qquad \alpha\!\left(\frac{t}{p}\right)\alpha\!\left(\frac{t}{q}\right) = \sum_{k=0}^{\infty} \left( \sum_{j=0}^{k} \frac{1}{p^j\,q^{k-j}}\, a_j\, a_{k-j} \right) t^k \,.$$

Now, (2.3) and (2.6) yield for $k \geq 1$

$$(2.7) \qquad \sum_{j=0}^{k} \left( \frac{1}{2^{k-j}} - \frac{3}{3^{k-j}} + \frac{3}{2^j\,3^{k-j}} \right) a_j\, a_{k-j} \;=\; 0\,.$$

Setting $k = 1$ we see that equation (2.7) has as solution any $a_1$. The assumption $F(0) = 0$ implies that there is $\lambda > 0$, such that $a_1 = \lambda > 0$. If $k = 2$, then (2.7) yields $a_2 = 0$. Assuming $a_j = 0$ for $2 \leq j \leq k - 1$, it follows from (2.7) that

$$\left( 1 - \frac{1}{2^{k-1}} \right) a_k \;=\; 0\,.$$

Thus, $a_k = 0$ for any $k \geq 3$. Therefore, (2.5) holds, which completes the proof. $\qquad\square$

Note that, conversely, if $X_i \sim \exp(\lambda)$ for $i = 1, 2, 3$, then (2.1) holds true. To show this, it is sufficient to verify (2.2). Indeed, assuming $X_1 \sim \exp(\lambda)$, we have $\varphi(t) = (1 + \lambda t)^{-1}$. Therefore,

$$
\begin{aligned}
\int_0^\infty e^{-tx} \left( \sum_{j=1}^{3} \binom{3}{j}(-1)^{j-1} j\, f(jx) \right) dx
&= 3\,\varphi(t) - 3\,\varphi\!\left(\frac{t}{2}\right) + \varphi\!\left(\frac{t}{3}\right) \\
&= \frac{3}{1 + \lambda t} - \frac{6}{2 + \lambda t} + \frac{3}{3 + \lambda t} \\
&= \varphi(t)\,\varphi\!\left(\frac{t}{2}\right)\varphi\!\left(\frac{t}{3}\right) \\
&= E\!\left[ e^{-t\left( X_1 + \frac{1}{2}X_2 + \frac{1}{3}X_3 \right)} \right],
\end{aligned}
$$

which is equivalent to (2.1).

## 3.  MAXIMUM OF THREE INDEPENDENT VARIABLES

In this section we will prove that, under some regularity assumptions on $F$, condition (1.2) is sufficient for $X_1, X_2$, and $X_3$ to be exponentially distributed. The proof will be based on the Taylor series expansion of $F$.

**Theorem 3.1.** *Assume the cdf $F$ has a power series representation for $x$ in a neighborhood of zero. If for $x > 0$*

$$(3.1) \qquad X_{3:3} \quad \text{has pdf} \quad \sum_{k=1}^{3} \binom{3}{k} (-1)^{k-1} k \, \bar{F}(kx),$$

*then $X_1 \sim \exp(1)$.*

**Proof:** The relation (3.1) implies

$$(3.2) \qquad F^2(x) f(x) + F(x) - 2 F(2x) + F(3x) = 0.$$

Since $F(x) = \sum_{k=0}^{\infty} c_k x^k$ and $f(x) = \sum_{k=0}^{\infty} (k+1) c_{k+1} x^k$, Cauchy formula for the product of three power series yields

$$(3.3) \qquad F^2(x) f(x) = \sum_{k=0}^{\infty} \left[ \sum_{i=0}^{k} \sum_{j=0}^{i} c_j c_{i-j} (k+1-i) c_{k+1-i} \right] x^k.$$

Using (3.2) and (3.3), we obtain for any $k \geq 0$

$$(3.4) \qquad \sum_{i=0}^{k} \sum_{j=0}^{i} c_j c_{i-j} (k+1-i) c_{k+1-i} + c_k (1 - 2^{k+1} + 3^k) = 0.$$

Since $F(0) = 0$, we have $c_0 = 0$. Also (3.4) with $k = 1$ yields $c_0^2 c_1 = 0$, which in turn implies that $c_1$ is undetermined. Let us set $c_1 = \delta$, where $-\infty < \delta < \infty$. Equation (3.4) with $k = 2$ yields $c_1^3 + 2 c_2 = 0$. Hence, $c_2 = \delta^3/2$. We will prove by induction that

$$(3.5) \qquad c_k = (-1)^{k-1} \frac{\delta^{2k-1}}{k!}, \qquad k = 1, 2, 3, \dots.$$

Indeed, assuming (3.5) holds true for $1, 2, ..., k$, we have

$$(3.6) \qquad \sum_{i=0}^{k+1} \sum_{j=0}^{i} c_j c_{i-j} (k+2-i) c_{k+2-i} = \sum_{i=2}^{k+1} \sum_{j=1}^{i-1} \frac{(-1)^{k+1} \delta^{2k+1}}{j! \, (i-j)! \, (k+1-i)!}.$$

Observe that

$$(3.7) \qquad \begin{aligned}
\sum_{i=2}^{k+1} \sum_{j=1}^{i-1} \frac{1}{j! \, (i-j)! \, (k+1-i)!} &= \sum_{i=2}^{k+1} \frac{1}{i! \, (k+1-i)!} \sum_{j=1}^{i-1} \frac{i!}{j! \, (i-j)!} \\
&= \frac{1}{(k+1)!} \sum_{i=2}^{k+1} \frac{(k+1)!}{i! \, (k+1-i)!} (2^i - 2) \\
&= \frac{1}{(k+1)!} \left[ \sum_{i=2}^{k+1} \binom{k+1}{i} 2^i - 2 \sum_{i=2}^{k+1} \binom{k+1}{i} \right] \\
&= \frac{1}{(k+1)!} \left( 3^{k+1} - 2^{k+2} + 1 \right).
\end{aligned}$$

It follows from (3.4), (3.6) and (3.7) that

$$(-1)^{k+1} \frac{\delta^{2k+1}}{(k+1)!} \left( 3^{k+1} - 2^{k+2} + 1 \right) + c_{k+1}(1 - 2^{k+2} + 3^{k+1}) = 0 .$$

Therefore,

$$c_{k+1} = (-1)^k \frac{\delta^{2k+1}}{(k+1)!} ,$$

which completes the induction and hence proves (3.5).

Now, we have

$$F(x) = \sum_{k=1}^{\infty} (-1)^{k+1} \frac{\delta^{2k-1}}{k!} x^k = \frac{1}{\delta} \left( 1 - e^{-\delta^2 x} \right) .$$

Since $\lim_{x \to \infty} F(x) = 1$, we obtain $\delta = 1$. The proof is complete. $\qquad \square$

It is not difficult to see that, conversely, if $X_1 \sim \exp(1)$, then (3.1) holds. Indeed, under the assumption of unit exponential parent variable, for the pdf of $X_{3:3}$ we obtain

$$3 F^2(x) f(x) = 3 (1 - e^{-x})^2 e^{-x} = 3 \bar{F}(x) - 6 \bar{F}(2x) + 3 \bar{F}(3x) ,$$

which is equivalent to (3.1).

## 4. SUMS OF DENSITY AND DISTRIBUTION FUNCTIONS

In this section we will prove that (1.3) is a sufficient condition for $X_1$ to be exponentially distributed. It is straightforward that (1.3) is a necessary condition as well.

**Theorem 4.1.** *Assume that $f$ is right-continuous at zero. If for $x > 0$*

$$(4.1) \qquad \sum_{j=1}^{3} \binom{3}{j} (-1)^{j-1} j f(j x) = \sum_{j=1}^{3} \binom{3}{j} (-1)^{j-1} j \bar{F}(j x) ,$$

*then $X_1 \sim \exp(1)$.*

**Proof:** The relation (4.1) leads to

$$(4.2) \qquad \left[ f(3x) - \bar{F}(3x) \right] - \left[ f(2x) - \bar{F}(2x) \right] = \left[ f(2x) - \bar{F}(2x) \right] - \left[ f(x) - \bar{F}(x) \right] .$$

Denoting $Q(y) = f(y) - \bar{F}(y)$, we rewrite (4.2) as

$$Q(y) - Q\left( \frac{2}{3} y \right) = Q\left( \frac{2}{3} y \right) - Q\left( \frac{1}{3} y \right) .$$

Iterating this equation $k$ times and taking limit as $k \to \infty$, we obtain

$$Q(y) - Q\left( \frac{2}{3} y \right) = Q\left( \frac{2}{3} y \right) - Q\left( \frac{1}{3} y \right) = \lim_{k \to \infty} Q\left( \left( \frac{2}{3} \right)^k y \right) - Q\left( \left( \frac{1}{3} \right)^k y \right) = 0 .$$

This implies $Q(y) = Q(2y/3)$ and thus,

$$(4.3) \qquad Q(y) = Q\left(\frac{2}{3}y\right) = \lim_{k\to\infty} Q\left(\left(\frac{2}{3}\right)^k y\right) = f(0+) - \bar{F}(0+) = f(0+) - 1\,.$$

On the other hand,

$$(4.4) \qquad \lim_{y\to\infty} f(y) = \lim_{y\to\infty} f(y) - \lim_{y\to\infty} \bar{F}(y) = \lim_{y\to\infty} Q(y) = f(0+) - 1\,.$$

But since $f$ is integrable, we have $\lim_{y\to\infty} f(y) = 0$, and therefore, by (4.3) and (4.4), $Q(x) = 0$. Thus, $f(x) = \bar{F}(x)$ for every $x \geq 0$. This, in turn, implies $X_1 \sim \exp(1)$. $\qquad\square$

---

## 5.   SUM AND MAXIMUM OF THREE VARIABLES

It is known (e.g., Arnold *et al.* (2008), p. 77) that if $X \sim \exp\{\lambda\}$, then

$$(5.1) \qquad \sum_{j=1}^{3} \frac{1}{j} X_j \overset{d}{=} X_{3:3} \qquad \text{and} \qquad X_{2:2} + \frac{1}{3} X_3 \overset{d}{=} X_{3:3}\,.$$

We will prove that both relations in (5.1) are also characterization properties of the exponential distribution. Next lemma provides the key argument in the proof of Theorem 1 in [4] and of the theorem below.

**Lemma 5.1.**   *If $F(0) = 0$, the pdf $f$ has a Taylor series expansion for $x > 0$, and*

$$(5.2) \qquad f^{(m)}(0) = \left[\frac{f'(0)}{f(0)}\right]^{m-1} f'(0)\,, \qquad m = 1, 2, \ldots,$$

*then $X_1 \sim \exp\{\lambda\}$ for some $\lambda > 0$.*

**Proof:**   For the Taylor series of $f(x)$, using (5.2), we have for $x > 0$

$$f(x) = \sum_{m=0}^{\infty} \frac{f^{(m)}(0)}{m!} x^m = f(0) + f(0) \sum_{m=1}^{\infty} \left[\frac{f'(0)}{f(0)}\right]^m \frac{x^m}{m!} = f(0)\exp\left\{\frac{f'(0)}{f(0)} x\right\}\,.$$

Since $f(x)$ is a pdf, we have $f'(0)/f(0) < 0$. Denoting $\lambda = -f'(0)/f(0) > 0$ and setting $\int_0^\infty f(x)\,dx = 1$, we obtain $\lambda = f(0)$. Therefore, $f(x) = \lambda e^{-\lambda x}$. $\qquad\square$

Next theorem can be obtained as a particular case of the results in [9]. We include it here since it complements the other results for samples of size three given in Theorems 2.1–4.1 and thus provides an easily reference.

**Theorem 5.1.**   *Assume the cdf $F$ admits a power series representation in a neighborhood of zero and $F(0) = 0$.*

   **(i)**   *If*

$$(5.3) \qquad X_{2:2} + \frac{1}{3} X_3 \overset{d}{=} X_{3:3}\,,$$

   *then $X_1 \sim \exp\{\lambda\}$ for some $\lambda > 0$.*

(**ii**)   *If*

(5.4)
$$\sum_{j=1}^{3} \frac{1}{j} X_j \overset{d}{=} X_{3:3},$$

*then $X_1 \sim \exp\{\lambda\}$ for some $\lambda > 0$.*

**Proof:**  (**i**). The pdf of the left-hand side of (5.3) is

(5.5)
$$
\begin{aligned}
f_{X_{2:2}+X_{3}/3}(x) &= \int_0^x f_{X_3/3}(y)\, f_{X_{2:2}}(x-y)\, dy \\
&= \int_0^x 3 f(3y) \frac{d}{dx} \big[F^2(x-y)\big]\, dy \\
&= 6 \int_0^x f(3y)\, F(x-y)\, f(x-y)\, dy\,.
\end{aligned}
$$

For the pdf of the right-hand side of (5.3), we have

(5.6)
$$f_{X_{3:3}}(x) = 3 F^2(x)\, f(x) = 6 f(x) \int_0^x F(y)\, f(y)\, dy\,.$$

Let $G(x) := F(x)\, f(x)$. It follows from (5.5) and (5.6) that (5.3) is equivalent to

(5.7)
$$\int_0^x f(3y)\, G(x-y)\, dy = f(x) \int_0^x G(y)\, dy\,.$$

Differentiating the left-hand side of (5.7) $n$ times with respect to $x$, we obtain

$$\frac{d^n}{dx^n} \int_0^x f(3y)\, G(x-y)\, dy = \sum_{i=1}^{n} f^{(n-i)}(3x)\, G^{(i-1)}(0) + \int_0^x f(3y)\, G^{(n)}(x-y)\, dy\,.$$

Applying the Leibniz rule for the $n$-th derivative of a product of two functions to the right-hand side of (5.7), we obtain

$$\frac{d^n}{dx^n}\left[f(x) \int_0^x G(y)\, dy\right] = \sum_{i=1}^{n} \binom{n}{i} f^{(n-i)}(x)\, G^{(i-1)}(x) + f^{(n)}(x) \int_0^x G(y)\, dy\,.$$

In the last two equations letting $x = 0$, we have

(5.8)
$$\sum_{i=1}^{n} 3^{n-i} f^{(n-i)}(0)\, G^{(i-1)}(0) = \sum_{i=1}^{n} \binom{n}{i} f^{(n-i)}(0)\, G^{(i-1)}(0)\,.$$

Since $G(0) = 0$ and $G'(0) = f^2(0)$, the above equation is equivalent to

(5.9)
$$\left[3^{n-2} - \binom{n}{2}\right] f^{(n-2)}(0)\, f^2(0) = \sum_{i=3}^{n} \left[\binom{n}{i} - 3^{n-i}\right] f^{(n-i)}(0)\, G^{(i-1)}(0)\,,$$

where $n \geq 4$. We will prove that (5.9) implies (5.2). Equation (5.2) is trivially true for $m = 1$. To proceed by induction, assume (5.2) holds true for all $1 \leq m \leq n - 3$, where $n \geq 4$. We need to prove it for $m = n - 2$. Using the induction assumption, it is not difficult to obtain for $j = 1, 2, ..., n - 2$

$$G^{(j)}(0) = \sum_{i=0}^{j} \binom{j}{i} F^{(i)}(0)\, f^{(j-i)}(0) = f^2(0) \left[\frac{f'(0)}{f(0)}\right]^{j-1} (2^j - 1)\,.$$

Therefore, using the induction assumption again, we have for $i = 3, 4, ..., n-1$

$$(5.10) \qquad f^{(n-i)}(0)\, G^{(i-1)}(0) \; = \; \left[\frac{f'(0)}{f(0)}\right]^{n-3} f'(0)\, f^2(0)\, (2^{i-1}-1)\,.$$

Substituting this in the right-hand side of (5.9) yields

$$\left[3^{n-2} - \binom{n}{2}\right] f^{(n-2)}(0) \; = \; \left[\frac{f'(0)}{f(0)}\right]^{n-3} f'(0) \sum_{i=3}^{n} \left[\binom{n}{i} - 3^{n-i}\right] (2^{i-1}-1)\,.$$

To complete the proof of (5.2), it is sufficient to show that

$$3^{n-2} - \binom{n}{2} \; = \; \sum_{i=3}^{n} \left[\binom{n}{i} - 3^{n-i}\right] (2^{i-1}-1)\,,$$

which can be easily verified. This proves (5.2). The claim in (i) follows from (5.2) and the lemma. $\qquad\square$

**Proof:** (ii). Equation (5.4) is equivalent to

$$(5.11) \qquad 6 \int_0^z f(y) \int_0^{z-y} f(2x)\, f\big(3(z-y-x)\big)\, dx\, dy \; = \; 6\, f(z) \int_0^z F(y)\, f(y)\, dy\,.$$

Denoting

$$(5.12) \qquad H(z-y) \; := \; \int_0^{z-y} f(2x)\, f\big(3(z-y-x)\big)\, dx\,,$$

we write (5.11) as

$$(5.13) \qquad \int_0^z f(y)\, H(z-y)\, dy \; = \; f(z) \int_0^z G(y)\, dy\,.$$

Similarly to the proof of (i), differentiating $n$ times both sides of (5.13) with respect to $z$ and setting $z = 0$, we have

$$\sum_{i=1}^{n-1} f^{(n-1-i)}(0)\, H^{(i)}(0) \; = \; \sum_{i=1}^{n-1} \binom{n}{i+1} f^{(n-1-i)}(0)\, G^{(i)}(0)\,.$$

Since $H'(0) = G'(0) = f^2(0)$, the last equation can be written for $k = n-1$ as

$$(5.14) \qquad \left[1 - \binom{k+1}{2}\right] f^{(k-1)}(0)\, f^2(0) \; = \; \sum_{i=2}^{k} \left[\binom{k+1}{i+1} G^{(i)}(0) - H^{(i)}(0)\right] f^{(k-i)}(0)\,.$$

Now we are in a position to prove (5.2) by induction. (5.2) holds true for $m = 1, 2, ..., k-2$. Differentiating (5.12) with respect to $z$ and setting $z = y$, we have

$$(5.15) \qquad H^{(n)}(0) \; = \; \sum_{i=1}^{n} 2^{n-i}\, f^{(n-i)}(0)\, 3^{i-1} f^{(i-1)}(0)\,.$$

Under the induction assumption, (5.15) implies for $j = 1, 2, ..., n-2$

$$H^{(j)}(0) = \left[\frac{f'(0)}{f(0)}\right]^{j-1} f^2(0) \left(3^j - 2^j\right).$$

Using the induction assumption again, we have for $i = 3, 4, ..., n-1$

$$f^{(n-i)}(0)\, H^{(i-1)}(0) = \left[\frac{f'(0)}{f(0)}\right]^{n-3} f'(0)\, f^2(0) \left(3^{i-1} - 2^{i-1}\right).$$

Recalling (5.10) from the proof of (i), we rewrite (5.14) as (note that $i = n$ corresponds to a 0 term)

$$\left[1 - \binom{n}{2}\right] f^{(n-2)}(0) = \left[\frac{f'(0)}{f(0)}\right]^{n-3} f'(0) \sum_{i=3}^{n} \left[\binom{n}{i}\left(2^{i-1} - 1\right) - \left(3^{i-1} - 2^{i-1}\right)\right].$$

Thus, to prove (5.2) for $k = n - 2$ it is sufficient to show that

$$1 - \binom{n}{2} = \sum_{i=3}^{n} \left[\binom{n}{i}\left(2^{i-1} - 1\right) - \left(3^{i-1} - 2^{i-1}\right)\right],$$

which verifies. This proves (5.2), which referring to the lemma, completes the proof of (ii).    $\square$

## 6.    EXAMPLE

We will illustrate a possible application of Theorem 5.1 with an example (see also [4]). Assume we have a simple random sample $X_1, X_2, ..., X_n$ for $n \geq 6$. Let us randomly divide the data set into six subsets, relabeled as

$$U_1, U_2, ..., U_{n/6}, \qquad V_1, V_2, ..., V_{n/6}, \qquad W_1, W_2, ..., W_{n/6},$$

$$X_1, X_2, ..., X_{n/6}, \qquad Y_1, Y_2, ..., Y_{n/6}, \qquad Z_1, Z_2, ..., Z_{n/6}.$$

Define for $i = 1, 2, ..., n/4$

$$R_i := U_i + \frac{1}{2}V_i + \frac{1}{3}W_i, \quad S_i := \max\{U_i, V_i\} + \frac{1}{3}W_i \quad \text{and} \quad T_i := \max\{X_i, Y_i, Z_i\}.$$

Then, according to Theorem 5.1, the $R$'s, the $S$'s, and the $T$'s will have a common distribution if and only if the original $X$'s follow an exponential distribution.

Let us simulate a sample of size $n = 180$ from a parent variable with $\exp(1)$ distribution. The values of $R_i$, $S_i$, and $T_i$ for $i = 1, 2, ..., 30$ are presented in Table 1.

Using the non-parametric two-sample Wicoxon rank test, we compare the sample distribution functions of the $R$'s and $T$'s on one hand and the $S$'s and $T$'s on another. The test results provide evidence supporting an exponential underlying distribution. Namely, the hypothesis that the distributions of the $R$'s and the $T$'s are the same cannot be rejected with $p$-value 0.7635 ($W = 471$). The hypothesis that the distributions of the $S$'s and the $T$'s are the same cannot be rejected with $p$-value 0.9357 ($W = 444$).

**Table 1**:   Values $R_i$, $S_i$, and $T_i$ for $i = 1, 2, ..., 30$.

| $R$ | 3.56 | 0.70 | 0.62 | 3.33 | 0.30 | 0.78 | 2.29 | 0.97 | 1.59 | 0.50 |
|---|---|---|---|---|---|---|---|---|---|---|
|  | 0.83 | 2.27 | 0.69 | 2.95 | 0.32 | 4.12 | 0.74 | 0.91 | 2.66 | 0.48 |
|  | 2.87 | 2.19 | 2.32 | 1.08 | 3.69 | 1.98 | 1.13 | 1.32 | 3.37 | 2.73 |
| $S$ | 2.98 | 1.23 | 0.77 | 2.75 | 0.44 | 0.75 | 1.97 | 1.08 | 1.43 | 0.50 |
|  | 0.76 | 1.65 | 0.58 | 2.39 | 0.27 | 3.41 | 0.73 | 0.89 | 2.63 | 0.40 |
|  | 2.22 | 1.87 | 4.25 | 1.07 | 2.72 | 1.74 | 1.07 | 1.11 | 2.71 | 3.87 |
| $T$ | 2.07 | 0.60 | 0.97 | 0.47 | 2.84 | 0.84 | 1.02 | 1.84 | 0.57 | 2.88 |
|  | 1.39 | 1.92 | 8.46 | 1.77 | 2.60 | 1.42 | 1.50 | 0.47 | 0.26 | 2.17 |
|  | 1.92 | 1.67 | 2.87 | 1.06 | 2.24 | 6.63 | 0.52 | 1.09 | 1.33 | 1.07 |

## 7.   CONCLUDING REMARKS

In this paper we proved characterizations of the exponential distribution conjectured by Arnold and Villaseñor in [4]. Furthermore, under the assumptions of Theorem 2.1 and using the same technique of proof, it can be seen that if $X_1 + \frac{1}{2}X_2 + \frac{1}{3}X_3$ has as its density any one of the following seven forms, then $X_i$'s are exponential:

$$3f(x) - 6f(2x) + 3\bar{F}(3x), \qquad 3f(x) - 6\bar{F}(2x) + 3f(3x),$$
$$3\bar{F}(x) - 6f(2x) + 3f(3x), \qquad 3f(x) - 6\bar{F}(2x) + 3\bar{F}(3x),$$
$$3F(x) - 6f(2x) + 3\bar{F}(3x), \qquad 3\bar{F}(x) - 6\bar{F}(2x) + 3f(3x),$$
$$3\bar{F}(x) - 6\bar{F}(2x) + 3\bar{F}(3x).$$

Likewise, under the assumptions of Theorem 3.1 and using the same technique of proof, it can be obtained that if $X_{3:3}$ has as its density any one of the preceding seven forms, then $X_i$'s are exponential.

The results presented here can be extended in several directions. Naturally, one would like to explore the general case of samples of size $n$ for any $n \geq 4$. As we mentioned earlier, generalizations of Theorem 5.1 for arbitrary sample size are proved in [9]. Here we would like to propose as open problems the following two characterizations, which would extend Theorem 2.1 and Theorem 3.1, respectively.

**Proposition 7.1.**   *Let $X_1, X_2, ..., X_n$ be i.i.d. random variables, where $n \geq 4$. Assume $\varphi(t)$ is finite for all $t$ in a neighbourhood of zero. If for $x > 0$*

$$\sum_{j=1}^{n} \frac{1}{j} X_j \quad \text{has pdf} \quad \sum_{j=1}^{n} \binom{n}{j} (-1)^{j-1} j\, f(j\,x),$$

*then $X_1 \sim \exp(\lambda)$ for some $\lambda > 0$.*

**Proposition 7.2.** *Let $X_1, X_2, ..., X_n$ be i.i.d. random variables, where $n \geq 4$. Assume the cdf $F$ has a power series representation in a neighborhood of zero. If for $x > 0$*

$$X_{n:n} \quad \text{has pdf} \quad \sum_{j=1}^{n} \binom{n}{j} (-1)^{j-1} j \, \bar{F}(j \, x),$$

*then $X_1 \sim \exp(1)$.*

## ACKNOWLEDGMENTS

## REFERENCES

[1] AHSANULLAH, M. (2017). *Characterizations of Univariate Continuous Distributions*, Atlantic Press, Amsterdam, the Netherlands.

[2] ARNOLD, B.C.; BALAKRISHNAN, N. and NAGARAJA, H.N. (2008). *A First Course in Order Statistics*, SIAM, Philadelphia, USA.

[3] ARNOLD, B.C. and HUANG, J.S. (1995). *Characterizations*. In "The Exponential Distribution: Theory, Methods and Applications" (N. Balakrishnan and A.P. Basu, Eds.), Gordon and Breach, Amsterdam, 79–95.

[4] ARNOLD, B.C. and VILLASEÑOR, J.A. (2013). Exponential characterizations motivated by the structure of order statistics in sample of size two, *Statistics and Probability Letters*, **83**, 596–601.

[5] AZLAROV, T. and VOLODIN, N.A. (1986). *Characterization Problems Associated with the Exponential Distribution*, Springer, Berlin.

[6] GALTHER, U.; KAMPS, U. and SCHWEITZER, N. (1998). *Characterizations of distributions via identically distributed functions of order statistics*. In "Order Statistics: Theory and Methods, Handbook of Statistics", Vol. 16 (N. Balakrishnan and C.R. Rao, Eds.), North–Holland, Amsterdam, 257–290.

[7] JOVANOVIC, M.; MILOSEVIC, B.; NIKITIN, YA.YU.; OBRADOVIC, M. and VOLKOVA, K.YU. (2015). Tests of exponentiality based on Arnold–Villasenor characterization and their efficiencies, *Comput. Statist. Data Analysis*, **90**, 100–113.

[8] NAGARAJA, H.N. (2006). *Characterizations of probability distributions*. In "Springer Handbook of Engineering Statistics" (H. Pham, Ed.), Springer, 395–402.

[9] YANEV, G.P. and CHAKRABORTY, S. (2016). Characterization of exponential distribution and Sakhatme–Renyi decomposition of exponential maxima, *Statistics and Probability Letters*, **110**, 94–102.

# TESTING FOR TRENDS
# IN EXCESSES OVER A THRESHOLD
# USING THE GENERALIZED PARETO DISTRIBUTION

Authors:  SERGIO JUÁREZ
– Universidad Veracruzana, Veracruz, México

JAVIER ROJO
– Korvis Professor of Statistics, Oregon State University, Oregon, USA
javier.rojo@oregonstate.edu

Abstract:

- The Generalized Pareto Distribution (GPD) is used for modeling exceedances over thresholds. The general form of the GPD depends on three parameters: the location parameter $\mu$; the scale parameter ($\beta > 0$); and the shape parameter ($-\infty < \xi < \infty$). This work restricts attention to the case where $\mu = 0$ and shows that, as $\xi$ decreases while $\beta$ is kept fixed, the family of GPD$(\xi, \beta)$ distributions increases in the usual stochastic order. This property is used for testing the significance of trends in the size of the exceedances over high thresholds in a time series consisting of ozone measurements.

Key-Words:

- *Stochastic Order; Peaks Over a Threshold; Ozone Concentrations; Likelihood Ratio Tests.*

## 1.    INTRODUCTION

Let $X$ be a random variable with continuous distribution function $F$ and corresponding survival function $\bar{F} = 1 - F$. Let $x^*$ be the right endpoint of the support of $F$ defined by $x^* = \sup\{x \in \mathbb{R} \colon F(x) < 1\}$. Given a real number $u < x^*$, referred to as the threshold, an *exceedance* over the threshold $u$ occurs when $X > u$. The residual life function of $F$ at time $u$, the probability that $X > u + x$ given that $X > u$, is

$$(1.1) \qquad \bar{F}_u(x) \,=\, P\big(X - u > x \mid X > u\big) \,=\, \frac{\bar{F}(x + u)}{\bar{F}(u)}\,, \qquad 0 < x < x^* - u\,.$$

The random variable $X - u$ is called the *excess* over the threshold $u$ and $\bar{F}_u$ is the *excess* survival function of $X$ over $u$. When $F$ belongs to the domain of attraction of one of the extreme value distributions, it follows that, for sufficiently large $u$, the distribution function of $X - u$ can be approximated by the Generalized Pareto Distribution (GPD). The distribution function of a GPD$(\xi, \beta)$ is

$$(1.2) \qquad F(x; \xi, \beta) = \begin{cases} 1 - (1 - \xi x/\beta)^{1/\xi}, & \xi \neq 0\,,\ \beta > 0\,, \\[2mm] 1 - \exp(-x/\beta)\,, & \xi = 0\,,\ \beta > 0\,, \end{cases}$$

where $\xi$ and $\beta$ are the shape and scale parameters, respectively. When $\xi < 0$ the support of $F(x; \xi, \beta)$ consists of the positive reals. When $\xi > 0$, the support is the interval $(0, \beta/\xi)$. The case $\xi = 0$ corresponds to the exponential distribution with mean $\beta$. When $\xi = 1$, the GPD distribution corresponds to the uniform distribution on $[0, \beta]$.

More precisely, let $X_1, ..., X_n$ be a sequence of independent and identically distributed random variables with continuos distribution $H$. Let $M_n = \max\{X_1, ..., X_n\}$. Suppose that there are sequences $a_n > 0$ and $b_n$ of real numbers such that

$$(1.3) \qquad P\big\{a_n(M_n - b_n) \leq z\big\} \to G(z)\,, \qquad \text{as} \ \ n \to \infty\,.$$

Then $G(z)$ is a member of the generalized extreme value distribution family defined by

$$G(z) \,=\, \exp\left\{ - \left\{ 1 - \xi\Big(\frac{z - \mu}{\sigma}\Big) \right\}^{1/\xi} \right\}.$$

The precise technical justification for modeling excesses using the GPD — expression (1.2) — was provided by Smith [32] and is based on the fact that

$$\lim_{u \to x^*} \sup_{0 < x < x^* - u} \big| F_u(x) - F\big(x; \xi, \beta(u)\big) \big| \,=\, 0\,,$$

for fixed $\xi$ and some positive function $\beta(u)$, if and only if $F$ is in the domain of attraction of some extreme value distribution. This result is from the parallel work done by Balkema and de Haan [1] and Pickands [23]. Since most of the common continuous distributions belong to the domain of attraction of one of the three extreme value distributions, this result makes the GPD the natural model for the excess distribution of the random variable $X$ when the threshold is high.

Starting with the early works by Smith [31] and Davison [6], the GPD has been used by many authors to model excesses over high thresholds in several fields such as river floods, air pollution, wind velocity, sea waves, insurance claims, etc. For the details of these applications see Hosking and Wallis [12], Smith [33], Dargahi-Noubary [5], Grimshaw [10], Rootzen and Tajvidi [29], Castillo and Hady [4], and Parisi and Lund [22]. Embrechts *et al.* [8], Falk *et al.* [9], and Reiss and Thomas [24] present detailed and elegant accounts of the theoretical underpinnings and the practical aspects of the modeling of extremes including discussions on the modeling of exceedances and excesses.

One of the main objectives of modeling excesses over high thresholds with the GPD is the estimation of tails of probability distributions — Smith [32]. But the GPD has also been used to detect and test for trends in the excesses. The papers by Smith [33], Davison and Smith [7], Smith and Huang [35] and Rootzen and Tajvidi [29] are some examples of such applications. Our interest in this article is also in testing for the existence of a long term trend in the excesses of a time series. The main difference with other works is our use of the concept of stochastic orderings of distribution functions. In Section 2 it is shown that given $k$ GPD distributions $F(\cdot; \xi_j, \beta)$, $(j = 1, ..., k)$, if $\xi_1 < \xi_2 < \cdots < \xi_k$, then $F(x; \xi_1, \beta) > F(x; \xi_2, \beta) > \cdots > F(x; \xi_k, \beta)$ for all $x$. That is, we give a sufficient condition for the GPD family to be stochastically ordered. This condition is used in Section 3 to develop a simple procedure based on a likelihood ratio statistic for testing $H_0 \colon \xi_1 = \xi_2 = \cdots = \xi_k$ vs. the isotonic alternative $H_a \colon \xi_1 \le \xi_2 \le \cdots \le \xi_k$. Our procedure is desirable when it is believed a priori that the GPDs satisfy the stochastic order restriction and, hence, it is desirable to have a test that is more powerful than an omnibus test.

The test being proposed here belongs to the field of restricted inference. There is a vast literature in this area. The literature consists of roughly two large subareas: shaped-restricted inference, and order-restricted inference. Barlow *et al.* [2] is a classic pioneering work based on isotonic regression ideas and the Pool-Adjacent-Violators-Algorithm. Robertson *et al.* [25] and the many references therein, summarize and extend the work of Barlow *et al.* and adopt the Nonparametric Maximum Likelihood paradigm proposed by Kiefer and Wolfowitz [14]. Kiefer and Wolfowitz [15] seem to have pioneered the area of shape-restricted inference. Wang [36, 37, 38], extended ideas of Kiefer and Wolfowitz to the estimation of distribution functions under the restriction of being star-shaped or being Increasing Failure Rate on Average. Lo [19], Rojo [26, 27], and Rojo and Ma [28], provide nonparametric estimators for distribution functions that are stochastically ordered. One recent monograph that examines shape-restricted inference is Groeneboom and Jongbloed [11]. Marshall and Olkin [20] and Shaked and Shanthikumar [30] provide excellent treatises on the topic of partial orders of distribution functions.

Finally, in Section 4, we apply our procedure to test for the existence of a monotonic trend in the size of the excesses of a time series of ozone measurements.

## 2.    STOCHASTIC ORDERING OF THE GPD

The concept of stochastic order permeates the theory and applications of statistics. The concept was introduced in the seminal paper by Lehmann [17] and was used to study the power properties of certain tests.

**Definition 2.1.** Let $X$ and $Y$ be random variables such that

$$P(X > x) \leq P(Y > x), \quad -\infty < x < \infty.$$

Then $X$ is said to be *stochastically smaller* than $Y$. This is denoted by $X <^{st} Y$.

We can also state that $Y$ is stochastically larger than $X$ and write $Y >^{st} X$. If $F$ and $G$ represent the cumulative distribution functions (*cdfs*) of $X$ and $Y$ respectively, then $X <^{st} Y$ if and only of $F(x) \geq G(x)$ for all $x \in \mathbb{R}$, and then we write $F <^{st} G$. As discussed by Lehmann [17], a convenient situation arises when the stochastic order is induced by the parameter as it varies monotonically in the parameter space. That is, a parametric family of cdfs $\{F(x; \theta) : \theta \in \Theta \subset \mathbb{R}\}$ is stochastically increasing in $\theta$ if $\theta_1 < \theta_2$ implies that $F(\cdot; \theta_1) <^{st} F(\cdot; \theta_2)$. Similarly, $\{F(x; \theta) : \theta \in \Theta \subset \mathbb{R}\}$ is stochastically decreasing in $\theta$ if $\theta_1 < \theta_2$ implies that $F(\cdot; \theta_2) <^{st} F(\cdot; \theta_1)$. Lehmann and Rojo [18] provided simple characterizations of this and other related orders.

Sufficient conditions are provided here for the family of GPD distribution functions $\mathcal{F} = \{F(x; \xi, \beta) : -\infty < \xi < \infty, \ \beta > 0\}$, to be stochastically ordered. Since $\beta$ is a scale parameter it is clear that the family $\mathcal{F}$ is stochastically ordered in $\beta$ for fixed $\xi$. The following Proposition states that the family $\mathcal{F}$ is stochastically decreasing in $\xi$ for fixed $\beta$.

**Proposition 2.1.** Let $F_1, F_2 \in \mathcal{F}$ with shape parameters $\xi_1$ and $\xi_2$, respectively and equal scale $\beta$. If $\xi_1 < \xi_2$ then $F_2 <^{st} F_1$.

**Proof:** The proof of Proposition 2.1 uses the following result.

**Proposition 2.2** (Mitrinovic [21], pp. 266, inequality 3.6.1). *If $a > 0$ and $x > 0$, then*

$$(2.1) \qquad\qquad e^{-x} \leq \left(\frac{a}{e\,x}\right)^a.$$

Setting $x = 1/u$ and $a = 1$ in (2.1) we obtain

$$(2.2) \qquad\qquad u \geq e^{1-1/u}, \qquad u > 0.$$

Now we prove Proposition 2.1. Let $F(\cdot; \xi, \beta) \in \mathcal{F}$ for $\beta$ fixed. From the definition of the usual stochastic order, it is enough to show that $F(\cdot; \xi, \beta)$ is an increasing function of the parameter $\xi \in \mathbb{R}$. This is true if and only if

$$h(\xi) = \log\big[1 - F(x; \xi, \beta)\big] = (1/\xi) \log(1 - \xi x/\beta)$$

is a decreasing function. First we analyze the case $\xi \neq 0$, for which the problem reduces to showing that

$$(2.3) \qquad h'(\xi) = -(1/\xi^2) \log(1 - \xi x/\beta) - \frac{x}{\xi\beta(1 - \xi x/\beta)} < 0.$$

Making the change of variable $u = 1 - \xi x/\beta$ we get $h'(\xi) = h'\big(\beta(1-u)/x\big) = g(u)$, where

$$g(u) = -\big[x/\beta(1-u)\big]^2 \big(\log u + (1/u) - 1\big),$$

for $0 < u < 1$ when $\xi > 0$, and $1 < u < \infty$ when $\xi < 0$. Then, $g(u) < 0$ if and only if $\log u + (1/u) - 1 > 0$, if and only if $u > e^{1-1/u}$, $u > 0$. But this is the strict inequality in (2.2). Hence (2.3) holds and therefore $F(x; \xi, \beta)$ is increasing in $\xi$ for $\xi \in \mathbb{R}\backslash\{0\}$. Now

$$\lim_{\xi \to 0} \{1 - (1 - \xi x/\beta)^{1/\xi}\} = 11 - e^{-x/\beta}.$$

This means that $F(x; \xi, \beta) \uparrow F(x; \xi = 0, \beta)$ as $\xi \uparrow 0$, and $F(x; \xi, \beta) \downarrow F(x; \xi = 0, \beta)$ as $\xi \downarrow 0$. Then, from the proved monotonicity of $F(x; \xi, \beta)$ in $\xi \in \mathbb{R}\backslash\{0\}$, the proposition follows.   □

Thus, the following result is obtained.

**Corollary 2.1.** *Let $F(x; \xi; \beta)$ denote the GPD distribution with scale parameter $\beta$ and shape parameter $\xi$ as defined by (1.2). Then,*

*If $\xi^* = \xi$ and $\beta < \beta^*$, $F(\cdot; \xi^*, \beta^*) \geq^{st} F(\cdot; \xi, \beta)$.*

*If $\xi > \xi^*$ and $\beta = \beta^*$, $F(\cdot; \xi^*, \beta^*) \geq^{st} F(\cdot; \xi, \beta)$.*

When $\xi > -1$, the expected value $\mu$ of a GPD$(\xi, \beta)$ is $\mu = \beta(1 + \xi)^{-1}$. Then $\xi = \xi(\mu) = (\beta/\mu) - 1$. Thus the shape parameter $\xi$ is a decreasing function of the mean $\mu$. So, if $X_1 \sim$ GPD$(\xi_1, \beta)$ and $X_2 \sim$ GPD$(\xi_2, \beta)$, with $\xi_1, \xi_2 > -1$, and we assume that the means $\mu_j = EX_j$ ($j = 1, 2$) are such that $\mu_2 \leq \mu_1$, then $\xi_1 < \xi_2$. Thus if $\mu_2 \leq \mu_1$ then $X_2 <^{st} X_1$. The converse is also true. To see this, let $F_j$ be the cdf of $X_j$ and assume $X_2 <^{st} X_1$, then we have $1 - F_2(x) \leq 1 - F_1(x)$ for all $x$, and since the GPD only takes positive values, it follows that

$$\mu_2 = \int_0^\infty [1 - F_2(x)]\,dx \leq \int_0^\infty [1 - F_1(x)]\,dx = \mu_1.$$

We can put together all these results in the following corollary.

**Corollary 2.2.** *Let $X_j \sim GPD(\xi_j, \beta)$, (or if $X_j \sim GPD(\xi, \beta_j)$), ($j = 1, ..., k$). Suppose that $E(X_j) = \mu_j$ exists for all $j$. Then the following propositions are equivalent.*

**a)** $X_1 {}^{st}> X_2 {}^{st}> \cdots {}^{st}> X_k$.

**b)** $\xi_1 < \xi_2 < \cdots < \xi_k$, $(\beta_1 > \beta_2 > \cdots > \beta_k)$.

**c)** $\mu_1 \geq \mu_2 \geq \cdots \geq \mu_k$.

---

## 3.   TESTING FOR A LINEAR TREND IN THE EXCESSES

---

Let $X_j \sim$ GPD$(\xi_j, \beta)$, ($j = 1, ..., k$), and denote equality in distribution by $\overset{\mathcal{D}}{=}$. Suppose we want to test the null hypothesis

$$H_0: \ X_1 \overset{\mathcal{D}}{=} X_2 \overset{\mathcal{D}}{=} \cdots \overset{\mathcal{D}}{=} X_k$$

vs. the alternative

$$H_a: \ X_1 >^{st} X_2 >^{st} \cdots >^{st} X_k.$$

From Corollary 2.2, we see that this would be equivalent to testing the null hypothesis

$$H_0: \quad \xi_1 = \xi_2 = \cdots = \xi_k$$

*vs.* the alternative hypothesis

$$H_a: \quad \xi_1 < \xi_2 < \cdots < \xi_k \,.$$

Similarly, the hypothesis $H_a: X_1 <^{\text{st}} X_2 <^{\text{st}} \cdots <^{\text{st}} X_k$ can be tested by using $H_a: \xi_1 > \xi_2 > \cdots > \xi_k$. From Corollary 2.2, a test for the stochastic order could also be based on the means of the GPD's. However the means do not always exist. Therefore we test the hypothesis of stochastic order on the basis of the shape parameter. Assume that for each $X_j$ we have a random sample of size $n_j$, $x_j = (x_{1j}, ..., x_{n_jj})'$ and let $x = (x_1, x_2, ..., x_k)$ be the full data vector. Furthermore, assume that we observe the $X_j$'s sequentially along time, and let $t_j$ be the epoch at which the random sample $x_j$ was observed. To detect a linear time trend, we introduce a third parameter $\theta$ by writing $\xi_j = \xi + \theta t_j$, $(j = 1, ..., k)$. When the $t_j$'s are equally spaced, $t_j$ can be set as $t_j = j$. Thus, we can test the hypothesis of order restriction by testing

$$H_0: \quad \theta = 0$$

*vs.* the alternative hypothesis

(3.1) $$H_a: \quad \theta \neq 0 \,.$$

Although other forms of monotonic trends could occur, e.g. $\xi_j = \xi \exp(\theta t_j)$, a test without assuming a particular form of the monotone trend would require a semiparametric model that would provide protection against misspecification of the functional form of the trend but would not perform as well as the current test for the specific alternative of a monotonic linear trend.

Modeling the parameters of the GPD in order to assess a trend is similar to the approach described in other works such as those by Smith [34], Smith and Huang [35] and Rootzen and Tajvidi [29]. For instance, Rootzen and Tajvidi model the scale parameter as $\beta = \exp(\alpha_0 + \alpha_1 t)$ where $t$ is time in years, and keep the shape parameter $\xi$ constant. In this work we reverse this procedure.

Let $\underset{\sim}{X}$ represent the data vector $X_1, X_2, ..., X_n$. For testing the hypothesis (3.1), we use the Likelihood Ratio Test (LRT) based on $\lambda(\underset{\sim}{X}) = L(\hat{\xi}, \hat{\beta})/L(\hat{\xi}, \hat{\theta}, \hat{\beta})$, where $L$ denotes the likelihood function and the estimators are maximum likelihood estimators (MLE). Then $-2 \log \lambda(\underset{\sim}{X})$ follows asymptotically a chi-square distribution with one degree of freedom. The detailed expression for $-2 \log \lambda(\underset{\sim}{X})$ is given in the Appendix.

## 4.    AN APPLICATION TO OZONE DATA

The data we analyze was collected in Yosemite National Park Wanona Valley and consists of hourly measurements of ozone (ppm) taken from April 1, 1987 to October 31, 1996. The time series contains 84,011 observations with 9412 missing values. The main concern is the detection of a long term trend in the extremal behavior of the time series.

More precisely, the problem is to detect either a decreasing or increasing trend in the size of the excesses over a certain high threshold, if in fact a trend exists. Table 1 displays the monthly number of exceedances over 0.08 ppm. The observations have a strong seasonal component with two periods: the exceedances period which extends from the month of April trough the month of October and the no-exceedances period in the remaining months. The frequency of exceedances increases in the summer months and then decreases in the fall months. Moreover, exploring the data we found that the ozone levels also tend to increase in the summer months and decrease in the fall months. Since the interest lies on the extremal behavior of the data, the analysis was based on the months from April to October.

**Table 1**:    Monthly exceedances over 8 ppm.

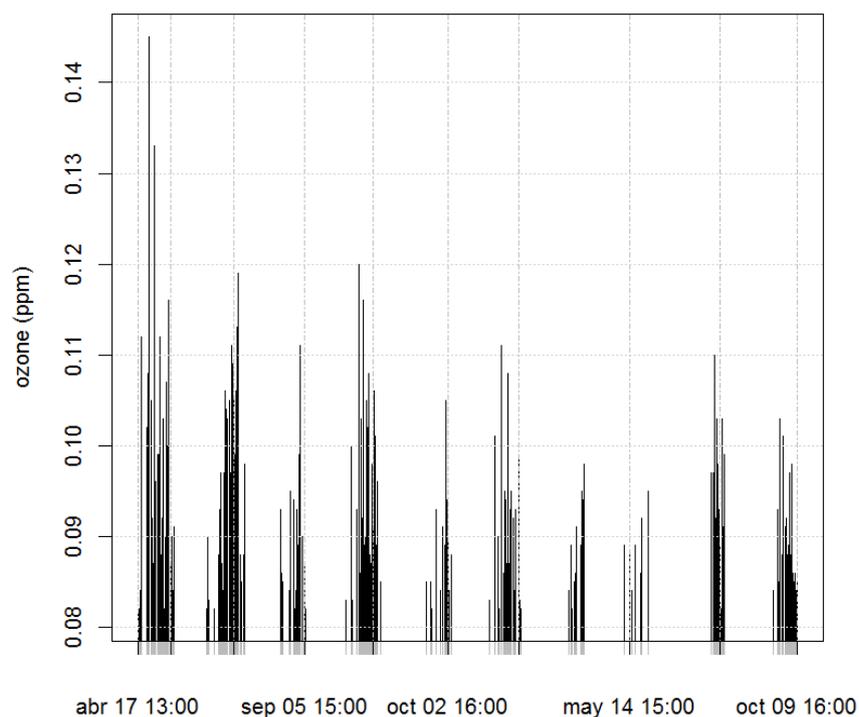| Year | Apr | May | Jun | Jul | Aug | Sep | Oct | Total ($N_u$) | $n$ |
|------|-----|-----|-----|-----|-----|-----|-----|---------------|-----|
| 1987 | 4 | 14 | 70 | 55 | 75 | 50 | 23 | 291 | 4742 |
| 1988 | 9 | 2 | 11 | 83 | 71 | 92 | 21 | 289 | 4856 |
| 1989 | 0 | 6 | 9 | 32 | 29 | 7 | 0 | 83 | 4913 |
| 1990 | 1 | 8 | 34 | 91 | 65 | 63 | 3 | 265 | 4630 |
| 1991 | 0 | 0 | 2 | 19 | 1 | 38 | 17 | 77 | 4463 |
| 1992 | 0 | 2 | 14 | 27 | 49 | 21 | 11 | 124 | 4736 |
| 1993 | 0 | 0 | 3 | 20 | 11 | 21 | 0 | 55 | 3860 |
| 1994 | 6 | 6 | 3 | 14 | 3 | 0 | 0 | 32 | 4720 |
| 1995 | 0 | 0 | 0 | 6 | 50 | 27 | 0 | 83 | 4804 |
| 1996 | 0 | 0 | 4 | 39 | 29 | 22 | 2 | 96 | 4636 |
| Total | 20 | 38 | 150 | 386 | 383 | 341 | 77 | 1395 | 46360 |



**Figure 1**:    Excesses over 0.08 ppm.

Figure 1 shows the empirical marked point processes of exceedances over 0.08 ppm. A clear decreasing trend in the size of the excesses appears. We assess the significance of this trend using the LRT from Section 3.

The LRT requires the excesses to be independent of one another. There is, however, a strong dependence between the exceedances because they tend to occur in clusters. That is, an exceedance tends to attract other exceedances. Several procedures to deal with dependent data have been proposed. One such procedure is to identify clusters of exceedances for which it can be assumed that the excesses within any cluster are independent of the excesses within any other cluster, and then select the maximum excesses within each cluster.

The practical problem with this approach is the identification of independent clusters. Two methods have been used. One is to select a time length $b$ (called block length) and then partition all the observations into consecutive blocks of length $b$. Then consider all the exceedances within a block as a cluster of exceedances. These are called block-clusters. See Leadbetter [16] for the formal justification of this approach as well as for some applications.

The second approach is to select a positive integer $r$ (called the run length) and then decide that any run of at least $r$ consecutive observations below the threshold separates two clusters of exceedances and then assume that such clusters are independent. These are called run-clusters. See Smith [33] for an application of this approach. In this work we use the run-cluster approach with 72 hours (three days) separation. This window of 72 hours is the common practice when analyzing ozone data. Once we have identified the run clusters, we take the maximum excess within each cluster. To distinguish from the *Exceedances over a Threshold* we call these values the *Peaks over a Threshold*, (POT's). Table 2 shows the POT's that we analyze in this work.

**Table 2**:   POT, run-clusters, 72 hours.

| Year | Peaks | | | | | | | | | |
|------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1987 | 0.002 | 0.004 | 0.032 | 0.022 | 0.065 | 0.025 | 0.053 | 0.032 | 0.027 | 0.036 |
|      | 0.010 | 0.011 | | | | | | | | |
| 1988 | 0.002 | 0.010 | 0.002 | 0.013 | 0.017 | 0.007 | 0.017 | 0.026 | 0.023 | 0.025 |
|      | 0.031 | 0.039 | 0.008 | 0.018 | | | | | | |
| 1989 | 0.013 | 0.005 | 0.015 | 0.014 | 0.004 | 0.013 | 0.009 | 0.019 | 0.031 | 0.010 |
|      | 0.007 | 0.002 | | | | | | | | |
| 1990 | 0.003 | 0.020 | 0.003 | 0.013 | 0.040 | 0.036 | 0.007 | 0.018 | 0.026 | 0.016 |
|      | 0.005 | | | | | | | | | |
| 1991 | 0.005 | 0.005 | 0.002 | 0.013 | 0.004 | 0.011 | 0.007 | 0.025 | 0.010 | 0.008 |
| 1992 | 0.003 | 0.021 | 0.010 | 0.031 | 0.006 | 0.015 | 0.007 | 0.028 | 0.012 | 0.013 |
|      | 0.019 | 0.003 | 0.002 | | | | | | | |
| 1993 | 0.004 | 0.009 | 0.011 | 0.009 | 0.015 | 0.018 | | | | |
| 1994 | 0.009 | 0.008 | 0.004 | 0.009 | 0.006 | 0.012 | 0.015 | | | |
| 1995 | 0.017 | 0.017 | 0.030 | 0.023 | 0.018 | 0.009 | 0.023 | 0.011 | 0.019 | |
| 1996 | 0.004 | 0.013 | 0.023 | 0.021 | 0.012 | 0.009 | 0.017 | 0.018 | 0.006 | 0.006 |
|      | 0.004 | 0.005 | | | | | | | | |

Figure 2 shows the POT's for all the years of the observation period. The decreasing trend in the POT's is evident. Under $H_0$ the estimates of the parameters are $\hat{\xi} = 0.2121$ and

$\hat{\beta} = 0.0179$. Under $H_a$ we have $\hat{\xi} = 0.164$, $\hat{\theta} = 0.0575$, and $\hat{\beta} = 0.0209$. The positive value of the estimate of $\xi$ is consistent with the observed decrease in the excesses of the ozone levels. The observed value of the LRT is $-2 \log \lambda(x) = 17.24$ which has a $p$-value of 0.000033. Thus we conclude that the observed decrease in the size of the excesses from 1987 to 1996 is statistically significant.
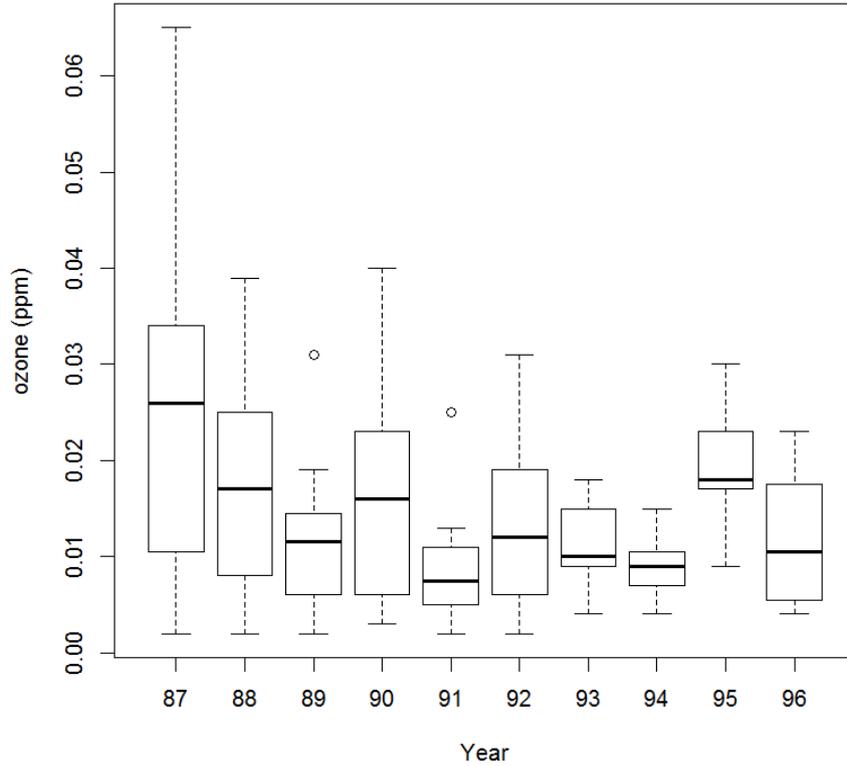


**Figure 2**: Maximum excesses within run-clusters grouped by years.

Once we have found statistical evidence for the decreasing trend in the excesses, we estimate the upper tail of the ozone levels as in Davison and Smith [7] or Embrechts *et al.* [8]. From (1.1) one gets

$$1 - F(u + x) = \gamma_u \big[ 1 - F_u(x) \big],$$

where $\gamma_u = \Pr(X > u) = 1 - F(u)$. Thus, if $N_u$ is the number of exceedances over $u$ and $n$ is the number of observations, then an estimator of $\gamma_u$ is $\hat{\gamma}_u = N_u/n$, and an estimator of the upper tail of $F_X$ is given by

$$(4.1) \qquad 1 - \hat{F}(u + x) = \hat{\gamma}_u \big[ 1 - \hat{F}_u(x) \big] = \frac{N_u}{n} \left( 1 - \hat{\xi} \frac{x}{\hat{\beta}} \right)^{1/\hat{\xi}}, \qquad x > 0.$$

Estimators of the quantiles of $F$ are obtained by solving $\hat{F}(x_p) = p$ for $x_p$ in (4.1), $0 \leq p \leq 1$. This yields

$$(4.2) \qquad \hat{x}_p = u + \frac{\hat{\beta}}{\hat{\xi}} \left[ 1 - \left( \frac{n(1-p)}{N_u} \right)^{\hat{\xi}} \right].$$

When $\hat{\xi} > 0$ by setting $p = 1$ we obtain the estimator of the right end point $\hat{x}^* = u + \hat{\beta}/\hat{\xi}$.

The ozone levels are not independent. So, to simplify, we assume that within the exceedances period in the year (from April to October) the ozone levels come from a strongly stationary process. Then, from the Ergodic Theorem — see Breiman [3], pp. 118 —, we have that $(1/n)\sum_{i=1}^n 1_{\{X_i>u\}} = N_u/n$ converges almost surely to $1 - F(u)$, where now $F$ is the marginal distribution of the ozone levels. Thus $N_u/n$ may be used as an estimator of $1 - F(u)$, and then we can use (4.2) to estimate the upper tail and high quantiles of the distribution of the ozone levels. Table 3 shows the estimates of the shape parameters and from Table 1 we get the number of observations (ozone measurements) and the number of exceedances per year. With this information we can estimate the extreme quantiles of the ozone levels. For instance, for 1987, we have

$$\hat{x}_p = 0.08 + (0.0209)\left(1 - \left[4742\,(1-p)/291\right]^{0.22}\right)/0.22\,, \qquad 0 \le p \le 1\,.$$

Figure 3 shows the estimated 0.99, 0.999 quantiles as well as the right endpoints of the marginal distribution of the ozone levels. The decreasing trend is evident.

**Table 3**:   Estimated shape parameters.

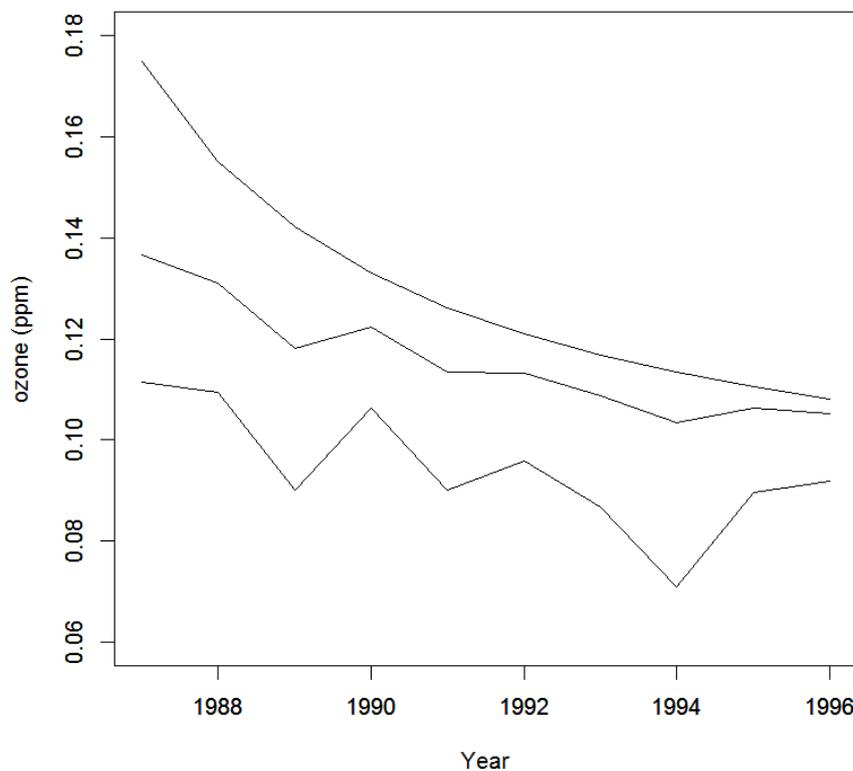| $t_j$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| $\hat{\xi}_j$ | .22 | .278 | .336 | .394 | .452 | .51 | .568 | .626 | .684 | .742 |



**Figure 3**:   Estimated 0.99 and 0.999 quantiles, and estimated right endpoints of the distribution of ozone levels for a threshold of 0.08 ppm.

## 5.    CONCLUSIONS

An exploratory data analysis of the extreme values of a time series of ozone levels made clear the existence of a decreasing linear trend in the size of the excesses over the threshold 8 ppm. We fitted the GPD to the POT's of the time series. By modeling the shape parameter of the GPD as a linear function of time in years, we were able to test the significance of a trend in the size of the excesses. More specifically, consider the years $s$ and $t$ with $s, t = 1987, ..., 1996$. Then we can say that the ozone excesses over 8 ppm for year $s$ were more likely to take larger values then the ozone excesses over 8 ppm for year $t$, when $s < t$.

## A.    APPENDIX – Maximum Likelihood Calculations

The density function of a $\text{GPD}(\xi, \beta)$ is

$$f(x; \xi, \beta) = \begin{cases} (1/\beta)\,(1 - \xi x/\beta)^{(1/\xi)-1}, & \xi \neq 0, \ \beta > 0, \\ \\ (1/\beta)\exp(-x/\beta), & \xi = 0, \ \beta > 0. \end{cases}$$

Let $X_j \sim \text{GPD}(\xi_j, \beta)$, and let $x_j = (x_{1j}, ..., x_{n_j j})'$ be a random sample from $X_j$, $(j = 1, ..., k)$. Write $\xi_j = \xi + \theta t_j$. Then the log-likelihood function under $H_0 \colon \theta = 0$ is

$$(\text{A.1}) \qquad l(\xi, \beta) = \sum_{j=1}^{k} \sum_{i=1}^{n_j} \log f(x_{ij}; \xi, \beta) = -n \log \beta + (\xi^{-1} - 1) \sum_{j=1}^{k} \sum_{i=1}^{n_j} \log(1 - \xi x_{ij}/\beta),$$

where $n = \sum_{j=1}^{k} n_j$, $(\xi, \beta) \in \Theta_0 = \big\{(\xi, \beta) \colon \xi < 0, \ \beta > 0\big\} \cup \big\{(\xi, \beta) \colon \xi > 0, \ \beta > 0, \text{ and } \beta/\xi > \max_{ij}(x_{ij})\big\}$. Making the reparametrization $(\xi, \beta) \mapsto (\xi, \tau)$, where $\tau = \xi/\beta$, the log-likelihood function becomes

$$l(\xi, \tau) = -n \log \xi + n \log \tau + (\xi^{-1} - 1) \sum_{j=1}^{k} \sum_{i=1}^{n_j} \log(1 - \tau x_{ij}),$$

where $\big\{\xi < 0, \ \tau > 0\big\} \cup \big\{0 < \xi \leq 1, \ \tau < 1/\max_{ij}(x_{ij})\big\}$. The log-likelihood equations are

$$(\text{A.2}) \qquad \frac{\partial l}{\partial \xi} = (n/\xi) - (1/\xi^2) \sum_{j=1}^{k} \sum_{i=1}^{n_j} \log(1 - \tau x_{ij}) = 0,$$

$$(\text{A.3}) \qquad \frac{\partial l}{\partial \tau} = (n/\tau) - (\xi^{-1} - 1) \sum_{j=1}^{k} \sum_{i=1}^{n_j} \frac{x_{ij}}{1 - \tau x_{ij}} = 0.$$

Solving equation (A.2) for $\xi$ we obtain

$$(\text{A.4}) \qquad \xi(\tau) = -(1/n) \sum_{j=1}^{k} \sum_{i=1}^{n_j} \log(1 - \tau x_{ij}).$$

Since equation (A.4) gives $\xi$ as an explicit function of $\tau$, we can substitute $\xi(\tau)$ of (A.4) in equation (A.3), and obtain

$$(n/\tau) - \left(\xi(\tau)^{-1} - 1\right) \sum_{j=1}^{k} \sum_{i=1}^{n_j} \frac{x_{ij}}{1 - \tau\,x_{ij}} \;=\; 0\,,$$

which can be solved numerically for $\tau$. If $\hat{\tau}$ is the solution, then the MLE's of $\xi$ and $\beta$ are given by $\hat{\xi} = \xi(\hat{\tau})$ and $\hat{\beta} = \hat{\xi}/\hat{\tau}$, respectively. This is the standard procedure to find the MLE's of the parameters of the GPD. For a detailed analysis of this procedure see Grimshaw [10]. Under $H_a\colon \theta > 0$ the log-likelihood function is

$$\begin{aligned}
l(\xi,\theta,\beta) &= \sum_{j=1}^{k} \sum_{i=1}^{n_j} \log f(x_{ij}; \xi, \theta, \beta) \\
&= -n \log\beta + \sum_{j=1}^{k} \left[(\xi + \theta t_j)^{-1} - 1\right] \sum_{i=1}^{n_j} \log\left[1 - (\xi + \theta t_j)\, x_{ij}/\beta\right],
\end{aligned}$$

where $(\xi,\theta,\beta) \in \Theta_a = \left\{(\xi,\theta,\beta)\colon \xi + \theta t_j < 0,\ j=1,...,k,\ \beta > 0,\ \theta > 0\right\} \cup \left\{(\xi,\theta,\beta)\colon \xi + \theta t_j > 0,\right.$ $j = 1,...,k,\ \beta > 0,\ \theta > 0$ and $\beta/(\xi + \theta t_j) > \max_i(x_{ij}),\ j=1,...,k\right\}$. Let $x_{(n_j)j} = \max_i(x_{ij})$, and note that the restriction $\beta/(\xi + \theta t_j) > x_{(n_j)j}$ is equivalent to $\beta - \xi x_{(n_j)j} - \theta x_{(n_j)j} t_j > 0$. So, the parameter space $\Theta_a \subset \mathbb{R}^3$ is given by all the $(\beta,\xi,\theta)'$ that satisfy the linear pointwise restrictions

$$\begin{bmatrix} 1 & -x_{(n_1)1} & -x_{(n_1)1}t_1 \\ 1 & -x_{(n_2)2} & -x_{(n_2)2}t_2 \\ \vdots & \vdots & \vdots \\ 1 & -x_{(n_k)k} & -x_{(n_k)k}t_k \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \beta \\ \xi \\ \theta \end{bmatrix} > \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix}.$$

Finding the MLE's of $\xi$, $\theta$, and $\beta$ becomes a problem of maximization with linear constraints. There are several numerical algorithms to deal with this type of problem. In this work we used the Price's controlled random search procedure. See Khuri [13], pp. 334–336, for the details of this algorithm. The calculations were performed with R. The test statistic is given by $-2\log\lambda(x) = 2\left[l(\hat{\xi},\hat{\theta},\hat{\beta}) - l(\hat{\xi},\hat{\beta})\right]$.

## REFERENCES

[1] BALKEMA, A.A. and DE HAAN, L. (1974). Residual life at great age, *Annals of Probability*, **2**, 792–804.

[2] BARLOW, R.E.; BARTHOLOMEW, D.J.; BREMNER, J.M. and BRUNK, H.D. (1972). *Statistical Inference Order Restrictions. The Theory and Application of Isotonic Regression*, London–New York–Sidney, Wiley.

[3] BREIMAN, L. (1968). *Probability*, Reading Massachussetts, Addisson Wesley.

[4] CASTILLO, E. and HADI, A.S. (1997). Fitting the generalized Pareto distribution to data, *Journal of the American Statistical Association*, **92**, 1609–1620.

[5] DARGAHI-NOUBARY, G.R. (1989). On tail estimation: an improved method, *Mathematical Geology*, **21**, 829–842.

[6] DAVISON, A.C. (1984). *Modeling excesses over high thresholds, with application.* In "Statistical Extremes and Applications" (J. Tiago de de Oliveira, Ed.), Dordrecht, Reidel Publishing Company.

[7] DAVISON, A.C. and SMITH, R.L. (1990). Models for exceedances over high threshold (with discussion), *J. R. Statistical Society B*, **52**, 393–442.

[8] EMBRECHTS, P.; KLUPPELBERG, C. and MIKOSCH, T. (1997). *Modelling Extremal Events*, Berlin, Springer Verlag.

[9] FALK, M.; HUSLER, J. and REISS, R.D. (1994). *Laws of Small Numbers: Extreme and Rare Events*, Basel, Birkhäuser Verlag.

[10] GRIMSHAW, S.D. (1993). Computing maximum likelihood estimates for the generalized Pareto distribution, *Technometrics*, **35**, 185–191.

[11] GROENEBOOM, P. and JONGBLOED, G. (2014). *Nonparametric Estimation under Shape Constraints*, New York, Cambridge University Press.

[12] HOSKING, J.R.M. and WALLIS, J.R. (1987). Parameter and quantile estimation for the Generalized Pareto distribution, *Technometrics*, **29**, 339–349.

[13] KHURI, A.I. (1993). *Advanced Calculus with Applications in Statistics*, New York, Wiley.

[14] KIEFER, J. and WOLFOWITZ, J. (1956). Consistency of the Maximum Likelihood Estimator in the Presence of Infinitely Many Incidental Parameters, *The Annals of Mathematical Statistics*, **27**, 887–906.

[15] KIEFER, J. and WOLFOWITZ, J. (1976). Asymptotically minimax estimation of concave and convex distribution function, *Z. Wahrsch. verw. Gebiete.*, **34**, 73–85.

[16] LEADBETTER, M.R. (1995). On high level exceedance modeling and tail inference, *Journal of Statistical Planning and Inference*, **45**, 247–260.

[17] LEHMANN, E.L. (1955). Ordered Families of Distributions, *Annals of Mathematical Statistics*, **26**, 399–419.

[18] LEHMANN, E.L. and ROJO, J. (1992). Invariant Directional Orderings, *Annals of Statistics*, **20**, 2100–2110.

[19] LO, S.H. (1987). Estimation of Distribution Functions Under Order Restrictions, *Statistics and Decisions*, **5**, 251–262.

[20] MARSHALL, A.W. and OLKIN, O. (2007). *Life Dsitributions — Structure of Nonparametric, Semiparametric, and Parametric Families*, Springer.

[21] MITRINOVIC, D.S. (1970). *Analytic Inequalities*, New York, Springer Verlag.

[22] PARISI, F. and LUND, R. (2000). Seasonality and Return Periods of Landfalling Atlantic Basin Hurricanes, *Australian New Zeland Journal of Statistics*, **42**(3), 271–282.

[23]  PICKANDS, J. (1975). Statistical inference using extreme order statistics, *Annals of Statistics*, **3**, 119–131.

[24]  REISS, R.D. and THOMAS, M. (1997). *Statistical Analysis of Extreme Values*, Basel, Birkhäuser Verlag.

[25]  ROBERTSON, T.; WRIGHT, F.T. and DYKSTRA, R.L. (1988). *Order Restricted Statistical Inference*, Chichister, Wiley.

[26]  ROJO, J. (1995). On the weak convergence of certain estimators of stochastically ordered survival functions, *Journal of Nonparametric Statistics*, **4**(4), 349–363.

[27]  ROJO, J. (2004). On the estimation of survival functions under stochastic order constraint, *Institute of Mathematical Statistics: Lecture Notes – Monograph Series*, **44**, 37–61, "The First Erich L. Lehmann Symposium – Optimality" (J. Rojo and V. Pérez-Abreu, Eds.).

[28]  ROJO, J. and MA, Z. (1988). On the Estimation of Stochastically Ordered Survival Functions, *Journal of Statistical Computation and Simulation*, **55**, 1–21.

[29]  ROOTZEN, H. and TAJVIDI, N. (1997). Extreme value statistics and wind storm losses: A case study, *Scandinavian Actuarial Journal*, **1**, 70–94.

[30]  SHAKED, M. and SHANTHIKUMAR, J.G. (1994). *Stochastic Orders and Their Applications*, Boston, Academic Press.

[31]  SMITH, R.L. (1984). *Threshold methods for sample extremes*. In "Statistical Extremes and Applications" (J. Tiago de Oliveira, Ed.), Dordrecht, Reidel Publishing Company.

[32]  SMITH, R.L. (1987). Estimating tails of probability distributions *Annals of Statistics*, **15**, 1174–1207.

[33]  SMITH, R.L. (1989). Extreme value analysis of environmental time series: an application to trend detection in ground level ozone, *Statistical Science*, **4**, 367–393.

[34]  SMITH, R.L. (1990). *Extreme value theory*, Chapter 14 of "Handbook of Applicable Mathematics Supplement" (W. Ledermann, E. Lloyd, S. Vajda and C. Alexander, Eds.), John Wiley, Chichester, 437–472.

[35]  SMITH, R.L. and HUANG, L.-S. (1993). *Modeling high threshold exceedances of urban ozone*, Technical Report # 6, National Institute of Statistical Sciences, Research Triangle Park.

[36]  WANG, J.L. (1987a). Estimating IFRA and NBU survival curves based on censored data, *Scandinavian Journal of Statistics*, **14**, 199–210.

[37]  WANG, J.L. (1987b). Estimators of a distribution function with increasing failure rate average, *Journal of Statistical Planning and Inference*, **16**, 415–427.

[38]  WANG, J.L. (1988). Optimal estimation of a star shaped distribution function, *Statisticcs and Decisions*, **6**, 21–32.

# ROBUST ESTIMATION OF REDUCED RANK MODELS TO LARGE SPATIAL DATASETS

Authors: Casey M. Jelsema
– Department of Biostatistics, West Virginia University,
Morgantown, West Virginia, USA
jelsema.casey@gmail.com

Rajib Paul
– Department of Public Health Sciences, University of North Carolina – Charlotte,
Charlotte, North Carolina, USA
Rajib.Paul@uncc.edu

Joseph W. McKean
– Department of Statistics, Western Michigan University,
Kalamazoo, Michigan, USA
joseph.mckean@wmich.edu

Abstract:

• For large datasets, spatial covariances are often modeled using basis functions and covariance of a reduced dimensional latent spatial process. For skewed data, likelihood based approaches with Gaussian assumption may not lead to faithful inference. Any $L_2$ norm based estimation is susceptible to long tails and outliers due to contamination. Our method is based on an empirical binned covariance matrix using the median absolute deviation and minimizes $L_1$ norm between empirical covariance and the model covariance. The consistency of the proposed estimate is established theoretically. The improvement is demonstrated using simulated data and cloud data obtained from NASA's Terra satellite.

## 1. INTRODUCTION

Analysis of geostatistical data is known to be computationally intense or infeasible when the number of observed locations, $n$, is large. This is due to the size of the covariance matrix, $\boldsymbol{\Sigma}$ (which is $n \times n$) and the computational demand of inverting or factoring it. Cressie and Johannesson [4] introduced Fixed Rank Kriging (FRK) to address the computational hurdle by modeling the spatial covariance through a fixed number of deterministic basis functions and a latent reduced rank spatial process. To introduce the parameters, we consider an observed spatial process $Z(\mathbf{s})$ to be made up of a hidden spatial process $Y(\mathbf{s})$ along with a white noise process $\varepsilon(\mathbf{s})$ which could represent, for example, measurement errors. So we write

$$(1.1) \qquad\qquad Z(\mathbf{s}) = Y(\mathbf{s}) + \varepsilon(\mathbf{s}) \,.$$

Typically $Y(\mathbf{s})$ and $\varepsilon(\mathbf{s})$ are assumed to be independent Gaussian distributions, with $\varepsilon(\mathbf{s})$ having mean of zero. In this work however we develop methods that are robust to departure from this assumption. Then, for $n$ observed locations, $Z(\mathbf{s}) \equiv \big\{Z(\mathbf{s}_1), ..., Z(\mathbf{s}_n)\big\}$ is an $n$-dimensional process with mean $E(Y(\mathbf{s})) = \boldsymbol{\mu}_Y$ and covariance matrix expressed as $\boldsymbol{\Sigma}_Z = \boldsymbol{\Sigma}_Y + \sigma^2 \mathbf{I}_n$, where $\boldsymbol{\Sigma}_Y$ is the covariance matrix of $Y(\mathbf{s}) \equiv \big\{Y(\mathbf{s}_1), ..., Y(\mathbf{s}_n)\big\}$ and $\mathbf{I}_n$ is the identity matrix of rank $n$. We then model $Y(\mathbf{s})$ using a mixed effects model such as

$$(1.2) \qquad\qquad Y(\mathbf{s}) = \mathbf{X}(\mathbf{s})\boldsymbol{\beta} + \mathbf{S}(\mathbf{s})\boldsymbol{\eta} + \delta(\mathbf{s}) \,.$$

In this model $\mathbf{X}(\mathbf{s})$ is a matrix of known covariates and $\boldsymbol{\beta}$ is the associated vector of regression coefficients; $\mathbf{S}(\mathbf{s})$ is a sparse $n \times r$ matrix of fixed, spatially varying basis functions which are centered at a set of $r$ knot locations. Dimension reduction is achieved by selecting $r \ll n$. Various classes of basis functions may be used, including wavelets (Shi and Cressie [18] and Zhu *et al.* [22]) and bisquare (Cressie and Johannesson [4] and Paul *et al.* [16]) functions. The latent process $\boldsymbol{\eta}$ is a zero-mean $r$-dimensional Gaussian process defined over the knot locations, with covariance matrix $\mathbf{V}$. Finally $\delta(\mathbf{s})$, the process error, is an *iid* zero-mean Gaussian process with variance $\tau^2$ which takes into account the variations unexplained by the large scale variations $\mathbf{X}(\mathbf{s})\boldsymbol{\beta}$ and spatial process $\mathbf{S}(\mathbf{s})\boldsymbol{\eta}$, and uncertainties arising from the dimension reduction. The process and measurement errors are usually assumed to be independent. When there is only one observation at each spatial location, $\tau^2$ and $\sigma^2$ are non-identifiable, instead their sum $\nu^2 = \sigma^2 + \tau^2$, called the nugget variance, is estimated (though indirect means exist to estimate these separately, see Katzfuss and Cressie [11]). Going forward, we suppress the dependence on $\mathbf{s}$ when possible by stacking scalers into vectors, and vectors into matrices (e.g., $Y(\mathbf{s})$ is replaced with $\mathbf{Y}$ and $\mathbf{X}(\mathbf{s})$ is replaced with $\mathbf{X}$).

With this framework, the covariance matrix $\boldsymbol{\Sigma}_Z$ can be written as $\boldsymbol{\Sigma}_Z = \mathbf{SVS}' + \nu^2 \mathbf{I}_n$. The objective is to estimate the model parameters: $\boldsymbol{\beta}, \mathbf{V}$ and $\nu^2$. Once this has been done one may obtain the inverse of $\boldsymbol{\Sigma}_Z$ easily using the Sherman–Morrison–Woodbury matrix identity. This model offers a large degree of flexibility. The only restriction on $\mathbf{V}$ is the positive-definiteness, hence the resulting covariance matrix may be both anisotropic and nonstationary.

A variety of approaches have been used to model or estimate $\mathbf{V}$. In introducing FRK, Cressie and Johannesson [4] used a Method of Moments (MoM) estimation scheme, while Katzfuss and Cressie [11] developed an expectation-maximization (EM) algorithm. Much attention has also been given to Bayesian hierarchical modeling (see, for example, Banerjee

*et al.* [1], Kang *et al.* [9] and Kang and Cressie [8]). To-date, little attention appears to have been given to robust estimation schemes. Zhu *et al.* [22] developed a method to reduce bias through improved basis function selection, but otherwise did not consider distributional assumptions. Paul *et al.* [16] developed a scale mixture model applicable to non-Gaussian datasets, but like many Bayesian methods it can be time-intensive to implement and run.

The basic FRK model we have described has been elaborated in various ways. For example, to obtain better representation of the spatial dependence some have used a tapering approach (Sang and Huang [17]) or multiple sets of knot locations with different resolutions (Cressie and Johannesson [4] and Kang *et al.* [10]). We demonstrate the latter approach in our data application in Section 5. Both the estimation and fitting stages in the existing MoM estimation use least-squares concepts, and therefore may suffer in the presence of skewed or contaminated data. In the present work we develop an alternative MoM estimator for the parameters of the RRSM. Our motivation in this is to provide an estimator that can model data containing outliers or exhibiting skewness, two features that are frequently encountered in geostatistical datasets, and which does not require significant computational resources.

MoM estimation of the model parameters is divided into two stages: an estimation stage and a fitting stage. In the estimation stage, the entire spatial domain is divided into $M$ bins such that $r < M \ll n$, and $\mathbf{\Sigma}_M$ is defined to be the covariance matrix over the bins. The bins are defined subjectively, though Cressie and Johannesson [4] and Katzfuss and Cressie [11] provide some recommendations. Then an empirical estimate $\hat{\mathbf{\Sigma}}_M$ is constructed using the *detail residuals*, $\mathbf{D} = \mathbf{Z} - \mathbf{X}\hat{\boldsymbol{\beta}}$, where $\hat{\boldsymbol{\beta}}$ is the ordinary least squares estimate of $\boldsymbol{\beta}$. Cressie and Johannesson [4] defined $\hat{\mathbf{\Sigma}}_M$ in the following manner: The $m^{\text{th}}$ diagonal elements $\hat{\mathbf{\Sigma}}_M(m, m) = \text{avg}(\mathbf{D}_m^2)$ and the $(m, m')$ off-diagonal element $\hat{\mathbf{\Sigma}}_M(m, m') = \text{avg}(\mathbf{D}_m) \times \text{avg}(\mathbf{D}_{m'})$. In these expressions, $\mathbf{D}_m$ is the vector of detail residuals in bin $m$, and $\text{avg}(\cdot)$ denotes the average.

Similarly $\mathbf{S}$ is binned into an $M \times r$ matrix by taking the column averages of the rows of $\mathbf{S}$ associated with the observed locations falling into each of the $M$ bins. Denoting this as $\overline{\mathbf{S}}$, one may then write

$$(1.3) \qquad \mathbf{\Sigma}_M = \overline{\mathbf{S}}\mathbf{V}\overline{\mathbf{S}}' + \nu^2 \mathbf{I}_M \,.$$

After estimation, the fitting stage obtains $\hat{\mathbf{V}}$ and $\hat{\nu}^2$ by minimizing the Frobenius norm between $\mathbf{\Sigma}_M$ and $\hat{\mathbf{\Sigma}}_M$, using the $QR$-decomposition on $\overline{\mathbf{S}}$. This is a two-step process resulting in the following estimates:

$$\hat{\nu}^2 \,=\, (\mathbf{F}'\mathbf{F})^{-1}\mathbf{F}'\big(\hat{\mathbf{\Sigma}}_M - \mathbf{Q}\mathbf{Q}'\hat{\mathbf{\Sigma}}_M\mathbf{Q}\mathbf{Q}'\big),$$

$$\hat{\mathbf{V}} \,=\, \mathbf{R}^{-1}\mathbf{Q}'\big(\hat{\mathbf{\Sigma}}_M - \hat{\nu}^2\mathbf{I}_M\big)\mathbf{Q}\mathbf{R}'^{-1},$$

where $\mathbf{F} = \mathbf{I}_M - \mathbf{Q}\mathbf{Q}'$. If $\hat{\mathbf{\Sigma}}_M$ is not positive-definite, the eigenvalues must be lifted to ensure that $\hat{\mathbf{V}}$ is positive-definite (see Kang *et al.* [9]). For further details on Fixed Rank Kriging, see Katzfuss and Cressie [11].

We redesign both the estimation and fitting stages for the MoM estimation using the Median Absolute Deviation and quantile regression (Section 2). Our work is novel in that we return to basic principles to redesign the estimation and fitting stages with a mind for resisting contaminated data. The consistency of our proposed estimate is shown (Section 3),

though the technical details are given in the Appendix. We describe and conduct a simulation study (Section 4) to investigate the performance of our proposed method. Finally, we provide a data example (Section 5) using a large remote sensing dataset and some concluding remarks (Section 6).

## 2.    ROBUST ESTIMATION AND FITTING

In this section we describe robust alternatives to both the estimation stage and fitting stage of MoM estimation for the FRK model. First we define $\hat{\boldsymbol{\Sigma}}_M^{(\mathrm{rob})}$ as an estimate empirical binned covariance matrix which is robust to contamination. Then we describe a robust strategy to fit the model parameters, which we call the *robust fit*. We denote the previous-described methods from Cressie and Johannesson [4] as $\hat{\boldsymbol{\Sigma}}_M^{(\mathrm{CJ})}$ and the Frobenius fit.

### 2.1.  Estimation stage

The diagonal elements of $\boldsymbol{\Sigma}_M$ represent the variance within a bin. We estimate this quantity using the median absolute deviation, $\mathrm{MAD}(X) = \mathrm{med}\big(|X - \mathrm{med}(X)|\big)$. A constant scale factor is applied to the MAD which causes it to be a consistent estimate for the standard deviation (see Hettmansperger and McKean [7], Eqn. 3.9.27). In the present work, we use the usual MAD which is consistent for $\sigma$ when the errors are normally distributed. Hence, the diagonal elements of our proposed estimate are given by

$$(2.1)\qquad\qquad \hat{\boldsymbol{\Sigma}}_M^{(\mathrm{rob})}(m, m) = \mathrm{MAD}^2(\mathbf{D}_m)\,,\qquad m = 1, ..., M\,.$$

Estimating the covariance between two bins is more challenging. First, recall that $\mathrm{cov}(A, B) = \frac{1}{4}\big[V(A + B) - V(A - B)\big]$. Estimating a covariance using this identity requires finding $\mathbf{D}_m \pm \mathbf{D}_{m'}$, however, these quantities are not well-defined. For example, two bins may not even have the same number of observations, much less any natural correspondence between observations. We therefore use the pairwise sums and pairwise differences, denoted using $\oplus$ and $\ominus$ respectively, to approximate $\mathbf{D}_m \pm \mathbf{D}_{m'}$. We again use the square of the MAD to estimate the variance, so the off-diagonal elements of our estimate are given by:

$$(2.2)\qquad \hat{\boldsymbol{\Sigma}}_M^{(\mathrm{rob})}(m, m') = \frac{1}{4}\Big[\mathrm{MAD}^2(\mathbf{D}_m \oplus \mathbf{D}_{m'}) - \mathrm{MAD}^2(\mathbf{D}_m \ominus \mathbf{D}_{m'})\Big]\,.$$

### 2.2.  Fitting stage

Given an empirical covariance matrix $\hat{\boldsymbol{\Sigma}}_M$, we fit $\mathbf{V}$ by minimizing some norm between $\hat{\boldsymbol{\Sigma}}_M$ and $\boldsymbol{\Sigma}_M$. To develop the robust fitting stage, we start from equation (1.3),

$$\hat{\boldsymbol{\Sigma}}_M = \overline{\mathbf{S}}\mathbf{V}\overline{\mathbf{S}}' + \hat{\nu}^2\mathbf{I}_M\,,$$

$$(2.3)\qquad\qquad \big(\hat{\boldsymbol{\Sigma}}_M - \hat{\nu}^2\mathbf{I}_M\big)\,\overline{\mathbf{S}}\big(\overline{\mathbf{S}}'\overline{\mathbf{S}}\big)^{-1} = \overline{\mathbf{S}}\mathbf{V}\,.$$

Then we may see equation (2.3) as a multivariate regression problem with $\overline{\mathbf{S}}$ as the design matrix and $\mathbf{V}$ as the matrix of regression coefficients. Any method of robust regression may then be implemented to obtain an estimate of $\mathbf{V}$. For this work, we use the popular least absolute deviations, $L_1$, estimator; see Koenker and Bassett [13] and Section 3.8 of Hettmansperger and McKean [7]. In comparison to least squares (LS), the least absolute deviation fit is obtained by replacing the squared Euclidean norm with the $L_1$ norm. Hence, the geometry and interpretation of the $L_1$ fit is quite similar to LS fit, but unlike the LS estimate, the $L_1$ estimate is robust. As discussed in Section 3.8 of Hettmansperger and McKean [7], the fit is also efficient. It attains efficiency 0.64 relative to LS for normal errors but is generally more efficient than LS for error distributions with tails heavier than the normal.

Each column of $\left(\hat{\boldsymbol{\Sigma}}_M - \hat{\nu}^2 \mathbf{I}_M\right) \overline{\mathbf{S}}\left(\overline{\mathbf{S}}'\overline{\mathbf{S}}\right)^{-1}$ is used as the response in a separate estimation. There are therefore $r$ estimates to obtain, each of which corresponds to a column of $\mathbf{V}$. As the final estimate $\mathbf{V}$ may not be numerically symmetric, we symmetrize $\hat{\mathbf{V}}$ by taking $\hat{\mathbf{V}} = 0.5\left(\hat{\mathbf{V}} + \hat{\mathbf{V}}'\right)$. We used the `quantreg R` package (Koenker [12]) for the computation of the $L_1$ fit.

Estimation of $\mathbf{V}$ requires an estimate of $\nu^2$. By substituting the left side of (2.3) for $\overline{\mathbf{S}}\mathbf{V}$ in (1.3) we obtain:

$$\hat{\boldsymbol{\Sigma}}_M = \left(\hat{\boldsymbol{\Sigma}}_M - \nu^2 \mathbf{I}_M\right) \overline{\mathbf{S}}\left(\overline{\mathbf{S}}'\overline{\mathbf{S}}\right)^{-1}\overline{\mathbf{S}}' + \nu^2 \mathbf{I}_M\,,$$

(2.4) $$\hat{\boldsymbol{\Sigma}}_M\left(\mathbf{I}_M - \overline{\mathbf{S}}\left(\overline{\mathbf{S}}'\overline{\mathbf{S}}\right)^{-1}\overline{\mathbf{S}}'\right) = \nu^2\left(\mathbf{I}_M - \overline{\mathbf{S}}\left(\overline{\mathbf{S}}'\overline{\mathbf{S}}\right)^{-1}\overline{\mathbf{S}}'\right)\,.$$

We then stack the columns of $\hat{\boldsymbol{\Sigma}}_M\left(\mathbf{I}_M - \overline{\mathbf{S}}\left(\overline{\mathbf{S}}'\overline{\mathbf{S}}\right)^{-1}\overline{\mathbf{S}}'\right)$ and the columns of $\left(\mathbf{I}_M - \overline{\mathbf{S}}\left(\overline{\mathbf{S}}'\overline{\mathbf{S}}\right)^{-1}\overline{\mathbf{S}}'\right)$. Doing this, we again cast the problem as a zero-intercept robust regression, where $\nu^2$ is the slope. This estimate is substituted into equation (2.3) to obtain an estimate of $\mathbf{V}$.

The estimate of $\mathbf{V}$ may not be positive-definite, so we may need to lift the eigenvalues (similar to Cressie and Johannesson [4]), while preserving the total variability. In our work, we compute the sum of the eigenvalues, $\Delta$, and proportionally redistribute this sum across the eigenvalues after shifting all eigenvalues to be non-negative.

---

## 3.    ASYMPTOTIC PROPERTIES

Here we discuss some of the infill asymptotic properties of our proposed estimator, $\hat{\boldsymbol{\Sigma}}_M^{(\text{rob})}$. Infill asymptotics is a common method of considering asymptotics related to geostatistical methodology in which the domain, $\mathcal{D}$, remains fixed but the density of observed locations is increased.

Recall that we obtain $\hat{\mathbf{V}}$ by minimizing some norm $\|\cdot\|$:

$$\hat{\mathbf{V}} = \operatorname{argmin}\left\|\hat{\boldsymbol{\Sigma}}_M - \boldsymbol{\Sigma}_M\right\|\,.$$

Hence, once $\hat{\boldsymbol{\Sigma}}_M$ is known, $\hat{\mathbf{V}}$ is fully determined by the fitting method. Therefore, a desirable property of the empirical binned covariance matrix $\hat{\boldsymbol{\Sigma}}_M^{(\text{rob})}$ is that it be consistent for $\boldsymbol{\Sigma}_M$, which we establish in this section.

There are two sets of assumptions that we need to make. From expressions (2.1) and (2.2), $\hat{\mathbf{\Sigma}}_M^{(\mathrm{rob})}$ is a function of MADs applied to the detail residuals. For each bin $m$, these residuals are obtained from ordinary least-squares regression, our proof requires that $\sqrt{n}(\hat{\beta} - \beta) = O(1)$ for each bin. For this, we assume the conditions in the paper by Lahiri *et al.* [14] for each bin.

Our process for bin $j$ (slightly abusing the notation to avoid double subscript), is $\{e_1, e_2, ..., e_{n_j}\}$ which we denote by $\{\mathbf{e}_j\}$. On this process we assume that

1.  $\{\mathbf{e}_j\}$ is stationary.
2.  $\{\mathbf{e}_j\}$ satisfies the strong mixing coefficients assumption given as follows. For $i \neq k$, let $A_i$ and $B_k$ be in the $\sigma$-fields generated by $e_i$ and $e_k$. Then

$$(3.1) \qquad \left| P[A_i \cap B_k] - P[A_i]P[B_k] \right| = O\big(\rho^{|i-k|}\big),$$

   where $0 \leq \rho < 1$.

Note that Assumption 2 implies that the spatial correlation between two locations exhibits exponential decay. This is a common feature in spatial modeling (e.g. the Matérn class of covariance models), and as such is not an unreasonable assumption.

For our proof, let $\mathbf{D}_m$ denote the random detail residual process within the $m^{\mathrm{th}}$ bin, and let $\mathbf{D}_m = \{\tilde{R}_{m_1}, ..., \tilde{R}_{m_k}\}$ be the $k$ observed detail residuals from that bin. We assume that $\mathbf{D}_m$ and, as will be seen, $|\mathbf{D}_m|$, exhibit strong mixing as described in conditions 1 and 2.

We now state the consistency result in theorem form. The proof is given in the Appendix.

**Theorem 3.1.**   *Under the above conditions, $\hat{\mathbf{\Sigma}}_M^{(\mathrm{rob})}$ is a consistent estimator of $\mathbf{\Sigma}_M$.*

Throughout we treat the number of bins, $M$, as fixed, and do not consider limits over that quantity. This is analogous to the work of Bliznyuk *et al.* [2]. In another context on binned estimation, they considered $m$ (the number of bins) as a radius to determine "adjacency" of locations, where $m$ does not depend on $n$, (the number of observations) and did not limit over $m$. The only restriction on $M$ is that it should be large enough to ensure that the assumption of stationary within bins holds for practical implementation.

## 4.  SIMULATION STUDY

To compare our proposed methods with the existing methods using simulated data, we generate a spatial process $Z$ according to the model:

$$\mathbf{Z} = \mathbf{X}\boldsymbol{\beta} + \mathbf{S}\boldsymbol{\eta} + \boldsymbol{\epsilon}.$$

First we select $n$ locations uniformly over a $100 \times 100$ domain, and $r_o = 1225$ knot locations on a $35 \times 35$ grid. These knot locations are used to simulate the data but not to fit the models (because reduced rank spatial models are designed as approximations of a more complex spatial process). Then we define $\mathbf{X}$ as an $n \times 3$ matrix where the columns correspond respectively to an intercept, the $x$-coordinate, and the $y$-coordinate.

To define $\mathbf{V}$ we first compute the pairwise distances between the knot locations, and generate a Matérn covariance matrix using these distances with sill and range parameters each set to 1, and smoothness set to 0.5. We use `cov.sp` in the R package `SpatialTools` (French [6]) to generate this matrix. We then obtain $\mathbf{V}$ as an observation from the inverse Wishart distribution using the Matérn covariance as a scale matrix and $2(r+1)$ degrees of freedom. In this way the covariance matrix used to simulate the data is not constrained to be either stationary or isotropic.

We construct $\mathbf{S}$ using the bisquare basis functions defined as

$$S_{i,j} = \begin{cases} \left(1 - \left(\|s_i - u_j\|/r_u\right)^2\right)^2 & \text{for } \|s_i - u_j\| \leq r_u\,, \\ 0 & \text{otherwise}\,, \end{cases}$$

where $r_u$ is 1.5 times the minimum distance between knots and $\|\cdot\|$ denotes the measure of distance appropriate to the data (e.g., in our simulations, we used Euclidean distance).

We used two methods to simulate the data, a Contaminated Normal distribution and an Exponential distribution. These simulate the presence of outliers or of skewness, respectively, in the resulting dataset. For either simulation method, we compare the model fits by splitting the simulated data into a training set and a held-out test set. The hold-out set was set as all of the locations in the square bounded by the points $(40, 40)$ and $(60, 60)$, which corresponds to approximately 4% of the observations. We use the estimated parameters to predict at the held-out locations and compute diagnostics to assess both the accuracy and uncertainty of the prediction, including the mean square error (MSE), mean square prediction error (MSPE), and the continuous ranked probability score (CRPS, Wilks [21]), a measure which incorporates both the prediction accuracy and the prediction uncertainty. Lower values are preferable for all of these measures.

## 4.1. Simulation 1: contaminated normal

For simulating datasets we first generate a $r_o$-dimensional process $\boldsymbol{\eta}$ from a zero-mean multivariate normal with covariance $\mathbf{V}$. To induce outliers, the measurement error process $\boldsymbol{\varepsilon}$ is generated from a contaminated normal distribution. We first draw a random sample from $\mathcal{N}(0, \nu^2)$, and then replace $\alpha n$ of the values with random draws from $\mathcal{N}(0, \nu_c^2)$. Finally, we obtain the simulated data by $\mathbf{Z} = \mathbf{X}\boldsymbol{\beta} + \mathbf{S}\boldsymbol{\eta} + \boldsymbol{\varepsilon}$. For each simulated dataset, estimate model parameters using both the method of Cressie and Johannesson [4] and the proposed robust method.

We considered three sample sizes, $n \in \{10000, 15000, 20000\}$ and five levels for the number of knots locations to fit the model, $r \in \{64, 100, 144, 196, 256\}$, intentionally chosen to much less than $r_o$, so that the "true" spatial process was more granular than the model. For the contamination level of $\boldsymbol{\varepsilon}$ we consider $\alpha \in \{0.00, 0.05, 0.10, 0.15, 0.20\}$. For the simulations shown, the values of $\boldsymbol{\beta} = (1, 0.01, 0.05)'$, $\nu^2 = 1$, and $\nu_c^2 = 100$ were held constant. These choices are not sensitive to our estimation technique except insofar as a larger or smaller $\nu_c^2$ would correspond to a larger or smaller effect from the contamination. For each combination of these parameters, we generated 50 replications of data. Hence, there were 75 settings of parameter levels, and 3750 replications in total.

## 4.2.  Simulation 2: exponential

As we have noted throughout, skewness can also be problematic for least-squares type estimators, and skewed data are not uncommon in geostatistics. Hence, we designed a second simulation in which we generate $\varepsilon$ from an Exponential distribution rather than from a contaminated Normal distribution. We use the same design as Simulation 1, but instead of $\alpha$, we consider the rate parameter of the Exponential distribution $\lambda \in \{0.10, 0.25, 0.50, 1.00\}$. Hence, for this simulation there were 60 settings and 3000 replications in total.

## 4.3.  Simulation results

The simulations suggest that the robust method is generally preferable to the CJ method. For brevity we present the results for the CRPS, but results for the MSE and MSPE were similar. We use two main values to compare the results: The median CRPS across the 50 replications, and the CRPS of the CJ method relative to that of the robust method (we refer to this as the CRPS ratio).

Results of Simulation 1 are shown in Figure 1, which plots the median CRPS over the 50 replications for each of the settings. In 67 of the 75 settings, the robust method produced a smaller median CRPS than the CJ method.
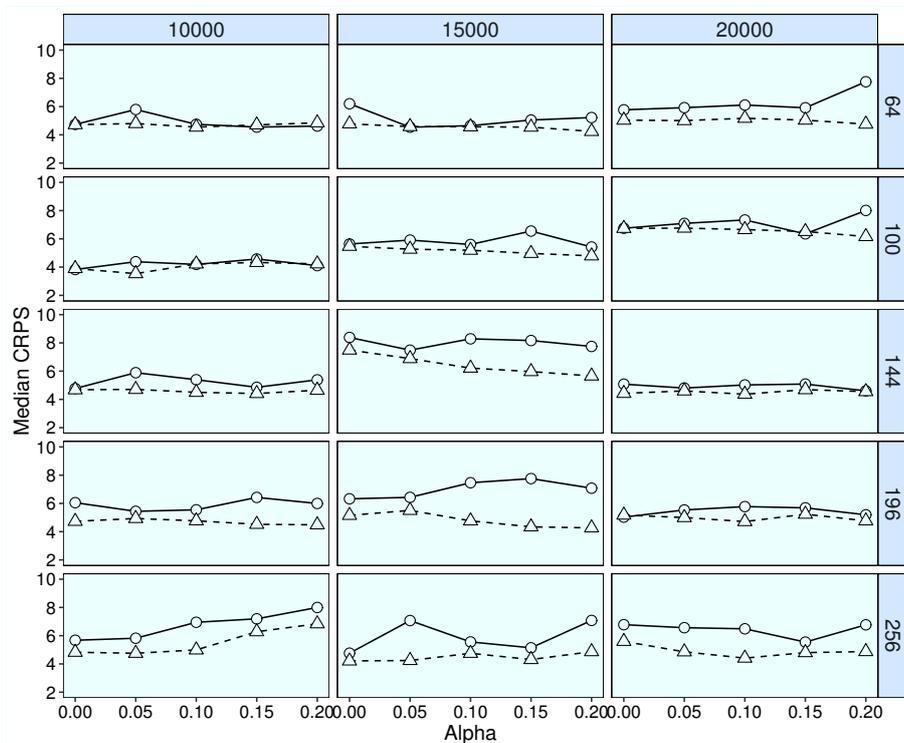


**Figure 1**:  Results for Simulation 1. Plotted points are median CRPS of the CJ method (circles) and the robust method (triangles) over the 50 replications.

In addition, the robust method produced a smaller CRPS (i.e. CRPS ratio greater than 1) in 68.8% of the replications, and the median of the CRPS ratio showed a 9% larger CRPS for the CJ method. When considering the CRPS ratio for each setting, the worst-performing setting for the robust method had a median CPRS ratio of 0.975 (near equivalence), while half of the settings had a median CRPS ratio showing an improvement of 10% or more.

The results for Simulation 2 were similar to those of Simulation 1, and are shown in Figure 2. In 55 of the 75 settings, the robust method produced a smaller median CRPS than the CJ method.
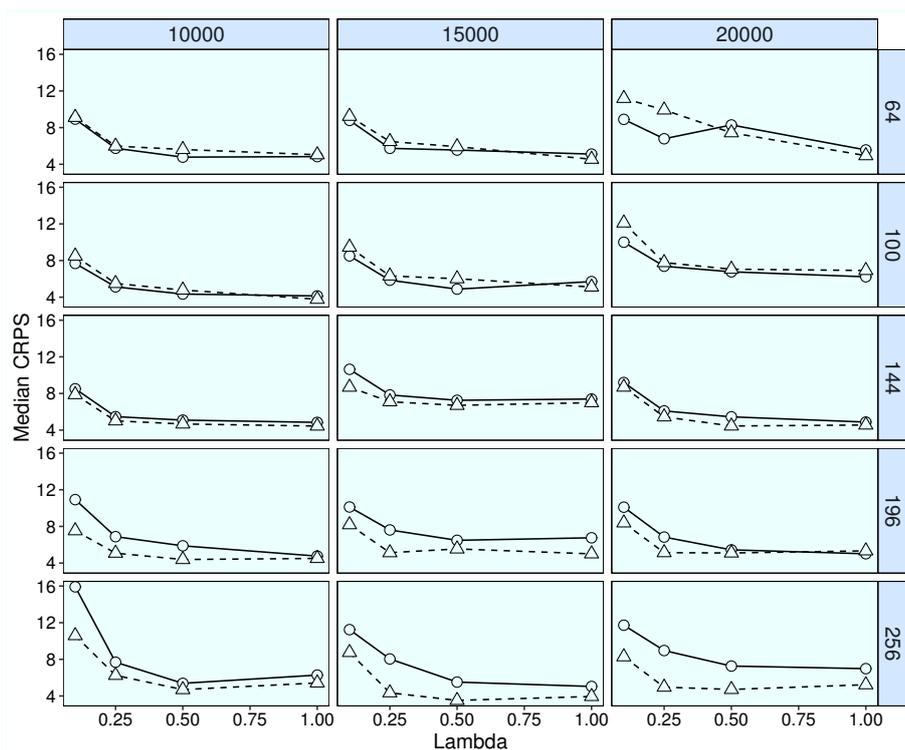


**Figure 2**: Results for Simulation 2. Plotted points are median CRPS of the CJ method (circles) and the robust method (triangles) over the 50 replications.

In addition, the robust method produced a smaller CRPS (i.e. CRPS ratio greater than 1) in 65.3% of the replications, and the median of the CRPS ratio showed an 8% larger CRPS for the CJ method. When considering the CRPS ratio for each setting, the worst-performing setting for the robust method had a median CPRS ratio of 0.957, which again shows minimal advantage for the CJ method, while half of the settings had a median CRPS ratio showing an improvement of at least 7%.

To provide an overall summary of our results, our findings suggest that the proposed robust method tends to be advantageous compared to the CJ method. While we acknowledge this is not uniformly the case, we note that in approximately two-thirds of cases, the proposed method resulted in smaller CRPS. It is unfortunately difficult to discern much in the way of a pattern across the simulation settings, to determine whether the robust or CJ method might be preferable in a specific setting. The main apparent pattern from these simulations is that the more knots, the better the robust method tended to perform against the CJ method.

This could potentially be a consequence of each bin from the estimation of $\boldsymbol{\Sigma}_M$ having fewer observations compared to a setting with the same sample size but smaller number of knots, in which case outliers would have an increased effect.

Since the number of knots is chosen by the modeler, one might be tempted to select a smaller value of $r$, so that any effect from the choice of method is minimized. However, fewer knots corresponds to a more coarse representation of the spatial variation, hence the general recommendation (e.g. Finley *et al.* [5]) is to use as many as possible (within any computational limits). Hence, the natural choice guiding the selection of $r$ will also tend to produce situations in which the robust method appears to perform better.

## 5.    APPLICATION TO NASA DATA

We use remote sensing data on daily cloud liquid water path (CWP), obtained through NASA's Moderate Resolution Imaging Spectroradiometer (MODIS) on the Terra satellite on April 22, 2012. Note that this date is an arbitrary choice, our interest here is to demonstrate our method outside of a fabricated example. Because the dataset is large ($n = 48552$), a reduced rank model is a reasonable choice for inference. The CWP data are right-skewed, so we restrict out focus to the log-scale.

### 5.1.  Original Data Analysis

The observed data are plotted in Figure 3. Due to a north-south trend (tending to smaller values closer to the equator), we model the large-scale variation using Legendre polynomials similar to Stein [19], though using only the latitude. Specifically, let $L$ denote the degrees latitudes and define $\ell = \pi(L/180)$.
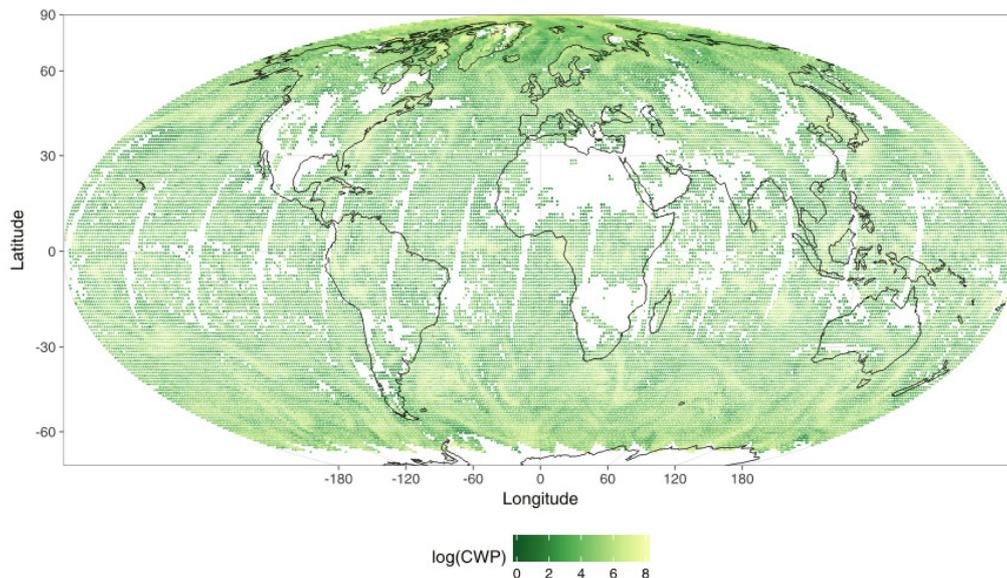


**Figure 3**:  Plot of observed Cloud Water Path over the spatial domain.

We compute Legendre polynomials $P_p^q(\sin(\ell))$ of degree $p = 80$ and order $q = 0, 1, ..., p$. This results in a design matrix consisting of 81 regressors of spherical harmonics. Stein [19] also included a cosine of the longitude. Since we observed primarily a trend over the latitudes, we do not include the cosine term on longitude. Since our focus is on the small-scale (spatial) variation rather than the large-scale variation, the main concern for us is that this model enables stationarity of the spatial process to be reasonable; visual inspection (figure not shown) of the predictions for each latitude show this to be the case.

For the MoM estimation described in the preceding sections we first compute the detailed residuals. The normal quantile-quantile plot of the detailed residuals in Figure 4 shows a heavy lower tail, which motivates the use of the proposed robust techniques. Initially we model the data as observed. Afterwards, we also induce outliers into the data and reanalyze the data.
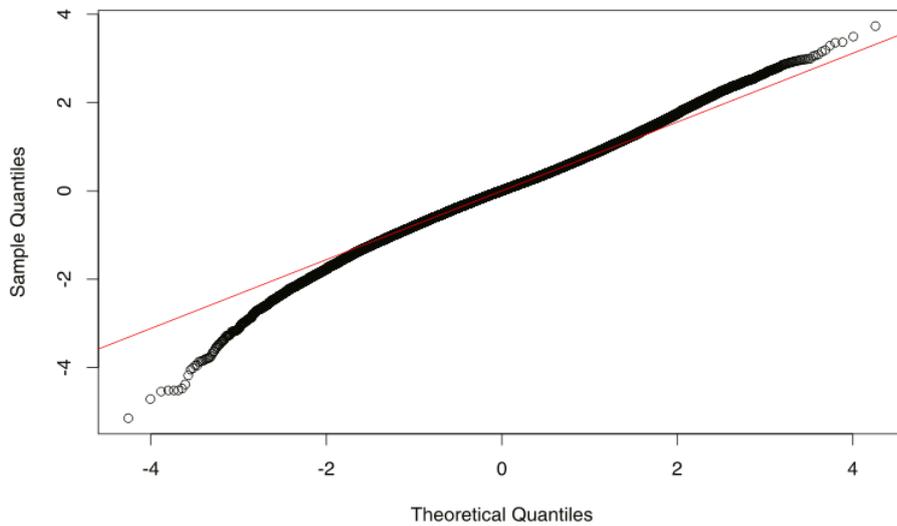


**Figure 4**: Normal quantile-quantile plot of the detailed residuals.

As recommended by Cressie and Johannesson [4], we use a multi-resolution model for CWP (see Nychka *et al.* [15]), to capture multiple scales of variation. We choose $r_1 = 38$ knot locations for the first resolution, and $r_2 = 97$ knot locations for the second resolution. Therefore the estimate of $\mathbf{V}$ is a $135 \times 135$ matrix. A map of these knot locations is given in Figure 5.

To construct the $\mathbf{S}$ matrix, we use the modified bisquare function, defined as:

$$\mathbf{S}_{i,j(l)} = \begin{cases} \left(1 - 0.25\, d^2\big(\mathbf{s}_i, \mathbf{u}_{j(l)}\big)\right) & \text{for } d\big(\mathbf{s}_i, \mathbf{u}_{j(l)}\big) \leq 2\,, \\ \\ 0 & \text{otherwise}\,, \end{cases}$$

where $\mathbf{u}_{j(l)}$ is the $j^{\text{th}}$ knot location of the $l^{\text{th}}$ resolution, $\mathbf{s}_i$ are the observed locations. The distance is given by:

$$d\big(\mathbf{s}_i, \mathbf{u}_{j(l)}\big) = \sqrt{d_{\text{long}}^2\big(\mathbf{s}_i, \mathbf{u}_{j(l)}\big)\big/r_{\text{long}(l)}^2 + d_{\text{lat}}^2\big(\mathbf{s}_i, \mathbf{u}_{j(l)}\big)\big/r_{\text{lat}(l)}^2}\,,$$

where $d_{\text{long}}\big(\mathbf{s}_i, \mathbf{u}_{j(l)}\big)$ and $d_{\text{lat}}\big(\mathbf{s}_i, \mathbf{u}_{j(l)}\big)$ denote the longitude (east-west) and latitude (north-south) distances, respectively, between the location $\mathbf{s}$ and the knot location $\mathbf{u}_{j(l)}$. The values

$r_{\text{long}(l)}$ and $r_{\text{lat}(l)}$ control the maximum distance between an observation and a knot such that there is non-zero weight associated between the two. We set these to be the minimum east-west distance and minimum north-south distance between two knot locations of the same resolution.
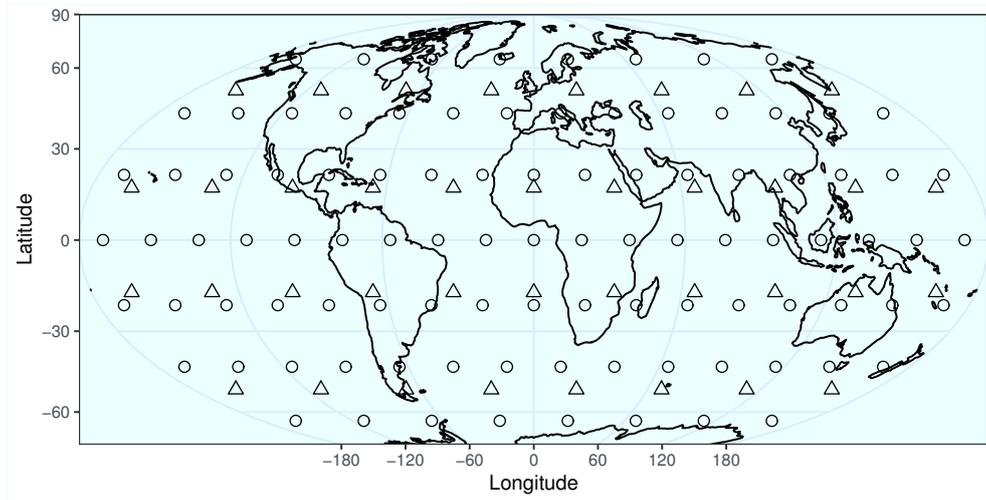


**Figure 5**:  Plot of the knot locations of the basis functions over the spatial domain. Triangles represent the 38 knot locations of the first resolution, and circles represent the 97 knot locations of the second resolution.

Figures of the predictions or prediction uncertainties are not particularly informative, as our focus is on comparing the robust method to the CJ method. The CJ method yielded larger RMSPEs by approximately 20%, and the CRPS tended to be larger as well. A plot of the CRPS ratio for each location is shown in Figure 6. On average, the CRPS ratio is 1.04, indicating better performance for the robust method.
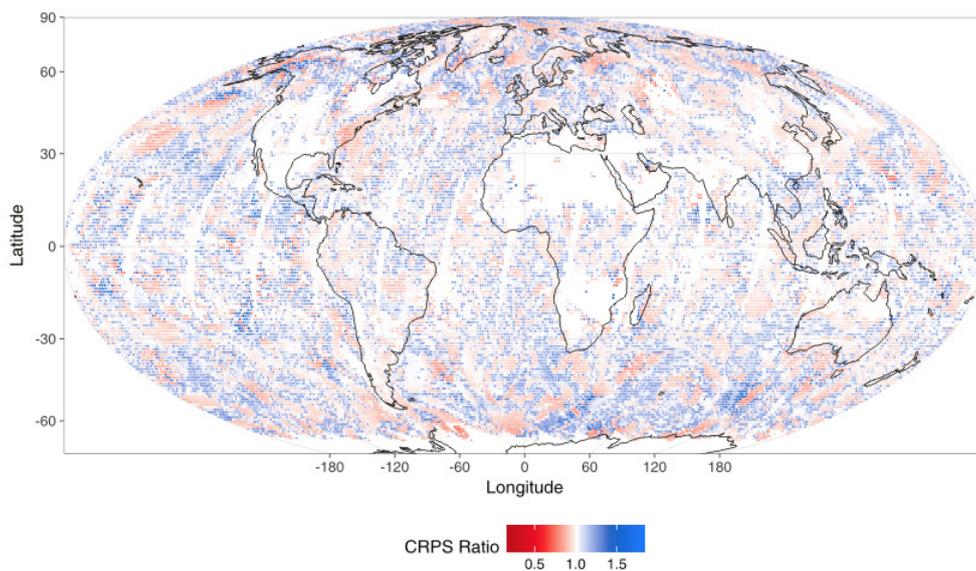


**Figure 6**:  Plot of the CRPS of predictions using the CJ method relative to those using the robust method. Larger values indicate the CJ method produced a larger CRPS at that location.

## 5.2. Analysis after inducing outliers

In addition to this analysis, we artificially contaminated the log CWP data by replacing the 2% of observed values $Z_i(\mathbf{s})$ with $1.5\,Z_i(\mathbf{s})$. Inspection of the normal quantile-quantile plot showed a heavy upper tail which also contained many outliers. The results followed the same pattern as those described above. The RMSPE were again uniformly larger for the CJ method, now averaging 78% larger, while the CRPS were, on average, 11% larger.

## 6.   CONCLUSIONS AND DISCUSSION

The Method of Moments is a flexible and powerful tool for estimating the parameters of a FRK model. Bayesian methods are more accurate than kriging (Kang and Cressie [8]), but they are also more time-consuming, and often come with some distributional assumptions. Kriging is typically a faster process, and kriging estimates are BLUP even in the face of non-normality, so kriging presents benefits of its own. However the typical parameter estimates using EM algorithm or MoM are susceptible to contaminated data. In this work we have provided robust alternatives to both stages of the MoM estimation.

Our results indicate that the proposed estimate and fitting scheme successfully capture the spatial covariance. In both our simulations and in our application to real data, the robust method tended to provide an advantage over the CJ method. At times the advantage was small, but in some cases the robust method showed substantial improvement, even when the data were neither contaminated or skewed.

Besides the $L_1$-fit, other robust fits can be used. For example, the Wilcoxon fit is a robust fit that minimizes the sum of the absolute differences of the residuals (see Hettmansperger and McKean [7], Section 3.8). The Wilcoxon fit is generally more efficient than the $L_1$-fit and it generalizes to fits for skewed-error distributions. We are currently investigating other robust norms which result in fits with higher efficiency than that of the $L_1$ fit for normal errors.

Again we emphasize that the kriging equations have been derived by minimizing the mean square prediction error. These predictions are then simply functions of $\mathbf{V}$ and $\nu$. In our work, we have provided robust methods of estimating these same parameters. Yet when using robust techniques, it may be desirable to derive predictions and measures of precision using a different loss function than the squared error loss, or such that the predictions are robust in addition to the parameter estimates (Cressie and Hawkins [3]). Our robust estimates perform well in spite of this.

## A.    APPENDIX – Proof of Theorem 3.1

The proof utilizes the consistency of a fit $\hat{\beta}$ such that $\sqrt{n}(\hat{\beta} - \beta) = O(1)$; the assumptions as discussed in Section 3, including $n_j \to \infty$, for $j = 1, ... M$; and the theory for the sign processes as discussed in Chapters 1 and 3 of Hettmansperger and McKean [7]. For the sign process theory, we assume that the pdf of the random errors is positive at its median. The proof is in two parts. Part 2 gives the desired result, while Part 1 establishes the consistency of the medians used in the second part.

**Part 1 of the Proof:**

Consider the $j$-th bin, for $j = 1, ... M$. Let $\{\mathbf{e}_j\}$ denote the process of random errors of the linear model $\mathbf{Z}_j = \mathbf{X}_j \boldsymbol{\beta}_j + \mathbf{e}_j$. Assume without loss of generality that $\boldsymbol{\beta}_j = \mathbf{0}$ and the median of $e_i$ is 0, where for ease of notation we have omitted the second subscript $j$ on $e_i$. Let $\hat{\mathbf{e}} = \mathbf{Z}_j - \mathbf{X}_j \hat{\boldsymbol{\beta}}_{j,ls}$ denote the residuals from the a fit such that $\sqrt{n}(\hat{\beta} - \beta) = \mathcal{O}(1)$. Let $F(t)$ and $f(t)$ denote the cdf and pdf of $e_i$, respectively.

Consider the sign process given by

$$(A.1) \qquad\qquad \overline{S}_j(\theta) = \frac{1}{n_j} \sum_{i=1}^{n_j} \text{sgn}(e_i - \theta),$$

where $\text{sgn}(u) = -1, 0,$ or $1$ for $u < 0$, $u = 0$, or $u > 0$. Denote the median of $e_1, ..., e_{n_j}$ by $\hat{\theta}_e$. Notice that $\hat{\theta}_e$ solves the equation $\overline{S}_j(\theta) = 0$. Our immediate goal is the asymptotic linearity of the process $\overline{S}_j(\theta)$ that is given in expression (A.3). We accomplish this by showing that the four sufficient conditions hold as given in Section 1.5 of Hettmansperger and McKean [7]. First note that $\overline{S}_j(\theta)$ is a nonincreasing function of $\theta$. Thus the first condition holds. For the second condition, by a simple shift theorem and stationarity, we have

$$\mu(\theta) = E_0\big[\overline{S}_j(\theta)\big] = E_\theta\big[\overline{S}_j(0)\big] = \frac{1}{n_j} \sum_{i=1}^{n_j} E_\theta\big[\text{sgn}(e_i)\big] = 1 - 2F(-\theta).$$

Hence, $\mu'(0) = 2f(0) > 0$ which establishes the second condition.

For the third condition, we need to show the variance of $\sqrt{n_j}\,\overline{S}_j(0)$ exists. This variance is

$$
\begin{aligned}
\sigma_{n_j}^2 &= V\big[\sqrt{n_j}\,\overline{S}_j(0)\big] \\
&= \frac{1}{n_j} \sum_{i=1}^{n_j} V\big(\text{sgn}(e_i)\big) + \frac{2}{n_j} \sum_{i=1}^{n_j-1} \sum_{k=i+1}^{n_j} \text{cov}\big[\text{sgn}(e_i), \text{sgn}(e_k)\big].
\end{aligned}
$$

The first term on the right is easily seen to be 1. Using $P[e_i < 0] = 1/2$ and expanding each covariance term into its expectation, we obtain four probability terms and, hence, the sum of four series. The absolute value of one of these four series is given next. As we show, we establish a bound on the series by invoking the assumption (3.1) and then applying properties

of the geometric series. A similar proof holds for the other three series.

$$\left| \frac{2}{n_j} \sum_{i=1}^{n_j-1} \sum_{k=i+1}^{n_j} \Big[ P(e_i<0, e_k<0) - P(e_i<0)\,P(e_k<0) \Big] \right| \leq$$

$$\leq \frac{2}{n_j} \sum_{i=1}^{n_j-1} \sum_{k=i+1}^{n_j} \Big| P(e_i<0, e_k<0) - P(e_i<0)\,P(e_k<0) \Big|$$

$$\leq K \frac{2}{n_j} \sum_{i=1}^{n_j-1} \sum_{k=i+1}^{n_j} \rho^{k-i}$$

$$= 2K \frac{\rho}{1-\rho} \frac{n_j-1}{n_j} - \left[ \frac{1}{n_j} \frac{\rho^2}{1-\rho^2} \left(1 - \rho^{n_j-1}\right) \right]$$

$$\leq 2K \frac{\rho}{1-\rho},$$

where the constants $K > 0$ and $0 \leq \rho < 1$ are given in expression (3.1). The last line follows because the term in brackets is nonnegative and the entire expression is nonnegative. Thus the above series is convergent. Since the other three series follow similarly and, since absolute convergence implies convergence, the series for the variance $\sigma^2_{n_j}$ converges. Let $\sigma^2(0)$ denote the value to which the series converges. The actual value is not needed in the proof but can be obtained from Wendler [20] as noted below.

The fourth condition requires that for all $b$, $\text{Var}_0\Big\{ \sqrt{n_j} \big[ \overline{S}(b/\sqrt{n_j}) - \overline{S}(0) \big] \Big\} \to 0$, as $n_j \to \infty$, where $I(x) = 1$ if $x$ is true, 0 otherwise. Based on the sign function, we have

$$V_{n_j,b} =_{\text{dfn}} \text{Var}\Big[ \sqrt{n_j} \big[ \overline{S}(b/\sqrt{n_j}) - \overline{S}(0) \big] \Big] = \text{Var}\left[ \frac{-2}{\sqrt{n_j}} \sum_{i=1}^{n_j} I\big(0 < e_i < b/\sqrt{n_j}\big) \right].$$

Thus,

(A.2)
$$V_{n_j,b} = \frac{4}{n_j} \sum_{i=1}^{n_j} \text{Var}\Big[ I\big(0 < e_i < b/\sqrt{n_j}\big) \Big]$$
$$+ \frac{8}{n_j} \sum_{i=1}^{n_j-1} \sum_{k=i+1}^{n_j} \text{cov}\Big[ I\big(0 < e_i < b/\sqrt{n_j}\big),\, I\big(0 < e_k < b/\sqrt{n_j}\big) \Big].$$

By stationarity and continuity of the cdf $F(t)$, $E\big[ I\big(0 < e_i < b/\sqrt{n_j}\big) \big] = F\big(b/\sqrt{n_j}\big) - \frac{1}{2} \to 0$, as $n_j \to \infty$; hence, the variance term on the right side of (A.2) goes to 0 as $n_j \to \infty$.

We can write the covariances as

$$c_{n_j,i,k} =_{\text{dfn}} \text{cov}\Big[ I\big(0 < e_i < b/\sqrt{n_j}\big),\, I\big(0 < e_k < b/\sqrt{n_j}\big) \Big]$$
$$= P\Big[ 0 < e_i < b/\sqrt{n_j},\; 0 < e_k < b/\sqrt{n_j} \Big]$$
$$- P\Big[ 0 < e_i < b/\sqrt{n_j} \Big] P\Big[ 0 < e_k < b/\sqrt{n_j} \Big].$$

Notice that this is similar to the above argument on the variance, except that the terms also go to zero as $n_j \to \infty$. Using mean value theorems it follows that the rate of this convergence is $1/n_j$. Using the assumptions from Section 3 and this rate we have $|c_{n_j,i,k}| \leq K \rho_{n_j}^{k-i}$, where $\rho_{n_j} = O(1/n_j)$. Following the same argument as used for the variance, the covariance term in (A.2) in absolute value is less than or equal to

$$2K \frac{\rho_{n_j}}{1-\rho_{n_j}} \leq O(1/n_j) \to 0, \qquad \text{as} \quad n_j \to \infty.$$

Thus $V_{n_j,b} \to 0$ as $n_j \to \infty$.

By these four conditions, as shown in Chapter 1 of Hettmansperger and McKean [7], the sign process satisfies the linearity result:

$$(A.3) \qquad \sqrt{n_j}\,\overline{S}_j(\theta) = \sqrt{n_j}\,\overline{S}_j(0) - 2\,f(0)\sqrt{n_j}\,\theta + o_p(1),$$

for $\sqrt{n_j}\,|\theta| \le B$, for all $B > 0$.

To obtain $\sigma^2(0)$, we can use Wendler [20]. He showed, under the mixing conditions above, that $\sqrt{n_j}\,|\hat{\theta}_e|$ converges in distribution and, hence, is tight. Since $\overline{S}_j(\theta) = 0$, we can use (A.3) and Wendler's asymptotic distribution to obtain the asymptotic normal distribution of $\sqrt{n_j}\,\overline{S}_j(0)$.

For our proof, we are interested in the residual process. Since for the proof the true parameters are 0, we can write the residuals as $\hat{e}_i = e_i - \mathbf{x}_i'\hat{\boldsymbol{\beta}}_{LS}$, $i = 1, ..., n_j$. The residual sign process is then given by

$$(A.4) \qquad \overline{S}_j^*(\theta) = \frac{1}{n_j}\sum_{i=1}^{n_j}\mathrm{sgn}(\hat{e}_i - \theta).$$

Let $\hat{\theta}^*$ denote median of the residuals. Notice that it solves $\overline{S}_j^*(\hat{\theta}^*) = 0$. In the independent error case, Hettmansperger and McKean [7] established the linearity of the residual process for any root-$n$ consistent estimate of $\boldsymbol{\beta}$; see their Section 3.5 and the associated parts of the Appendix. A key result used in their proof was the linearity for the single sample case, i.e., in the current proof, the result (A.3). See Lemma A.3.2 of Hettmansperger and McKean [7]. The remainder of the proof for the linearity of $\overline{S}_j^*(\theta)$ follows using similar reasoning as above. The result is

$$(A.5) \qquad \sqrt{n_j}\,\overline{S}_j^*(\theta) = \sqrt{n_j}\,\overline{S}_j^*(0) - 2\,f(0)\sqrt{n_j}\,\theta + o_p(1),$$

for $\sqrt{n_j}\,|\theta| \le B$, for all $B > 0$. Using this and $\overline{S}_j^*(\hat{\theta}^*) = 0$, we obtain the asymptotic distribution of $\hat{\theta}^*$ and, hence, its consistency.

The second part of our proof requires the consistency of three other estimators. The first is the median of the absolute value of the residuals. This is easily obtained by replacing $e_i$ with $|e_i|$ in the above processes. Since the pdf of $|e_i|$ is strictly positive at the true median, the proof holds in this case too. The second estimator is a function of the residuals from two bins, say, $j$ and $j'$. More specifically, it is a function of the residuals

$$\hat{e}_{j,i} + \hat{e}_{j',i'} = e_{j,i} + e_{j',i'} - \begin{bmatrix} \mathbf{x}_{j,i}^{\mathsf{T}} & \mathbf{x}_{j',i'}^{\mathsf{T}} \end{bmatrix}\begin{bmatrix} \hat{\boldsymbol{\beta}}_j \\ \hat{\boldsymbol{\beta}}_{j'} \end{bmatrix},$$

where $\hat{\boldsymbol{\beta}}_j$ and $\hat{\boldsymbol{\beta}}_{j'}$ denote the LS estimates from bins $j$ and $j'$, respectively. Because the vector $(\hat{\boldsymbol{\beta}}_j^{\mathsf{T}}, \hat{\boldsymbol{\beta}}_{j'}^{\mathsf{T}})^{\mathsf{T}}$ is root-$n$ consistent and the convolution of identical pdfs is positive at its median when each pdf is positive at its median, nothing in the above proof precludes the use of random errors of the form $e_{j,i} + e_{j',i'}$. Thus the theory holds in this case also. These comments apply to the third estimator also because it is based on the residuals $\hat{e}_{j,i} - \hat{e}_{j',i'}$.

**Part 2 of the Proof:**

This part of the proof makes use of the standard inequality $|a| = |a - b + b| \leq |a - b| + |b|$. It suffices to show consistency of $\hat{\boldsymbol{\Sigma}}_M^{(\text{rob})}$ element-wise. We first show the consistency of the diagonal elements. The statistic and functional of the $m^{\text{th}}$ diagonal of $\hat{\boldsymbol{\Sigma}}_M^{(\text{rob})}$ are given by:

$$\text{MAD}\{\hat{\mathbf{e}}_m\} = \text{med}_i \big|\hat{\mathbf{e}}_{m_i} - \text{med}_j\{\hat{\mathbf{e}}_{m_j}\}\big| \quad \text{with functional} \quad \xi_m = \text{med}\big|\mathbf{e}_m - \text{med}\{\mathbf{e}_m\}\big|.$$

Without loss of generality, assume that $\text{med}\{\mathbf{e}_m\} = 0$. From Part 1, $\text{med}_i\{\hat{\mathbf{e}}_{m_i}\} \xrightarrow{P} 0$, in probability. Next, assume that $\text{med}\{|\mathbf{e}_m|\} = \xi$. Then also from Part 1, $\text{med}_i|\hat{\mathbf{e}}_{m_i}| \xrightarrow{P} \xi$. Choose $N_0$ sufficiently large so that, given $\varepsilon > 0$,

$$(A.6) \qquad\qquad k \geq N_0 \implies \big|\text{med}_{1 \leq i \leq k}\{\hat{\mathbf{e}}_{m_i}\}\big| < \varepsilon$$

with probability greater than $(1 - (\varepsilon/2))$. Let $A_n$ denote the event where (A.6) occurs. Then, on $A_n$ we have

$$\begin{aligned} |\hat{\mathbf{e}}_{m_i}| &= \big|\hat{\mathbf{e}}_{m_i} - \text{med}_j\{\hat{\mathbf{e}}_{m_j}\} + \text{med}_j\{\hat{\mathbf{e}}_{m_j}\}\big| \\ &\leq \big|\hat{\mathbf{e}}_{m_i} - \text{med}_j\{\hat{\mathbf{e}}_{m_j}\}\big| + \big|\text{med}_j\{\hat{\mathbf{e}}_{m_j}\}\big| \\ &< \big|\hat{\mathbf{e}}_{m_i} - \text{med}_j\{\hat{\mathbf{e}}_{m_j}\}\big| + \varepsilon. \end{aligned}$$

So, on $A_n$,

$$(A.7) \qquad\qquad \text{med}_i|\hat{\mathbf{e}}_{m_i}| < \text{med}_i\big|\hat{\mathbf{e}}_{m_i} - \text{med}_j\{\hat{\mathbf{e}}_{m_j}\}\big| + \varepsilon,$$

and

$$\begin{aligned} \big|\hat{\mathbf{e}}_{m_i} - \text{med}_j\{\hat{\mathbf{e}}_{m_j}\}\big| &= \big|\hat{\mathbf{e}}_{m_i} - \text{med}_j\{\hat{\mathbf{e}}_{m_j}\} - \hat{\mathbf{e}}_{m_i} + \hat{\mathbf{e}}_{m_i}\big| \\ &\leq \big|\text{med}_j\{\hat{\mathbf{e}}_{m_j}\}\big| + |\hat{\mathbf{e}}_{m_i}| \\ &< |\hat{\mathbf{e}}_{m_i}| + \varepsilon. \end{aligned}$$

Hence, on $A_n$,

$$(A.8) \qquad\qquad \text{med}_i\big|\hat{\mathbf{e}}_{m_i} - \text{med}_j\{\hat{\mathbf{e}}_{m_j}\}\big| < \text{med}_i|\hat{\mathbf{e}}_{m_j}| + \varepsilon.$$

Putting (A.7) and (A.8) together, we have on $A_n$,

$$(A.9) \qquad\qquad \Big|\text{med}_i\big|\hat{\mathbf{e}}_{m_i} - \text{med}_j\{\hat{\mathbf{e}}_{m_j}\}\big| - \text{med}_i|\hat{\mathbf{e}}_{m_i}|\Big| < \varepsilon.$$

Since this occurs with probability of at least $(1 - (\varepsilon/2))$, the difference on the left-side goes to 0 in probability. As noted above, from Part 1, $\text{med}_i|\hat{\mathbf{e}}_{m_i}| \xrightarrow{P} \xi$; hence, $\text{med}_i\big|\hat{\mathbf{e}}_{m_i} - \text{med}_j\{\hat{\mathbf{e}}_{m_j}\}\big| \xrightarrow{P} \xi$.

For the off-diagonal elements, let $m \neq m'$ be given. Recall that the off-diagonal elements of $\hat{\boldsymbol{\Sigma}}_M^{(\text{rob})}$ are given by equation (2.2), which can be expressed as follows:

$$(A.10) \qquad \hat{\boldsymbol{\Sigma}}_M^{(\text{rob})}(m, m') = \left(\text{MAD}\left\{\frac{\hat{\mathbf{e}}_m \oplus \hat{\mathbf{e}}_{m'}}{2}\right\}\right)^2 - \left(\text{MAD}\left\{\frac{\hat{\mathbf{e}}_m \ominus \hat{\mathbf{e}}_{m'}}{2}\right\}\right)^2.$$

It suffices to show consistency for each of the terms on the right-side. Define $\mathbf{t} = \frac{1}{2}(\mathbf{e}_m \oplus \mathbf{e}_{m'})$. Then the statistic and its functional, respectively, for the off-diagonal elements are:

$$\text{MAD}\{\hat{\mathbf{t}}\} = \text{med}_i\big|\hat{\mathbf{t}}_i - \text{med}_j\{\hat{\mathbf{t}}_j\}\big| \quad \text{with functional} \quad \xi_{m,m'} = \text{med}\big|\mathbf{t} - \text{med}\{\mathbf{t}\}\big|.$$

Without loss of generality let $\text{med}\{\mathbf{t}\} = 0$. From Part 1, $\text{med}_i\{\hat{\mathbf{t}}_i\} \xrightarrow{P} 0$. Then the proof follows in the same manner as for the diagonal elements. So each of the MADs in equation (A.10) is consistent. Therefore, the entire expression is consistent. Thus, the diagonal and off-diagonal entries of $\hat{\boldsymbol{\Sigma}}_M^{(\text{rob})}$ are consistent. Hence, $\hat{\boldsymbol{\Sigma}}_M^{(\text{rob})}$ is a consistent estimator of $\boldsymbol{\Sigma}_M$. $\qquad\square$

## ACKNOWLEDGMENTS

## REFERENCES

[1]    BANERJEE, S.; GELFAND, A.; FINLEY, A. and SANG, H. (2008). Gaussian predictive process models for large spatial datasets, *Journal of the Royal Statistical Society: Series B*, **70**(4), 825–844.

[2]    BLIZNYUK, N.; CARROLL, R.; GENTON, M. and WANG, Y. (2012). Variogram estimation in the presence of trend, *Statistics and Its Interface*, **5**, 159–168.

[3]    CRESSIE, N. and HAWKINS, D. (1984). Robust kriging – a proposal, *Mathematical Geology*, **16**, 3–18.

[4]    CRESSIE, N. and JOHANNESSON, G. (2008). Fixed rank kriging for very large spatial data sets, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **70**, 209–226.

[5]    FINLEY, A.O.; SANG, H.; BANERJEE, S. and GELFAND, A.E. (2009). Improving the performance of predictive process modeling for large datasets, *Computational Statistics & Data Analysis*, **53**(8), 2873–2884.

[6]    FRENCH, J. (2018). *SpatialTools: Tools for Spatial Data Analysis*, R package version 1.0.4.

[7]    HETTMANSPERGER, T. and MCKEAN, J. (2011). *Robust Nonparametric Statistical Methods*, Chapman Hall, New York, 2nd edition.

[8]    KANG, E. and CRESSIE, N. (2011). Bayesian inference for the spatial random effects model, *Journal of the American Statistical Association*, **106**, 972–983.

[9]    KANG, E.; CRESSIE, N. and SHI, T. (2010). Using temporal variability to improve spatial mapping with application to satellite data, *The Canadian Journal of Statistics*, **38**, 271–289.

[10]   KANG, E.L.; CRESSIE, N. and SAIN, S.R. (2012). Combining outputs from the north american regional climate change assessment program by using a bayesian hierarchical model, *Journal of the Royal Statistical Society C*, **61**(2), 291–313.

[11]   KATZFUSS, M. and CRESSIE, N. (2011). *Tutorial on fixed rank kriging (frk) of co2 data*, Technical Report, The Ohio State University, 858.

[12]   KOENKER, R. (2018). *quantreg: Quantile Regression*, R package version 5.35.

[13]   KOENKER, R. and BASSETT, G. (1978). Regression quantiles, *Econormetrica*, **46**, 33–50.

[14]   LAHIRI, S.; LEE, Y. and CRESSIE, N. (2002). Asymptotic distribution and asymptotic efficiency of least squares estimators of variogram parameters, *Journal of Statistical Planning and Inference*, **103**, 65–85.

[15]   NYCHKA, D.; WIKLE, C. and ROYLE, J.A. (2002). Multiresolution models for nonstationary spatial covariance functions, *Statistical Modelling*, **2**, 315–331.

[16]    PAUL, R.; JELSEMA, C.M. and LAU, K.W. (2015). *A flexible class of reduced rank spatial models for large non-gaussian datasets.* In "Current Trends in Bayesian Methodology with Applications" (S.K. Upadhyay, U. Singh, D.K. Dey and A. Loganathan, Eds.), Chapman & Hall/CRC Press.

[17]    SANG, H. and HUANG, J.Z. (2011). A full scale approximation of covariance functions for large spatial data sets, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **74**(1), 111–132.

[18]    SHI, T. and CRESSIE, N. (2007). Global statistical analysis of misr aerosol data: a massive data product from nasa's terra satellite, *Environmetrics*, **18**, 665–680.

[19]    STEIN, M.L. (2007). Spatial variation of totla column ozone on a global scale, *The Annals of Applied Statistics*, **1**, 191–210.

[20]    WENDLER, M. (2011). Bahadur representation for $U$-quantiles of dependent data, *Journal of Multivariate Analysis*, **102**, 1064–1079.

[21]    WILKS, D. (2006). *Statistical Methods in the Atmospheric Sciences*, Academic Press, San Diego, USA, 2nd edition.

[22]    ZHU, Y.; KANG, E.L.; BO, Y.; TANG, Q.; CHENG, J. and HE, Y. (2015). A robust fixed rank kriging method for improving the spatial completeness and accuracy of satellite sst products, *IEEE Transactions on Geoscience and Remote Sensing*, **53**, 5021–5035.

# ESTIMATION OF SMALL AREA TOTAL WITH RANDOMIZED DATA

Authors:  Shakeel Ahmed
– Department of Statistics, Quaid-i-Azam University Islamabad,
Islamabad, Pakistan
sahmed@qau.edu.pk

Javid Shabbir
– Department of Statistics, Quaid-i-Azam University Islamabad,
Islamabad, Pakistan
javidshabbir@gmail.com

Sat Gupta
– Department of Mathematics and Statistics, University of North Carolina at
Greensboro, North Carolina, USA
sngupta@uncg.edu

Frank Coolen
– Department of Mathematical Sciences, Durham University,
Durham, DH1 3LE, UK
frank.coolen@durham.ac.uk

Abstract:

• In social surveys involving questions that are sensitive or personal in nature, respondents may
not provide correct answers to certain questions asked by the interviewer. The impact of this non-
response or inaccurate response becomes even more acute in the case of small area estimation (SAE)
where we already have the problem of small sample size coming from the small area. To obtain a
truthful response, we use randomized response techniques in each small area. We assume that a
non-sensitive auxiliary variable, highly correlated with the study variable, is available. We use the
word model in two senses — one in the context of population models, i.e. the relationship between
the study variable and the auxiliary variable; and second, the scrambled response model. We focus
on the problem of estimating small area total and examine its performance both theoretically and
numerically.

## 1.    INTRODUCTION

In social sciences, responses on some stigmatizing variables are often needed to make inference about the behavior of some human populations. Examples of such situations are where questions are asked that are related to topics like tax evasion, use of illegal drugs, extra marital affairs, ethical issues, political affiliation, etc. In the case of stigmatizing study variables, non-sampling error may increase due to missing or false responses, which leads to biased estimates of population parameters such as mean, total or proportion. To reduce such bias in sample surveys, [34] proposed a randomized response technique (RRT) for obtaining more accurate estimates. A lot of research has been done for improving the original RRT model of [34]. Authors contributing in this area include [17], [18], [35], [6], [12], [22], [7], [3], [19, 20], [21] and [9, 10, 11]. In RRT literature, much more attention has been paid to design-based approach which assumes the population to consist of fixed constants. But in many real-life situations, population values are generated as realizations of a set of stochastic variables. Such population is called a superpopulation and the statistical models for such type of populations are called superpopulation models. Superpopulation models help in sample selection, constructing estimators for population parameters of interest, and enhancing the precision of estimates. A superpopulation model uses the relationship between the study variable and the auxiliary variable(s) to predict the population values for the non-sampled units assuming non-informative sampling approach. Under the framework of model-based inference, [14] dealt with the problem of estimation of a finite population mean or total. [27] and [8] attempted to obtain optimal model-unbiased estimators of the population mean and total using least squares estimation methods and the well-known Gauss–Markov theorem. Some discussion on model-based approach can be found in [2], [15], [16], [30, 31], [29], [28], and [33]. A detailed review of model-based estimation is also available in [32].

[13] and [24] have suggested post-censal estimates (estimates obtained immediately after census using the census results) for small areas and called it small area estimation (SAE). [23] dealt with labor force trend estimation for small areas. Work related to such methods can also be found in [25, 26] and [38]. More recently, [36] have considered estimation of uncertainty in spatial micro-simulation approaches for SAE. The main purpose of SAE is to overcome the problem of small sample when separate estimates for domains are needed. In this article, we develop some model-based estimators for small area totals assuming the study variable in each domain is sensitive. A generalized randomized response model has been used to collect information about the study variable. The rest of the article is structured as follows: an overview of SAE under direct response is considered in Section 2 with some superpopulation models. Section 3 extends the SAE given in Section 2 to randomized response models, assuming a sensitive quantitative study variable and non-sensitive auxiliary variable. Section 4 presents a numerical study based on two real life data sets. Some concluding remarks are provided in Section 5.

## 2.   SAE UNDER DIRECT RESPONSE

Consider a finite population $U = \{U_1, U_2, ..., U_N\}$ of $N$ units as a realization of a super-population with variable of interest $y$, and auxiliary variable $x$. For a specific sup-population $A_k$, also known as "small area", let $d_{ki}$ be an area specific binary variable, for $k = 1, 2, 3, ...m$ and $i = 1, 2, ...N$, such that $d_{ki} = 1$ if $U_i$ belongs to $A_k$, and zero otherwise. Further, let $N_k = \sum_U d_{ki}$ be the size of the $k$-th sub-population or $k$-th small area (usually unknown), $T_{yk} = \sum_U d_{ki}\, y_i$ and $T_{xk} = \sum_U d_{ki}\, x_i$ be the population totals, $\mu_{yk} = \frac{T_{yk}}{N_k}$ and $\mu_{xk} = \frac{T_{xk}}{N_k}$ be the population means, and $\sigma_{yk}^2 = \frac{1}{N_k} \sum_U d_{ki}(y_i - \mu_{yk})^2$ and $\sigma_{xk}^2 = \frac{1}{N_k} \sum_U d_{ki}(x_i - \mu_{xk})^2$ be the population variances of the study variable and the auxiliary variable respectively in the $k$-th area. The notation $\sum_U$ is used for summing the values over $U$. Also, let the covariance between the study variable and the auxiliary variable in the $k$-th area be $\sigma_{yxk} = \frac{1}{N_k} \sum_U d_{ik}(y_i - \mu_{yk})(x_i - \mu_{xk})$. Suppose that $s$ is a member of the set $S$ of all possible samples that can be drawn from $U$ using simple random sampling without replacement (SRSWOR) scheme with size $n$, and $\bar{s}$ consists of all those elements of $U$ that are not selected in sample $s$. The population total for the study variable, quantity of interest or estimand, in $k$-th area can then be expressed as $T_{yk} = \sum_s d_{ki}\, y_i + \sum_{\bar{s}} d_{ki}\, y_i$. A predictor for $T_{yk}$ is obtained as follows:

$$(2.1) \qquad \hat{T}_{yk} = \sum_s d_{ki}\, y_i + \sum_{\bar{s}} d_{ki}\, \hat{y}_i.$$

The main problem is to find $\hat{y}_i$ for $U_i \in \bar{s}$. The predictor $\hat{y}_i$ is obtained assuming different superpopulation models. We consider three most widely used population models:

1.   Homogenous Population Model (HPM):  $y = \mu_{yk} + \varepsilon$,
2.   Linear Population Model (LPM):  $y = \alpha_k + \beta x + \varepsilon$,
3.   Ratio Population Model (RPM):  $y = \gamma x + x^{1/2} \varepsilon$,

for $k = 1, 2, ..., m$, where $\varepsilon$ is the stochastic error term which has mean 0 and a constant variance $\sigma^2$. Also, $\mu_{yk}$ and $\alpha_k$ are mean effects in $k$-th area and $\beta$ and $\gamma$ are the coefficients of the regression line of $y$ on $x$ for the whole population for the cases with and without intercepts. In model based approach, these parameters are termed as superpopulation parameters.

## 2.1.  Homogeneous Population Model (HPM)

In case of HPM, a BLUP for $\mu_{yk}$, obtained by minimizing the residual sum of square $\sum_s d_{ki}(y_i - \mu_{yk})^2$ is $\bar{y}_k = \frac{1}{n_k} \sum_s d_{ki}\, y_i$, which yields an estimator for $T_{yk}$ given by

$$(2.2) \qquad \hat{t}_{kh} = \sum_s d_{ki}\, y_i + \sum_{\bar{s}} d_{ki}\, \bar{y}_k = \frac{N}{n} \sum_s d_{ki}\, y_i.$$

The sub-script 'h' is used to indicate that the superpopulation model is homogeneous. It is straight forward to show that $\hat{t}_{kh}$ is an unbiased estimator of population total $T_{yk}$ with variance given by

$$(2.3) \qquad \mathrm{Var}(\hat{t}_{kh}) = \lambda \left[ \theta_k\, \sigma_{yk}^2 + \theta_k(1 - \theta_k)\, \mu_{yk}^2 \right],$$

where $\theta_k = \frac{N_k}{N}$ is the population proportion of the units belonging to $k$-th small area, and $\lambda = \frac{N(N-N)}{n}$. For proof readers can see [5, p. 156–160].

## 2.2. Linear Population Model (LPM)

Now consider LPM for finding $\hat{y}_i$, $U_i \in \bar{s}$. The BLUP for $\alpha_k$ and $\beta$ are obtained by minimizing the sum of squared prediction errors for specific areas, i.e.

$$\text{SSPE} = \sum_s d_{ki}(y_i - \alpha_k - x_i \beta)^2 \,.$$

These are given by $\hat{\alpha}_k = \bar{y}_k - \hat{\beta}\bar{x}_k$ and $\hat{\beta} = \frac{\sum_s d_{ki}(y_i - \bar{y}_k)(x_i - \bar{x}_k)}{\sum_s d_{ki}(x_i - \bar{x}_k)^2}$, where $\bar{y}_k$ and $\bar{x}_k$ are the sample means corresponding to $k$-th small area. The estimator of $T_{yk}$ under LPM is given by

$$\hat{t}_{k\text{lr}} = \sum_s d_{ki}\, y_i + \sum_{\bar{s}} d_{ki}\big(\hat{\alpha}_k + \hat{\beta}\, x_i\big)\,.$$

After some simplifications and using assumption from [5], i.e. $\frac{N_k}{N} \approx \frac{n_k}{n}$, we get

$$(2.4) \qquad\qquad \hat{t}_{k\text{lr}} = \frac{N}{n} t_{yk} + \hat{\beta}\left(T_{xk} - \frac{N}{n} t_{xk}\right),$$

where $t_{yk} = \sum_s d_{ki}\, y_i$ and $t_{xk} = \sum_s d_{ki}\, x_i$ are the sample totals for $k$-th small area. Further, $\hat{\beta}$ given in (2.4) is based on local (area specific) observations only, which do not account for relationship between the variables for the entire population. To overcome this deficiency, different area level models have been proposed in literature. For simplicity, we assume that the regression coefficient $\beta$ of $y$ on $x$ is known for the whole population. For known $\beta$, we have

$$(2.5) \qquad\qquad \hat{t}_{k\text{lr}} = \frac{N}{n} t_{yk} + \beta\left(T_{xk} - \frac{N}{n} t_{xk}\right).$$

The sub-script 'lr' is used to denote that the underlying model is linear. For known $\beta$, $\hat{t}_{k\text{lr}}$ is unbiased for $T_{yk}$ with variance given by

$$(2.6) \qquad\qquad \text{Var}(\hat{t}_{k\text{lr}}) = \lambda\big(\sigma_{yk}^{*2} + \beta^2 \sigma_{xk}^{*2} - 2\beta \sigma_{yxk}^{*}\big)\,,$$

where $\sigma_{yk}^{*2} = \theta_k \sigma_{yk}^2 + \theta_k(1-\theta_k)\mu_{yk}^2$, $\sigma_{xk}^{*2} = \theta_k \sigma_{xk}^2 + \theta_k(1-\theta_k)\mu_{xk}^2$ and $\sigma_{yxk}^{*} = \theta_k \sigma_{yxk} + \theta_k(1-\theta_k)\mu_{yk}\mu_{xk}$. The value of $\beta$ that minimizes the variance is $\beta_{\text{opt}} = \frac{\sigma_{yxk}^{*}}{\sigma_{xk}^{*2}}$. The corresponding minimum variance of $\hat{t}_{k\text{lr}}$ is given by

$$(2.7) \qquad\qquad \text{Var}(\hat{t}_{k\text{lr}})_{\text{opt}} = \lambda\big(1 - \rho_{yxk}^{*2}\big)\sigma_{yk}^{*2}\,,$$

where $\rho_{yxk}^{*} = \frac{\sigma_{yxk}^{*}}{\sigma_{yk}^{*}\sigma_{xk}^{*}}$. From Equations (2.7) and (2.3), it is obvious that $\hat{t}_{k\text{lr}}$ is always more efficient than $\hat{t}_{k\text{h}}$ for any linear relationship between $y$ and $x$.

## 2.3. Ratio Population Model (RPM)

For situations when there is a proportional relationship between the survey variable and the auxiliary variables, the RPM [32] is often preferred as the working model. RPM is given by

$$(2.8) \qquad y = \gamma x + x^{1/2}\varepsilon \,.$$

The estimator for $\gamma$ which minimizes the sum of squared errors, i.e. $\text{SSE}^* = \sum_s d_{ki}\left(\frac{y_i - x_i\gamma}{x_i^{1/2}}\right)^2$, is given by $\hat{\gamma} = \frac{\sum_s d_{ki}\, y_i}{\sum_s d_{ki}\, x_i}$. Now consider

$$(2.9) \qquad \hat{t}_{kr} = \sum_s d_{ki}\, y_i + \sum_{\bar{s}} d_{ki}(\hat{\gamma}\, x_i)$$

as an estimator of $T_{yk}$. The sub-script 'r' is used to denote the ratio population model for the response variable. After simplification and assuming $\frac{N_k}{N} \approx \frac{n_k}{n}$, we get

$$(2.10) \qquad \hat{t}_{kr} = \frac{\sum_s d_{ki}\, y_i}{\sum_s d_{ki}\, x_i} \sum_{\bar{s}} d_{ki}\, x_i = \frac{N}{n}\left[t_{yk}\frac{n\,\mu_{xk}}{t_{xk}}\right].$$

The bias and MSE respectively, of $\hat{t}_{kr}$, are given by

$$(2.11) \qquad \text{Bias}(\hat{t}_{klr}) \cong \frac{\lambda}{N}\,\mu_{yk}\left(C_{xk}^{*2} - C_{yxk}^{*}\right)$$

and

$$(2.12) \qquad \text{MSE}(\hat{t}_{kr}) \cong \lambda\mu_{yk}^2\left(C_{yk}^{*2} + C_{xk}^{*2} - 2C_{yxk}^{*}\right),$$

where $C_{yk}^{*2} = \frac{\sigma_{yk}^{*2}}{\mu_{yk}^2}$, $C_{xk}^{*} = \frac{\sigma_{xk}^{*2}}{\mu_{xk}^2}$ and $C_{yxk}^{*} = \frac{\sigma_{yxk}^{*}}{\mu_{yk}\,\mu_{yk}}$. From (2.3) and (2.12), it can be inferred that $\text{MSE}(\hat{t}_{kr}) \leq \text{Var}(\hat{t}_{kh})$ if $\rho_{yxk}^{*} \geq \frac{1}{2}\frac{C_{xk}^{*}}{C_{yk}^{*}}$.

## 3. SAE UNDER RANDOMIZED RESPONSE TECHNIQUE

When the study variable is of sensitive nature, it is difficult to obtain 100% response through direct response method. For improved response rate in such situations, survey statisticians prefer to use RRT. Assuming quantitative study variable, and following [11], we use the following scrambled response model

$$(3.1) \qquad z = a\,y + b \,,$$

where $y$ is the sensitive study variable which follows one of the population models given in Section 2, $a$ and $b$ are two uncorrelated scrambling variables with means $\mu_a$ and $\mu_b$, and variances $\sigma_a^2$ and $\sigma_b^2$ respectively. Further, $a$ and $b$ are independent of the study variable $y$. Note that respondents from each small area use the same scrambling variables $a$ and $b$ whose distributions are unknown to the interviewer while the means and variances are known. Taking expectation of Equation (3.1) with respect to randomization mechanism, we have $E_R(z) = \mu_a y + \mu_b$. The transformed scrambled response is obtained as $y = \frac{E_R(z) - \mu_b}{\mu_a}$. A sample unbiased estimate for $y$ is $\tilde{y} = \frac{z - \mu_b}{\mu_a}$.

---

### 3.1.　Homogeneous Population Model (HPM)

---

When the underlying population model is homogeneous, i.e. when there is no covariate affecting the outcome variable, a BLUP for the superpopulation parameter $\mu_{yk}$ is $\tilde{\bar{y}}_k = \tilde{t}_{yk}/n_k$ which yields an estimator for $T_{yk}$ given by

$$(3.2) \qquad \tilde{t}_{k\mathrm{h}} \;=\; \sum_s d_{ki}\,\tilde{y}_i + \sum_{\bar{s}} d_{ki}\,\tilde{\bar{y}}_k \;=\; n_k\,\tilde{\bar{y}}_k + (N_k - n_k)\,\tilde{\bar{y}}_k \;=\; \frac{N}{n}\,\tilde{t}_{yk}\,,$$

where $\tilde{t}_{yk} = \sum_s d_{ki}\,\tilde{y}_i$. We assume that the sampling weights for the whole sample and the sample within $k$-th domain are same, i.e. $\frac{N_k}{N} \approx \frac{n_k}{n}$. It is easy to show that $\tilde{t}_{k\mathrm{h}}$ is an unbiased estimator of population total $T_{yk}$ with variance

$$(3.3) \qquad \mathrm{Var}(\tilde{t}_{k\mathrm{h}}) \;=\; \lambda\left(\theta_k\,\tilde{\sigma}_{yk}^2 + \theta_k(1-\theta_k)\,\mu_{yk}^2\right),$$

where $\tilde{\sigma}_{yk}^2 = \mathrm{Var}(\tilde{y}_i \mid d_{ki}{=}1) = \frac{1}{\mu_a^2}\,\mathrm{Var}(z_i \mid d_{ki}{=}1)$, and

$$(3.4) \qquad \begin{aligned} \mathrm{Var}(z_i \mid d_{ki}{=}1) &= V_s\{E_R(z_i \mid d_{ki}{=}1)\} + V_R\{E_S(z_i \mid d_{ki}{=}1)\} \\ &= E_s\left(\sigma_a^2\,y_i^2 + \sigma_b^2 \mid d_{ki}{=}1\right) + V_s\left(\mu_a\,y_i + \mu_b \mid d_{ki}{=}1\right) \\ &= \sigma_a^2\,\mu_{2,yk} + \sigma_b^2 + \mu_a^2\,\sigma_{yk}^2\,, \end{aligned}$$

where $E_s$ and $V_s$ are the expectation and variance with respect to the data generating mechanism. Also $\mu_{2,yk}$ is the second order raw moment for $k$-th area. Using value of $\tilde{\sigma}_{yk}^2$ from (3.3), we get

$$\mathrm{Var}(\tilde{t}_{k\mathrm{h}}) \;=\; \lambda\left(\theta_k\,\sigma_{yk}^2 + \theta_k(1-\theta_k)\,\mu_{yk}^2 + \theta_k\,\psi_{yk}^2\right),$$

$$(3.5) \qquad \mathrm{Var}(\tilde{t}_{k\mathrm{h}}) \;=\; \mathrm{Var}(\hat{t}_{k\mathrm{h}}) + \lambda\left(\theta_k\,\psi_{yk}^2\right),$$

where $\psi_{yk}^2 = \frac{1}{\mu_a^2}\left(\sigma_a^2\,\mu_{2,yk} + \sigma_b^2\right)$. It is observed from (3.5) that $\mathrm{Var}(\tilde{t}_{k\mathrm{h}})$ is always larger than $\mathrm{Var}(\hat{t}_{k\mathrm{h}})$ as the second term is positive. For detailed derivation, see [1]. The $\mathrm{Var}(\tilde{t}_{k\mathrm{h}})$ decreases with decrease in variance of the scrambled variables but this leads to reduction in respondent's privacy as well. Hence, the variance of the scrambled response models should be of a reasonable size resulting in a proper tradeoff between respondent's privacy and the efficiency of the proposed estimators.

To improve efficiency for a fixed level of privacy protection, we use model relationship between the available auxiliary variable and the study variable. Subsections 3.2 and 3.3 cover linear and ratio population models respectively that utilize the relationship between the variables at unit level to increase efficiency.

---

### 3.2.　Linear Population Model (LPM)

---

Assuming LPM, we find the predicted transformed scrambled response $\tilde{y}_i$, $U_i \in \bar{s}$. The BLUP for $\alpha_k$ and $\beta$ are obtained by minimizing the sum of squared errors for the $k$-th area as follows:

$$\mathrm{SSE} \;=\; \sum_s d_{ki}\,\tilde{e}_i^2 \;=\; \sum_s d_{ki}(\tilde{y}_i - \alpha_k - x_i\beta)^2\,,$$

where $\tilde{\tilde{\alpha}}_k = \tilde{\bar{y}}_k - \tilde{\tilde{\beta}}\bar{x}_k$ and $\tilde{\tilde{\beta}} = \frac{\sum_s d_{ki}(\tilde{y}_i - \tilde{\bar{y}}_k)(x_i - \bar{x}_k)}{\sum_s d_{ki}(x_i - \bar{x}_k)^2}$. The predictive estimator under LPM using transformed scrambled response is given by

$$(3.6) \qquad \tilde{t}_{k\mathrm{lr}} = \sum_s d_{ki}\,\tilde{y}_i + \sum_{\bar{s}} d_{ki}\big(\tilde{\tilde{\alpha}}_k + \tilde{\tilde{\beta}}x_i\big).$$

After some simplification, we get

$$\tilde{t}_{k\mathrm{lr}} = \frac{N}{n}\,\tilde{t}_{yk} + \tilde{\tilde{\beta}}\left(T_{xk} - \frac{N}{n}\,t_{xk}\right).$$

By same argument as given in Subsection 2.2, we have

$$(3.7) \qquad \tilde{t}_{k\mathrm{lr}} = \frac{N}{n}\,\tilde{t}_{yk} + \beta\left(T_{xk} - \frac{N}{n}\,t_{xk}\right).$$

For known $\beta$, $\tilde{t}_{k\mathrm{lr}}$ is unbiased for $T_{yk}$, with variance given by

$$(3.8) \qquad \mathrm{Var}(\tilde{t}_{k\mathrm{lr}}) = \lambda\big(\tilde{\sigma}_{yk}^{*2} + \beta^2\sigma_{xk}^{*2} - 2\,\beta\,\sigma_{yxk}^{*}\big).$$

The optimum value of $\beta$ is $\beta_{\mathrm{opt}} = \frac{\sigma_{yxk}^{*}}{\sigma_{xk}^{*2}}$ with corresponding design optimum variance

$$(3.9) \qquad \mathrm{Var}(\tilde{t}_{k\mathrm{lr}})_{\mathrm{opt}} = \lambda\big(1 - \tilde{\rho}_{yxk}^{*2}\big)\,\tilde{\sigma}_{yk}^{*2},$$

where $\tilde{\rho}_{yxk}^{*} = \frac{\sigma_{yxk}^{*}}{\tilde{\sigma}_{yk}^{*}\sigma_{xk}^{*}}$. Equation (3.9) shows that $\tilde{t}_{k\mathrm{lr}}$ is always more efficient than $\tilde{t}_{k\mathrm{h}}$ for any correlation between $y$ and $x$.

## 3.3. Ratio Population Model (RPM)

For the situation when there is a proportional relationship between the sensitive study variable, and the auxiliary variable whose values are available for all population units and the variance of the survey variable is also proportional to the auxiliary variable, the RPM is often preferred. Consider (3.1), where $y$ follows the ratio population model. The estimator for $\gamma$ which minimizes the sum of squared errors, i.e. $\mathrm{SSE}^{*} = \sum_s d_{ki}\big(\frac{\tilde{y}_i - x_i\gamma}{x_i^{1/2}}\big)^2$, is given by $\tilde{\tilde{\gamma}} = \frac{\sum_s d_{ki}\tilde{y}_i}{\sum_s d_{ki}x_i}$. Consider the prediction problem as follows

$$(3.10) \qquad \tilde{t}_{k\mathrm{r}} = \sum_s \tilde{d}_{ki}\,\tilde{y}_i + \sum_{\bar{s}} d_{ki}\big(\tilde{\tilde{\gamma}}x_i\big).$$

After simplification, we get

$$(3.11) \qquad \tilde{t}_{k\mathrm{r}} = \frac{\sum_s d_{ki}\,\tilde{y}_i}{\sum_s d_{ki}x_i}\sum_{\bar{s}} d_{ki}x_i = \frac{N}{n}\left[\tilde{t}_{yk}\frac{n\,\mu_{xk}}{t_{xk}}\right].$$

The bias and MSE of $\tilde{t}_{k\mathrm{r}}$ are given by

$$(3.12) \qquad \mathrm{Bias}(\tilde{t}_{k\mathrm{lr}}) \cong \frac{\lambda}{N}\,\mu_{yk}\big(C_{xk}^{*2} - C_{yxk}^{*}\big)$$

and

$$(3.13) \qquad \mathrm{MSE}(\tilde{t}_{k\mathrm{r}}) \cong \lambda\mu_{yk}^2\big(\tilde{C}_{yk}^{*} + \tilde{C}_{xk}^{*} - 2\,C_{yxk}^{*}\big),$$

where $\tilde{C}_{yk}^{*2} = \frac{\tilde{\sigma}_{yk}^{*2}}{\mu_{yk}^2}$ and $C_{yxk}^{*} = \frac{\sigma_{yxk}^{*}}{\mu_{yk}\,\mu_{xk}}$. Equation (3.12) shows that the use of RRT to collect response on the dependent variable does not affect the bias of ratio estimator. From (3.5) and (3.13), it can be inferred that $\mathrm{MSE}(\tilde{t}_{k\mathrm{r}}) \leq \mathrm{Var}(\tilde{t}_{k\mathrm{h}})$ if $\tilde{\rho}_{yxk}^{*} \geq \frac{1}{2}\frac{C_{xk}^{*}}{\tilde{C}_{yk}^{*}}$.

## 4.    NUMERICAL STUDY

For numerical validation of our proposed estimators, two real life data sets, one with two small areas and the other with three small areas, are used. The detailed descriptions along with summary statistics of the populations are given in following subsections.

### Blood transfusion data

The data are taken from [37], where $F$, the frequency of donations, is the study variable, $T$ (Time in months since first donation) is taken as the covariate, and a binary variable representing whether he/she donated blood in March 2007 (1 stands for donating blood; 0 stands for not donating blood) is taken as the area membership variable.

### Players head circumference data

This data is taken from [4] which contains physical measures of $N = 90$ players forming three groups, i.e. high school football players (Group 1), college football players (Group 2) and Non-football players (Group 3), each having 30 students. The three groups represent the small areas. The study variable $y$ and the auxiliary variable $x$ respectively are jaw width and ear-to-top-of-head measurement of players. The scrambling variables $a$ and $b$ are generated from Uniform distributions with different ranges.

**Table 1**:    Summary statistics.

| Parameter | Data 1 | | Data 2 | | |
|---|---|---|---|---|---|
| $k$ | 1 | 2 | 1 | 2 | 3 |
| $\theta_k$ | 0.7620 | 0.2380 | 0.3333 | 0.3333 | 0.3333 |
| $\mu_{yk}$ | 4.8018 | 7.7978 | 13.0833 | 10.0800 | 10.9467 |
| $\mu_{xk}$ | 4.8018 | 7.7978 | 14.7333 | 13.4533 | 13.6967 |
| $\sigma_{yk}^2$ | 22.5318 | 64.5916 | 1.0876 | 1.1520 | 1.4577 |
| $\sigma_{xk}^2$ | 605.4251 | 558.3500 | 0.8920 | 0.5702 | 0.3921 |
| $\sigma_{yxk}$ | 76.3885 | 140.5756 | 0.5402 | 0.0870 | 0.0870 |
| $\rho_{yxk}$ | 0.6540 | 0.7402 | 0.3333 | 0.3333 | 0.3333 |

Table 1 provides the summary statistics for the data sets. The theoretical results (TR) are obtained using Variance/MSE expressions given in Section 2. The simulated results (SR) are obtained using following algorithm:

1.    Select a simple random sample of size $n$ (100 and 30 for Populations I and II respectively) without replacement from the populations described above and stratify the populations according to the domain membership variable $d_k$.

2. Record information $y$ and $x$ for all small areas after generating values of scrambling variables $a$ and $b$ from uniform distribution with different ranges.

3. Calculate the values of small area estimators under direct and randomized response technique.

4. Repeat Steps 1–3 50000 times and obtain the simulated Variance, MSE and PRE.

The PRE in Table 2 are computed as $\text{PRE}_\text{r} = \frac{\text{Var}(\hat{t}_{k\text{h}})}{\text{MSE}(\hat{t}_{k\text{r}})}$ and $\text{PRE}_\text{lr} = \frac{\text{Var}(\hat{t}_{k\text{h}})}{\text{Var}(\hat{t}_{k\text{lr}})}$ for $\hat{t}_{k\text{r}}$ and $\tilde{t}_{k\text{lr}}$ are respectively while $\text{PRE}_\text{h}$ is 100 for $\hat{t}_{k\text{h}}$. Table 2 gives the theoretical and simulated PREs of the small area total estimators for different domains under direct response (without using randomized response techniques) with both data sets. PREs in Tables 3 and 4 are obtained in similar manner using the Variances and MSEs under RRT. The theoretical and simulated values of PRE are reported in Tables 2–4 with notations TR and SR respectively.

**Table 2**: PREs of the SAE under direct response.

|         |         | Type | $\text{PRE}_\text{h}$ | $\text{PRE}_\text{r}$ | $\text{PRE}_\text{lr}$ |
|---------|---------|------|------|-----------|-----------|
| Data I  | $k = 1$ | TR   | 100  | 215.864   | 216.839   |
|         |         | SR   | 100  | 217.230   | 218.106   |
|         | $k = 2$ | TR   | 100  | 378.592   | 379.123   |
|         |         | SR   | 100  | 370.443   | 375.775   |
| Data II | $k = 1$ | TR   | 100  | 13853.214 | 13862.216 |
|         |         | SR   | 100  | 12993.771 | 14382.960 |
|         | $k = 2$ | TR   | 100  | 5134.249  | 5137.867  |
|         |         | SR   | 100  | 4855.831  | 5352.376  |
|         | $k = 3$ | TR   | 100  | 6770.974  | 6770.974  |
|         |         | SR   | 100  | 6371.760  | 7076.799  |

From Table 2, one can infer that for both data sets, total estimators under RPM and LPM (see the last two columns) which utilize auxiliary information provide smaller variance than the MSE of Total estimator under HPM. Further, estimator obtained through LPM outperforms the other two competitors in all cases.

Tables 3 and Table 4 give a comparison of the three competing population models in term of PREs for Data I and Data II respectively under randomized response. Going from top to bottom in Tables 3 and 4, we observe that the PREs decrease with increase in variability in the scrambling variables. Also, comparing Table 2 with Tables 3 and 4, we can infer that the efficiency of the domain estimators decreases when using randomized response technique. But that is expected given that RRT introduces noise in the data. Without RRT, the real loss of efficiency will be much larger due to "invisible" response bias.

**Table 3**:   PRE of the SAE under randomized response for Data I.

|        | $a$     | $b$      | Type | $\mathrm{PRE_h}$ | $\mathrm{PRE_r}$ | $\mathrm{PRE_{lr}}$ |
|--------|---------|----------|------|-------|---------|---------|
| $k=1$  | $U(2,3)$ | $U(0,1)$ | TR   | 100   | 210.572 | 211.480 |
|        |         |          | SR   | 100   | 210.726 | 211.501 |
|        |         | $U(0,5)$ | TR   | 100   | 208.031 | 208.907 |
|        |         |          | SR   | 100   | 207.984 | 208.938 |
|        | $U(1,4)$ | $U(0,1)$ | TR   | 100   | 181.451 | 182.026 |
|        |         |          | SR   | 100   | 179.096 | 179.553 |
|        |         | $U(0,5)$ | TR   | 100   | 180.063 | 180.625 |
|        |         |          | SR   | 100   | 177.685 | 178.235 |
| $k=2$  | $U(2,3)$ | $U(0,1)$ | TR   | 100   | 363.439 | 363.921 |
|        |         |          | SR   | 100   | 355.997 | 360.266 |
|        |         | $U(0,5)$ | TR   | 100   | 360.746 | 361.220 |
|        |         |          | SR   | 100   | 351.889 | 357.175 |
|        | $U(1,4)$ | $U(0,1)$ | TR   | 100   | 284.007 | 284.270 |
|        |         |          | SR   | 100   | 275.869 | 277.921 |
|        |         | $U(0,5)$ | TR   | 100   | 282.689 | 282.949 |
|        |         |          | SR   | 100   | 273.902 | 276.442 |

**Table 4**:   PRE of the SAE under randomized response for Data II.

|        | $a$     | $b$      | Type | $\mathrm{PRE_h}$ | $\mathrm{PRE_r}$ | $\mathrm{PRE_{lr}}$ |
|--------|---------|----------|------|-------|---------|---------|
| $k=1$  | $U(2,3)$ | $U(0,1)$ | TR   | 100   | 3740.45 | 3740.97 |
|        |         |          | SR   | 100   | 2624.11 | 2797.29 |
|        |         | $U(0,5)$ | TR   | 100   | 3403.93 | 3404.35 |
|        |         |          | SR   | 100   | 2368.88 | 2524.28 |
|        | $U(1,4)$ | $U(0,1)$ | TR   | 100   | 631.56  | 631.57  |
|        |         |          | SR   | 100   | 435.21  | 460.29  |
|        |         | $U(0,5)$ | TR   | 100   | 623.77  | 623.78  |
|        |         |          | SR   | 100   | 429.72  | 454.59  |
| $k=2$  | $U(2,3)$ | $U(0,1)$ | TR   | 100   | 2578.63 | 2579.54 |
|        |         |          | SR   | 100   | 1991.16 | 2127.92 |
|        |         | $U(0,5)$ | TR   | 100   | 2318.17 | 2318.90 |
|        |         |          | SR   | 100   | 1758.21 | 1875.70 |
|        | $U(1,4)$ | $U(0,1)$ | TR   | 100   | 593.55  | 593.59  |
|        |         |          | SR   | 100   | 424.67  | 447.19  |
|        |         | $U(0,5)$ | TR   | 100   | 582.27  | 582.31  |
|        |         |          | SR   | 100   | 416.73  | 438.87  |
| $k=3$  | $U(2,3)$ | $U(0,1)$ | TR   | 100   | 2930.02 | 2930.02 |
|        |         |          | SR   | 100   | 2203.50 | 2354.74 |
|        |         | $U(0,5)$ | TR   | 100   | 2642.70 | 2642.70 |
|        |         |          | SR   | 100   | 1967.09 | 2095.90 |
|        | $U(1,4)$ | $U(0,1)$ | TR   | 100   | 608.22  | 608.22  |
|        |         |          | SR   | 100   | 432.67  | 455.85  |
|        |         | $U(0,5)$ | TR   | 100   | 598.12  | 598.12  |
|        |         |          | SR   | 100   | 426.19  | 448.82  |

## 5. CONCLUSION

In this study, an attempt for obtaining separate total estimates for the sensitive study variable in each domain (small area) is made using the model relationship between the sensitive study variable and the auxiliary variable. It is observed that the small area total estimators under randomized response techniques possess larger variance (as they should) as compared to the estimators obtained through direct responses. As the privacy and efficiency move in opposite directions, one can't improve both at the same time. Our proposed estimators provide greater efficiency in estimating small area totals when an appropriate model relationship between the study variable and the auxiliary variable is used. Our numerical study with two real life data sets supports the theoretical findings. This is clear from the fact that both $PRE_r$ and $PRE_{lr}$ are greater than $PRE_h$.

## REFERENCES

[1] AHMED, S.; SHABBIR, J. and GUPTA, S. (2017). Use of scrambled response model in estimating the finite population mean in presence of non response when coefficient of variation is known, *Communication in Statistics – Theory and Methods*, **46**(17), 8435–8449.

[2] BASU, D. (1971). *An Essay on the Logical Foundations of Survey Sampling, Part I*. In "Foundations of Statistical Inference" (V.P. Godambe and D. Sprott, Eds.), Holt, Rinehart and Winston, Toronto, 203–233.

[3] BAR-LEV, S.K.; BOBOVITCH, E. and BOUKAI, B. (2004). A note on randomized response models for quantitative data, *Metrika*, **60**, 255–260.

[4] BRYCE, G.R. (1980). "Some observations on the analysis of growth curves", Paper No. SD-025-R, Brigham Young University, Department of Statistics, **41**, 627–636.

[5] CHAMBERS, R. and CLARK, R. (2012). An introduction to model-based survey sampling with applications, *OUP Oxford*, **37**.

[6] CHAUDHURI, A. and MUKERJEE, R. (1988). *Randomized Response: Theory and Techniques*, Marcel Dekker, New York.

[7] CHAUDHURI, A. and ROY, D. (1997). Model assisted survey sampling strategies with randomized response, *Journal of Statistical Planning Inference*, **60**, 61–68.

[8] CUMBERLAND, W.G. and ROYALL, R.M. (1981). Prediction models and unequal probability sampling, *Journal of the Royal Statistical Society: Series B (Methodological)*, **43**(3), 353–367.

[9]   DIANA, G. and PERRI, P.F. (2009). Estimating a sensitive proportion through randomized response procedures based on auxiliary information, *Statistical Papers*, **50**, 661–672.

[10]  DIANA, G. and PERRI, P.F. (2010). New scrambled response models for estimating the mean of a sensitive quantitative character, *Journal of Applied Statistics*, **37**, 1875–1890.

[11]  DIANA, G. and PERRI, P.F. (2011). A class of estimators for quantitative sensitive data, *Statistical Papers*, **52**, 633–650.

[12]  EICHHORN, B.H. and HAYRE, L.S. (1983). Scrambled randomized response methods for obtaining sensitive quantitative data, *Journal of Statistical Planning Inference*, **7**, 307–316.

[13]  FAY, R.E. and HERRIOT, R.A. (1979). Estimation of income for small places: An application of James Stein procedures to census data, *Journal of the American Statistical Association*, **74**, 268–277.

[14]  FULLER, W.A. (1970). "Simple Estimators for the Mean of Skewed Populations", Technical report, Iowa State University, Dept. of Statistics.

[15]  GODAMBE, V.P. (1955). A unified theory of sampling from finite populations, *Journal of the Royal Statistical Society: Series B*, **17**, 269–278.

[16]  GODAMBE, V.P. and JOSHI, V.M. (1965). Admissibility and Bayes estimation in sampling finite populations, I, *The Annals of Mathematical Statistics*, **36**, 1707–1722.

[17]  GREENBERG, B.G.; ABUL-ELA, A.L.A.; SIMMONS, W.R. and HORVITZ, D.G. (1969). The unrelated question randomized response model: theoretical framework, *Journal of the American Statistical Association*, **64**, 520–539.

[18]  GREENBERG, B.G.; KUEBLER, R.R.; ABERNATHY, J.R. and HORVITZ, D.G. (1971). Application of the randomized response technique in obtaining quantitative data, *Journal of the American Statistical Association*, **66**, 243–250.

[19]  GUPTA, S.; GUPTA, B. and SINGH, S. (2002). Estimation of sensitive level of personal interview survey questions, *Journal of Statistical Planning Inference*, **100**, 239–247.

[20]  GUPTA, S.N.; SHABBIR, J. and SEHRA, S. (2010). Mean and sensitivity estimation in optional randomized response models, *Journal of Statistical Planning Inference*, **140**, 2870–2874.

[21]  HUANG, K.C. (2010). Unbiased estimators of mean, variance and sensitivity level for quantitative characteristics in finite population sampling, *Metrika*, **71**, 341–352.

[22]  MANGAT, N.S. and SINGH, R. (1990). An alternative randomized response procedure, *Biometrika*, **77**, 439–442.

[23]  PFEFFERMANN, D.; BELL, P. and SIGNORELLI, D. (1996). Labour force trend estimation in small areas, *Proceedings of the Annual Research Conference, US Bureau of the Census*, 407–431.

[24]  PURCELL, N.I. and KISH, L. (1980). Postcensal estimates for local areas (or domains), *International Statistical Review*, **48**, 3–18.

[25]  RAO, J. (1994). Small area estimation by combining time series and cross sectional data, *Canadian Journal of Statistics*, **22**, 511–528.

[26]  RAO, J. (2003). *Small Area Estimation*, New York, Wiley.

[27]  ROYALL, R.M. (1970). An old approach to finite population sampling theory, *Journal of the American Statistical Association*, **63**, 1269–1279.

[28]  ROYALL, R.M. (1992). The model based (prediction) approach to finite population sampling theory, *Lecture Notes – Monograph Series*, **17**, 225–240.

[29]  SARNDAL, C.E.; THOMSEN, I.; HOEM, J.M.; LINDLEY, D.V.; BARNDORFF-NIELSEN, O. and DALENIUS, T. (1978). Design-based and model-based inference in survey sampling [with discussion and reply], *Scandinavian Journal of Statistics*, 27–52.

[30]  SMITH, T.M.F. (1976). The foundations of survey sampling: A review, *Journal of the Royal Statistical Society: Series A*, **139**, 183–195.

[31]  SMITH, T.M.F. (1983). On the validity of inferences from non-random samples, *Journal of the Royal Statistical Society: Series A (General)*, **146**(4), 394–403.

[32]  VALLIANT, R. (2009). Model-Based Prediction of Finite Population Totals, *Sample Surveys: Inference and Analysis*, **29B**, 23–31.

[33]  VALLIANT, R.; DORFMAN, A.H. and ROYALL, M. (2000). *Finite Population Sampling and Inference: A Prediction Approach*, John Wiley & Sons, New York.

[34]  WARNER, S.L. (1965). Randomized response: A survey technique for eliminating evasive answer bias, *Journal of the American Statistical Association*, **60**, 63–69.

[35]  WARNER, S.L. (1971). The linear randomized response model, *Journal of the American Statistical Association*, **66**(336), 884–888.

[36]  WHITWORTH, A.; CARTER, E.; BALLAS, D. and MOON, G. (2017). Estimating uncertainty in spatial microsimulation approaches to small area estimation: A new approach to solving an old problem. Computers, *Environment and Urban Systems*, **63**, 50–57.

[37]  YEH, I.-C.; YANG, KING-JANG and TING, T.-M. (2008). Knowledge discovery on RFM model using Bernoulli sequence, *Expert Systems with Applications*, **36**, 5866–5871.

[38]  YOU, Y. (2008). Small area estimation using area level models with model checking and applications, *Proceedings of the Survey Methods Section*, Statistical Society of Canada.

# REVSTAT – STATISTICAL JOURNAL

**Background**

Statistics Portugal (INE, I.P.), well aware of how vital a statistical culture is in understanding most phenomena in the present-day world, and of its responsibility in disseminating statistical knowledge, started the publication of a scientific statistical journal called Revista de Estatística. The original language used in this publication was Portuguese and the idea behind it was to publish it, three times a year, containing original research results, and application studies, namely in the economic, social and demographic fields.

In 1998 it was decided that the publication should also include papers in English. This step was taken to achieve a broader dissemination, and to encourage foreign contributors to submit their work for publication.

At the time, the Editorial Board was mainly comprised of Portuguese university professors. It is now comprised of international university faculties and this has been the first step aimed at changing the character of Revista de Estatística from a national to an international scientific journal.

We have also initiated a policy of publishing special volumes that may be thematic highlighting areas of interest or associated with scientific events in Statistics. For example, in 2001, a special issue of Revista de Estatística was published containing three volumes of extended abstracts of the invited contributed papers presented at the 23rd European Meeting of Statisticians.

In 2003, the name of the Journal has been changed to REVSTAT - STATISTICAL JOURNAL, now fully published in English, with a prestigious international editorial board, aiming to become a reference scientific journal that promotes the dissemination of relevant research results in Statistics.

The editorial policy of REVSTAT Statistical Journal is mainly placed on the originality and importance of the research.

All articles consistent with REVSTAT aims and scope will undergo scientific evaluation by at least two reviewers, one from the Editorial Board and another external.

The only working language allowed is English.

## Abstract and Indexing Services

The REVSTAT is covered by the following abstracting/indexing services:

- Current Index to Statistics

- Google Scholar

- Mathematical Reviews® (MathSciNet®)

- Science Citation Index Expanded

- Zentralblatt für Mathematic

- Scimago Journal & Country Rank

- Scopus

## Instructions to Authors

### Articles must be written in English and will be submitted according to the following guidelines:

The corresponding author sends the manuscript in PDF format to the Executive Editor (revstat@ine.pt) with the Subject "New Submission to REVSTAT"; a MS#REVSTAT reference will be assigned later.

Optionally, in a mail cover letter, authors are welcome to suggest one of the Editors or Associate Editors, whose opinion may be considered suitable to be taken into account.

The submitted manuscript should be original and not have been previously published nor about to be published elsewhere in any form or language, avoiding concerns about self-plagiarism'.

Content published in this journal is peer-reviewed (Single Blind).

All research articles will be refereed by at least two researchers, including one from the Editorial Board unless the submitted manuscript is judged unsuitable for REVSTAT or does not contain substantial methodological novelty, in which case is desk rejected.

Manuscripts should be typed only in black, in double-spacing, with a left margin of at least 3 cm, with numbered lines, and with less than 25 pages. Figures (minimum of 300dpi) will be reproduced online in colours, if produced this way; however, authors should take into account that the printed version is always in black and grey tones.

The first page should include the name, ORCID iD (optional), Institution, country, and mail-address of the author(s) and a summary of fewer than one hundred words, followed by a maximum of six keywords and the AMS 2000 subject classification.

Authors are encouraged to submit articles using LaTeX, in the REVSTAT style, which is available at the LaTeX2e MACROS webpage.

References about the format and other useful information on the submission are available in the LaTeX2e Templates page.

Acknowledgments of people, grants or funds should be placed in a short section before the References title page. Note that religious beliefs, ethnic background, citizenship and political orientations of the author(s) are not allowed in the text.

Supplementary files (in REVSTAT style) may be published online along with an article, containing data, programming code, extra figures, or extra proofs, etc; however, REVSTAT is not responsible for any supporting information supplied by the author(s).

Any contact with REVSTAT must always contain the assigned REVSTAT reference number.

## Accepted papers

Authors of accepted papers are requested to provide the LaTex files to the Secretary of the REVSTAT revstat@ine.pt. The authors should also mention if figure files were included, and submit electronic figures separately in .gif, .jpg, .png or .pdf format. Figures must be a minimum of 300dpi.

## Copyright and reprints

Upon acceptance of an article, the author(s) will be asked to transfer copyright of the article to the publisher, Statistics Portugal, in order to ensure the widest possible dissemination of information, namely through the Statistics Portugal website (http://www.ine.pt).

After assigning copyright, authors may use their own material in other publications provided that REVSTAT is acknowledged as the original place of publication. The Executive Editor of the Journal must be notified in writing in advance.