



INSTITUTO NACIONAL DE ESTATÍSTICA
STATISTICS PORTUGAL

REVSTAT

Statistical Journal



REVSTAT

Statistical Journal

Catálogo Recomendada

REVSTAT. Lisboa, 2003-
Revstat : statistical journal / ed. Instituto Nacional
de Estatística. - Vol. 1, 2003- . - Lisboa I.N.E.,
2003- . - 30 cm
Trimestral. - Continuação de : Revista de Estatística =
ISSN 0873-4275. - edição exclusivamente em inglês
ISSN 1645-6726 ; e-ISSN 2183-0371

CREDITS

- EDITOR-IN-CHIEF

- *Isabel Fraga Alves*

- CO-EDITOR

- *Giovani L. Silva*

- ASSOCIATE EDITORS

- *Marília Antunes*
- *Barry Arnold*
- *Narayanaswamy Balakrishnan*
- *Jan Beirlant*
- *Graciela Boente (2019-2020)*
- *Paula Brito*
- *Vanda Inácio de Carvalho*
- *Arthur Charpentier*
- *Valérie Chavez-Demoulin*
- *David Conesa*
- *Charmaine Dean*
- *Jorge Milhazes Freitas*
- *Alan Gelfand*
- *Stéphane Girard*
- *Wenceslao Gonzalez-Manteiga*
- *Marie Kratz*
- *Victor Leiva*
- *Maria Nazaré Mendes-Lopes*
- *Fernando Moura*
- *John Nolan*
- *Paulo Eduardo Oliveira*
- *Pedro Oliveira*
- *Carlos Daniel Paulino (2019-2021)*
- *Arthur Pewsey*
- *Gilbert Saporta*
- *Alexandra M. Schmidt*
- *Julio Singer*

- *Manuel Scotto*

- *Lisete Sousa*

- *Milan Stehlik*

- *María Dolores Ugarte*

- FORMER EDITOR-IN-CHIEF

- *M. Ivette Gomes*

- FORMER CO-EDITOR

- *M. Antónia Amaral Turkman*

- EXECUTIVE EDITOR

- *José A. Pinto Martins*

- FORMER EXECUTIVE EDITOR

- *Maria José Carrilho*

- *Ferreira da Cunha*

- SECRETARIAT

- *José Cordeiro*

- *Olga Bessa Mendes*

- PUBLISHER

- *Instituto Nacional de Estatística, I.P. (INE, I.P.)*

Web site: <http://www.ine.pt>

- COVER DESIGN

- *Mário Bouçadas, designed on the stain glass window at INE by the painter Abel Manta*

- LAYOUT AND GRAPHIC DESIGN

- *Carlos Perpétuo*

- PRINTING

- *Instituto Nacional de Estatística, I.P.*

- EDITION

- *140 copies*

- LEGAL DEPOSIT REGISTRATION

- *N.º 191915/03*

- PRICE [VAT included]

- *€ 9,00*



INDEX

Nonparametric Regression Based on Discretely Sampled Curves	
<i>Liliana Forzani, Ricardo Fraiman and Pamela Llop</i>	1
A Transition Model for Analysis of Zero-Inflated Longitudinal Count Data Using Generalized Poisson Regression Model	
<i>Taban Baghfalaki and Mojtaba Ganjali</i>	27
Compound Power Series Distribution with Negative Multinomial Sums: Characterisation and Risk Process	
<i>Pavlina Jordanova, Monika Petkova and Milan Stehlik</i>	47
Characterization of the Maximum Probability Fixed Marginals $r \times c$ Contingency Tables	
<i>Francisco Requena</i>	71
On the Occurrence of Boundary Solutions in Two-Way Incomplete Tables	
<i>Sayan Ghosh and Palaniappan Vellaisamy</i>	89
Depth-Based Signed-Rank Tests for Bivariate Central Symmetry	
<i>Sakineh Dehghan and Mohammad Reza Faridrohani</i>	109
Prediction Intervals for Time Series and Their Applications to Portfolio Selection	
<i>Shih-Feng Huang and Hsiang-Ling Hsu</i>	131

NONPARAMETRIC REGRESSION BASED ON DISCRETELY SAMPLED CURVES

Authors: LILIANA FORZANI
– Facultad de Ingeniería Química, UNL and researcher of CONICET, Argentina
liliana.forzani@gmail.com

RICARDO FRAIMAN
– Centro de Matemática, Facultad de Ciencias (UdelaR), Uruguay
fraimanricardo@gmail.com

PAMELA LLOP
– Facultad de Ingeniería Química, UNL and researcher of CONICET, Argentina
lloppamela@gmail.com

Received: January 2017

Revised: August 2017

Accepted: September 2017

Abstract:

- In the context of nonparametric regression, we study conditions under which the consistency (and rates of convergence) of estimators built from discretely sampled curves can be derived from the consistency of estimators based on the unobserved whole trajectories. As a consequence, we derive asymptotic results for most of the regularization techniques used in functional data analysis, including smoothing and basis representation.

Key-Words:

- *nonparametric regression; functional data; discrete curves.*

AMS Subject Classification:

- 62G08, 62M99.

1. INTRODUCTION

Technological progress in collecting and storing data provides datasets recorded at finite grids of points that become denser and denser over time. Although in practice data always comes in the form of finite dimensional vectors, from the theoretical point of view, the classic multivariate techniques are not well suited to deal with data which, essentially, is infinite dimensional and whose observations within the same curve are highly correlated.

From a practical point of view, a commonly used technique to treat this kind of data is to transform the (observed) discrete values into a function via smoothing or a series approximations (see [5], [21], [24, 25, 26], or chapter 9 of [13] and the references therein). For the analysis, we can use the intrinsic infinite dimensional nature of the data and assume the existence of continuous underlying stochastic processes which are observed ideally at every point. In this context, the theoretical analysis is performed on the functional space where they take values (see [15]). In what follows, we will refer to this last setting as the *full model*.

Nonparametric regression is an important tool in functional data analysis (FDA) which has received considerable attention from different authors in both settings. For the full model, consistency results have been obtained by, among others, [1], [3], [4], [7], [10], [15], [22], and [23]. In particular, [16] (see also the Corrigendum [17]) prove a consistency result close to universality for the kernel (with random bandwidth) estimator. The first contribution of the present paper will be to prove the consistency of the k -nearest neighbor with kernel regression estimator (Proposition 2.2) when the full trajectories are observed. This family, considered by [12], combines the smoothness properties of the kernel function with the locality properties of the k -nearest neighbors distances.

Regarding regression when discretized curves are available, [19] study the mean square consistency of the kernel estimator when the sample size as well as the grid size discretization go to infinity. More precisely, from independent realizations of a random process with continuous covariance structure, they estimate the regression function, assuming its smoothness. Under the same assumptions, but using interpolation of the data, [27], in a mainly practical approach, propose a method to estimate the regression function via smoothing splines (see also [20]). More recently, [8] establish minimax rates of convergence of estimators of the mean based on discretized sampled data while [9] establish the minimax rates of convergence for the covariance operator when data are observed on a lattice (see also [18] for the problem of principal components analysis for longitudinal data). In this context it is natural to assess the relation between the *ideal* nonparametric regression estimator constructed with the entire set of curves and the one computed with the discretized sample. In this direction, we are interested in addressing the following question:

- Under what conditions can the consistency (and rates of convergence) of the estimate computed with the discretized trajectories be derived from the consistency of the estimate based on the full curves?

Clearly, the asymptotic results for estimates computed with the discretized sample will not be a direct consequence of those for the full model. However, we provide reasonable conditions in order to still get the consistency and find rates of convergence of the estimator.

In this context we state the results for the well known kernel and k -nearest neighbor with kernel estimators. These results are a consequence of a more general result, which, besides discretization, also includes the cases of regularization via smoothing and basis representation.

This paper is organized as follows: In Section 2 we state the consistency of the k -nearest neighbor with kernel estimator in the infinite dimensional setting (for the full model). This result is not only interesting by itself but also, it will be used to prove consistency results when discretely sample data are available. In Section 3 we provide conditions for the consistency of the kernel and k -nearest neighbor with kernel estimators when we do not observe the whole trajectories but only a function of them (Theorems 3.1 and 3.2). In Section 4 the results for discretization, smoothing and basis representation are obtained as a consequence of Theorems 3.1 and 3.2. Finally, in Section 5 we perform a small simulation study where we compare the behaviour of the estimators computed with the discretized trajectories and with the full curves. Proofs are given in Appendices A and B.

2. CONSISTENCY RESULTS FOR FULLY OBSERVED CURVES

In this section we provide two L^2 -consistency results for the full model, i.e., when ideally all trajectories are observed at every point of the interval $[0, 1]$. The first one corresponds to kernel estimates, and was obtained in [16], while the second one for k -NN with kernel estimates is derived in the present paper. Both results will be used, in Section 3, to prove the consistency of that estimators when only discretely sampled curves in $[0, 1]$ are observed.

We will use the notation $f \lesssim g$ when there exists a constant $C > 0$ such that $f \leq Cg$ and $f \approx g$ if there exists a constant $C > 0$ such that $f = Cg$.

Let (\mathcal{H}, d) be a separable metric space and let $(\mathcal{X}_1, Y_1), \dots, (\mathcal{X}_n, Y_n)$ be independent identically distributed (i.i.d.) random elements in $\mathcal{H} \times \mathbb{R}$ with the same law as the pair (\mathcal{X}, Y) fulfilling the model:

$$(2.1) \quad Y = \eta(\mathcal{X}) + e,$$

where the error e satisfies $\mathbb{E}_{e|\mathcal{X}}(e|\mathcal{X}) = 0$ and $\text{var}_{e|\mathcal{X}}(e|\mathcal{X}) = \sigma^2 < \infty$. In this context, the regression function $E(Y|\mathcal{X}) = \eta(\mathcal{X})$ can be estimated by

$$(2.2) \quad \hat{\eta}_n(\mathcal{X}) = \sum_{i=1}^n W_{ni}(\mathcal{X}) Y_i,$$

where the weights $W_{ni}(\mathcal{X}) = W_{ni}(\mathcal{X}, \mathcal{X}_1, \dots, \mathcal{X}_n) \geq 0$ and $\sum_{i=1}^n W_{ni}(\mathcal{X}) = 1$. In this paper, we first consider the weights corresponding to the family of kernel estimators given by

$$(2.3) \quad W_{ni}(\mathcal{X}) = \frac{K\left(\frac{d(\mathcal{X}, \mathcal{X}_i)}{h_n(\mathcal{X})}\right)}{\sum_{j=1}^n K\left(\frac{d(\mathcal{X}, \mathcal{X}_j)}{h_n(\mathcal{X})}\right)},$$

where K is a regular kernel, i.e., there are constants $0 < c_1 < c_2 < \infty$ such that $c_1 \mathbb{I}_{[0,1]}(u) \leq K(u) \leq c_2 \mathbb{I}_{[0,1]}(u)$. Here $0/0$ is assumed to be 0. In this general setting, [16] proved the following result.

Proposition 2.1 (Theorem 5.1 in [16]). *Assume that*

K1) K is a regular and Lipschitz kernel;

F1) (\mathcal{H}, d) is a separable metric space;

F2) $\{(\mathcal{X}_i, Y_i)\}_{i \geq 1}$ are i.i.d. random elements with the same law as the pair $(\mathcal{X}, Y) \in \mathcal{H} \times \mathbb{R}$ fulfilling model (2.1) with, for each $i = 1, \dots, n$, joint distribution $\mathbb{P}_{\mathcal{X}, \mathcal{X}_i}$;

F3) μ is a Borel probability measure of \mathcal{X} and $\eta \in L^2(\mathcal{H}, \mu) = \{f: \mathcal{H} \rightarrow \mathbb{R}: \int_{\mathcal{H}} f^2(z) d\mu(z) < \infty\}$ is a bounded function which satisfies the Besicovitch condition:

$$(2.4) \quad \lim_{\delta \rightarrow 0} \frac{1}{\mu(\mathcal{B}(\mathcal{X}, \delta))} \int_{\mathcal{B}(\mathcal{X}, \delta)} |\eta(z) - \eta(\mathcal{X})| d\mu(z) = 0,$$

in probability, where $\mathcal{B}(\mathcal{X}, \delta)$ is the closed ball of center \mathcal{X} and radius δ with respect to d .

For any $x \in \text{supp}(\mu)$ and any sequence $h_n(x) \rightarrow 0$ such that $\frac{n\mu(\mathcal{B}(x, h_n(x)))}{\log n} \rightarrow \infty$, the estimator given in (2.2) with weights given in (2.3) satisfies

$$\lim_{n \rightarrow \infty} \mathbb{E}((\hat{\eta}_n(\mathcal{X}) - \eta(\mathcal{X}))^2) = 0.$$

Remark 2.1. The Besicovitch condition in F3 is a differentiation type condition which, as is well known, in finite dimensional spaces automatically holds for any integrable function η . Unfortunately, it is no longer true in infinite dimensional spaces and it can be proved, for instance, that it is necessary in order to get the L_1 -consistency of uniform kernel estimates (see Proposition 5.1 in [16]). However, it holds in a general setting if, for instance, the function η is continuous. For a deeper reading on this topic see [10] or [16].

Remark 2.2. Note that for $x \in \text{supp}(\mu)$ the consistency of this estimator holds for every sequence $\tilde{h}_n(x) \rightarrow 0$ such that $\tilde{h}_n(x) \geq h_n(x)$, where $h_n(x)$ is given in Proposition 2.1, since if $\tilde{h}_n(x) \geq h_n(x)$, then $\frac{n\mu(\mathcal{B}(x, \tilde{h}_n(x)))}{\log n} \geq \frac{n\mu(\mathcal{B}(x, h_n(x)))}{\log n} \rightarrow \infty$.

The existence of a sequence verifying $\frac{n\mu(\mathcal{B}(x, h_n(x)))}{\log n} \rightarrow \infty$ in Proposition 2.1 follows from the next lemma.

Lemma 2.1 (Lemma A.5 in [16]). *For any $x \in \text{supp}(\mu)$, there exists a sequence of positive real numbers $h_n(x) \rightarrow 0$ such that $\frac{n\mu(\mathcal{B}(x, h_n(x)))}{\log n} \rightarrow \infty$.*

Let $H_n(x)$ be the distance from x to its k_n -nearest neighbor among $\{\mathcal{X}_1, \dots, \mathcal{X}_n\}$. Recall that the k_n -nearest neighbor of x among $\{\mathcal{X}_1, \dots, \mathcal{X}_n\}$ is the sample point \mathcal{X}_i reaching the k_n -th smallest distance to x in the sample. Then, when the bandwidth in (2.3) is given by $H_n(x)$, we obtain the family of k_n -nearest neighbor (k -NN) with kernel estimates. For the uniform kernel, the consistency of the estimator was proven in [16], Theorem 4.1. For more general kernels, the consistency could be a consequence of Proposition 2.1 if we can prove that $H_n(x) \rightarrow 0$ and $\frac{n\mu(\mathcal{B}(x, H_n(x)))}{\log n} \rightarrow \infty$. Although it can be proved that $H_n(x) \rightarrow 0$ (see [16], Lemma A.4 stated below) the condition $\frac{n\mu(\mathcal{B}(x, H_n(x)))}{\log n} \rightarrow \infty$ is not necessary true for $H_n(x)$. However, as we will see in Proposition 2.2, we can still prove the mean square consistency of this estimator under the same weak conditions as in Proposition 2.1.

Lemma 2.2 (Lemma A.4 in [16]). *Let \mathcal{H} be a separable metric space, μ a Borel probability measure, and $\{\mathcal{X}_i\}_{i=1}^n$ a random sample of \mathcal{X} . If $x \in \text{supp}(\mu)$ and k_n is a sequence of positive real numbers such that $k_n \rightarrow \infty$ and $k_n/n \rightarrow 0$, then $H_n(x) \rightarrow 0$.*

Proposition 2.2. *Assume K1, F1–F3 hold. Let k_n be a sequence of positive real numbers such that $k_n \rightarrow \infty$, $k_n/n \rightarrow 0$ and let $H_n(x)$ be the distance from x to its k_n -nearest neighbor among $\{\mathcal{X}_1, \dots, \mathcal{X}_n\}$. Then, the estimator given by (2.2) with weights given in (2.3) is mean square consistent for any sequence $h_n(x) \rightarrow 0$ such that $h_n(x) \geq H_n(x)$, $x \in \text{supp}(\mu)$.*

Remark 2.3. Observe that, unlike [15] or [7], we ask d to be a metric not a semi-metric (which is a milder condition). Nevertheless, we do not ask for conditions neither on small ball probabilities nor on the smoothness of the regression function as in the cited papers. Further study is needed to extend our results to the case of semi-metrics.

3. CONSISTENCY RESULTS FOR DISCRETELY SAMPLED CURVES

In this section we will assume that we are not able to observe the whole trajectories \mathcal{X}_i in \mathcal{H} given in F2, but only a function of them. As we will see in Section 4, different choices of that function will correspond to discretizations, eigenfunction expansions, or smoothing. In this context, the weights of the estimator given in (2.3) cannot be computed because we have not a distance d defined for the discretized sample curves (as a consequence, we do not have the validity of the Besicovitch condition (2.4) for the discretized data) or a bandwidth h_n .

We are interested in defining an estimator and proving its consistency in this setting. For that, let us consider the following assumptions:

H1) (\mathcal{H}, d) is a separable (metric) Hilbert space and $F: \mathcal{H} \rightarrow \mathcal{H}$ is a function such that, for each $i = 1, \dots, n$, $F(\mathcal{X}_i) = \mathcal{X}_i^p$;

H2) $d_p: \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$ is a semi-metric in \mathcal{H} defined by $d_p(\mathcal{X}, \mathcal{Y}) = d(\mathcal{X}^p, \mathcal{Y}^p)$ such that there exists a sequence $c_{n,p} \rightarrow 0$ as $n, p \rightarrow \infty$ satisfying, for each $i = 1, \dots, n$,

$$(3.1) \quad n^2 \mathbb{E}_{\mathcal{X}} \left(\mathbb{P}_{\mathcal{X}_i|\mathcal{X}}^2 \left(|d(\mathcal{X}, \mathcal{X}_i) - d_p(\mathcal{X}, \mathcal{X}_i)| \geq c_{n,p} \mid \mathcal{X} \in \text{supp}(\mu) \right) \right) \rightarrow 0.$$

Here, $\mathbb{P}_{\mathcal{Y}|\mathcal{X}}^2(\cdot)$ means the square of $\mathbb{P}_{\mathcal{Y}|\mathcal{X}}(\cdot)$.

Remark 3.1. Observe that in H1 neither \mathcal{H} nor F change with the sample. This implies that in this case, the functional data falls into the category of sparsely and regularly sampled data.

The estimator of η based on $\{(\mathcal{X}_i^p, Y_i)\}_{i=1}^n$ will be defined as in (2.2) and (2.3) but with the semi-metric d_p instead of the metric d . More precisely, for $h_{n,p}(\mathcal{X}) > 0$, we define

$$(3.2) \quad \hat{\eta}_{n,p}(\mathcal{X}) = \frac{\sum_{i=1}^n K\left(\frac{d_p(\mathcal{X}, \mathcal{X}_i)}{h_{n,p}(\mathcal{X})}\right) Y_i}{\sum_{j=1}^n K\left(\frac{d_p(\mathcal{X}, \mathcal{X}_j)}{h_{n,p}(\mathcal{X})}\right)}.$$

For this estimator, we state the following two asymptotic results.

Theorem 3.1. *Assume K1, F2, F3, H1 and H2 hold.*

- (a) *(Kernel estimator) For any $x \in \text{supp}(\mu)$, let $h_n^*(x) \rightarrow 0$ be a sequence of positive real numbers such that $\frac{n\mu(\mathcal{B}(x, h_n^*(x)))}{\log n} \rightarrow \infty$. Then, for $c_{n,p}$ given in H2 and $h_{n,p}(x) \rightarrow 0$ such that there exists a sequence $h_n(x) \rightarrow 0$, $h_n(x) \geq h_n^*(x)$ satisfying:*

$$(H3.1) \quad \mathbb{E}_{\mathcal{X}}(c_{n,p}^2/h_n^2(\mathcal{X})) \rightarrow 0 \text{ as } n, p \rightarrow \infty;$$

$$(H3.2) \quad c_{n,p} \leq h_{n,p}(x) - h_n(x) \leq C_2 c_{n,p} \text{ for } C_2 \geq 1;$$

we have

$$(3.3) \quad \lim_{n,p \rightarrow \infty} \mathbb{E}((\hat{\eta}_{n,p}(\mathcal{X}) - \eta(\mathcal{X}))^2) = 0.$$

- (b) *(k_n -NN with kernel estimator) Let $c_{n,p}$ given in H2 and $H_n(x)$ the distance from x to its k_n -nearest neighbor among $\{\mathcal{X}_1, \dots, \mathcal{X}_n\}$. For any $x \in \text{supp}(\mu)$, let $h_{n,p}(x) \rightarrow 0$ be such that there exists a sequence $h_n(x) \rightarrow 0$, $h_n(x) \geq H_n(x)$ satisfying assumptions (H3.1) and (H3.2). Then, for $k_n \rightarrow \infty$ and $k_n/n \rightarrow 0$ we have (3.3).*

Remark 3.2. Observe that the sequence $h_n^*(x)$ in Theorem 3.1 always exists by Lemma 2.1. In addition, under H2, it is always possible to choose a sequence $h_{n,p}(x) \rightarrow 0$ fulfilling the conditions in Theorem 3.1. Indeed, taking $h_n(x) = \max\{h_n^*(x), \sqrt{c_{n,p}}\}$ and $h_{n,p}(x) = h_n(x) + Cc_{n,p}$, with $C \geq 1$, we have that $h_n(x) \rightarrow 0$, $h_{n,p}(x) \rightarrow 0$, $h_n(x) \geq h_n^*(x)$, (H3.1) holds since $h_n(x) \geq \sqrt{c_{n,p}}$ and (H3.2) holds by definition of $h_{n,p}(x)$. The same happens if instead of taking $h_n^*(x)$ we take $H_n(x)$.

Theorem 3.2. *Under the assumptions of Theorem 3.1, let $\gamma_n \rightarrow \infty$ as $n \rightarrow \infty$ be such that, as $n, p \rightarrow \infty$,*

$$(a) \quad \mathbb{E}_{\mathcal{X}}\left(\gamma_n \left(\frac{c_{n,p}}{h_n(\mathcal{X})}\right)^2\right) \rightarrow 0;$$

$$(b) \quad \gamma_n n^2 \mathbb{E}_{\mathcal{X}}\left(\mathbb{P}_{\mathcal{X}_i|\mathcal{X}}^2\left(|d(\mathcal{X}, \mathcal{X}_i) - d_p(\mathcal{X}, \mathcal{X}_i)| \geq c_{n,p} \mid \mathcal{X} \in \text{supp}(\mu)\right)\right) \rightarrow 0, \quad \text{for each } i = 1, \dots, n.$$

Then

$$\lim_{n \rightarrow \infty} \mathbb{E}(\gamma_n (\hat{\eta}_n(\mathcal{X}) - \eta(\mathcal{X}))^2) = 0$$

implies

$$\lim_{n,p \rightarrow \infty} \mathbb{E}(\gamma_n (\hat{\eta}_{n,p}(\mathcal{X}) - \eta(\mathcal{X}))^2) = 0.$$

4. PARTICULAR CASES

In this section we provide definitions of \mathcal{H} and d_p for discretization, smoothing, and eigenfunction expansions, which satisfy conditions H1 and H2. Then, for any sequence $h_{n,p}(x) \rightarrow 0$ satisfying (H3.1) and (H3.2) in Theorem 3.1, we get the consistency of $\hat{\eta}_{n,p}$ as a consequence of the consistency results for $\hat{\eta}_n$ in the full model.

Consider the case where the elements of the dataset are curves in $L^2([0, 1])$ that are only observed at a discrete set of points in the interval $[0, 1]$. More precisely, let us assume that $\{\mathcal{X}_i\}_{i=1}^n$ are observed only at some points: $(\mathcal{X}_i(t_1), \dots, \mathcal{X}_i(t_{p+1}))$ where $0 = t_1 < t_2 < \dots < t_{p+1} = 1$, which for simplicity we will assume are equally spaced, i.e., $\Delta t = t_{i+1} - t_i = 1/p$. In this case, we will need to require the trajectories to satisfy some regularity condition. More precisely, we will assume that \mathcal{X} is a random element of $\mathcal{H} \doteq H^1([0, 1])$, the Sobolev space defined as

$$H^1([0, 1]) = \left\{ f: [0, 1] \rightarrow \mathbb{R}: f \text{ and } Df \in L^2([0, 1]) \right\},$$

where Df is the weak derivative of f , i.e., Df is a function in $L^2([0, 1])$ which satisfies

$$\int_0^1 f(t) D\phi(t) dt = - \int_0^1 Df(t) \phi(t) dt, \quad \forall \phi \in C_0^\infty.$$

In this space, the norm is defined by

$$\|f\|_{H^1([0,1])} = \|f\|_{L^2([0,1])} + \|Df\|_{L^2([0,1])}.$$

In this setting, we will prove consistency for the semi-metrics d_p given below.

4.1. Discretization

Consider the semi-metric

$$d_p(\mathcal{X}, \mathcal{X}_1) = d(\mathcal{X}^p, \mathcal{X}_1^p) = \left(\frac{1}{p} \sum_{j=1}^p |\mathcal{X}(t_j) - \mathcal{X}_1(t_j)|^2 \right)^{1/2},$$

where $\mathcal{X}^p(t) = F(\mathcal{X})(t) = \sum_{j=1}^p \phi_j(t) \mathcal{X}(t_j)$ with $\phi_j(t) = \mathbb{I}_{[t_j, t_{j+1})}(t)$. In this case, consistency will hold for any sequence $c_{n,p} \rightarrow 0$ as $n, p \rightarrow \infty$ such that $n^2 \mathbb{P}_{\mathcal{X}, \mathcal{X}_1}(\|\mathcal{X}\|_{\mathcal{H}} + \|\mathcal{X}_1\|_{\mathcal{H}} \geq p c_{n,p}) \rightarrow 0$.

4.2. Kernel smoothing

Let us consider now the semi-metric

$$d_p(\mathcal{X}, \mathcal{X}_1) = d(\mathcal{X}^p, \mathcal{X}_1^p) = \left(\int_0^1 |\mathcal{X}^p(t) - \mathcal{X}_1^p(t)|^2 dt \right)^{1/2},$$

where $\mathcal{X}^p(t) = F(\mathcal{X})(t) = \sum_{j=1}^p \phi_j(t) \mathcal{X}(t_j)$ with $\phi_j(t) = \frac{K(|t-t_j|/h)}{\sum_{i=1}^p K(|t-t_i|/h)}$ and K is a regular kernel supported in $[0, 1]$. In this case, consistency will be true for any sequence $c_{n,p} \rightarrow 0$ as $n, p \rightarrow \infty$ satisfying $n^2 \mathbb{P}_{\mathcal{X}, \mathcal{X}_1}(\|\mathcal{X}\|_{\mathcal{H}} + \|\mathcal{X}_1\|_{\mathcal{H}} \geq p c_{n,p}) \rightarrow 0$.

Let us note that if $\mathbb{E}_{\mathcal{X}}(\|\mathcal{X}\|_{\mathcal{H}}^2) < \infty$, the consistency for the cases given in Sections 4.1 and 4.2 will hold for any sequence $c_{n,p}$ such that $\frac{n}{p c_{n,p}} \rightarrow 0$.

4.3. Eigenfunction expansions

Let $\mathcal{X}, \mathcal{X}_1$ be i.i.d. random elements on $\mathcal{H} = L^2[0, 1]$. Let v_1, v_2, \dots be the orthonormal eigenfunctions of the covariance operator $\mathbb{E}_{\mathcal{X}}(\mathcal{X}(t)\mathcal{X}(s))$ (without loss of generality we have assumed that $\mathbb{E}(\mathcal{X}(t)) = 0$) associated with the eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots$ such that

$$\mathbb{E}_{\mathcal{X}}(\mathcal{X}(t)\mathcal{X}(s)) = \sum_{k=1}^{\infty} \lambda_k v_k(t) v_k(s).$$

If $\mathbb{E}(\int \mathcal{X}^2(s) ds) < \infty$ is finite, using the Karhunen–Loève representation, we can write \mathcal{X} as

$$(4.1) \quad \mathcal{X}(t) = \sum_{k=1}^{\infty} \left(\int \mathcal{X}(s) v_k(s) ds \right) v_k(t) \doteq \sum_{k=1}^{\infty} \xi_k v_k(t),$$

with $\mathbb{E}(\xi_k) = 0$, $\mathbb{E}(\xi_k \xi_j) = 0$ (i.e., ξ_1, ξ_2, \dots uncorrelated) and $\text{var}(\xi_k) = \mathbb{E}(\xi_k^2) = \lambda_k = \mathbb{E}\left(\left(\int \mathcal{X}(s) v_k(s) ds\right)^2\right)$. The classical L^2 -norm in \mathcal{H} can be written as

$$(4.2) \quad d(\mathcal{X}, \mathcal{X}_1) = \sqrt{\sum_{k=1}^{\infty} \left(\int (\mathcal{X}(t) - \mathcal{X}_1(t)) v_k(t) dt \right)^2}.$$

If we consider the truncated expansion of \mathcal{X} as given in [15],

$$(4.3) \quad \mathcal{X}^p(t) = \sum_{k=1}^p \left(\int \mathcal{X}(s) v_k(s) ds \right) v_k(t),$$

we can define the parametrized class of seminorms from the classical L^2 -norm given by

$$\|\mathcal{X}\|_p = \sqrt{\int (\mathcal{X}^p(t))^2 dt} = \sqrt{\sum_{k=1}^p \left(\int \mathcal{X}(t) v_k(t) dt \right)^2},$$

which leads to the semi-metric

$$(4.4) \quad d_p(\mathcal{X}, \mathcal{X}_1) = d(\mathcal{X}^p, \mathcal{X}_1^p) = \sqrt{\sum_{k=1}^p \left(\int (\mathcal{X}(t) - \mathcal{X}_1(t)) v_k(t) dt \right)^2}.$$

In this case, the consistency will hold for any sequence $c_{n,p} \rightarrow 0$ such that $\frac{n^2}{c_{n,p}^2} \sum_{k=p+1}^{\infty} \lambda_k \rightarrow 0$ as $n, p \rightarrow \infty$.

5. SIMULATION STUDY

In order to illustrate the results given in Theorems 3.1 and 3.2, we perform a small simulation study where we compare the behaviour of the estimators, $\hat{\eta}_n$ and $\hat{\eta}_{n,p}$ for finite sample sizes settings. Following [7], we simulate n pairs $\{(\mathcal{X}_i(t), Y_i)\}_{i=1}^n$ where, for $t \in [0, \pi]$, and for each $i = 1, \dots, n$,

$$\mathcal{X}_i(t) = a_i \cos(2t), \quad a_i \sim N(0, \sigma = 0.1).$$

The plot of $n = 100$ curves is shown in Figure 1.

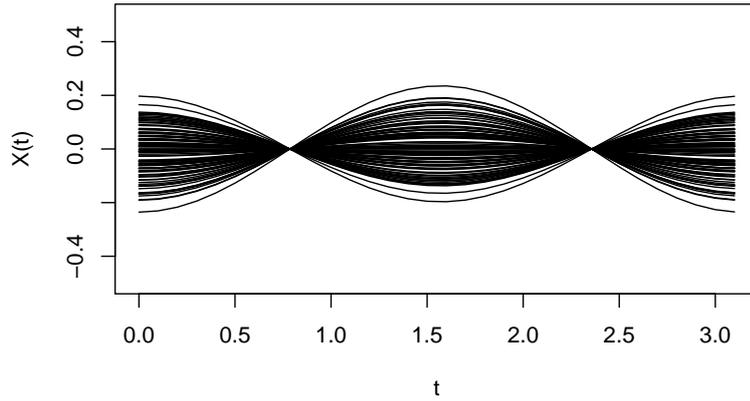


Figure 1: Simulated curves for $n = 100$.

The responses were generated following the model

$$Y_i = \eta(\mathcal{X}_i) + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma = 0.4),$$

for different regression functions η as listed below:

Setting 1: $\eta(\mathcal{X}_i) = a_i^2$ (see [7]);

Setting 2: $\eta(\mathcal{X}_i) = \left(\int_0^\pi \sin(4\pi t) \mathcal{X}_i(t) dt\right)^2$ (see [11]);

Setting 3: $\eta(\mathcal{X}_i) = \int_0^\pi |\mathcal{X}_i(t)| \log(|\mathcal{X}_i(t)|) dt$ (see [14]);

Setting 4: $\eta(\mathcal{X}_i) = \int_0^\pi \mathcal{X}_i^2(t) dt$ (see [2]).

For the *full model* we used the classical L^2 -metric which in this case gives

$$\begin{aligned} d(\mathcal{X}_i, \mathcal{X}_j) &= \left(\int_0^\pi (\mathcal{X}_i(t) - \mathcal{X}_j(t))^2 dt \right)^{1/2} = \left(\int_0^\pi (a_i - a_j)^2 \cos^2(2t) dt \right)^{1/2} \\ &= \left(\int_0^\pi \cos^2(2t) dt \right)^{1/2} |a_i - a_j| = \sqrt{\frac{\pi}{2}} |a_i - a_j|. \end{aligned}$$

For the discretized model, we divided the interval of time $[0, \pi]$ in $p+1$ subintervals of length $\frac{\pi}{p}$. The semimetric in this case is given by

$$\begin{aligned} d_p(\mathcal{X}, \mathcal{X}_1) &= d(\mathcal{X}^p, \mathcal{X}_1^p) = \left(\int_0^\pi |\mathcal{X}^p(t) - \mathcal{X}_1^p(t)|^2 dt \right)^{1/2} \\ &\approx \left(\frac{1}{p} \sum_{k=1}^p (\mathcal{X}_i(t_k) - \mathcal{X}_j(t_k))^2 \right)^{1/2}. \end{aligned}$$

For both estimators $\hat{\eta}_n$ and $\hat{\eta}_{n,p}$, we used the Epanechnikov kernel $K(u) = \frac{3}{4}(1 - u^2)\mathbb{I}_{[0,1]}(u)$ and the bandwidths h_n and $h_{n,p}$ were chosen via cross validation.

In both cases the sample of size n was divided in two samples of the same size, the learning sample, used to compute the optimal smoothing parameter and the testing sample, used to measure the power of both methods by the Mean Square Error (MSE). For different combination of n and p we repeated 250 times the procedure of building $n/2$ learning samples and $n/2$ testing samples and computing the MSE's for the full and discretized models.

The following tables show the mean over the 250 MSE's for all estimators. As we can see, the simulations confirm our theoretical results since, for the four different settings we can see the consistency as $n, p \rightarrow \infty$ stated in Theorem 3.1 and also the equal order or convergence stated in Theorem 3.2.

Table 1: MSE's for Setting 1.

n	Discretized model				Full model
	20	40	60	80	
50	0.1871725	0.1829381	0.1819154	0.1817674	0.1818614
100	0.1784129	0.1661579	0.1661309	0.1660854	0.1659922
150	0.1727869	0.1674195	0.1675846	0.1674071	0.1672996
200	0.1671014	0.1629972	0.1629855	0.1630360	0.1631458
250	0.1646048	0.1631582	0.1631817	0.1632266	0.1632193
300	0.1653583	0.1638297	0.1637960	0.1638118	0.1637993

Table 2: MSE's for Setting 2.

n	Discretized model				Full model
	20	40	60	80	
50	0.1919580	0.1796157	0.1795600	0.1789984	0.1789860
100	0.1787471	0.1684685	0.1684097	0.1684710	0.1685058
150	0.1731875	0.1661859	0.1661971	0.1663508	0.1663451
200	0.1695872	0.1646054	0.1646025	0.1646861	0.1646566
250	0.1658714	0.1622371	0.1621559	0.1621067	0.1621016
300	0.1655437	0.1633919	0.1634236	0.1634164	0.1634100

Table 3: MSE's for Setting 3.

n	Discretized model				Full model
	20	40	60	80	
50	0.1875816	0.1752962	0.1744660	0.1751941	0.1748388
100	0.1797477	0.1672346	0.1671503	0.1671671	0.1671481
150	0.1706658	0.1662048	0.1661369	0.1661024	0.1660888
200	0.1696802	0.1683357	0.1681568	0.1681344	0.1681435
250	0.1666817	0.1651802	0.1652298	0.1652369	0.1652162
300	0.1626991	0.1622967	0.1623146	0.1622935	0.1623169

Table 4: MSE's for Setting 4.

n	Discretized model				Full model
	20	40	60	80	
50	0.1951465	0.1867710	0.1872990	0.1870323	0.1869950
100	0.1824836	0.1694453	0.1694464	0.1695669	0.1695569
150	0.1717909	0.1655053	0.1656256	0.1657503	0.1657367
200	0.1692647	0.1657557	0.1655030	0.1655163	0.1655050
250	0.1651644	0.1630851	0.1631351	0.1630439	0.1630378
300	0.1665684	0.1655066	0.1655070	0.1654343	0.1654715

APPENDIX A – Proofs of auxiliary results

To prove the consistency of the examples given in sections 4.1 and 4.2 we need the following result.

Proposition A.1. *Let $\mathcal{X}^p(t) = \sum_{j=1}^p \phi_j(t)\mathcal{X}(t_j)$ with ϕ_j satisfying:*

- (a) *for each $t \in [0, 1]$, $\sum_{j=1}^p \phi_j(t) = 1$;*
- (b) *for each $t \in [0, 1]$, $\sum_{j=i}^p \phi_j^2(t) \leq C_3$ for some constant C_3 ;*
- (c) *$\text{supp}(\phi_j) \subset [t_{(j-m)}, t_{(j+m)}]$ with m independent of p .*

If $c_{n,p} \rightarrow 0$ as $n, p \rightarrow \infty$ is such that $n^2 \mathbb{P}_{\mathcal{X}, \mathcal{X}_1}(\|\mathcal{X}\|_{\mathcal{H}} + \|\mathcal{X}_1\|_{\mathcal{H}} \geq pc_{n,p}) \rightarrow 0$, H2 is fulfilled.

Proof of Proposition A.1: Using the Fundamental Theorem of Calculus (FTC) (see Theorem 8.2 in [6]) for $H^1([0, 1])$, we get

$$\begin{aligned}
d^2(\mathcal{X}^p, \mathcal{X}) &= \int_0^1 \left| \sum_{j=1}^p \mathcal{X}(t_j)\phi_j(t) - \mathcal{X}(t) \right|^2 dt \\
&= \int_0^1 \left| \sum_{j=1}^p (\mathcal{X}(t_j) - \mathcal{X}(t))\phi_j(t) \right|^2 dt && \text{(by (a))} \\
&= \int_0^1 \left| \sum_{j=1}^p \left(\int_{t_j}^t D\mathcal{X}(s) ds \right) \phi_j(t) \right|^2 dt && \text{(from FTC)} \\
&\leq \int_0^1 \left(\sum_{j=1}^p \left(\int_{t_j}^t D\mathcal{X}(s) ds \right)^2 \mathbb{I}_{\{\text{supp}(\phi_j)\}}(t) \right) \left(\sum_{j=1}^p \phi_j^2(t) \right) dt && \text{(by C-S Ineq.)} \\
&\lesssim \int_0^1 \sum_{j=1}^p \left(\int_{t_j}^t D\mathcal{X}(s) ds \right)^2 \mathbb{I}_{\{\text{supp}(\phi_j)\}}(t) dt && \text{(by (b))} \\
&\lesssim \int_0^1 \sum_{j=1}^p \left(\int_{t_j}^t (D\mathcal{X}(s))^2 ds \right) |t - t_j| \mathbb{I}_{\{\text{supp}(\phi_j)\}}(t) dt && \text{(by C-S Ineq.)} \\
&= \sum_{i=1}^p \int_{t_i}^{t_{i+1}} \sum_{\substack{j=1 \\ \hat{j}|j-i| \leq m}}^p \left(\int_{t_j}^t (D\mathcal{X}(s))^2 ds \right) |t - t_j| dt && \text{(by (c))} \\
&\lesssim \sum_{i=1}^p \sum_{\substack{j=1 \\ \hat{j}|j-i| \leq m}}^p \int_{t_{i-m}}^{t_{i+m}} (D\mathcal{X}(s))^2 \left(\int_{t_j}^{t_{j+1}} |t - t_j| dt \right) ds \\
&\lesssim \frac{m}{p^2} \sum_{i=1}^p \sum_{\substack{j=1 \\ \hat{j}|j-i| \leq m}}^p \int_{t_{i-m}}^{t_{i+m}} (D\mathcal{X}(s))^2 ds
\end{aligned}$$

$$\begin{aligned}
&\lesssim \frac{m^2}{p^2} \sum_{i=1}^p \int_{t_{i-m}}^{t_{i+m}} (D\mathcal{X}(s))^2 ds \\
&= \frac{m^2}{p^2} \int_0^1 \sum_{i=1}^p \mathbb{I}_{[t_{i-m}, t_{i+m}]}(s) (D\mathcal{X}(s))^2 ds \lesssim \frac{1}{p^2} \|\mathcal{X}\|_{\mathcal{H}}^2,
\end{aligned}$$

from where we get $d(\mathcal{X}^p, \mathcal{X}) \lesssim \frac{1}{p} \|\mathcal{X}\|_{\mathcal{H}}$. Analogously we can prove that $d(\mathcal{X}_1^p, \mathcal{X}_1) \lesssim \frac{1}{p} \|\mathcal{X}_1\|_{\mathcal{H}}$. By triangular inequality,

$$\begin{aligned}
n^2 \mathbb{E}_{\mathcal{X}} \left(\mathbb{P}_{\mathcal{X}_1 | \mathcal{X}}^2 \left(|d(\mathcal{X}, \mathcal{X}_1) - d_p(\mathcal{X}, \mathcal{X}_1)| \geq c_{n,p} \mid \mathcal{X} \in \text{supp}(\mu) \right) \right) \\
\leq n^2 \mathbb{P}_{\mathcal{X}, \mathcal{X}_1} (\|\mathcal{X}\|_{\mathcal{H}} + \|\mathcal{X}_1\|_{\mathcal{H}} \geq pc_{n,p}),
\end{aligned}$$

and therefore, for any $c_{n,p} \rightarrow 0$ such that $n^2 \mathbb{P}_{\mathcal{X}, \mathcal{X}_1} (\|\mathcal{X}\|_{\mathcal{H}} + \|\mathcal{X}_1\|_{\mathcal{H}} \geq pc_{n,p}) \rightarrow 0$ H2 is fulfilled. \square

A.1. Consistency for the example in Section 4.1

Since the functions $\phi_j(t) = \mathbb{I}_{[t_j, t_{j+1})}(t)$ satisfy trivially conditions (a)–(c) of Proposition A.1, H2 is fulfilled and therefore, for any sequence $h_{n,p}(x) \rightarrow 0$ satisfying (H3.1) and (H3.2) in Theorem 3.1, we get the consistency of $\hat{\eta}_{n,p}$.

A.2. Consistency for the example in Section 4.2

Observe that $\phi_j(t) = \frac{K(|t-t_j|/h)}{\sum_{i=1}^p K(|t-t_i|/h)}$ satisfies conditions (a)–(c) in Proposition A.1:

- (a) for each $t \in [0, 1]$, $\sum_{j=1}^p \phi_j(t) = \sum_{j=1}^p \frac{K(|t-t_j|/h)}{\sum_{i=1}^p K(|t-t_i|/h)} = 1$;
- (b) since K is nonnegative and $\frac{K(|t-t_j|/h)}{\sum_{i=1}^p K(|t-t_i|/h)} \leq 1$, for each $t \in [0, 1]$, there exists $C_3 = 1$ such that

$$\sum_{j=1}^p \phi_j^2(t) = \sum_{j=1}^p \left(\frac{K(|t-t_j|/h)}{\sum_{i=1}^p K(|t-t_i|/h)} \right)^2 \leq \sum_{j=1}^p \frac{K(|t-t_j|/h)}{\sum_{i=1}^p K(|t-t_i|/h)} = 1;$$

- (c) $\text{supp}(\phi_j) = \text{supp}(K(|t-t_j|/h)) = [t_j - h, t_j + h]$, which implies that, for $h \leq m/p$, $\text{supp}(\phi_j) \subset [t_{(j-m)}, t_{(j+m)}]$.

This implies that H2 is fulfilled then, for any sequence $h_{n,p}(x) \rightarrow 0$ satisfying (H3.1) and (H3.2) in Theorem 3.1, we get the consistency of $\hat{\eta}_{n,p}$.

A.3. Consistency for the example in Section 4.3

Let us consider the truncated expansion of \mathcal{X} , $\mathcal{X}^p(t)$, given by (4.3) and the pseudo-metric $d_p(\mathcal{X}, \mathcal{X}_1) = d(\mathcal{X}^p, \mathcal{X}_1^p)$ given by (4.4). In order to prove H2, let us consider $c_{n,p}$ such that $\frac{n^2}{c_{n,p}^2} \sum_{k=p+1}^{\infty} \lambda_k \rightarrow 0$. Using Chebyshev's Inequality in (3.1) followed by Cauchy Schwartz, we get

$$(A.1) \quad n^2 \mathbb{E}_{\mathcal{X}} \left(\mathbb{P}_{\mathcal{X}_1 | \mathcal{X}}^2 (|d(\mathcal{X}, \mathcal{X}_1) - d_p(\mathcal{X}, \mathcal{X}_1)| \geq c_{n,p} \mid \mathcal{X} \in \text{supp}(\mu)) \right) \\ \leq \frac{n^2}{c_{n,p}^2} \mathbb{E}_{\mathcal{X}, \mathcal{X}_1} ((d(\mathcal{X}, \mathcal{X}_1) - d_p(\mathcal{X}, \mathcal{X}_1))^2).$$

Now, since $d(\mathcal{X}, \mathcal{X}_1) \geq d_p(\mathcal{X}, \mathcal{X}_1)$ we have that $0 \leq d(\mathcal{X}, \mathcal{X}_1) - d_p(\mathcal{X}, \mathcal{X}_1) = d(\mathcal{X}, \mathcal{X}_1) - d(\mathcal{X}^p, \mathcal{X}_1^p)$ and, by triangular inequality $d(\mathcal{X}, \mathcal{X}_1) \leq d(\mathcal{X}, \mathcal{X}^p) + d(\mathcal{X}^p, \mathcal{X}_1^p) + d(\mathcal{X}_1^p, \mathcal{X}_1)$ which implies that

$$(A.2) \quad 0 \leq d(\mathcal{X}, \mathcal{X}_1) - d_p(\mathcal{X}, \mathcal{X}_1) \leq d(\mathcal{X}, \mathcal{X}^p) + d(\mathcal{X}_1^p, \mathcal{X}_1)$$

and, taking squares,

$$0 \leq (d(\mathcal{X}, \mathcal{X}_1) - d_p(\mathcal{X}, \mathcal{X}_1))^2 \leq (d(\mathcal{X}, \mathcal{X}^p) + d(\mathcal{X}_1^p, \mathcal{X}_1))^2 \leq 2 (d^2(\mathcal{X}, \mathcal{X}^p) + d^2(\mathcal{X}_1^p, \mathcal{X}_1)).$$

As a consequence, to proof this proposition it will sufficient to bound $\mathbb{E}_{\mathcal{X}} (d^2(\mathcal{X}, \mathcal{X}^p))$ (equivalently, $\mathbb{E}_{\mathcal{X}_1} (d^2(\mathcal{X}_1, \mathcal{X}_1^p))$). Since v_k are orthonormal,

$$d^2(\mathcal{X}, \mathcal{X}^p) = \int \left(\mathcal{X}(s) - \sum_{k=1}^p \left(\int \mathcal{X}(t) v_k(t) dt \right) v_k(s) \right)^2 ds \\ = \sum_{k=p+1}^{\infty} \left(\int \mathcal{X}(t) v_k(t) dt \right)^2.$$

Then, we have

$$\mathbb{E}_{\mathcal{X}} (d^2(\mathcal{X}, \mathcal{X}^p)) = \mathbb{E}_{\mathcal{X}} \left(\sum_{k=p+1}^{\infty} \left(\int \mathcal{X}(t) v_k(t) dt \right)^2 \right) \\ = \sum_{k=p+1}^{\infty} \lambda_k \quad (\text{from (4.1)}).$$

Analogously we can prove that $\mathbb{E}_{\mathcal{X}_1} (d^2(\mathcal{X}_1, \mathcal{X}_1^p)) = \sum_{k=p+1}^{\infty} \lambda_k$. Therefore, in (A.1) we get

$$n^2 \mathbb{E}_{\mathcal{X}} \left(\mathbb{P}_{\mathcal{X}_1 | \mathcal{X}}^2 (|d(\mathcal{X}, \mathcal{X}_1) - d_p(\mathcal{X}, \mathcal{X}_1)| \geq c_{n,p} \mid \mathcal{X} \in \text{supp}(\mu)) \right) \lesssim \frac{n^2}{c_{n,p}^2} \sum_{k=p+1}^{\infty} \lambda_k \rightarrow 0.$$

This implies that H2 is fulfilled then, for any sequence $h_{n,p}(x) \rightarrow 0$ satisfying (H3.1) and (H3.2) in Theorem 3.1, we get the consistency of $\hat{\eta}_{n,p}$.

APPENDIX B – Proof of Proposition 2.2 and Theorems 3.1 and 3.2

To prove Proposition 2.2 we need some preliminary results whose proofs can be found in [16].

Theorem B.1 (Theorem 3.4). *If $\eta \in L^2(\mathcal{H}, \mu)$ and $\hat{\eta}_n$ is the estimator given in (2.2) with weights $W_n(\mathcal{X}) = \{W_{ni}(\mathcal{X})\}_{i=1}^n$ satisfying the following conditions:*

(i) *there is a sequence of nonnegative random variables $a_n(\mathcal{X}) \rightarrow 0$ a.s. such that*

$$\lim_{n \rightarrow \infty} \mathbb{E} \left(\sum_{i=1}^n W_{ni}(\mathcal{X}) \mathbb{I}_{\{d(\mathcal{X}, \mathcal{X}_i) > a_n(\mathcal{X})\}} \right) = 0;$$

(ii)

$$\lim_{n \rightarrow \infty} \mathbb{E} \left(\max_{1 \leq i \leq n} W_{ni}(\mathcal{X}) \right) = 0;$$

(iii) *for all $\epsilon > 0$ there exists $\delta > 0$ such that for any η^* bounded and continuous function fulfilling $\mathbb{E}_{\mathcal{X}}((\eta(\mathcal{X}) - \eta^*(\mathcal{X}))^2) < \delta$ we have that*

$$\mathbb{E} \left(\sum_{i=1}^n W_{ni}(\mathcal{X}) (\eta^*(\mathcal{X}_i) - \eta(\mathcal{X}_i))^2 \right) < \epsilon;$$

then $\hat{\eta}_n$ is mean square consistent.

Corollary B.1 (Corollary 3.3). *Let U_n be a sequence of probability weights satisfying conditions (i), (ii) and (iii) of Theorem B.1. If W_n is a sequence of weights such that $\sum_{i=1}^n W_{ni}(\mathcal{X}) = 1$ and, for each $n \geq 1$, $|W_n| \leq MU_n$ for some constant $M \geq 1$, then the estimator given in (2.2) with weights $W_n(\mathcal{X})$ is mean square consistent.*

Lemma B.1 (Lemma A.1). *Let \mathcal{H} be a separable metric space. If $A = \text{supp}(\mu) = \{x \in \mathcal{H} : \mu(\mathcal{B}(x, \epsilon)) > 0, \forall \epsilon > 0\}$ then $\mu(A) = 1$.*

Proof of Proposition 2.2: Let $x \in \text{supp}(\mu)$ be fixed. Let us observe that, since K is regular, there exist constants $0 < c_1 < c_2 < \infty$ such that, for each i ,

$$(B.1) \quad W_{ni}(x) = \frac{K\left(\frac{d(\mathcal{X}_i, x)}{h_n(x)}\right)}{\sum_{j=1}^n K\left(\frac{d(\mathcal{X}_j, x)}{h_n(x)}\right)} \leq \frac{c_2}{c_1} \frac{\mathbb{I}_{\{d(\mathcal{X}_i, x) \leq h_n(x)\}}}{\sum_{j=1}^n \mathbb{I}_{\{d(\mathcal{X}_j, x) \leq h_n(x)\}}} \doteq \frac{c_2}{c_1} U_{ni}(x).$$

Let $h_n(x) \rightarrow 0$ such that $h_n(x) \geq H_n(x)$ ($H_n(x) \rightarrow 0$ by Lemma 2.2, for $x \in \text{supp}(\mu)$). From (B.1) and Corollary B.1, it suffices to prove that the weights U_{ni} satisfy conditions (i), (ii) and (iii) of Theorem B.1. To prove (i) let us take $a_n(x) = h_n^{1/2}(x) \rightarrow 0$. Then, by Lemma B.1,

$$\begin{aligned} & \mathbb{E} \left(\sum_{i=1}^n U_{ni}(\mathcal{X}) \mathbb{I}_{\{d(\mathcal{X}_i, \mathcal{X}) > h_n(\mathcal{X})^{1/2}\}} \right) \\ &= \mathbb{E}_{\mathcal{X}} \left(\mathbb{E}_{\mathcal{D}_n | \mathcal{X}} \left(\mathbb{I}_{\{\mathcal{X} \in \text{supp}(\mu)\}} \sum_{i=1}^n U_{ni}(\mathcal{X}) \mathbb{I}_{\{d(\mathcal{X}_i, \mathcal{X}) > h_n(\mathcal{X})^{1/2}\}} \middle| \mathcal{X} \in \text{supp}(\mu) \right) \right). \end{aligned}$$

Given $\epsilon > 0$, let $x \in \text{supp}(\mu)$ be fixed. Since $h_n(x) \rightarrow 0$, there exists $N_1 = N_1(x)$ such that if $n \geq N_1$, $\mathbb{I}_{\{h_n(x)^{1/2} < d(x_i, x) \leq h_n(x)\}} = 0$ for all i and, consequently,

$$\mathbb{E}_{\mathcal{D}_n} \left(\frac{1}{\sum_{j=1}^n \mathbb{I}_{\{d(x_j, x) \leq h_n(x)\}}} \sum_{i=1}^n \mathbb{I}_{\{h_n(x)^{1/2} < d(x_i, x) \leq h_n(x)\}} \right) < \epsilon.$$

In addition, $\frac{\sum_{i=1}^n \mathbb{I}_{\{h_n(x)^{1/2} < d(x_i, x) \leq h_n(x)\}}}{\sum_{j=1}^n \mathbb{I}_{\{d(x_j, x) \leq h_n(x)\}}} \leq 1$, from what follows that

$$\mathbb{E}_{\mathcal{D}_n} \left(\frac{1}{\sum_{j=1}^n \mathbb{I}_{\{d(x_j, x) \leq h_n(x)\}}} \sum_{i=1}^n \mathbb{I}_{\{h_n(x)^{1/2} < d(x_i, x) \leq h_n(x)\}} \right) \leq 1.$$

Therefore, by the dominated convergence theorem we have that condition (i) is satisfied. Now, since $h_n(x) \geq H_n(x)$,

$$\sum_{j=1}^n \mathbb{I}_{\{d(\mathcal{X}_j, x) \leq h_n(x)\}} \geq \sum_{j=1}^n \mathbb{I}_{\{d(\mathcal{X}_j, x) \leq H_n(x)\}} = k_n \rightarrow \infty.$$

Therefore,

$$\max_{1 \leq i \leq n} U_{ni}(x) \leq \max_{1 \leq i \leq n} \frac{1}{\sum_{j=1}^n \mathbb{I}_{\{d(\mathcal{X}_j, x) \leq h_n(x)\}}} \leq \frac{1}{k_n} \rightarrow 0,$$

from what we derive (ii) using the dominated convergence theorem. It remains to verify that condition (iii) holds. Since $\eta \in L^2(\mathcal{H}, \mu)$ which is separable and complete, there exists η^* continuous and bounded such that, for all $\delta > 0$, $\mathbb{E}_{\mathcal{X}}((\eta(\mathcal{X}) - \eta^*(\mathcal{X}))^2) < \delta$. Then,

$$\begin{aligned} & \mathbb{E} \left(\sum_{i=1}^n U_{ni}(\mathcal{X})(\eta^*(\mathcal{X}_i) - \eta(\mathcal{X}_i))^2 \right) \\ &= \mathbb{E}_{\mathcal{X}} \left(\mathbb{E}_{\mathcal{D}_n | \mathcal{X}} \left(\mathbb{I}_{\{\mathcal{X} \in \text{supp}(\mu)\}} \sum_{i=1}^n U_{ni}(\mathcal{X})(\eta^*(\mathcal{X}_i) - \eta(\mathcal{X}_i))^2 | \mathcal{X} \in \text{supp}(\mu) \right) \right). \end{aligned}$$

Let $x \in \text{supp}(\mu)$ be fixed. From [16], Lemma A.7, for any nonnegative bounded measurable function f , we have

$$\mathbb{E}_{\mathcal{D}_n} \left(\sum_{i=1}^n U_{ni}(x) f(\mathcal{X}_i) \right) \leq 12 \frac{1}{\mu(\mathcal{B}(x, h_n(x)))} \int_{\mathcal{B}(x, h_n(x))} f(y) d\mu(y).$$

Then, applying the inequality to $f(\mathcal{X}_i) = (\eta^*(\mathcal{X}_i) - \eta(\mathcal{X}_i))^2$, we get

$$\begin{aligned} & \mathbb{E}_{\mathcal{D}_n} \left(\sum_{i=1}^n U_{ni}(x) (\eta^*(\mathcal{X}_i) - \eta(\mathcal{X}_i))^2 \right) \\ & \lesssim \frac{1}{\mu(\mathcal{B}(x, h_n(x)))} \int_{\mathcal{B}(x, h_n(x))} (\eta^*(y) - \eta(y))^2 d\mu(y) \\ & \leq \frac{1}{\mu(\mathcal{B}(x, h_n(x)))} \int_{\mathcal{B}(x, h_n(x))} (\eta^*(y) - \eta^*(x))^2 d\mu(y) \\ & \quad + \frac{1}{\mu(\mathcal{B}(x, h_n(x)))} \int_{\mathcal{B}(x, h_n(x))} (\eta^*(x) - \eta(x))^2 d\mu(y) \\ & \quad + \frac{1}{\mu(\mathcal{B}(x, h_n(x)))} \int_{\mathcal{B}(x, h_n(x))} (\eta(x) - \eta(y))^2 d\mu(y) \\ & \doteq f_{1,n}(x) + f_{2,n}(x) + f_{3,n}(x). \end{aligned}$$

This part will be complete if we show that the expectation with respect to \mathcal{X} of these three functions converges to zero. For this, let $\epsilon > 0$ and $\delta \leq \epsilon$. Since η^* is continuous, there exists $r = r(x, \epsilon) > 0$ such that if $d(x, y) < r$ then $|\eta^*(x) - \eta^*(y)| < \epsilon$. On the other hand, since $h_n(x) \rightarrow 0$, for that $r(x, \epsilon) > 0$, there exists $N_2 = N_2(x, r(x, \epsilon))$ such that if $n \geq N_2$, $h_n(x) < r$. Then, $f_{1,n}(x) = \frac{1}{\mu(\mathcal{B}(x, h_n(x)))} \int_{\mathcal{B}(x, h_n(x))} (\eta^*(y) - \eta^*(x))^2 d\mu(y) < \epsilon$ for $n \geq N_2$ and in addition it is bounded so, by the dominated convergence theorem we have that

$$\mathbb{E}_{\mathcal{X}}(f_{1,n}(\mathcal{X})) \rightarrow 0.$$

For the second term, since $\delta \leq \epsilon$, we have that

$$\mathbb{E}_{\mathcal{X}}(f_{2,n}(\mathcal{X})) = \mathbb{E}_{\mathcal{X}}((\eta(\mathcal{X}) - \eta^*(\mathcal{X}))^2) < \epsilon.$$

Finally, since η is bounded,

$$\mathbb{E}_{\mathcal{X}}(f_{3,n}(\mathcal{X})) \lesssim \mathbb{E}_{\mathcal{X}}\left(\frac{1}{\mu(\mathcal{B}(\mathcal{X}, h_n(\mathcal{X})))} \int_{\mathcal{B}(\mathcal{X}, h_n(\mathcal{X}))} |\eta(\mathcal{X}) - \eta(y)| d\mu(y)\right),$$

which converge to zero if the bounded random variables

$$\frac{1}{\mu(\mathcal{B}(\mathcal{X}, h_n(\mathcal{X})))} \int_{\mathcal{B}(\mathcal{X}, h_n(\mathcal{X}))} |\eta(\mathcal{X}) - \eta(y)| d\mu(y)$$

converge to zero in probability. To see this, let $\lambda > 0$ be fixed. For every $\delta_0 > 0$,

$$\begin{aligned} \mathbb{P}_{\mathcal{X}}\left(\frac{1}{\mu(\mathcal{B}(\mathcal{X}, h_n(\mathcal{X})))} \int_{\mathcal{B}(\mathcal{X}, h_n(\mathcal{X}))} |\eta(\mathcal{X}) - \eta(y)| d\mu(y) > \lambda\right) \\ \leq \mathbb{P}_{\mathcal{X}}(h_n(\mathcal{X}) > \delta_0) + \sup_{\delta \leq \delta_0} \mathbb{P}_{\mathcal{X}}\left(\frac{1}{\mu(\mathcal{B}(\mathcal{X}, \delta))} \int_{\mathcal{B}(\mathcal{X}, \delta)} |\eta(\mathcal{X}) - \eta(y)| d\mu(y) > \lambda\right). \end{aligned}$$

Since $h_n(\mathcal{X}) \rightarrow 0$ a.s. the first term converges to zero while the second term does thanks to the truth of the Besicovitch condition (2.4). \square

Proof of Theorem 3.1:

Proof of (a): Let us define $\mathcal{D}_n = \{\mathcal{X}_1, \dots, \mathcal{X}_n\}$ and $\mathcal{C}_n = \{Y_1, \dots, Y_n\}$. In order to prove the mean square consistency, we consider

$$\mathbb{E}((\hat{\eta}_{n,p}(\mathcal{X}) - \eta(\mathcal{X}))^2) = \mathbb{E}_{\mathcal{X}}(\mathbb{E}_{\mathcal{D}_n, \mathcal{C}_n | \mathcal{X}}((\hat{\eta}_{n,p}(\mathcal{X}) - \eta(\mathcal{X}))^2) | \mathcal{X}).$$

Let $x \in \text{supp}(\mu)$ be fixed. To simplify the notation, we set $\mathbb{E}(\cdot) = \mathbb{E}_{\mathcal{D}_n, \mathcal{C}_n | \mathcal{X}}(\cdot)$. Then, for a particular $h_n(x) \geq h_n^*(x)$ to be defined later, let us define the *theoretical quantities*

$$K\left(\frac{d(x, \mathcal{X}_i)}{h_n(x)}\right) \doteq K_i(x) \doteq K_i \quad \text{and} \quad K\left(\frac{d_p(x, \mathcal{X}_i)}{h_{n,p}(x)}\right) \doteq K_{i,p}(x) \doteq K_{i,p},$$

and, as in (2.3),

$$\frac{K_i}{\sum_{j=1}^n K_j} \doteq W_i \quad \text{and} \quad \frac{K_{i,p}}{\sum_{j=1}^n K_{j,p}} \doteq W_{i,p}.$$

Let us consider the following auxiliary unobservable quantities:

$$\hat{\eta}_n(x) = \sum_{i=1}^n W_i Y_i, \quad \eta_n(x) = \sum_{i=1}^n W_i \eta(\mathcal{X}_i) \quad \text{and} \quad \eta_{n,p}(x) = \sum_{i=1}^n W_{i,p} \eta(\mathcal{X}_i).$$

Then, we have

$$\begin{aligned}
\widehat{\eta}_{n,p}(x) - \eta(x) &= [\widehat{\eta}_{n,p}(x) - \eta_{n,p}(x)] + [\eta_{n,p}(x) - \eta_n(x)] + [\eta_n(x) - \widehat{\eta}_n(x)] + [\widehat{\eta}_n(x) - \eta(x)] \\
&= \sum_{i=1}^n W_{i,p}(Y_i - \eta(\mathcal{X}_i)) + \sum_{i=1}^n (W_{i,p} - W_i)\eta(\mathcal{X}_i) + \sum_{i=1}^n W_i(\eta(\mathcal{X}_i) - Y_i) \\
&\quad + [\widehat{\eta}_n(x) - \eta(x)] \\
&= \sum_{i=1}^n (W_{i,p} - W_i)(Y_i - \eta(\mathcal{X}_i)) + \sum_{i=1}^n (W_{i,p} - W_i)\eta(\mathcal{X}_i) \\
&\quad + [\widehat{\eta}_n(x) - \eta(x)].
\end{aligned}$$

Taking squares and expectation in $\mathcal{D}_n, \mathcal{C}_n$, we have

$$\begin{aligned}
\mathbb{E}((\widehat{\eta}_{n,p}(x) - \eta(x))^2) &\lesssim \mathbb{E}\left(\left(\sum_{i=1}^n (W_{i,p} - W_i)(Y_i - \eta(\mathcal{X}_i))\right)^2\right) \\
&\quad + \mathbb{E}\left(\left(\sum_{i=1}^n (W_{i,p} - W_i)\eta(\mathcal{X}_i)\right)^2\right) \\
&\quad + \mathbb{E}\left([\widehat{\eta}_n(x) - \eta(x)]^2\right) \\
&\doteq I + II + III.
\end{aligned}$$

By Proposition 2.1 and Remark 2.2 (since $h_n(x) \rightarrow 0$ and $h_n(x) \geq h_n^*(x)$), taking expectation on \mathcal{X} we have that term *III* converges to zero. For the first term, we have

$$\begin{aligned}
I &\approx \mathbb{E}\left(\left(\sum_{i=1}^n (W_{i,p} - W_i)(Y_i - \eta(\mathcal{X}_i))\right)^2\right) \\
&= \mathbb{E}\left(\sum_{i=1}^n \sum_{j=1}^n (W_{i,p} - W_i)(W_{j,p} - W_j)e_i e_j\right) \quad (Y_i - \eta(\mathcal{X}_i) = e_i) \\
&= \mathbb{E}\left(\sum_{i=1}^n \sum_{j=1}^n (W_{i,p} - W_i)(W_{j,p} - W_j)\mathbb{E}_{\mathcal{C}_n|\mathcal{D}_n}(e_i e_j|\mathcal{D}_n)\right) \\
&= \mathbb{E}\left(\sum_{i=1}^n |W_{i,p} - W_i|^2 \mathbb{E}_{\mathcal{C}_n|\mathcal{D}_n}(e_i^2|\mathcal{D}_n)\right) \quad (\text{cond. ind.}) \\
&= \sigma^2 \mathbb{E}\left(\sum_{i=1}^n |W_{i,p} - W_i|^2\right).
\end{aligned}$$

On the other hand, since η is bounded, in *II* we have

$$II = \mathbb{E}\left(\left(\sum_{i=1}^n (W_{i,p} - W_i)\eta(\mathcal{X}_i)\right)^2\right) \lesssim \mathbb{E}\left(\left(\sum_{i=1}^n |W_{i,p} - W_i|\right)^2\right).$$

We will see that terms I and II converge to zero by splitting the sum in different pieces:

- (1) $A_1 \doteq \{i: d_p(x, \mathcal{X}_i) > h_{n,p}(x), d(x, \mathcal{X}_i) > h_n(x)\};$
- (2) $A_2 \doteq \{i: d_p(x, \mathcal{X}_i) > h_{n,p}(x), d(x, \mathcal{X}_i) \leq h_n(x)\};$
- (3) $A_3 \doteq \{i: d_p(x, \mathcal{X}_i) \leq h_{n,p}(x), d(x, \mathcal{X}_i) > 3h_n(x)\};$
- (4) $A_4 \doteq \{i: d_p(x, \mathcal{X}_i) \leq h_{n,p}(x), d(x, \mathcal{X}_i) \leq 3h_n(x)\}.$

Case (1) is trivial since in this case K is supported in $[0, 1]$ which implies that $W_{i,p} = W_i = 0$. Let us start, therefore, with case (2).

(2) Let $A_2 \doteq \{i: d_p(x, \mathcal{X}_i) > h_{n,p}(x), d(x, \mathcal{X}_i) \leq h_n(x)\}$. Observe that in this case $W_{i,p} = 0$ since K is supported in $[0, 1]$. Therefore, since $|W_i| \leq 1$ we get

$$I_{A_2} \doteq \mathbb{E} \left(\sum_{i=1}^n |W_i|^2 \mathbb{I}_{\{i \in A_2\}} \right) \leq \mathbb{E} \left(\sum_{i=1}^n \mathbb{I}_{\{i \in A_2\}} \right)$$

and

$$(B.2) \quad II_{A_2} \doteq \mathbb{E} \left(\left(\sum_{i=1}^n |W_i| \mathbb{I}_{\{i \in A_2\}} \right)^2 \right) \leq \mathbb{E} \left(\left(\sum_{i=1}^n \mathbb{I}_{\{i \in A_2\}} \right)^2 \right) \doteq C_{A_2}.$$

Observe that the i.i.d. random variables $\mathbb{I}_{\{i \in A_2\}}$ have a Bernoulli distribution with parameter

$$\begin{aligned} p &= \mathbb{P}_{\mathcal{X}_1} (d_p(x, \mathcal{X}_1) > h_{n,p}(x), d(x, \mathcal{X}_1) \leq h_n(x)) \\ &\leq \mathbb{P}_{\mathcal{X}_1} (d_p(x, \mathcal{X}_1) - d(x, \mathcal{X}_1) \geq h_{n,p}(x) - h_n(x)) \\ &\leq \mathbb{P}_{\mathcal{X}_1} (|d_p(x, \mathcal{X}_1) - d(x, \mathcal{X}_1)| \geq c_{n,p}) \end{aligned} \quad (\text{by H3.2}).$$

As a consequence, the random variable $Z \doteq \sum_{i=1}^n \mathbb{I}_{\{i \in A_2\}}$ has Binomial distribution with parameters n and p and expectation $\mathbb{E}(Z) = np$. This implies that

$$(B.3) \quad I_{A_2} \lesssim \mathbb{E}(Z) \leq n \mathbb{P}_{\mathcal{X}_1} (|d_p(x, \mathcal{X}_1) - d(x, \mathcal{X}_1)| \geq c_{n,p}),$$

and, since $\mathbb{E}(Z^2) = np(1-p) + n^2p^2 \leq np + (np)^2$,

$$(B.4) \quad II_{A_2} \leq C_{A_2} \lesssim \mathbb{E}(Z^2) \leq n \mathbb{P}_{\mathcal{X}_1} (|d_p(x, \mathcal{X}_1) - d(x, \mathcal{X}_1)| \geq c_{n,p}) + \left(n \mathbb{P}_{\mathcal{X}_1} (|d_p(x, \mathcal{X}_1) - d(x, \mathcal{X}_1)| \geq c_{n,p}) \right)^2.$$

(3) Let $A_3 \doteq \{i: d_p(x, \mathcal{X}_i) \leq h_{n,p}(x), d(x, \mathcal{X}_i) > 3h_n(x)\}$. Observe that in this case $W_i = 0$ since K is supported in $[0, 1]$. Then, since $\forall i, |W_{i,p}| \leq 1$, we get

$$I_{A_3} \doteq \mathbb{E} \left(\sum_{i=1}^n |W_{i,p}|^2 \mathbb{I}_{\{i \in A_3\}} \right) \leq \mathbb{E} \left(\sum_{i=1}^n \mathbb{I}_{\{i \in A_3\}} \right),$$

and

$$(B.5) \quad II_{A_3} \doteq \mathbb{E} \left(\left(\sum_{i=1}^n |W_{i,p}| \mathbb{I}_{\{i \in A_3\}} \right)^2 \right) \leq \mathbb{E} \left(\left(\sum_{i=1}^n \mathbb{I}_{\{i \in A_3\}} \right)^2 \right).$$

Now, the i.i.d. random variables $\mathbb{I}_{\{i \in A_3\}}$ have Bernoulli distribution with parameter

$$\begin{aligned} p &= \mathbb{P}_{\mathcal{X}_1}(d_p(x, \mathcal{X}_1) \leq h_{n,p}(x), d(x, \mathcal{X}_1) > 3h_n(x)) \\ &\leq \mathbb{P}_{\mathcal{X}_1}(d(x, \mathcal{X}_1) - d_p(x, \mathcal{X}_1) \geq 3h_n(x) - h_{n,p}(x)). \end{aligned}$$

As a consequence, the random variable $Z \doteq \sum_{i=1}^n \mathbb{I}_{\{i \in A_3\}}$ has Binomial distribution with parameters n and p . But from (H3.1), for n large enough, $h_n(x) \geq \left(\frac{1+C_2}{2}\right) c_{n,p}$ which, together with H3.2 implies that

$$3h_n(x) - h_{n,p}(x) \geq 2h_n(x) - C_2 c_{n,p} \geq c_{n,p},$$

and then, for n large enough,

$$p \leq \mathbb{P}_{\mathcal{X}_1}(|d_p(x, \mathcal{X}_1) - d(x, \mathcal{X}_1)| \geq c_{n,p}).$$

Therefore, since $\mathbb{E}(Z) = np$ we have

$$(B.6) \quad I_{A_3} \lesssim \mathbb{E}(Z) \leq n\mathbb{P}_{\mathcal{X}_1}(|d_p(x, \mathcal{X}_1) - d(x, \mathcal{X}_1)| \geq c_{n,p}),$$

and since $\mathbb{E}(Z^2) = np(1-p) + n^2p^2 \leq np + (np)^2$,

$$(B.7) \quad \begin{aligned} II_{A_3} &\lesssim \mathbb{E}(Z^2) \leq n\mathbb{P}_{\mathcal{X}_1}(|d_p(x, \mathcal{X}_1) - d(x, \mathcal{X}_1)| \geq c_{n,p}) \\ &\quad + \left(n\mathbb{P}_{\mathcal{X}_1}(|d_p(x, \mathcal{X}_1) - d(x, \mathcal{X}_1)| \geq c_{n,p})\right)^2. \end{aligned}$$

(4) Let $A_4 \doteq \{i: d_p(x, \mathcal{X}_i) \leq h_{n,p}(x), d(x, \mathcal{X}_i) \leq 3h_n(x)\}$. In this case we write,

$$\begin{aligned} W_{i,p} - W_i &= \frac{K_{i,p}}{\sum_{j=1}^n K_{j,p}} - \frac{K_i}{\sum_{j=1}^n K_j} \\ &= \frac{K_{i,p}}{\sum_{j=1}^n K_{j,p}} - \frac{K_i}{\sum_{j=1}^n K_{j,p}} + \frac{K_i}{\sum_{j=1}^n K_{j,p}} - \frac{K_i}{\sum_{j=1}^n K_j} \\ &= (K_{i,p} - K_i) \frac{1}{\sum_{j=1}^n K_{j,p}} + K_i \frac{\sum_{j=1}^n (K_j - K_{j,p})}{\sum_{j=1}^n K_j \sum_{j=1}^n K_{j,p}} \\ &= (K_{i,p} - K_i) \frac{1}{\sum_{j=1}^n K_{j,p}} + W_i \frac{\sum_{j=1}^n (K_j - K_{j,p})}{\sum_{j=1}^n K_{j,p}}. \end{aligned}$$

Then,

$$(B.8) \quad \begin{aligned} I_{A_4} &\doteq \mathbb{E} \left(\sum_{i=1}^n |W_{i,p} - W_i|^2 \mathbb{I}_{\{i \in A_4\}} \right) \\ &\lesssim \mathbb{E} \left(\sum_{i=1}^n |K_{i,p} - K_i|^2 \frac{\mathbb{I}_{\{i \in A_4\}}}{(\sum_{j=1}^n K_{j,p})^2} \right) \\ &\quad + \mathbb{E} \left(\sum_{i=1}^n W_i^2 \mathbb{I}_{\{i \in A_4\}} \left(\frac{\sum_{j=1}^n (K_j - K_{j,p})}{\sum_{j=1}^n K_{j,p}} \right)^2 \right) \\ &\lesssim \mathbb{E} \left(\sum_{i=1}^n |K_{i,p} - K_i|^2 \frac{\mathbb{I}_{\{i \in A_4\}}}{(\sum_{j=1}^n \mathbb{I}_{\{j: d_p(x, \mathcal{X}_j) \leq h_{n,p}(x)\}})^2} \right) \quad (K \text{ regular}) \\ &\quad + \mathbb{E} \left(\left(\frac{\sum_{j=1}^n |K_j - K_{j,p}|}{\sum_{j=1}^n K_{j,p}} \right)^2 \right) \quad \left(|W_i| \leq 1, \sum_{i=1}^n W_i = 1 \right) \\ &\doteq I_{A_4}^1 + I_{A_4}^2 \end{aligned}$$

and

$$\begin{aligned}
(B.9) \quad II_{A_4} &\doteq \mathbb{E} \left(\left(\sum_{i=1}^n |W_{i,p} - W_i| \mathbb{I}_{\{i \in A_4\}} \right)^2 \right) \\
&\lesssim \mathbb{E} \left(\left(\sum_{i=1}^n |K_{i,p} - K_i| \frac{\mathbb{I}_{\{i \in A_4\}}}{\sum_{j=1}^n K_{j,p}} \right)^2 \right) \\
&\quad + \mathbb{E} \left(\left(\sum_{i=1}^n W_i \mathbb{I}_{\{i \in A_4\}} \frac{\sum_{j=1}^n (K_j - K_{j,p})}{\sum_{j=1}^n K_{j,p}} \right)^2 \right) \\
&\lesssim \mathbb{E} \left(\left(\sum_{i=1}^n |K_{i,p} - K_i| \frac{\mathbb{I}_{\{i \in A_4\}}}{\sum_{j=1}^n \mathbb{I}_{\{j: d_p(x, \mathcal{X}_j) \leq h_{n,p}(x)\}}} \right)^2 \right) \quad (K \text{ regular}) \\
&\quad + \mathbb{E} \left(\left(\frac{\sum_{j=1}^n |K_j - K_{j,p}|}{\sum_{j=1}^n K_{j,p}} \right)^2 \right) \quad (|W_i| \leq 1) \\
&\doteq II_{A_4}^1 + II_{A_4}^2.
\end{aligned}$$

Observe that if $\sum_{j=1}^n \mathbb{I}_{\{j: d_p(x, \mathcal{X}_j) \leq h_{n,p}(x)\}} = 0$ then $\forall j, \mathbb{I}_{\{j \in A_4\}} = 0$ so in this case, $I_{A_4}^1$ and $II_{A_4}^1$ are zero. Then, in what follows we will assume that $\sum_{j=1}^n \mathbb{I}_{\{j: d_p(x, \mathcal{X}_j) \leq h_{n,p}(x)\}} \neq 0$. Since K is Lipschitz and we are only considering the indexes i such that $d_p(x, \mathcal{X}_i) \leq h_{n,p}(x)$, we get

$$\begin{aligned}
|K_{i,p} - K_i| &= \left| K \left(\frac{d_p(x, \mathcal{X}_i)}{h_{n,p}(x)} \right) - K \left(\frac{d(x, \mathcal{X}_i)}{h_n(x)} \right) \right| \\
&\lesssim \left| \frac{d_p(x, \mathcal{X}_i)}{h_{n,p}(x)} - \frac{d(x, \mathcal{X}_i)}{h_n(x)} \right| \\
&= \frac{|d_p(x, \mathcal{X}_i)h_n(x) - d(x, \mathcal{X}_i)h_{n,p}(x)|}{h_{n,p}(x)h_n(x)} \\
&\leq \frac{|d_p(x, \mathcal{X}_i) - d(x, \mathcal{X}_i)|}{h_n(x)} + \frac{d_p(x, \mathcal{X}_i)|h_n(x) - h_{n,p}(x)|}{h_n(x)h_{n,p}(x)} \\
&\lesssim \frac{|d_p(x, \mathcal{X}_i) - d(x, \mathcal{X}_i)|}{h_n(x)} + \frac{c_{n,p}}{h_n(x)} \quad (\text{by H3.2}).
\end{aligned}$$

Therefore,

$$\begin{aligned}
(B.10) \quad I_{A_4}^1 &\lesssim \frac{1}{h_n^2(x)} \mathbb{E} \left(\sum_{i=1}^n |d_p(x, \mathcal{X}_i) - d(x, \mathcal{X}_i)|^2 \frac{\mathbb{I}_{\{i \in A_4\}}}{\left(\sum_{j=1}^n \mathbb{I}_{\{j: d_p(x, \mathcal{X}_j) \leq h_{n,p}(x)\}} \right)^2} \right) \\
&\quad + \left(\frac{c_{n,p}}{h_n(x)} \right)^2 \mathbb{E} \left(\sum_{i=1}^n \frac{\mathbb{I}_{\{i \in A_4\}}}{\left(\sum_{j=1}^n \mathbb{I}_{\{j: d_p(x, \mathcal{X}_j) \leq h_{n,p}(x)\}} \right)^2} \right) \\
&\lesssim \frac{1}{h_n^2(x)} \mathbb{E} \left(\sum_{i=1}^n |d_p(x, \mathcal{X}_i) - d(x, \mathcal{X}_i)|^2 \frac{\mathbb{I}_{\{j \in A_4\}}}{\left(\sum_{j=1}^n \mathbb{I}_{\{j: d_p(x, \mathcal{X}_j) \leq h_{n,p}(x)\}} \right)^2} \right) \\
&\quad + \left(\frac{c_{n,p}}{h_n(x)} \right)^2
\end{aligned}$$

and

$$\begin{aligned}
(B.11) \quad II_{A_4}^1 &\lesssim \frac{1}{h_n^2(x)} \mathbb{E} \left(\left(\sum_{i=1}^n |d_p(x, \mathcal{X}_i) - d(x, \mathcal{X}_i)| \frac{\mathbb{I}_{\{i \in A_4\}}}{\sum_{j=1}^n \mathbb{I}_{\{j d_p(x, \mathcal{X}_j) \leq h_{n,p}(x)\}}} \right)^2 \right) \\
&\quad + \left(\frac{c_{n,p}}{h_n(x)} \right)^2 \mathbb{E} \left(\left(\sum_{i=1}^n \frac{\mathbb{I}_{\{i \in A_4\}}}{\sum_{j=1}^n \mathbb{I}_{\{j d_p(x, \mathcal{X}_j) \leq h_{n,p}(x)\}}} \right)^2 \right) \\
&\lesssim \frac{1}{h_n^2(x)} \mathbb{E} \left(\left(\sum_{i=1}^n |d_p(x, \mathcal{X}_i) - d(x, \mathcal{X}_i)| \frac{\mathbb{I}_{\{i \in A_4\}}}{\sum_{j=1}^n \mathbb{I}_{\{j d_p(x, \mathcal{X}_j) \leq h_{n,p}(x)\}}} \right)^2 \right) \\
&\quad + \left(\frac{c_{n,p}}{h_n(x)} \right)^2.
\end{aligned}$$

(4.1) Let $\mathbf{A}_{41} \doteq \mathbf{A}_4 \cap \{i: |d_p(x, \mathcal{X}_i) - d(x, \mathcal{X}_i)| \leq c_{n,p}\}$. In this case, by (H3.1) we get

$$(B.12) \quad I_{A_{41}}^1 \doteq \frac{c_{n,p}^2}{h_n^2(x)} \mathbb{E} \left(\frac{\sum_{i=1}^n \mathbb{I}_{\{i \in A_4\}}}{\left(\sum_{j=1}^n \mathbb{I}_{\{j d_p(x, \mathcal{X}_j) \leq h_{n,p}(x)\}} \right)^2} \right) + \left(\frac{c_{n,p}}{h_n(x)} \right)^2 \lesssim \left(\frac{c_{n,p}}{h_n(x)} \right)^2$$

and

$$(B.13) \quad II_{A_{41}}^1 \doteq \frac{c_{n,p}^2}{h_n^2(x)} \mathbb{E} \left(\left(\frac{\sum_{i=1}^n \mathbb{I}_{\{i \in A_4\}}}{\sum_{j=1}^n \mathbb{I}_{\{j d_p(x, \mathcal{X}_j) \leq h_{n,p}(x)\}}} \right)^2 \right) + \left(\frac{c_{n,p}}{h_n(x)} \right)^2 \lesssim \left(\frac{c_{n,p}}{h_n(x)} \right)^2.$$

(4.2) Let $\mathbf{A}_{42} \doteq \mathbf{A}_4 \cap \{i: |d_p(x, \mathcal{X}_i) - d(x, \mathcal{X}_i)| > c_{n,p}\}$. Let us define the i.i.d. random variables $Z_i \doteq d_p(x, \mathcal{X}_i) - d(x, \mathcal{X}_i)$, $i = 1, \dots, n$. Since $d_p(x, \mathcal{X}_i) \leq h_{n,p}(x)$ and $d(x, \mathcal{X}_i) \leq 3h_n(x)$ we have that $|Z_i| \leq h_{n,p}(x) + 3h_n(x)$. Observe that, from (H3.2) and (H3.1), respectively, for n large enough we have

$$h_{n,p} \leq h_n(x) + C_2 c_{n,p} \leq C h_n(x).$$

Which implies that, for n large enough, $|Z_i| \leq C h_n(x)$. Therefore,

$$\begin{aligned}
(B.14) \quad I_{A_{42}}^1 &\doteq \frac{1}{h_n^2(x)} \mathbb{E} \left(\sum_{i=1}^n |Z_i|^2 \mathbb{I}_{\{i c_{n,p} \leq |Z_i| \leq C h_n(x)\}} \right) + \left(\frac{c_{n,p}}{h_n(x)} \right)^2 \\
&\leq \frac{1}{h_n^2(x)} \mathbb{E} \left(\sum_{i=1}^n |Z_i|^2 \mathbb{I}_{\{i c_{n,p} \leq |Z_i| \leq C h_n(x)\}} \right) + \left(\frac{c_{n,p}}{h_n(x)} \right)^2 \\
&\leq \frac{n}{h_n^2(x)} \mathbb{E} (|Z_1|^2 \mathbb{I}_{\{c_{n,p} \leq |Z_1| \leq C h_n(x)\}}) + \left(\frac{c_{n,p}}{h_n(x)} \right)^2 \quad (\#A_{42} \leq n) \\
&\lesssim \frac{n}{h_n(x)} \mathbb{E} (|Z_1| \mathbb{I}_{\{c_{n,p} \leq |Z_1| \leq C h_n(x)\}}) + \left(\frac{c_{n,p}}{h_n(x)} \right)^2 \quad (|Z_1| \lesssim h_n(x)).
\end{aligned}$$

On the other hand,

$$\begin{aligned}
(B.15) \quad II_{A_{42}}^1 &\doteq \frac{1}{h_n^2(x)} \mathbb{E} \left(\left(\sum_{i=1}^n |Z_i| \mathbb{I}_{\{i c_{n,p} \leq |Z_i| \leq C h_n(x)\}} \right)^2 \right) + \left(\frac{c_{n,p}}{h_n(x)} \right)^2 \\
&\leq \frac{1}{h_n^2(x)} \mathbb{E} \left(\left(\sum_{i=1}^n |Z_i| \mathbb{I}_{\{i c_{n,p} \leq |Z_i| \leq C h_n(x)\}} \right)^2 \right) + \left(\frac{c_{n,p}}{h_n(x)} \right)^2.
\end{aligned}$$

Observe that, for $i \neq j$, Z_i is independent of Z_j , then

$$\begin{aligned}
& \mathbb{E} \left(\left(\sum_{i=1}^n |Z_i| \mathbb{I}_{\{i c_{n,p} \leq |Z_i| \leq C h_n(x)\}} \right)^2 \right) \\
&= \mathbb{E} \left(\sum_{i=1}^n \sum_{j=1}^n |Z_i| |Z_j| \mathbb{I}_{\{i c_{n,p} \leq |Z_i| \leq C h_n(x)\}} \mathbb{I}_{\{j c_{n,p} \leq |Z_j| \leq C h_n(x)\}} \right) \\
&= \mathbb{E} \left(\sum_{i=1}^n |Z_i|^2 \mathbb{I}_{\{i c_{n,p} \leq |Z_i| \leq C h_n(x)\}} \right) \\
&\quad + \mathbb{E} \left(\sum_{i=1}^n \sum_{\substack{j=1 \\ j \neq i}}^n |Z_i| |Z_j| \mathbb{I}_{\{i c_{n,p} \leq |Z_i| \leq C h_n(x)\}} \mathbb{I}_{\{j c_{n,p} \leq |Z_j| \leq C h_n(x)\}} \right) \\
&\leq n \mathbb{E} (|Z_1|^2 \mathbb{I}_{\{c_{n,p} \leq |Z_1| \leq C h_n(x)\}}) + n^2 \mathbb{E} (|Z_1| \mathbb{I}_{\{c_{n,p} \leq |Z_1| \leq C h_n(x)\}}) \mathbb{E} (|Z_1| \mathbb{I}_{\{c_{n,p} \leq |Z_1| \leq C h_n(x)\}}) \\
&\lesssim n h_n(x) \mathbb{E} (|Z_1| \mathbb{I}_{\{c_{n,p} \leq |Z_1| \leq C h_n(x)\}}) + n^2 (\mathbb{E} (|Z_1| \mathbb{I}_{\{c_{n,p} \leq |Z_1| \leq C h_n(x)\}}))^2 \quad (|Z_1| \lesssim h_n(x)).
\end{aligned}$$

Using this bound in (B.15), we get

$$\begin{aligned}
\text{(B.16)} \quad II_{A_{42}}^1 &\lesssim \frac{n}{h_n(x)} \mathbb{E} (|Z_1| \mathbb{I}_{\{c_{n,p} \leq |Z_1| \leq C h_n(x)\}}) \\
&\quad + \frac{n^2}{h_n^2(x)} (\mathbb{E} (|Z_1| \mathbb{I}_{\{c_{n,p} \leq |Z_1| \leq C h_n(x)\}}))^2 + \left(\frac{c_{n,p}}{h_n(x)} \right)^2.
\end{aligned}$$

We need to compute the expectation $\mathbb{E} (|Z_1| \mathbb{I}_{\{c_{n,p} \leq |Z_1| \leq C h_n(x)\}})$, which is

$$\begin{aligned}
\mathbb{E} (|Z_1| \mathbb{I}_{\{c_{n,p} \leq |Z_1| \leq C h_n(x)\}}) &= \int_{c_{n,p}}^{h_n(x)} \mathbb{P} (|Z_1| > t) dt \\
&\leq \mathbb{P} (|Z_1| > c_{n,p}) \int_{c_{n,p}}^{h_n(x)} dt \\
&\leq \mathbb{P} (|Z_1| > c_{n,p}) h_n(x).
\end{aligned}$$

Therefore, with this inequality in (B.14), we have

$$\begin{aligned}
\text{(B.17)} \quad I_{A_{42}}^1 &\lesssim n \mathbb{P} (|Z_1| > c_{n,p}) + \left(\frac{c_{n,p}}{h_n(x)} \right)^2 \\
&= n \mathbb{P} (|d_p(x, \mathcal{X}_1) - d(x, \mathcal{X}_1)| > c_{n,p}) + \left(\frac{c_{n,p}}{h_n(x)} \right)^2
\end{aligned}$$

and, with the same inequality in (B.16),

$$\begin{aligned}
\text{(B.18)} \quad II_{A_{42}}^1 &\lesssim n \mathbb{P} (|Z_1| > c_{n,p}) + (n \mathbb{P} (|Z_1| > c_{n,p}))^2 + \left(\frac{c_{n,p}}{h_n(x)} \right)^2 \\
&= n \mathbb{P} (|d_p(x, \mathcal{X}_1) - d(x, \mathcal{X}_1)| > c_{n,p}) \\
&\quad + \left(n \mathbb{P} (|d_p(x, \mathcal{X}_1) - d(x, \mathcal{X}_1)| > c_{n,p}) \right)^2 + \left(\frac{c_{n,p}}{h_n(x)} \right)^2.
\end{aligned}$$

Then, with (B.12) and (B.17) in (B.10) we get

$$(B.19) \quad I_{A_4}^1 \lesssim \left(\frac{c_{n,p}}{h_n(x)} \right)^2 + n\mathbb{P}(|d_p(x, \mathcal{X}_1) - d(x, \mathcal{X}_1)| > c_{n,p})$$

and, with (B.13) and (B.18) in (B.11),

$$(B.20) \quad II_{A_4}^1 \lesssim \left(\frac{c_{n,p}}{h_n(x)} \right)^2 + n\mathbb{P}(|d_p(x, \mathcal{X}_1) - d(x, \mathcal{X}_1)| > c_{n,p}) \\ + \left(n\mathbb{P}(|d_p(x, \mathcal{X}_1) - d(x, \mathcal{X}_1)| > c_{n,p}) \right)^2.$$

On the other hand, observe that $I_{A_4}^2 = \mathbb{E} \left(\left(\frac{\sum_{j=1}^n |K_j - K_{j,p}|}{\sum_{j=1}^n K_{j,p}} \right)^2 \right)$. Since $A_4^c = \{j : d(x, \mathcal{X}_j) > 3h_n(x)\} \cup \{j : d_p(x, \mathcal{X}_j) > h_{n,p}(x)\}$, we can write

$$\frac{\sum_{j=1}^n |K_j - K_{j,p}|}{\sum_{j=1}^n K_{j,p}} \leq \frac{\sum_{j=1}^n |K_j - K_{j,p}| \mathbb{I}_{\{j \in A_4\}}}{\sum_{j=1}^n K_{j,p}} \\ + \frac{\sum_{j=1}^n |K_j - K_{j,p}| \mathbb{I}_{\{j : d(x, \mathcal{X}_j) > 3h_n(x)\}}}{\sum_{j=1}^n K_{j,p}} \\ + \frac{\sum_{j=1}^n |K_j - K_{j,p}| \mathbb{I}_{\{j : d_p(x, \mathcal{X}_j) > h_{n,p}(x)\}}}{\sum_{j=1}^n K_{j,p}}.$$

Using that K is regular and that $\sum_{j=1}^n K_{j,p} \geq 1$ (this is since $\{j : d_p(x, \mathcal{X}_j) \leq h_{n,p}(x)\} \neq \emptyset$), we get

$$I_{A_4}^2 = \mathbb{E} \left(\left(\frac{\sum_{j=1}^n |K_j - K_{j,p}|}{\sum_{j=1}^n K_{j,p}} \right)^2 \right) \\ \lesssim II_{A_4}^1 + \mathbb{E} \left(\left(\sum_{j=1}^n |W_{j,p}| \mathbb{I}_{\{j : d_p(x, \mathcal{X}_j) \leq h_{n,p}(x), d(x, \mathcal{X}_j) > 3h_n(x)\}} \right)^2 \right) \\ + \frac{\sum_{j=1}^n K_j \mathbb{I}_{\{j : d_p(x, \mathcal{X}_j) > h_{n,p}(x)\}}}{\sum_{j=1}^n K_{j,p}} \\ \lesssim II_{A_4}^1 + II_{A_3} + \mathbb{E} \left(\left(\sum_{j=1}^n \mathbb{I}_{\{j : d_p(x, \mathcal{X}_j) > h_{n,p}(x), d(x, \mathcal{X}_j) \leq h_n(x)\}} \right)^2 \right) \\ \leq II_{A_4}^1 + II_{A_3} + C_{A_2},$$

where $II_{A_4}^1$ was defined in (B.9), II_{A_3} in (B.5), and C_{A_2} in (B.2). Then, from (B.20), (B.7) and (B.4), we have

$$(B.21) \quad I_{A_4}^2 \lesssim \left(\frac{c_{n,p}}{h_n(x)} \right)^2 + n\mathbb{P}(|d_p(x, \mathcal{X}_1) - d(x, \mathcal{X}_1)| > c_{n,p}) \\ + \left(n\mathbb{P}(|d_p(x, \mathcal{X}_1) - d(x, \mathcal{X}_1)| > c_{n,p}) \right)^2.$$

Therefore, with (B.19) and (B.21) in (B.8) we have

$$(B.22) \quad I_{A_4} \lesssim \left(\frac{c_{n,p}}{h_n(x)} \right)^2 + n\mathbb{P}(|d_p(x, \mathcal{X}_1) - d(x, \mathcal{X}_1)| > c_{n,p}) \\ + \left(n\mathbb{P}(|d_p(x, \mathcal{X}_1) - d(x, \mathcal{X}_1)| > c_{n,p}) \right)^2,$$

and with (B.20) and (B.21) in (B.9),

$$(B.23) \quad II_{A_4} \lesssim \left(\frac{c_{n,p}}{h_n(x)} \right)^2 + n\mathbb{P}(|d_p(x, \mathcal{X}_1) - d(x, \mathcal{X}_1)| > c_{n,p}) \\ + \left(n\mathbb{P}(|d_p(x, \mathcal{X}_1) - d(x, \mathcal{X}_1)| > c_{n,p}) \right)^2.$$

Finally, to complete the proof of this result (i.e. that I and II converge to zero) we need to show that the expectation on \mathcal{X} of

$$\left(\frac{c_{n,p}}{h_n(x)} \right)^2 + n\mathbb{P}_{\mathcal{X}_1}(|d_p(x, \mathcal{X}_1) - d(x, \mathcal{X}_1)| > c_{n,p}) + (n\mathbb{P}_{\mathcal{X}_1}^2(|d_p(x, \mathcal{X}_1) - d(x, \mathcal{X}_1)| > c_{n,p}))$$

converges to zero. In order to show it, recall that from H2 we have

$$n^2\mathbb{E}_{\mathcal{X}} \left(\mathbb{P}_{\mathcal{X}_1|\mathcal{X}}^2(|d_p(\mathcal{X}, \mathcal{X}_1) - d(\mathcal{X}, \mathcal{X}_1)| \geq c_{n,p}) \mid \mathcal{X} \in \text{supp}(\mu) \right) \rightarrow 0,$$

and consequently, by Cauchy Schwartz inequality,

$$n\mathbb{E}_{\mathcal{X}} \left(\mathbb{P}_{\mathcal{X}_1|\mathcal{X}}(|d_p(\mathcal{X}, \mathcal{X}_1) - d(\mathcal{X}, \mathcal{X}_1)| \geq c_{n,p}) \mid \mathcal{X} \in \text{supp}(\mu) \right) \rightarrow 0.$$

In addition, from (H3.1) we have

$$\mathbb{E}_{\mathcal{X}} \left(\left(\frac{c_{n,p}}{h_n(\mathcal{X})} \right)^2 \right) \rightarrow 0.$$

Therefore, taking expectation with respect to \mathcal{X} in (B.3), (B.4), (B.6), (B.7), (B.22) and (B.23), we prove Part (a) of the Theorem.

Proof of (b): The only difference with item (a) is the convergence of term III to zero which is ensured by Proposition 2.2. \square

Proof of Theorem 3.2: Let $\gamma_n \rightarrow \infty$ as $n \rightarrow \infty$ a sequence such that, as $n, p \rightarrow \infty$, $\mathbb{E}_{\mathcal{X}} \left(\gamma_n \left(\frac{c_{n,p}}{h_n(\mathcal{X})} \right)^2 \right) \rightarrow 0$ and, for each $i = 1, \dots, n$,

$$\gamma_n n^2 \mathbb{E}_{\mathcal{X}} \left(\mathbb{P}_{\mathcal{X}_i|\mathcal{X}}^2(|d(\mathcal{X}, \mathcal{X}_i) - d_p(\mathcal{X}, \mathcal{X}_i)| \geq c_{n,p}) \mid \mathcal{X} \in \text{supp}(\mu) \right) \rightarrow 0.$$

From proof of Theorem 3.1 we get

$$\mathbb{E}(\gamma_n(\widehat{\eta}_{n,p}(\mathcal{X}) - \eta(\mathcal{X}))^2) \lesssim \gamma_n n \mathbb{E}_{\mathcal{X}}(\mathbb{P}_{\mathcal{X}_1}(|d_p(x, \mathcal{X}_1) - d(x, \mathcal{X}_1)| \geq c_{n,p})) \\ + \mathbb{E}_{\mathcal{X}} \left(\gamma_n \left(\frac{c_{n,p}}{h_n(\mathcal{X})} \right)^2 \right) + \mathbb{E}(\gamma_n(\widehat{\eta}_n(\mathcal{X}) - \eta(\mathcal{X}))^2),$$

from what follows that

$$\lim_{n,p \rightarrow \infty} \mathbb{E}(\gamma_n(\widehat{\eta}_{n,p}(\mathcal{X}) - \eta(\mathcal{X}))^2) = 0. \quad \square$$

ACKNOWLEDGMENTS

The authors would like to thank one of the anonymous referees and the Associate Editor for their constructive comments which improved the present version of the paper.

REFERENCES

- [1] ABRAHAM, C.; BIAU, G. and CADRE, B. (2006). On the kernel rule for function classification, *Ann. Inst. Statist. Math.*, **58**, 619–633.
- [2] AMIRI, A.; CRAMBES, C. and THIAM, B. (2014). Recursive estimation of nonparametric regression with functional covariate, *Comput. Statist. Data Anal.*, **69**, 154–172.
- [3] BIAU, G.; BUNEA, F. and WEGKAMP, M.H. (2005). Functional classification in Hilbert spaces, *IEEE Trans. Inf. Theory*, **51**, 2163–2172.
- [4] BIAU, G.; CÉROU, F. and GUYADER, A. (2010). Rates of convergence of the functional k -nearest neighbor estimate, *IEEE Trans. Inform. Theory*, **56**, 2034–2040.
- [5] BIGOT, J. (2006). Landmark-based registration of curves via the continuous wavelet transform, *J. Comput. Graph. Statist.*, **15**, 542–564.
- [6] BREZIS, H. (2010). *Functional Analysis, Sobolev Spaces and Partial Differential Equations*, Springer-Verlag, Berlin.
- [7] BURBA, F.; FERRATY, F. and VIEU, P. (2009). k -nearest neighbor method in functional nonparametric regression, *Journal of Nonparametric Statistics*, **21**, 453–469.
- [8] CAI, T. and YUAN, M. (2011). Optimal estimation of the mean function based on discretely sampled functional data: Phase transition, *The Annals of Statistics*, **39**, 2330–2355.
- [9] CAI, T. and YUAN, M. (2016). Minimax and Adaptive Estimation of Covariance Operator for Random Variables Observed on a Lattice Graph, *J. Amer. Statist. Assoc.*, **39**, 2330–2355.
- [10] CÉROU, F. and GUYADER, A. (2006). Nearest neighbor classification in infinite dimension, *ESAIM Probab. Stat.*, **10**, 340–355.
- [11] CHAGNY, G. and ROCHE, A. (2014). Adaptive and minimax estimation of the cumulative distribution function given a functional covariate, *Electron. J. Stat.*, **8**, 2352–2404.
- [12] COLLOMB, G. (1980). Estimation de la regression par la méthode des k points les plus proches avec noyau: Quelques propriétés de convergence ponctuelle, *Lectures Notes in Mathematics*, **821**, 159–175, Springer-Verlag, Berlin.
- [13] FERRATY, F. and ROMAIN, Y., Eds. (2011). *The Oxford Handbook of Functional Data Analysis*, Oxford University Press.
- [14] FERRATY, F. and VIEU, P. (2002). The functional nonparametric model and application to spectrometric data, *Comput. Statist.*, **17**(4), 545–564.
- [15] FERRATY, F. and VIEU, P. (2006). *Nonparametric Functional Data Analysis. Theory and Practice*. Springer-Verlag, Berlin.
- [16] FORZANI, L.; FRAIMAN, R. and LLOP, P. (2012). Consistent nonparametric regression for functional data under the Stone–Besicovitch conditions, *IEEE Trans. Inform. Theory*, **58**, 6697–6708.
- [17] FORZANI, L.; FRAIMAN, R. and LLOP, P. (2014). Corrigendum to consistent nonparametric regression for functional data under the Stone–Besicovitch conditions, *IEEE Trans. Inform. Theory*, **60**, 3069.

- [18] HALL, P.; MÜLLER, H.G. and WANG, J.L. (2006). Properties of principal component methods for functional and longitudinal data analysis, *The Annals of Statistics*, **34**(3), 1493–1517.
- [19] HART, J.D. and WHERLY, T.E. (1986). Kernel regression estimation using repeated measurements data, *J. Amer. Statist. Assoc.*, **81**, 1080–1088.
- [20] HASTIE, T.J. and TIBSHIRANI, R.J. (1990). *Generalized Additive Models*, Chapman and Hall, London.
- [21] KNEIP, A. and RAMSAY, J.O. (2008). Combining registration and fitting for functional models, *J. Amer. Statist. Assoc.*, **103**, 1155–1165.
- [22] LIAN, H. (2011). Convergence of functional k -nearest neighbor regression estimate with functional responses, *Electronic Journal of Statistics*, **5**, 31–40.
- [23] MÜLLER, S. (2011). *Consistency and bandwidth selection for dependent data in non-parametric functional data analysis*, Ph.D. Thesis, Faculty of Mathematics and Physics, University of Stuttgart, Germany.
- [24] RAMSAY, J.O. and SILVERMAN, B.W. (1997). *Functional Data Analysis*, McGraw-Hill, New York.
- [25] RAMSAY, J.O. and SILVERMAN, B.W. (2002). *Applied Functional Data Analysis. Methods and Case Studies*. Springer-Verlag, New York.
- [26] RAMSAY, J.O. and SILVERMAN, B.W. (2005). *Functional Data Analysis*, 2nd ed., Springer-Verlag, New York.
- [27] RICE, J.A. and SILVERMAN, B.W. (1991). Estimating the mean and covariance structure nonparametrically when the data are curves, *J. Roy. Statist. Soc. Ser. B*, **53**, 233–243.

A TRANSITION MODEL FOR ANALYSIS OF ZERO-INFLATED LONGITUDINAL COUNT DATA USING GENERALIZED POISSON REGRESSION MODEL

Authors: TABAN BAGHFALAKI
– Department of Statistics, Faculty of Mathematical Sciences,
Tarbiat Modares University, Tehran, Iran
t.baghfalaki@modares.ac.ir

MOJTABA GANJALI
– Department of Statistics, Faculty of Mathematical Sciences,
Shahid Beheshti University, Tehran, Iran
m-ganjali@sbu.ac.ir

Received: September 2016

Revised: May 2017

Accepted: September 2017

Abstract:

- In most of the longitudinal studies, involving count responses, excess zeros are common in practice. Usually, the current response measurement in a longitudinal sequence is a function of previous outcomes. For example, in a study about acute renal allograft rejection, the number of acute rejection episodes for a patient in current time is a function of this outcome at previous follow-up times. In this paper, we consider a transition model for accounting the dependence of current outcome on the previous outcomes in the presence of excess zeros. We propose the use of the generalized Poisson distribution as a flexible distribution for considering overdispersion (or underdispersion). The maximum likelihood estimates of the parameters are obtained using the EM algorithm. Some simulation studies are performed for illustration of the proposed methods. Also, analysis of a real data set of a kidney allograft rejection study illustrates the application of the proposed model.

Key-Words:

- *count data; EM algorithm; generalized Poisson distribution; longitudinal data; transition models; zero-inflated models.*

AMS Subject Classification:

- 62J99, 62P10.

1. INTRODUCTION

In modeling many count longitudinal clinical studies, the excess of zero is a common problem. For example, in a study about acute renal allograft rejection, many patients may have no acute rejection episodes at some follow-up times or in an asthma-related study, if the response variable is the number of asthma-related hospitalizations at each follow-up time, many patients may report no asthma-related hospitalizations. In these examples, the response variable for the patient can be considered as a count variable which may be recorded with extra zeros. Useful models for describing these kinds of data sets are zero-inflated models. In these models a special probability is allocated to zero observations (see Section 2 for more details).

Several approaches are proposed for analyzing these data sets. For example, hurdle model [25, 15, 16] and zero-inflated Poisson (ZIP) model [19, 12] are two well-known approaches for analysing zero-inflated (ZI) count data. Also, zero-inflated generalized Poisson (ZIGP) and zero-inflated negative binomial (ZINB) models are two other well-known approaches for considering overdispersion of which ZIGP model can also consider underdispersion to analyse inflated count data. [7] proposed a ZIGP model to analyse the data set of outsourcing of patent applications.

The analysis of longitudinal ZI count data are discussed frequently in literature. [4] proposed ZIP and ZINB models for analysing data of a study of growth. They describe their approaches as mixture models with a proportion P of subjects not at risk, and a proportion of $1-P$ at risk subjects who take on outcome values following a Poisson or negative binomial distribution. [21] used the ZIP and ZINB models to analyze longitudinal studies in epidemiology. [23] proposed a random effect model to analysis the ZI longitudinal count data. [28] discussed application of the ZI and hurdle models for longitudinal studies concerning vaccination safety. [14] used ZIP regression for analysing longitudinal data. [2] proposed a two-part regression model for analysing ZI longitudinal count data. They used their proposed approach for analysing an healthcare utilization data set. [26] discussed a Bayesian paradigm for ZIP and ZINB model for analysing data set of a study of psychiatric outpatient service. [27] provided a review of the literature and tests the Poisson, the ZIP, the negative binomial (NB) and the ZINB models in the context of longitudinal count data. [3] give many examples of the use of ZI distributions to model longitudinal data and consider this approach as a conventional one. [22] described a mixed-effect hurdle model for ZI longitudinal count data, where a baseline variable is included in the model specification. They used their proposed approach to analyse a healthcare utilization data.

A common problem in the practice of studying count data is overdispersion or underdispersion. The use of Poisson distribution to analyze count data has a lack of fit because of ignoring to consider these problems. To deal with overdispersion the use of NB distribution is proposed. But, this distribution has a lack of fit for considering the possible underdispersion. A distribution function which considers both the overdispersion and underdispersion is the generalized Poisson distribution [6, 5]. Note that the zero-inflation generally involve overdispersion or underdispersion. Here, the use of ZIGP distribution is recommended to consider both problems of underdispersion and overdispersion. Underdispersion is rarely occurred in practice. Therefore, the most concern of this paper is on the overdispersion in zero-inflated longitudinal data.

Three main modeling families are introduced to model longitudinal data: marginal models, subject-specified models and conditionally specified models [9, 24]. In a marginal model, marginal distributions are used to describe the longitudinal outcomes vector given a set of predictor variables. The correlation among the components of the longitudinal measurements can be captured by a fully parametric approach or by modeling a limited number of lower-order moments such as generalized estimating equations (GEE). In random effects or subject-specified models the longitudinal outcome vector is modeled by a vector of random effects. Several software and programs, for instance SAS and Mplus, make it possible to fit ZIP and ZINB distributions to longitudinal ZI data using random effects models. Finally in a conditionally specified model any response within the sequence of longitudinal measurements is modeled conditional upon the outcome on the previous time or a subset of previous outcomes. A particular relevant class of conditional models is the so-called autoregressive or transition models. In a transition model a current measurement in a longitudinal study is described as a function of the previous outcomes [9]. In this paper, our focus is on transition models. For some applications of the transition models in repeated measurement outcomes see [1, 18, 11]. Also, for reviews of transition models for analyzing the longitudinal data see [9], [30] and [10].

In this paper, we use the ZIGP transition models to analyze longitudinal count data with extra zeros. We use the usual EM algorithm for parameters estimation. The proposed model is illustrated using some simulation studies, where the performance of the proposed distributional assumption for transition model is compared with ZIP, ZINB, NB and GP distributional assumptions. Also, the proposed method is used for analyzing a real data set of a kidney allograft rejection study in application section where the best fitting model is selected by using Akaike information criterion (AIC), Bayesian information criterion (BIC) and Hannan–Quinn criterion (HQC).

This paper is organized as follows: Section 2 is a review on generalized Poisson and zero-inflated generalized Poisson distributions and the relation of these distributions with Poisson and zero-inflated Poisson distributions. Section 3 includes some notation, definitions of models, likelihood functions, the EM algorithm and our illustration of the proposed transition model for analyzing zero-inflated longitudinal data. In Section 4, some simulation studies are performed. In this section four different structures are considered for generating data and performance of ZIGP, ZINB, ZIP, NB and GP transition models are compared for each structure. The description and the analysis of a real data set using the proposed model are given and comparison of the performance of our approach with some other distributional assumptions is given in Section 5. The last section includes some conclusions and discussions.

2. ZERO-INFLATED GENERALIZED POISSON DISTRIBUTION

The random variable Y is said to have a generalized Poisson distribution, if its probability mass function is given by

$$(2.1) \quad f(y; \xi, \omega) = \frac{\xi(\xi + \omega y)^{y-1}}{y!} e^{-(\xi + \omega y)}, \quad y = 0, 1, 2, \dots$$

where $\xi > 0$ and $\max(-1, -\xi/4) < \omega < 1$ [13]. The mean and variance of this distribution are given by

$$E(Y) = \frac{\xi}{1-\omega}, \quad \text{Var}(Y) = \frac{\xi}{(1-\omega)^3} = \frac{1}{(1-\omega)^2} E(Y),$$

therefore, the term $\frac{1}{(1-\omega)^2}$ plays the role of a dispersion factor. Clearly, when $\omega = 0$, the generalized Poisson distribution reduces to the usual Poisson distribution with parameter ξ . Further, when $\omega > 0$, we have overdispersion in the model; when $\omega < 0$, we have underdispersion.

A parameterization of this distribution is given by setting $\lambda = \frac{\xi}{1-\omega}$ and $\phi = \frac{\omega}{\xi}$, denoted by $Y \sim GP(\lambda, \phi)$, and its probability mass function is given by

$$(2.2) \quad f_{GP}(y; \lambda, \phi) = \left(\frac{\lambda}{1 + \phi\lambda} \right)^y \frac{(1 + \phi y)^{y-1}}{y!} \exp\left(\frac{-\lambda(1 + \phi y)}{1 + \phi\lambda} \right), \quad y = 0, 1, 2, \dots, \quad \lambda > 0,$$

where ϕ is a real value parameter such that for all y , $1 + \phi y > 0$ and $1 + \phi\lambda > 0$. These restrictions are confirmed by the restriction on the distribution (2.1). The generalized Poisson distribution (2.2) is a natural extension of the Poisson distribution. If $\phi = 0$, then the probability function (2.2) reduces to the Poisson distribution, denoted by $Y \sim P(\lambda)$. By the above mentioned parameterization, the mean of Y is given by $E(Y) = \lambda$ and the variance of Y is given by $\text{Var}(Y) = \lambda(1 + \phi\lambda)^2$. In the generalized Poisson distribution, the ϕ parameter is called dispersion parameter. When $\phi > 0$, the overdispersion is presented in the model, whereas when $\phi < 0$, the underdispersion is included in the model. The generalized Poisson distribution is a more flexible distribution than the negative binomial distribution for considering possibility of underdispersion or overdispersion. This property is one of the well-known properties of generalized Poisson distribution. [17] proved that the generalized Poisson distribution, the same as negative binomial distribution, can be considered as a mixture of the Poisson distribution. [17] show that there are some differences between the fits of the generalized Poisson and negative binomial distributions. When the first two moments are fixed, the negative binomial distribution have larger mass at zero than the generalized Poisson distribution. This means their zero-inflated variations tend to have larger discrepancy. However, the fits of their zero-inflated variations may differ when there is a large zero fraction [17]. For more details about generalized Poisson distribution see [6] and [5]. Also, *VGAM* and *HMMpa* packages of R can be applied to use the generalized Poisson distribution.

A zero-inflated generalized Poisson distribution for a positive value π ($0 \leq \pi \leq 1$) is defined as follows:

$$(2.3) \quad f_{ZIGP}(y; \lambda, \phi, \pi) = \begin{cases} \pi + (1 - \pi)f_{GP}(0; \lambda, \phi), & y = 0, \\ (1 - \pi)f_{GP}(y; \lambda, \phi), & y > 0, \end{cases}$$

where $f_{GP}(\cdot; \lambda, \phi)$ is the probability mass function of (2.2). We will use the notation $Y \sim ZIGP(\lambda, \phi, \pi)$ to denote the distribution of (2.3). The mean and variance of this distribution are given by $E(Y) = (1 - \pi)\lambda$ and $\text{var}(Y) = E(Y)[(1 + \phi\lambda)^2 + \pi\lambda]$, respectively. The variance of this distribution shows that for $\pi > 0$ and $\phi > 0$ the distribution of Y exhibits overdispersion. The distribution (2.3) reduced to the generalized Poisson distribution when $\pi = 0$ and it reduced to zero-inflated Poisson distribution when $\phi = 0$, denoted by $Y \sim ZIP(\lambda, \pi)$. When π is allowed to be negative, the distribution (2.3) presents a zero-deflated generalized Poisson distribution which rarely occurs in practice.

3. ZERO-INFLATED TRANSITION MODELS FOR COUNT RESPONSES

Suppose N individuals are participated in a longitudinal study and for each individual n_i ($i = 1, 2, \dots, N$) repeated measurements are recorded as response variables. Also, let Y_{ij} , $i = 1, 2, \dots, N$ and $j = 1, 2, \dots, n_i$ be the longitudinal measurements for the i^{th} individual at j^{th} time point and let W_{ij} , $i = 1, 2, \dots, N$ and $j = 1, 2, \dots, n_i$ be indicator variables as follows:

$$W_{ij} = \begin{cases} 1, & Y_{ij} \text{ is from the perfect state,} \\ 0, & Y_{ij} \text{ is from the Poisson state.} \end{cases}$$

where by perfect we means that the sample is from a degenerated distribution at 0. It is clear that W_{ij} is a latent variable. Also, let $\mathbf{h}_{ij} = (Y_{i1}, \dots, Y_{i,j-1})$ be the previous outcomes up to time j or in other words history of outcomes for the i^{th} individual.

In a transition model, the outcome Y_{ij} is modeled in term of \mathbf{h}_{ij} [9]. The order of a transition model is the number of the previous measurements that are considered for modeling the measurement of the current time. We consider a first order zero-inflated transition model as follows:

$$\begin{aligned} P_{ZI}(Y_{ij} = y_{ij} | \pi_{ij}, \lambda_{ij}, \phi, \mathbf{x}_{ij}, \mathbf{z}_{ij}, y_{i,j-1}) &= \\ (3.1) \quad &= \begin{cases} \pi_{ij} + (1 - \pi_{ij}) P(Y_{ij} = y_{ij} | \lambda_{ij}, \phi, \mathbf{x}_{ij}, y_{i,j-1}), & y_{ij} = 0, \\ (1 - \pi_{ij}) P(Y_{ij} = y_{ij} | \lambda_{ij}, \phi, \mathbf{x}_{ij}, y_{i,j-1}), & y_{ij} \neq 0, \end{cases} \end{aligned}$$

where

$$(3.2) \quad \log(\lambda_{i1}) = \mathbf{x}'_{i1} \boldsymbol{\beta},$$

$$(3.3) \quad \text{logit}(\pi_{i1}) = \mathbf{z}'_{i1} \boldsymbol{\alpha},$$

and, for $j = 2, 3, \dots, n_i$,

$$(3.4) \quad \log(\lambda_{ij}) = \mathbf{x}'_{ij} \boldsymbol{\beta} + \gamma_1 I_{\{0\}}(Y_{i,j-1}) + \gamma_2 y_{i,j-1} (1 - I_{\{0\}}(Y_{i,j-1})),$$

$$(3.5) \quad \text{logit}(\pi_{ij}) = \mathbf{z}'_{ij} \boldsymbol{\alpha} + \tau_1 I_{\{0\}}(Y_{i,j-1}) + \tau_2 y_{i,j-1} (1 - I_{\{0\}}(Y_{i,j-1})), \quad j = 2, \dots, n_i,$$

where $\pi_{ij} = P(Y_{ij} = 0 | \boldsymbol{\alpha}, \mathbf{z}_{ij}, \mathbf{h}_{ij}) = P(Y_{ij} = 0 | \boldsymbol{\alpha}, \mathbf{z}_{ij}, y_{i,j-1})$ is the rate of zeros given some covariates and the history of outcomes. In this model the effect of the previous zero response on current measurement (γ_1) and the effect of the non-zero previous response on current mean (γ_2) are separately considered. This is due to the fact that one expects to have the current mean to be close to the previous values of responses.

We will use the notation $Y \sim ZIGP(\lambda_{ij}, \pi_{ij}, \phi)$ to denote model (3.1). Note that (3.1) is reduced to zero-inflated Poisson model when $\phi = 0$. We will use the notation $Y \sim ZIP(\lambda_{ij}, \pi_{ij})$ to denote model (3.1) when $\phi = 0$. Let $\boldsymbol{\theta} = (\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\tau}, \phi)$ be the vector of all the unknown parameters in the model where $\boldsymbol{\alpha} = (\alpha_1, \alpha_2)'$ and $\boldsymbol{\gamma} = (\gamma_1, \gamma_2)'$. The likelihood

function of the model can be written as:

$$\begin{aligned} L(\boldsymbol{\theta}|\mathbf{y}, \mathbf{x}, \mathbf{z}) &= \prod_{i=1}^N \left\{ f(y_{i1}) \times \prod_{j=2}^{n_i} f(y_{ij}|y_{i,j-1}) \right\} \\ &= \prod_{i=1}^N \prod_{j=1}^{n_i} \left(\pi_{ij} + (1 - \pi_{ij}) P(Y_{ij} = 0 | \lambda_{ij}, \phi, \mathbf{x}_{ij}, \mathbf{h}_{ij}) \right)^{I(y_{ij}=0)} \\ &\quad \times \left((1 - \pi_{ij}) P(Y_{ij} \neq 0 | \lambda_{ij}, \phi, \mathbf{x}_{ij}, \mathbf{h}_{ij}) \right)^{1-I(y_{ij}=0)}, \end{aligned}$$

where $h_{i1} = 0$ and it will not be considered in the model. This likelihood function can be maximized using some numerical methods such as Newton–Raphson [20].

Another approach for obtaining parameter estimates is the use of the Expectation-Maximization (EM) [8] algorithm. To obtain the EM estimates of the parameters, we use the indicator variable, W_{ij} , $i = 1, 2, \dots, N$, $j = 1, 2, \dots, n_i$. As mentioned earlier W_{ij} is a latent variable for indicating the perfect state versus the Poisson state outcome. Therefore, the log-likelihood function of (\mathbf{Y}, \mathbf{W}) as complete data is given by

$$\begin{aligned} \ell_c(\boldsymbol{\theta}|\mathbf{y}, \mathbf{w}, \mathbf{x}, \mathbf{z}) &= \sum_{i=1}^N \sum_{j=1}^{n_i} w_{ij} \log(\pi_{ij}) + \sum_{i=1}^N \sum_{j=1}^{n_i} (1 - w_{ij}) \log(1 - \pi_{ij}) \\ &\quad + \sum_{i=1}^N \sum_{j=1}^{n_i} (1 - w_{ij}) \left\{ y_{ij} \log(\lambda_{ij}) - y_{ij} \log(1 + \phi \lambda_{ij}) \right. \\ &\quad \left. + (y_{ij} - 1) \log(1 + \phi y_{ij}) - \log(y_{ij}!) - \lambda_{ij} \frac{1 + \phi y_{ij}}{1 + \phi \lambda_{ij}} \right\}. \end{aligned}$$

The EM algorithm contains two steps: in the first step (E-step), the expectation of the complete likelihood function (here $\ell_c(\boldsymbol{\theta}|\mathbf{y}, \mathbf{w}, \mathbf{x}, \mathbf{z})$) given the observed data (here \mathbf{Y}) and the current value of the parameters in the r^{th} step (called $\boldsymbol{\theta}^{(r)}$) is calculated, by defining $Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(r)}) = E[\ell_c(\boldsymbol{\theta}|\mathbf{y}, \mathbf{w}, \mathbf{x}, \mathbf{z}) | \mathbf{y}, \mathbf{x}, \mathbf{z}, \boldsymbol{\theta}^{(r)}]$. We have

$$\begin{aligned} Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(r)}) &= \sum_{i=1}^N \sum_{j=1}^{n_i} E[W_{ij} | \mathbf{y}, \mathbf{x}, \mathbf{z}, \boldsymbol{\theta}^{(r)}] \log(\pi_{ij}) \\ &\quad + \sum_{i=1}^N \sum_{j=1}^{n_i} (1 - E[W_{ij} | \mathbf{y}, \mathbf{x}, \mathbf{z}, \boldsymbol{\theta}^{(r)}]) \log(1 - \pi_{ij}) \\ &\quad + \sum_{i=1}^N \sum_{j=1}^{n_i} (1 - E[W_{ij} | \mathbf{y}, \mathbf{x}, \mathbf{z}, \boldsymbol{\theta}^{(r)}]) \left\{ y_{ij} \log(\lambda_{ij}) - y_{ij} \log(1 + \phi \lambda_{ij}) \right. \\ &\quad \left. + (y_{ij} - 1) \log(1 + \phi y_{ij}) - \log(y_{ij}!) - \lambda_{ij} \frac{1 + \phi y_{ij}}{1 + \phi \lambda_{ij}} \right\}. \end{aligned}$$

For computing the EM algorithm, the following expectation is needed:

$$\begin{aligned} E[W_{ij} | \mathbf{y}, \mathbf{w}, \mathbf{x}, \mathbf{z}, \boldsymbol{\theta}^{(r)}] &= P(W_{ij} = 1 | y_{i,j-1}, \mathbf{x}, \mathbf{z}, \boldsymbol{\theta}^{(r)}) \\ &= \begin{cases} \frac{\pi_{ij}^{(r)}}{\pi_{ij}^{(r)} + (1 - \pi_{ij}^{(r)}) P(Y_{ij} = 0 | \lambda_{ij}^{(r)}, \phi^{(r)}, \mathbf{x}_{ij}, y_{i,j-1})}, & y_{ij} = 0, \\ 0, & y_{ij} \neq 0, \end{cases} \end{aligned}$$

where $\pi_{ij}^{(r)} = P(Y_{ij} = 0 | \mathbf{z}_{ij}, \boldsymbol{\theta}^{(r)}, y_{i,j-1})$ and $\lambda_{ij}^{(r)}$ is the current Poisson rate at the r^{th} iteration.

In the second step (M-step), we define

$$\boldsymbol{\theta}^{(r+1)} = \arg \max_{\boldsymbol{\theta} \in \Theta} Q(\boldsymbol{\theta} | \boldsymbol{\theta}^{(r)}).$$

The algorithm is converged and is stopped when

$$\left\| \boldsymbol{\theta}^{(r)} - \boldsymbol{\theta}^{(r+1)} \right\| < \varepsilon,$$

where $\|\cdot\|$ is a pre-specified measure.

4. SIMULATION STUDIES

In this section some simulation studies are performed for investigating the performance of the proposed approach. At first, the data are generated from ZIGP and the performance of ZIGP, GP, ZINB, NB and ZIP are compared on analyzing these data. Two other simulated data are generated under ZINB and ZIP where the performance of analyzing ZIGP, ZINB and ZIP are investigated in each case. Note that ZIP model is a ZIGP model with $\phi = 0$. The last simulation study is used to examine the performance of ZIGP, ZINB and ZIP in the presence of underdispersion.

4.1. Zero-inflated generalized Poisson model

In this simulation study the data set is generated from a transition model under ZIGP. The simulation study contains two sample sizes $N = 100$ and 500 where $M = 1000$ iterations are performed. For generating data, we consider a ZIGP model as follows:

$$(4.1) \quad Y_{ij} | \lambda_{ij}, \pi_{ij} \sim ZIGP(\lambda_{ij}, \pi_{ij}, \phi),$$

where

$$(4.2) \quad \begin{aligned} \log(\lambda_{i1}) &= \beta_0 + \beta_1 x_i + \beta_2 t_1, \\ \text{logit}(\pi_{i1}) &= \alpha_0 + \alpha_1 x_i + \alpha_2 t_1, \\ \log(\lambda_{ij}) &= \beta_0 + \beta_1 x_i + \beta_2 t_j + \gamma_1 I_{\{0\}}(Y_{i,j-1}) + \gamma_2 y_{i,j-1} (1 - I_{\{0\}}(Y_{i,j-1})), \quad j = 2, 3, 4, \\ \text{logit}(\pi_{ij}) &= \alpha_0 + \alpha_1 x_i + \alpha_2 t_j + \tau_1 I_{\{0\}}(Y_{i,j-1}) + \tau_2 y_{i,j-1} (1 - I_{\{0\}}(Y_{i,j-1})), \quad j = 2, 3, 4. \end{aligned}$$

For this simulation study, two sets of real values are considered as follows:

- 1) $\alpha_0 = -1, \alpha_1 = 1, \alpha_2 = 0, \beta_0 = -3, \beta_1 = \beta_2 = 1, \gamma_1 = -1, \gamma_2 = 0, \tau_1 = 0, \tau_2 = 1$ and $\phi = 1$.
- 2) $\alpha_0 = -1, \alpha_1 = -1, \alpha_2 = 0, \beta_0 = -3, \beta_1 = \beta_2 = 1, \gamma_1 = 1, \gamma_2 = -1, \tau_1 = 1, \tau_2 = -1$ and $\phi = 0.5$.

The results of these simulation studies are summarized in Tables 1 and 2, respectively.

Table 1: Results of simulation study for generated data under ZIGP model, estimate (Est.), standard error (S.E.), relative bias (Bias) and mean square error (MSE) for $M = 1000$ simulated data with sample sizes 100 and 500 and the first set of real values.

N	Para.	Real	ZIGP			GP			ZINB			NB			ZIP			
			Est. (S.E.)	Bias	MSE	Est. (S.E.)	Bias	MSE	Est. (S.E.)	Bias	MSE	Est. (S.E.)	Bias	MSE	Est. (S.E.)	Bias	MSE	
100	α_0	-1.00	-1.26 (0.39)	0.26	0.46	—	—	—	-15.35 (28.71)	14.35	102.37	—	—	—	0.53 (0.71)	-1.53	2.84	
	α_1	1.00	1.18 (0.97)	0.18	0.91	—	—	—	8.51 (18.13)	7.51	3.82	—	—	—	0.05 (0.12)	-0.94	0.90	
	α_2	0.00	-0.04 (0.39)	*	0.14	—	—	—	0.48 (2.81)	*	8.11	—	—	—	0.02 (0.03)	*	0.00	
	τ_1	0.00	0.14 (0.35)	*	0.48	—	—	—	1.19 (23.82)	*	564.39	—	—	—	0.16 (0.13)	*	0.21	
	τ_2	-1.00	-0.93 (0.69)	-0.06	0.85	—	—	—	-3.42 (9.77)	2.42	100.75	—	—	—	-0.02 (0.04)	-0.97	0.95	
	β_0	-3.00	-2.98 (0.59)	-0.00	0.32	-3.35 (0.46)	0.11	0.33	-3.33 (0.64)	0.11	0.53	-3.16 (0.50)	0.05	0.26	-0.05 (0.29)	-0.98	8.76	
	β_1	1.00	1.03 (0.42)	0.03	0.17	0.62 (0.28)	-0.37	0.21	1.02 (0.41)	0.02	0.16	0.64 (0.27)	-0.35	0.19	0.11 (0.18)	-0.88	0.81	
	β_2	1.00	0.99 (0.21)	-0.00	0.04	1.02 (0.21)	0.02	0.04	1.04 (0.18)	0.04	0.03	1.04 (0.14)	0.04	0.02	0.36 (0.15)	-0.63	0.42	
	γ_1	-1.00	-0.96 (0.45)	-0.03	0.19	-1.17 (0.55)	0.17	0.33	-1.07 (0.48)	0.07	0.23	-1.34 (0.31)	0.34	0.21	-0.05 (0.20)	-0.94	0.93	
	γ_2	0.00	0.09 (0.39)	*	0.15	-0.00 (0.02)	*	0.00	-0.02 (0.09)	*	0.10	0.01 (0.08)	*	0.00	0.01 (0.07)	*	0.00	
	ϕ	1.00	0.85 (0.16)	-0.14	0.04	1.86 (0.27)	0.86	0.86	0.27 (0.10)	-0.72	0.54	0.16 (0.02)	-0.83	0.69	—	—	—	
	500	α_0	-1.00	-0.95 (0.23)	-0.05	0.39	—	—	—	-11.92 (7.01)	10.92	166.61	—	—	—	0.09 (0.21)	-1.09	1.23
		α_1	1.00	1.06 (0.40)	0.06	0.16	—	—	—	10.01 (7.39)	9.01	133.75	—	—	—	0.08 (0.14)	-0.91	0.84
		α_2	0.00	-0.03 (0.10)	*	0.01	—	—	—	0.41 (0.24)	*	0.22	—	—	—	-0.08 (0.17)	*	0.03
τ_1		0.00	0.09 (0.40)	*	0.16	—	—	—	-0.51 (1.07)	*	1.37	—	—	—	0.24 (0.31)	*	0.15	
τ_2		-1.00	-0.97 (0.45)	-0.03	0.25	—	—	—	-1.89 (2.22)	0.89	5.53	—	—	—	-0.02 (0.03)	-0.97	0.95	
β_0		-3.00	-2.99 (0.31)	-0.00	0.09	-3.21 (0.28)	0.07	0.26	-3.35 (0.19)	0.11	0.16	-3.33 (0.26)	0.11	0.18	-0.40 (0.65)	-0.86	7.16	
β_1		1.00	1.00 (0.17)	0.00	0.03	0.65 (0.26)	-0.34	0.18	0.92 (0.14)	-0.07	0.02	0.71 (0.16)	-0.28	0.10	0.27 (0.33)	-0.72	0.62	
β_2		1.00	1.01 (0.08)	0.01	0.00	1.02 (0.21)	0.02	0.04	1.08 (0.06)	0.08	0.01	1.07 (0.06)	0.07	0.01	0.49 (0.21)	-0.50	0.30	
γ_1		-1.00	-1.02 (0.18)	0.02	0.03	-1.17 (0.57)	0.17	0.34	-1.23 (0.16)	0.23	0.08	-1.28 (0.17)	0.28	0.10	-0.25 (0.39)	-0.74	0.70	
γ_2		0.00	-0.00 (0.02)	*	0.00	0.27 (0.52)	*	0.34	-0.00 (0.02)	*	0.00	0.00 (0.03)	*	0.00	0.01 (0.06)	*	0.00	
ϕ		1.00	0.93 (0.11)	-0.06	0.01	1.76 (0.28)	0.76	0.66	0.21 (0.02)	-0.78	0.62	0.45 (0.01)	-0.55	0.52	—	—	—	

Table 2: Results of simulation study for generated data under ZIGP model, estimate (Est.), standard error (S.E.), relative bias (Bias) and mean square error (MSE) for $M = 1000$ simulated data with sample sizes 100 and 500 and the second set of real values.

N	Para.	Real	ZIGP			GP			ZINB			NB			ZIP			
			Est. (S.E.)	Bias	MSE	Est. (S.E.)	Bias	MSE	Est. (S.E.)	Bias	MSE	Est. (S.E.)	Bias	MSE	Est. (S.E.)	Bias	MSE	
100	α_0	-1.00	-1.19 (0.72)	0.09	0.81	—	—	—	65.76 (198.29)	-66.76	37221.00	—	—	—	0.11 (0.33)	-1.11	1.31	
	α_1	-1.00	-0.98 (0.53)	-0.02	0.22	—	—	—	-82.51 (200.41)	81.51	40113.72	—	—	—	-0.58 (0.71)	-0.42	0.56	
	α_2	0.00	0.03 (0.17)	*	0.02	—	—	—	0.93 (0.51)	*	1.08	—	—	—	-0.06 (0.12)	*	0.01	
	τ_1	1.00	1.14 (0.40)	0.14	0.55	—	—	—	-70.94 (199.79)	-71.94	38439.76	—	—	—	0.49 (0.59)	-0.51	0.52	
	τ_2	-1.00	-1.01 (0.21)	-0.01	0.16	—	—	—	-82.87 (193.31)	81.87	37844.46	—	—	—	0.31 (0.34)	-1.31	1.81	
	β_0	-3.00	-3.11 (0.40)	0.04	0.14	-3.40 (0.35)	0.13	0.27	-3.62 (0.68)	0.21	0.76	-3.57 (0.65)	0.19	0.53	-0.79 (1.15)	-0.74	5.87	
	β_1	1.00	1.08 (0.23)	0.08	0.05	1.22 (0.18)	0.22	0.08	1.22 (0.48)	0.22	0.24	1.40 (0.56)	0.40	0.32	0.20 (0.18)	-0.80	0.66	
	β_2	1.00	0.98 (0.15)	-0.02	0.02	1.03 (0.12)	0.03	0.01	1.15 (0.08)	0.15	0.03	0.95 (0.11)	-0.05	0.01	0.55 (0.17)	-0.45	0.23	
	γ_1	1.00	1.14 (0.15)	0.14	0.04	0.62 (0.39)	-0.38	0.28	0.75 (0.31)	-0.25	0.14	1.06 (0.22)	0.06	0.03	0.72 (0.76)	-0.28	0.51	
	γ_2	-1.00	-1.05 (0.26)	0.05	0.06	-0.98 (0.28)	-0.02	0.07	-1.22 (0.33)	0.22	0.14	-0.84 (0.25)	-0.16	0.06	-0.14 (0.19)	-0.86	0.77	
	ϕ	0.50	0.47 (0.03)	-0.06	0.00	0.91 (0.11)	0.82	0.18	0.34 (0.06)	-0.33	0.03	0.28 (0.01)	-0.44	0.05	—	—	—	
	500	α_0	-1.00	-1.08 (0.32)	0.08	0.37	—	—	—	-20.95 (5.89)	19.95	426.92	—	—	—	1.45 (0.29)	-2.45	6.06
		α_1	-1.00	-1.12 (0.03)	0.12	0.01	—	—	—	-0.84 (0.28)	-0.16	0.09	—	—	—	-1.02 (0.09)	0.02	0.00
		α_2	0.00	-0.10 (0.06)	*	0.01	—	—	—	0.44 (0.12)	*	0.21	—	—	—	-0.25 (0.00)	*	0.06
τ_1		1.00	1.01 (0.10)	0.01	0.11	—	—	—	18.64 (5.94)	17.64	340.38	—	—	—	0.17 (0.28)	-0.83	0.73	
τ_2		-1.00	-1.04 (0.13)	-0.04	0.11	—	—	—	-0.42 (3.05)	-0.58	8.07	—	—	—	0.23 (0.17)	-1.23	1.53	
β_0		-3.00	-3.00 (0.31)	0.00	0.05	-3.46 (0.04)	0.15	0.21	-3.25 (0.13)	0.08	0.08	-3.39 (0.19)	0.13	0.18	-1.17 (0.03)	-0.61	3.34	
β_1		1.00	1.01 (0.01)	0.01	0.00	1.29 (0.04)	0.29	0.09	1.14 (0.11)	0.14	0.03	1.32 (0.16)	0.32	0.12	0.32 (0.30)	-0.68	0.50	
β_2		1.00	0.98 (0.01)	-0.02	0.00	1.05 (0.01)	0.05	0.00	1.02 (0.07)	0.02	0.00	0.97 (0.04)	-0.03	0.00	0.75 (0.07)	-0.25	0.07	
γ_1		1.00	1.11 (0.25)	0.11	0.04	0.61 (0.06)	-0.39	0.15	0.82 (0.15)	-0.18	0.05	0.79 (0.08)	-0.21	0.05	0.68 (0.05)	-0.32	0.10	
γ_2		-1.00	-0.98 (0.10)	0.02	0.02	-0.98 (0.06)	-0.02	0.00	-0.98 (0.11)	-0.02	0.01	-0.89 (0.08)	-0.11	0.02	-0.72 (0.01)	-0.28	0.08	
ϕ		0.50	0.47 (0.02)	-0.05	0.00	0.92 (0.03)	0.85	0.18	0.38 (0.02)	-0.25	0.02	0.26 (0.01)	-0.49	0.06	—	—	—	

The simulated data set are analyzed using NB, GP, ZIP, ZINB and ZIGP models, such that

$$\begin{aligned}
 Y_{ij}|\lambda_{ij}, \phi &\sim NB\left(\phi, \frac{\phi}{\phi + \lambda_{ij}}\right), \\
 Y_{ij}|\lambda_{ij}, \phi &\sim GP(\lambda_{ij}, \phi), \\
 (4.3) \quad Y_{ij}|\lambda_{ij}, \pi_{ij} &\sim ZIP(\lambda_{ij}, \pi_{ij}), \\
 Y_{ij}|\lambda_{ij}, \pi_{ij}, \phi &\sim ZINB\left(\phi, \frac{\phi}{\phi + \lambda_{ij}}, \pi_{ij}\right), \\
 Y_{ij}|\lambda_{ij}, \pi_{ij}, \phi &\sim ZIGP(\lambda_{ij}, \phi, \pi_{ij}).
 \end{aligned}$$

Note that $Y \sim NB(\phi, \kappa)$ if the probability mass function is given by $f_{NB}(y; \phi, \kappa) = \frac{\Gamma(y+\phi)}{\Gamma(\phi)y!} \kappa^\phi (1-\kappa)^y$, $y = 0, 1, \dots, r$ and $r > 0$. Also, $Y \sim ZINB(\phi, \kappa, \pi)$ is a zero-inflated negative binomial distribution which can be obtained by (2.3) by replacing $f_{GP}(\cdot; \lambda, \phi)$ by $f_{NB}(\cdot; \phi, \kappa)$. In order to compare the results, the mean of the estimated values, the standard errors, relative biases and mean of square errors (MSEs) are used. The latter two criteria are defined as follows:

$$\begin{aligned}
 Bias(\theta) &= \frac{1}{M} \sum_{k=1}^M \left(\frac{\hat{\theta}_k}{\theta} - 1 \right), \\
 MSE(\theta) &= \frac{1}{M} \sum_{k=1}^M \left(\hat{\theta}_k - \theta \right)^2,
 \end{aligned}$$

where $\hat{\theta}_k$ is the estimate of θ for the k^{th} sample, $k = 1, 2, \dots, M$.

The results of Tables 1 and 2 show that the performance of the ZIGP in parameter estimation is better than those of the other models. The performance of ZINB in estimating parameters of the logistic model is not well while in estimating the other parameters is almost good. The GP and NB models do not have good performances in this simulation study and the ZIP has a good performance in estimating some parameters. The results of the simulation study for ZIGP show that the increase in the sample size is an effective way of decreasing biases and standard deviations of parameters estimates. As shown in these tables, relative biases and MSEs are reduced by increasing the sample size. This suggests that the method in finding estimates is consistent.

4.2. Zero-inflated Poisson model

In this simulation study, we simulate data from the following model:

$$(4.4) \quad Y_{ij}|\lambda_{ij}, \pi_{ij} \sim ZIP(\lambda_{ij}, \pi_{ij}),$$

such that the parameterizations and real values of parameters in λ_{ij} and π_{ij} are the same as the first set of real values and those described in equation (4.2). The results of this simulation study are summarized in Table 3. The results show the well performance of the ZIP model. Also, the results show that the performance of ZIGP is as good as ZIP model. The ZINB model dose not have a good performance when the sample size is 100 while for N=500 has a performance which is as good as the other two models. The overdispersion parameter ϕ

is estimated zero in ZIGP model but it has a large value in ZINB (note that in negative binomial distribution the dispersion index is proportion to ϕ^{-1} and overdispersion presents in the data when the value of ϕ is very large). As a conclusion, this simulation study shows that the use of ZIGP model is preferred to the use of ZINB model. The ZIGP has a similar performance to ZIP and, for moderate sample size, a much better performance than ZINB model.

Table 3: Results of simulation study for generated data under ZIP model, estimate (Est.), standard error (S.E.), relative bias (Bias) and mean square error (MSE) for $M=1000$ simulated data with sample sizes 100 and 500.

N	Para.	Real	ZIP			ZIGP			ZINB		
			Est. (S.E.)	Bias	MSE	Est. (S.E.)	Bias	MSE	Est. (S.E.)	Bias	MSE
100	α_0	-1.00	-1.06 (0.86)	0.06	0.74	-1.20 (0.98)	0.20	0.95	$<-10^3$ ($>10^3$)	$>10^3$	$>10^3$
	α_1	1.00	1.03 (0.50)	0.03	0.24	1.23 (0.57)	0.23	0.55	$>10^3$ ($>10^3$)	$>10^3$	$>10^3$
	α_2	0.00	0.02 (0.20)	*	0.04	-0.00 (0.19)	*	0.00	0.02 (0.20)	*	0.04
	τ_1	0.00	-0.07 (0.56)	*	0.32	-0.04 (0.61)	*	0.03	-0.09 (0.64)	*	0.42
	τ_2	1.00	-1.26 (0.91)	0.26	0.96	-1.05 (0.36)	0.05	0.01	-1.07 (0.37)	0.07	0.14
	β_0	-3.00	-3.00 (0.22)	0.00	0.05	-2.98 (0.33)	-0.00	0.00	-3.04 (0.25)	0.01	0.06
	β_1	1.00	1.00 (0.08)	0.00	0.00	0.99 (0.14)	-0.01	0.00	1.01 (0.08)	0.01	0.01
	β_2	-1.00	0.99 (0.05)	-0.00	0.00	0.99 (0.06)	-0.00	0.00	1.00 (0.05)	0.00	0.00
	γ_1	-1.00	-0.99 (0.08)	-0.00	0.00	-0.99 (0.09)	-0.00	0.00	-1.00 (0.09)	0.00	0.00
	γ_2	0.00	-0.00 (0.01)	*	0.00	0.00 (0.01)	*	0.00	-0.00 (0.01)	*	0.00
	ϕ	0.00	—	—	—	-0.00 (0.00)	*	0.00	$>10^3$ ($>10^3$)	$>10^3$	$>10^3$
500	α_0	-1.00	-1.02 (0.36)	0.02	0.13	-1.02 (0.36)	0.02	0.13	-1.03 (0.36)	0.03	0.13
	α_1	1.00	1.02 (0.21)	0.02	0.04	1.02 (0.21)	0.02	0.04	1.02 (0.21)	0.02	0.04
	α_2	0.00	0.00 (0.07)	*	0.00	0.00 (0.07)	*	0.00	0.00 (0.07)	*	0.00
	τ_1	0.00	-0.02 (0.24)	*	0.06	-0.02 (0.24)	*	0.06	-0.02 (0.24)	*	0.06
	τ_2	1.00	-1.01 (0.17)	0.01	0.02	-1.01 (0.17)	0.01	0.02	-1.02 (0.17)	0.02	0.06
	β_0	-3.00	-3.00 (0.11)	0.00	0.01	-3.00 (0.11)	0.00	0.01	-3.00 (0.11)	0.00	0.01
	β_1	1.00	1.00 (0.04)	0.00	0.00	1.00 (0.04)	0.00	0.00	1.00 (0.04)	0.00	0.00
	β_2	-1.00	0.99 (0.02)	-0.00	0.00	0.99 (0.02)	0.00	0.00	1.00 (0.02)	0.00	0.00
	γ_1	-1.00	-1.00 (0.03)	0.00	0.00	-0.99 (0.03)	0.00	0.00	-1.00 (0.03)	0.00	0.00
	γ_2	0.00	-0.00 (0.00)	*	0.00	-0.00 (0.00)	*	0.00	-0.00 (0.00)	*	0.00
	ϕ	0.00	—	—	—	-0.00 (0.00)	*	0.00	$>10^3$ ($>10^3$)	$>10^3$	$>10^3$

4.3. Zero-inflated negative binomial model

In this simulation study, we simulate data from the following model:

$$(4.5) \quad Y_{ij} | \lambda_{ij}, \pi_{ij} \sim ZINB\left(\phi, \frac{\phi}{\phi + \lambda_{ij}}, \pi_{ij}\right),$$

such that the parameterizations and real values of parameters in λ_{ij} and π_{ij} are the same as the first set of real values and those described in equation (4.2), also, we consider $\phi = 1$. The results of this simulation study are summarized in Table 4. The results show the well performance of the ZINB model in large sample size. Also, the results show that the performance of ZIGP is as good as ZINB model expect for estimating intercept and the overdispersion parameters. Also, in moderate sample size the performance of ZIGP model is better

than those in ZINB model. The results show that the ZIP model dose not have a good performance.

Table 4: Results of simulation study for generated data under ZINB model, estimate (Est.), standard error (S.E.), relative bias (Bias) and mean square error (MSE) for $M = 1000$ simulated data with sample sizes 100 and 500.

N	Para.	Real	ZINB			ZIGP			ZIP		
			Est. (S.E.)	Bias	MSE	Est. (S.E.)	Bias	MSE	Est. (S.E.)	Bias	MSE
100	α_0	-1.00	-1.47 (1.29)	0.47	1.87	-0.20 (0.89)	-0.79	1.41	0.69 (0.77)	-1.69	3.47
	α_1	1.00	1.24 (0.91)	0.24	0.88	0.67 (0.51)	-0.32	0.36	0.42 (0.37)	-0.57	0.47
	α_2	0.00	0.06 (0.27)	*	0.08	-0.09 (0.17)	*	0.04	-0.32 (0.19)	*	0.14
	τ_1	0.00	-0.03 (1.10)	*	1.20	0.10 (0.64)	*	0.42	0.59 (0.51)	*	0.61
	τ_2	-1.00	-1.61 (3.28)	0.61	11.05	-0.82 (0.48)	-0.17	0.26	-0.37 (0.26)	-0.62	0.46
	β_0	-3.00	-3.13 (0.45)	0.04	0.21	-2.73 (0.50)	-0.08	0.32	-2.13 (0.63)	-0.28	1.14
	β_1	1.00	1.03 (0.26)	0.03	0.06	0.98 (0.22)	-0.01	0.05	0.87 (0.28)	-0.12	0.09
	β_2	1.00	1.02 (0.09)	0.02	0.00	0.95 (0.10)	-0.04	0.01	0.84 (0.14)	-0.15	0.04
	γ_1	-1.00	-0.97 (0.21)	-0.02	0.04	-0.95 (0.20)	-0.04	0.04	-0.85 (0.28)	-0.14	0.09
	γ_2	0.00	-0.00 (0.03)	*	0.00	-0.00 (0.03)	*	0.00	0.00 (0.04)	*	0.00
	ϕ	1.00	1.11 (0.25)	0.11	0.07	0.22 (0.03)	-0.77	0.60	—	—	—
500	α_0	-1.00	-1.00 (0.50)	0.00	0.24	-0.07 (0.37)	-0.92	1.00	0.72 (0.34)	-1.72	3.10
	α_1	1.00	1.03 (0.33)	0.00	0.11	0.66 (0.22)	-0.33	0.16	0.34 (0.18)	-0.65	0.46
	α_2	0.00	-0.01 (0.10)	*	0.01	-0.15 (0.08)	*	0.03	-0.33 (0.07)	*	0.11
	τ_1	0.00	0.02 (0.34)	*	0.11	-0.25 (0.24)	*	0.12	0.69 (0.20)	*	0.52
	τ_2	-1.00	-0.99 (0.29)	0.00	0.08	-0.73 (0.20)	-0.26	0.11	-0.29 (0.10)	-0.70	0.50
	β_0	-3.00	-3.00 (0.20)	0.00	0.04	-2.65 (0.20)	-0.11	0.16	-2.20 (0.33)	-0.26	0.75
	β_1	1.00	1.00 (0.07)	0.00	0.00	0.95 (0.08)	-0.04	0.00	0.88 (0.13)	-0.11	0.03
	β_2	1.00	0.99 (0.04)	0.00	0.00	0.93 (0.04)	-0.06	0.00	0.86 (0.08)	-0.13	0.02
	γ_1	-1.00	-0.99 (0.09)	-0.00	0.00	-0.91 (0.09)	-0.08	0.01	-0.86 (0.13)	-0.13	0.03
	γ_2	0.00	-0.00 (0.01)	*	0.00	-0.00 (0.01)	*	0.00	0.00 (0.03)	*	0.00
	ϕ	1.00	1.01 (0.10)	0.00	0.01	0.23 (0.01)	-0.76	0.57	—	—	—

4.4. Zero-inflated underdispersion generalized Poisson model

For investigating the performance of the proposed transition model, the data set of this subsection are generated from a zero-inflated underdispersed generalized Poisson model and the performance of the ZIGP, ZINB and ZIP models are compared. The data set are generated from a $ZIGP(\lambda_{ij}, \pi_{ij}, \phi)$ such that $\log(\lambda_{i1}) = \beta_0$, $\text{logit}(p_{i1}) = \alpha_0$, $\log(\lambda_{ij}) = \beta_0 + \gamma_1 I_{\{0\}}(Y_{i,j-1}) + \gamma_2 y_{i,j-1}(1 - I_{\{0\}}(Y_{i,j-1}))$, $j = 2, 3, 4$, $\text{logit}(p_{ij}) = \alpha_0 + \tau_1 I_{\{0\}}(Y_{i,j-1}) + \tau_2 y_{i,j-1}(1 - I_{\{0\}}(Y_{i,j-1}))$, $j = 2, 3, 4$, where $\alpha_0 = -1$, $\tau_1 = -1$, $\tau_2 = 1$, $\beta_0 = 1$, $\gamma_1 = 0$, $\gamma_2 = -1$ and $\phi = -0.3$. Also, two sample sizes $N=500$ and 1000 are selected where $M = 1000$ iterations are performed. The results of this simulation study are summarized in Table 5. These results show the well performance of the ZIGP model as the best fitting model while the performance of ZINB model is poor. Also, the results show that the performance of ZIP is better than those of ZINB model. Note that the underdispersion rarely occur in practice. The well performance of ZIGP model are only satisfied in large sample size as described in this simulation study.

Table 5: Results of simulation study for generated data under ZIGP model in the presence of underdispersion, estimate (Est.), standard error (S.E.), relative bias (Bias) and mean square error (MSE) for $M=1000$ simulated data with sample sizes 500 and 1000.

N	Para.	Real	ZIGP			ZINB			ZIP		
			Est. (S.E.)	Bias	MSE	Est. (S.E.)	Bias	MSE	Est. (S.E.)	Bias	MSE
500	α_0	-1.00	-0.97 (0.11)	-0.03	0.01	$< -10^3$ ($> 10^3$)	3618.73	$> 10^3$	-4.47 (4.31)	3.47	29.19
	τ_1	-1.00	-1.03 (0.15)	0.03	0.02	$< -10^3$ ($> 10^3$)	12854.74	$> 10^3$	-18.20 (9.03)	17.20	371.03
	τ_2	1.00	1.15 (0.17)	0.15	0.15	$< -10^3$ ($> 10^3$)	-63370.51	$> 10^3$	-11.15 (6.46)	-12.15	186.23
	β_0	1.00	0.93 (0.01)	-0.07	0.14	0.84 (0.05)	-0.18	0.04	0.84 (0.10)	-0.18	0.15
	γ_1	0.00	0.00 (0.01)	*	0.00	0.16 (0.04)	*	0.03	0.13 (0.05)	*	0.02
	γ_2	-1.00	-1.20 (0.75)	0.20	0.65	-1.55 (0.33)	0.55	0.50	-1.39 (0.45)	0.39	0.46
	ϕ	-0.30	-0.37 (0.00)	0.24	0.03	$> 10^3$ ($> 10^3$)	$< -10^3$	$> 10^3$	—	—	—
1000	α_0	-1.00	-0.97 (0.05)	-0.03	0.00	$< -10^3$ ($> 10^3$)	3537.27	$> 10^3$	-3.78 (3.78)	2.78	21.80
	τ_1	-1.00	-1.03 (0.10)	0.03	0.01	$< -10^3$ ($> 10^3$)	12912.70	$> 10^3$	-18.35 (4.68)	17.35	322.60
	τ_2	1.00	1.06 (0.09)	0.06	0.08	$< -10^3$ ($> 10^3$)	-83788.10	$> 10^3$	-13.00 (6.35)	-14.00	235.61
	β_0	1.00	0.98 (0.01)	-0.02	0.04	0.92 (0.06)	-0.08	0.06	0.95 (0.05)	-0.05	0.06
	γ_1	0.00	0.00 (0.01)	*	0.00	0.17 (0.03)	*	0.03	0.14 (0.02)	*	0.02
	γ_2	-1.00	-1.12 (0.32)	0.12	0.29	-1.47 (0.40)	0.47	0.31	-1.52 (0.35)	0.52	0.44
	ϕ	-0.30	-0.32 (0.00)	0.07	0.03	$> 10^3$ ($> 10^3$)	$< -10^3$	$> 10^3$	—	—	—

5. APPLICATION

The data set of this paper is extracted from a longitudinal study on kidney transplant patients in Imam Khomeini hospital of Urmia in Iran. The data set contains some information about $N = 129$ patients who have kidney transplant in this hospital. The response variable in this study is the number of acute rejections which is count response with extra zeros. The data are recorded in one year period which contain the number of acute rejection each four months. The barchart of the response variable for each time point (month 4, 8 and 12) is showed in Figure 1. In this figure, Y_k , $k = 1, 2, 3$, is used for indicating the response variable at the k^{th} time point. The number of extra zeros is clear in these charts.

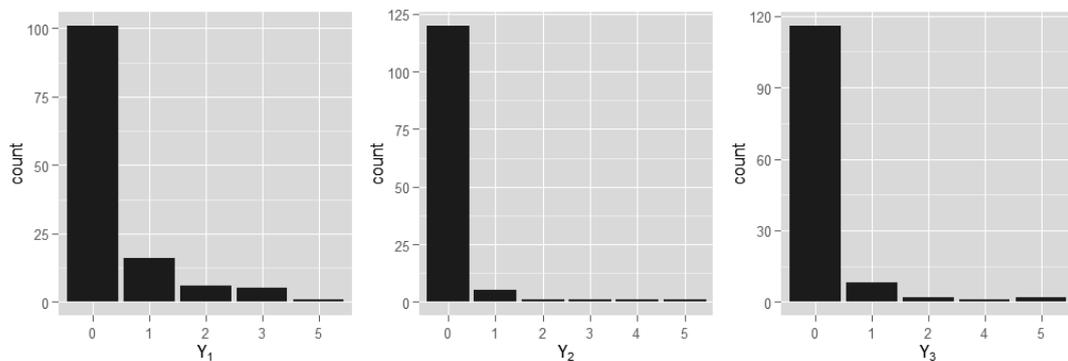


Figure 1: Barcharts of the number of acute rejections for time point at month 4 (first panel), month 8 (middle panel) and month 12 (third panel).

The collected explanatory variables which are considered in our analysis are creatinine index as a continuous covariate and having hyperacute rejection of kidney (rejection in the first 24 hours after surgery) as a categorical covariate. Figure 2 presents the boxplots of the creatinine index versus the number of acute rejections for each time. Also, Table 6 summarizes frequency of the number of acute rejections for each category of this variable for each time point.

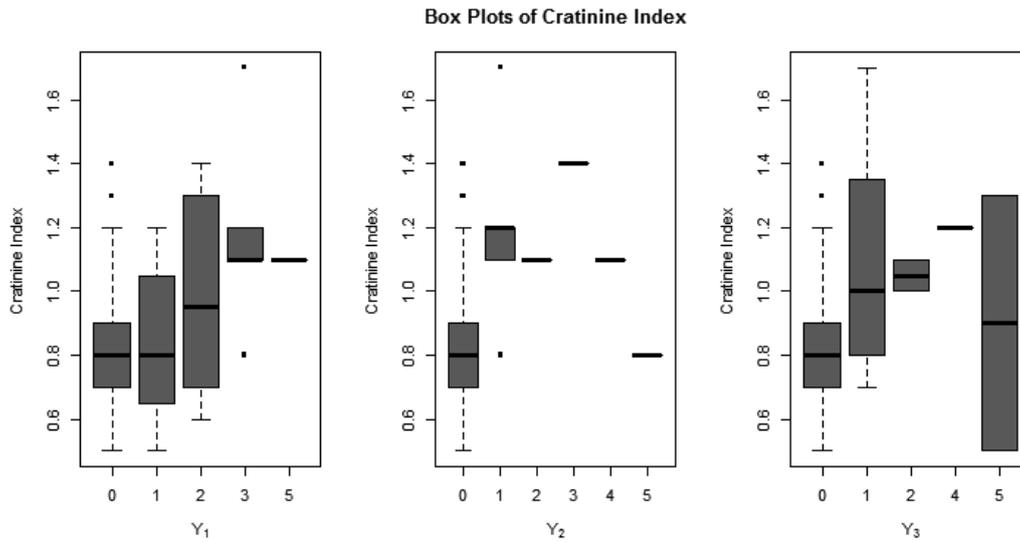


Figure 2: Boxplots of creatinine index versus the number of acute rejections for all time points.

For analyzing this data set, we use the proposed zero-inflated generalized Poisson transition model, also, Poisson (PM), negative binomial (NBM), generalized Poisson (GPM), zero-inflated Poisson (ZIPM), zero-inflated generalized Poisson (ZIGPM) and zero-inflated negative binomial (ZINBM) models under the transition structure are used for analyzing the data set. The explanatory variables which are considered for analysing the data are creatinine index (CRAT), having early acute rejection (EAR) and time ($t = 4, 8, 12$).

Table 6: Frequency of early acute rejection of kidney on the total number of acute rejection at each time point. “Yes” is used for having early acute rejection and “No” is used for not having early acute rejection.

Number	Early acute rejection					
	1st time point		2nd time point		3rd time point	
	Yes	No	Yes	No	Yes	No
0	19	82	31	89	28	88
1	7	9	1	4	5	3
2	5	1	1	0	1	1
3	3	2	1	0	0	0
4	0	0	1	0	0	1
5	1	0	0	1	1	1

We consider models (4.3) for analysing this data, where

$$(5.1) \quad \begin{aligned} \log(\lambda_{ij}) = & \beta_0 + \beta_1 CRAT_i + \beta_2 Time_j + \beta_3 EAR_i \\ & + \gamma_1 I_{\{0\}}(Y_{i,j-1}) + \gamma_2 (1 - I_{\{0\}}(Y_{i,j-1})) y_{i,j-1}, \end{aligned}$$

and

$$(5.2) \quad \begin{aligned} \text{logit}(\pi_{ij}) = & \alpha_0 + \alpha_1 CRAT_i + \alpha_2 Time_j + \alpha_3 EAR_i \\ & + \tau_1 I_{\{0\}}(Y_{i,j-1}) + \tau_2 (1 - I_{\{0\}}(Y_{i,j-1})) y_{i,j-1}. \end{aligned}$$

For model comparison, we evaluate different model fits by considering some information criteria. These criteria are AIC, BIC and HQC, which are defined as follows:

Let $\boldsymbol{\theta}$ be the vector of unknown parameters, then

$$\begin{aligned} AIC &= -2\ell(\hat{\boldsymbol{\theta}}|\mathbf{Y}) + 2|\boldsymbol{\theta}|, \\ BIC &= -2\ell(\hat{\boldsymbol{\theta}}|\mathbf{Y}) + |\boldsymbol{\theta}| \ln(N), \\ HQC &= -2\ell(\hat{\boldsymbol{\theta}}|\mathbf{Y}) + 2 \ln(\ln(N)), \end{aligned}$$

where $|\boldsymbol{\theta}|$ is the number of unknown parameters in vector $\boldsymbol{\theta}$, N is the number of subjects and $\hat{\boldsymbol{\theta}}$ is the vector of parameters estimates. The smaller values of AIC, BIC and HQC indicate a better fitting model.

We use the EM algorithm, as described in Section 3, for parameters estimation of zero-inflated models, also, the usual maximum likelihood approach is used for parameter estimation of other models. In the EM algorithm, the initial values for unknown parameters were set equal to the estimates obtained by analysing separate models. The results of the above mentioned models are summarized in Table 7. This table contains parameter estimates and their standard errors for the first order transition model where standard deviations for zero-inflated models are estimated using a Bootstrap approach with 10000 iterations and for the others we use inverse of the Hessian matrix. The results show, based on the values of different criteria, that for this data set, the performance of ZIGP and ZINB models are similar and the difference between them is negligible. After them ZIP has the best fitting model and the worst fitting model based on these criteria is the PM. The results show some evidence for existence of mild overdispersion.

The results show that for zero-inflated models creatinine index (CRAT), having early acute rejection (EAR) and time are significant variables such that the more the creatinine index is, the larger is the estimated probability of nonzeros. Also, two covariates time and early acute rejection are positively significant, i. e. by increasing them the probability of zero increases. The results of zero-inflated models also show that only transition parameter τ_1 is significant. The results show that significant covariates in non-inflated models are similar to those in modeling zero probability in zero-inflated models, that is, the significant parameters in modeling zero probability of zero-inflated models have similar interpretation to those in modeling the rate of distributions in non-inflated models. Also, ϕ and τ_1 are the other significant parameters in these models.

Note that in a first order transition model the first response of each individual should be modeled given its previous response which is not recorded. How to face this issue, called

the initial condition problem [11, 10]. This problem does not exist in this study, because the number of acute rejections before the time of study is zero. In other words, the patients have been entered in the study from the time of kidney transplant and they have been followed for one year. Also, in this paper, we consider the first order transition model for modeling the data set, because the number of replications in our real data is three and a first order transition model for considering between-group dependence in data is adequate.

Table 7: Results of fitting (parameter estimations and standard errors in parenthesis) the Poisson model (PM), negative binomial model (NBM), generalized Poisson model (GPM), zero-inflated Poisson model (ZIPM), zero-inflated generalized Poisson model (ZIGPM) and zero-inflated negative binomial model (ZINBM) to kidney transplant study (significant parameters are highlighted in bold).

Parameter	ZIGPM	ZINBM	ZIPM	GPM	NBM	PM
	Est. (S.E.)					
α_0	2.66 (1.13)	2.62 (1.16)	3.10 (0.91)	—	—	—
α_1 (CRAT)	-3.43 (1.09)	-3.42 (1.09)	-3.35 (0.92)	—	—	—
α_2 (Time)	0.11 (0.04)	0.12 (0.05)	0.11 (0.05)	—	—	—
α_3 (EAR)	1.10 (0.58)	1.12 (0.59)	1.01 (0.44)	—	—	—
β_0	-0.71 (0.82)	-0.72 (0.83)	-0.18 (0.54)	-3.57 (0.94)	-3.19 (0.77)	-2.49 (0.46)
β_1 (CRAT)	0.55 (0.67)	0.58 (0.69)	0.32 (0.52)	2.19 (0.73)	2.08 (0.66)	1.99 (0.37)
β_2 (Time)	-0.25 (0.40)	-0.24 (0.39)	-0.21 (0.29)	-0.93 (0.34)	-0.87 (0.32)	-0.79 (0.23)
β_3 (EAR)	0.79 (1.35)	0.70 (1.28)	0.60 (0.98)	3.03 (1.93)	2.10 (1.27)	0.41 (0.65)
τ_1	1.83 (0.62)	1.85 (0.63)	1.69 (0.49)	—	—	—
τ_2	0.04 (0.31)	0.04 (0.32)	0.02 (0.25)	—	—	—
γ_1	-0.35 (1.01)	-0.29 (0.96)	-0.24 (0.72)	-2.82 (1.18)	-2.27 (0.78)	-1.34 (0.41)
γ_2	0.00 (0.19)	0.00 (0.19)	-0.01 (0.14)	-0.17 (0.28)	-0.13 (0.23)	-0.02 (0.14)
ϕ	0.21 (0.07)	2.11 (0.99)	—	1.22 (0.36)	0.34 (0.10)	—
AIC	384.04	384.98	386.19	392.96	392.51	441.87
BIC	419.16	419.20	417.64	412.98	412.53	459.03
HQC	364.00	364.05	367.35	382.13	381.67	433.03

6. CONCLUSION AND DISCUSSION

In this paper, we have discussed a new transition model for analysing longitudinal outcomes with extra zeros. We compare the performance of different distributional assumptions: zero-inflated generalized Poisson, zero-inflated negative binomial and zero-inflated Poisson and we conclude that zero-inflated generalized Poisson is a flexible distributional assumption.

We have used the EM algorithm for parameter estimation. For illustration of the proposed models some simulation studies have been conducted. Also, a real data set of a kidney allograft rejection study has been analyzed as an illustrative example. Based on the results the creatinine index, having early acute rejection and time are significant covariates such that the more the creatinine index is, the larger is the estimated probability of nonzeros acute rejection. Also, two covariates time and early acute rejection are positively significant, i. e. by increasing them the probability of zero acute rejection increases. The results show

that the significant parameters in modeling zero probability of zero-inflated models have similar effect to parameters in the modeling rate of distributions in non-inflated models. We have considered a first order transition model for considering within-group dependence in longitudinal measurements, because the number of repeated longitudinal measurements has been small in our real data set. As a future work, illustration of the proposed approach for higher order of transition model for analyzing data set with larger number of repeated measures and comparison of the performance of it with that of the first order transition model may be performed. For this purpose (3.2) and (3.4) can be improve to be $\log(\lambda_{i1}) = \mathbf{x}'_{i1}\boldsymbol{\beta}$, $\text{logit}(\pi_{i1}) = \mathbf{z}'_{i1}\boldsymbol{\alpha}$, $\log(\lambda_{ij}) = \mathbf{x}'_{ij}\boldsymbol{\beta} + \boldsymbol{\gamma}'\mathbf{h}_{i,j}$ and $\text{logit}(\pi_{ij}) = \mathbf{z}'_{ij}\boldsymbol{\alpha} + \boldsymbol{\tau}'\mathbf{h}_{i,j}$, $j = 2, \dots, n_i$. Another parameterizations for λ_{ij} and π_{ij} of (3.2) and (3.4) may be the use of the first order transition model along with some random effects, that is, $\log(\lambda_{i1}) = \mathbf{x}'_{i1}\boldsymbol{\beta} + b_{i1}$, $\text{logit}(\pi_{i1}) = \mathbf{z}'_{i1}\boldsymbol{\alpha} + b_{i2}$, $\log(\lambda_{ij}) = \mathbf{x}'_{ij}\boldsymbol{\beta} + \gamma_1 I_{\{0\}}(Y_{i,j-1}) + \gamma_2 y_{i,j-1}(1 - I_{\{0\}}(Y_{i,j-1})) + b_{i1}$ and $\text{logit}(\pi_{ij}) = \mathbf{z}'_{ij}\boldsymbol{\alpha} + \tau_1 I_{\{0\}}(Y_{i,j-1}) + \tau_2 y_{i,j-1}(1 - I_{\{0\}}(Y_{i,j-1})) + b_{i2}$, $j = 2, \dots, n_i$. where $\mathbf{b}_i = (b_{i1}, b_{i2})'$ is a bivariate random effects. As a parameterization for the random effects, one can write $b_{1i} \sim N(0, \sigma_1^2)$, $b_{2i}|b_{1i} \sim N(\psi b_{1i}, \sigma_2^2)$. We have used the EM algorithm for parameter estimation, one can use a Bayesian paradigm using MCMC for parameter estimation [29]. The priors elicitation are an important issue for performing this paradigm. The data set which analyzed in this paper has not had any missing values. The proposed method can be extended for modeling data sets in the presence of missing values as a future work. For this purpose, an ignorable or non-ignorable missing mechanism should be selected. The modeling of missing data mechanism for modeling non-ignorable missing data mechanism is necessary and these a sensitivity analysis is commonly suggested.

ACKNOWLEDGMENTS

This work has been supported by the grant number 96000139 from Iranian National Science Foundation (INSF). The authors would like to thank the INSF. The authors also acknowledge the valuable suggestions from the referees.

REFERENCES

- [1] AGRESTI, A. (1999). Modeling ordered categorical data: recent advances and future challenges, *Statistics in medicine*, **18**, 2191–2207.
- [2] ALFO, M. and MARUOTTI, A. (2014). Two-part regression models for longitudinal zero-inflated count data, *The Canadian Journal of Statistics*, **38**, 197–216.
- [3] BUU, A.; LI, R.; TAN, X. and ZUCKER, R.A. (2012). Statistical models for longitudinal zero-inflated count data with applications to the substance abuse field, *Statistics in medicine*, **31**(29), 4074–4086.
- [4] CHEUNG, Y.B. (2002). Zero-inflated models for regression analysis of count data: a study of growth and development, *Statistics in Medicine*, **21**, 1461–1469.
- [5] CONSUL, P.C. (1989). *Generalized Poisson distribution: Properties and Applications*, Marcel Dekker, New York.

- [6] CONSUL, P.C. and FAMOYE, F. (1992). Generalized Poisson regression model, *Communications in Statistics – Theory and Methods*, **21**, 81–109.
- [7] CZADO, C.; ERHARDT, V.; MIN, A. and WAGNER, S. (2007). Zero-inflated generalized Poisson models with regression effects on the mean dispersion and zero-inflation level applied to patent outsourcing rates, *Statistical Modelling*, **7**(2), 125–153.
- [8] DEMPSTER, A.P.; LAIRD, N.M. and RUBIN, D.B. (1977). Maximum likelihood for incomplete data via the EM algorithm, *Journal of the Royal Statistical Society B*, **39**, 1–38.
- [9] DIGGLE, P.J.; HEAGERTY, P.; LIANG, K.Y. and ZEGER, S.L. (2002). *Analysis of Longitudinal Data*, Oxford: Oxford University Press.
- [10] GANJALI, M.; BAGHFALAKI, T. and GHAHRODI, Z.R. (2017). Transitional Ordinal Modeling, *Wiley StatsRef: Statistics Reference Online*, Free Trial, 1–13.
- [11] GANJALI, M. and REZAEI, Z. (2007). A Transition Model for Analysis of Repeated Measure Ordinal Response Data to Identify the Effects of Different Treatments, *Drug Information Journal*, **41**, 527–534.
- [12] HALL, D.B. (2000). Zero-Inflated Poisson and Binomial Regression with Random Effects: A Case Study, *Biometrics*, **56**, 1030–1039.
- [13] HARRIS, T.; YANG, Z. and HARDIN, J.W. (2012). Modeling underdispersed count data with generalized Poisson regression, *Stata Journal*, **12**(1), 736–747.
- [14] HASAN, M.T. and SNEDDON, G. (2009). Zero-Inflated Poisson Regression for Longitudinal Data, *Communications in Statistics – Simulation and Computation*, **38**(3), 638–653.
- [15] HEILBRON, D.C. (1994). Zero-altered and other regression models for count data with added zeros, *Biometrical Journal*, **36**, 531–547.
- [16] HU, M.C.; PAVLICOVA, M. and NUNES, E.V. (2011). Zero-inflated and hurdle models of count data with extra zeros: examples from an HIV-risk reduction intervention trial, *The American Journal of Drug and Alcohol Abuse*, **37**, 367–375.
- [17] JOE, H. and ZHU, R. (2005). Generalized Poisson distribution: the property of mixture of Poisson and comparison with negative binomial distribution, *Biometrical Journal*, **47**(2), 219–229.
- [18] LALL, R.; CAMPBELL, M.J.; WALTERS, S.J. and MORGAN, K. (2002). A review of ordinal regression models applied on health-related quality of life assessments, *Statistical Methods in Medical Research*, **11**, 49–67.
- [19] LAMBERT, D. (1992). Zero-inflated Poisson regression, with an application to defects in manufacturing, *Technometrics*, **34**, 1–14.
- [20] LANGE, K. (2004). *Optimization*, Springer-Verlag, New York.
- [21] LEWSEY, J.D. and THOMSON, W.M. (2004). The utility of the zero-inflated Poisson and zero-inflated negative binomial models: a case study of cross-sectional and longitudinal DMF data examining the effect of socio-economic status, *Community of Dentistry and Oral Epidemiology*, **32**, 183–189.
- [22] MARUOTTIAB, A. and RAPONIC, V. (2014). On Baseline Conditions for Zero-Inflated Longitudinal Count Data, *Communications in Statistics – Simulation and Computation*, **43**, 743–760.
- [23] MIN, Y. and AGRESTI, A. (2005). Random effect models for repeated measures of zero-inflated count data, *Statistical Modeling*, **5**, 1–19.
- [24] MOLENBERGHS, G. and VERBEKE, G. (2005). *Models for Discrete Longitudinal Data*, Springer-Verlag.
- [25] MULLAHDY, J. (1986). Specification and testing of some modified count data models, *Journal of Econometrics*, **33**, 341–365.
- [26] NEELON, B.H.; OMALLEY, A.J. and NORMAND, S.L. (2010). A Bayesian model for repeated measures zero-inflated count data with application to outpatient psychiatric service use, *Statistical Modelling*, **10**, 421–439.

- [27] REINECKE, J. and SEDDIG, D. (2011). Growth mixture models in longitudinal research, *AStA Advances in Statistical Analysis*, **95**(4), 415–434.
- [28] ROSE, C.E.; MARTIN, S.W.; WANNEMUEHLER, K.A. and PLIKAYTIS, B.D. (2006). On the use of zero-inflated and hurdle models for modeling vaccine adverse event count data, *Journal of Biopharmaceutical Statistics*, **16**, 463–481.
- [29] SILVA, G.L.; JUAREZ-COLUNGA, E. and DEAN, C. (2015). A joint analysis of counts and severity with zero-inflated longitudinal data, *CEB-EIB 2015*, Bilbao, 23–25 September.
- [30] SONG, P.X.K. (2007). *Correlated Data Analysis*, Springer-Verlag, New York.

COMPOUND POWER SERIES DISTRIBUTION WITH NEGATIVE MULTINOMIAL SUMMANDS: CHARACTERISATION AND RISK PROCESS

Authors: PAVLINA JORDANOVA
– Faculty of Mathematics and Informatics, Shumen University, Bulgaria
pavlina_kj@abv.bg

MONIKA PETKOVA
– Faculty of Mathematics and Informatics, Sofia University, Bulgaria
monikapetevapetkova@abv.bg

MILAN STEHLÍK
– Department of Applied Statistics and Linz Institute of Technology,
Johannes Kepler University in Linz, Austria
and
Institute of Statistics, University of Valparaíso, Chile
Milan.Stehlik@jku.at

Received: March 2017

Revised: October 2017

Accepted: October 2017

Abstract:

- The paper considers a multivariate distribution whose coordinates are compounds. The number of the summands is itself also a multivariate compound with one and the same univariate Power series distributed number of summands and negative multinomially distributed summands. In the total claims amount process the summands are independent identically distributed random vectors. We provide the first full characterization of this distribution. We show that considered as a mixture this distribution would be Mixed Negative multinomial distribution having the possibly scale changed power series distributed THE first parameter. We provide an interesting application to risk theory.

Key-Words:

- *power series distributions; Multivariate Compound distribution; Negative multinomial distribution; risk process.*

AMS Subject Classification:

- 60E05, 62P05.

1. INTRODUCTION AND PRELIMINARY RESULTS

It seems that Bates and Neyman [3] were first to introduce Negative multinomial (NMn) distribution in 1952. They obtained it by considering a mixture of independent Poisson distributed random variables (r.v.s) with one and the same Gamma distributed mixing variable. Their first parameter could be a real number. Wishart [25] considers the case when the first parameter could be only integer. He calls this distribution Pascal multinomial distribution. At the same time Tweedie [24] obtained estimators of the parameters. Sibuya *et al.* [18] make a systematic investigation of this distribution and note that the relation between Binomial distribution and Negative binomial (NBi) distribution is quite similar to that between the Multinomial distribution and NMn distribution. The latter clarifies the probability structure of the individual distributions. The bivariate case of the compound power series distribution with geometric summands (i.e. $n = 1$ and $k = 2$) is partially investigated in [12]. Another related work is [10].

A version of k -variate negative binomial distribution with respect to risk theory is considered in [2, 26]. The authors show that it can be obtained by mixing of iid Poisson random variables with a multivariate finite mixture of Erlang distributions with one and the same second parameter. Further on they interpret it as the loss frequencies and obtain the main characteristics. Due to covariance invariance property, the corresponding counting processes can be useful to model a wide range of dependence structures. See [2, 26] for examples. Using probability generating functions, the authors present a general result on calculating the corresponding compound, when the loss severities follow a general discrete distribution. The similarity of our paper and papers [2, 26] is that both consider the aggregate losses of an insurer that runs through several correlated lines of business. In (2.1) and (2.2) [2] consider Mixed k -variate Poisson distribution (with independent coordinates, given the mixing variable) and the mixing variable is Mixed Erlang distributed. More precisely the first parameter in the Erlang distribution is replaced with a random variable. The mixing variable is multivariate and the coordinates of the compounding vector are independent. In our case the mixing variable is one and the same and the coordinates of the counting vector are dependent.

Usually Negative Multinomial (NMn) distribution is interpreted as the one of the numbers of outcomes A_i , $i = 1, 2, \dots, k$ before the n -th B , in series of independent repetitions, where A_i , $i = 1, 2, \dots, k$ and B form a partition of the sample space. See e.g. Johnson *et al.* [7]. Let us recall the definition.

Definition 1.1. Let $n \in \mathbb{N}$, $0 < p_i$, $i = 1, 2, \dots, k$ and $p_1 + p_2 + \dots + p_k < 1$. A vector $(\xi_1, \xi_2, \dots, \xi_k)$ is called Negative multinomially distributed with parameters n, p_1, p_2, \dots, p_k , if its probability mass function (p.m.f.) is

$$\begin{aligned} P(\xi_1 = i_1, \xi_2 = i_2, \dots, \xi_k = i_k) &= \\ &= \binom{n + i_1 + i_2 + \dots + i_k - 1}{i_1, i_2, \dots, i_k, n - 1} p_1^{i_1} p_2^{i_2} \dots p_k^{i_k} (1 - p_1 - p_2 - \dots - p_k)^n, \\ & \qquad \qquad \qquad i_s = 0, 1, \dots, \quad s = 1, 2, \dots, k. \end{aligned}$$

Briefly $(\xi_1, \xi_2, \dots, \xi_k) \sim NMn(n; p_1, p_2, \dots, p_k)$.

If A_1, A_2, \dots, A_k describe all possible mutually exclusive “successes” and the event $\bar{A}_1 \cap \bar{A}_2 \cap \dots \cap \bar{A}_k$ presents the “failure”, then the coordinates ξ_i of the above vector can be interpreted as the number of “successes” of type A_i , $i = 1, 2, \dots, k$ until n -th “failure”.

This distribution is a particular case of Multivariate Power series distribution. (The definition is recalled below.) Considering this distribution for $k = 1$, we obtain a version of NBi distribution used in this paper. We denote the membership of a random variable ξ_1 to this class of distributions by $\xi_1 \sim \text{NBi}(n; 1 - p_1)$.

Notice that the marginal distributions of NMn distributed random vector are $\text{NBi}(n, 1 - \rho_i)$, $\rho_i = \frac{p_i}{1 - \sum_{j \neq i} p_j}$. More precisely their probability generating function (p.g.f.) is $G_{\xi}(z) = E z^{\xi_i} = \left(\frac{1 - \rho_i}{1 - \rho_i z} \right)^n$, $|z| < \frac{1}{\rho_i}$, $i = 1, 2, \dots, k$.

The distribution in Definition 1.1 is sometimes called Multivariate Negative Binomial distribution.

For $n = 1$ the NMn distribution is a Multivariate geometric distribution. Some properties of the bivariate version of this distribution are considered e.g. by Phatak *et al.* [15]. A systematic investigation of multivariate version could be found e.g. in Srivastava *et al.* [21].

If $(\xi_1, \xi_2, \dots, \xi_k) \sim \text{NMn}(n; p_1, p_2, \dots, p_k)$, its probability generating function (p.g.f.) is

$$(1.1) \quad G_{\xi_1, \xi_2, \dots, \xi_k}(z_1, z_2, \dots, z_k) = \left\{ \frac{1 - p_1 - p_2 - \dots - p_k}{1 - (p_1 z_1 + p_2 z_2 + \dots + p_k z_k)} \right\}^n,$$

$$|p_1 z_1 + p_2 z_2 + \dots + p_k z_k| < 1.$$

For $m = 2, 3, \dots, k - 1$, its finite dimensional distributions (f.d.ds) are, $(\xi_{i_1}, \xi_{i_2}, \dots, \xi_{i_m}) \sim \text{NMn}(n; \rho_{i_1}, \rho_{i_2}, \dots, \rho_{i_m})$, with

$$(1.2) \quad \rho_{i_s} = \frac{p_{i_s}}{1 - \sum_{j \notin \{i_1, i_2, \dots, i_m\}} p_j}, \quad s = 1, 2, \dots, m,$$

and for the set of indexes $\bar{i}_1, \bar{i}_2, \dots, \bar{i}_{k-m}$ that complements i_1, i_2, \dots, i_m to the set $1, 2, \dots, k$ its conditional distributions are

$$(1.3) \quad (\xi_{\bar{i}_1}, \xi_{\bar{i}_2}, \dots, \xi_{\bar{i}_{k-m}} \mid \xi_{i_1} = n_1, \xi_{i_2} = n_2, \dots, \xi_{i_m} = n_m) \sim \\ \sim \text{NMn}(n + n_1 + n_2 + \dots + n_m; p_{\bar{i}_1}, p_{\bar{i}_2}, \dots, p_{\bar{i}_{k-m}}).$$

More properties of NMn distribution can be found in Bates and Neyman [3] or Johnson *et al.* [7].

The set of all NMn distributions with one and the same p_1, p_2, \dots, p_k is closed with respect to convolution.

Lemma 1.1. *If r vs $\mathbf{S}_i \sim \text{NMn}(n_i; p_1, p_2, \dots, p_k)$, $i = 1, 2, \dots, m$ are independent, then the random vector*

$$(1.4) \quad \mathbf{S}_1 + \mathbf{S}_2 + \dots + \mathbf{S}_m \sim \text{NMn}(n_1 + n_2 + \dots + n_m; p_1, p_2, \dots, p_k).$$

One of the most comprehensive treatments with a very good list of references on Multivariate discrete distributions is the book of Johnson *et al.* [7].

The class of Power Series (PS) Distributions seems to be introduced by Noack (1950) [13] and Khatri (1959) [11]. A systematic approach on its properties could be found e.g. in Johnson *et al.* [8]. We will recall now only the most important for our work.

Definition 1.2. Let $\vec{a} = (a_0, a_1, \dots)$, where $a_i \geq 0$, $i = 0, 1, \dots$ and $\theta \in \mathbb{R}$ is such that

$$(1.5) \quad 0 < g_{\vec{a}}(\theta) = \sum_{n=0}^{\infty} a_n \theta^n < \infty.$$

A random variable (r.v.) X is Power series distributed, associated with the function $g_{\vec{a}}$ and the parameter θ (or equivalently associated with the sequence \vec{a} and the parameter θ), if it has p.m.f.

$$(1.6) \quad P(X=n) = \frac{a_n \theta^n}{g_{\vec{a}}(\theta)}, \quad n = 0, 1, \dots$$

Briefly $X \sim PS(a_1, a_2, \dots; \theta)$ or $X \sim PS(g_{\vec{a}}(x); \theta)$. The radius of convergence of the series (1.5) determines the parametric space Θ for θ . Further on we suppose that $\theta \in \Theta$.

Notice that given a PS distribution and the function $g_{\vec{a}}$ the constants θ and a_1, a_2, \dots are not uniquely determined, i.e. it is an ill-posed inverse problem. However, given the constants θ and a_0, a_1, \dots (or the function $g_{\vec{a}}(x)$ and θ) the corresponding PS distribution is uniquely determined. In this case, it is well known that:

- The p.g.f. of X is

$$(1.7) \quad \mathbb{E}z^X = \frac{g_{\vec{a}}(\theta z)}{g_{\vec{a}}(\theta)}, \quad z\theta \in \Theta.$$

- The type of all PS distributions is closed under convolution and more precisely if $X_1 \sim PS(g_1(x); \theta)$ and $X_2 \sim PS(g_2(x); \theta)$ are independent and $\theta \in \Theta_1 \cap \Theta_2$, then

$$(1.8) \quad X_1 + X_2 \sim PS(g_1(x)g_2(x); \theta).$$

- The mean is given by

$$(1.9) \quad \mathbb{E}X = \theta \frac{g'_{\vec{a}}(\theta)}{g_{\vec{a}}(\theta)} = \theta [\log(g_{\vec{a}}(\theta))]'.$$

From now on we denote the first and the second derivative of $g(x)$ with respect to x briefly by $g'(x)$ and $g''(x)$.

- The variance of X has the form

$$(1.10) \quad \text{Var } X = \theta^2 [\log(g_{\vec{a}}(\theta))]'' + \mathbb{E}X;$$

- The Fisher index is given by

$$FI X = 1 + \theta \frac{[\log(g_{\vec{a}}(\theta))]''}{[\log(g_{\vec{a}}(\theta))]'}.$$

We show that the class of Compound Power Series Distributions with Negative Multinomial Summands is a particular case of Multivariate Power series distribution (MPSD) considered by Johnson *et al.* [7]. Therefore let us remind the definition and its main properties.

Definition 1.3. Let $\theta_j > 0$, $j = 1, 2, \dots, k$ be positive real numbers and $a_{(i_1, i_2, \dots, i_k)}$, $i_j = 0, 1, \dots$, be non-negative constants such that

$$(1.11) \quad A_{\vec{a}}(\theta_1, \theta_2, \dots, \theta_k) = \sum_{i_1=0}^{\infty} \cdots \sum_{i_k=0}^{\infty} a_{(i_1, i_2, \dots, i_k)} \theta_1^{i_1} \theta_2^{i_2} \cdots \theta_k^{i_k} < \infty.$$

The distribution of the random vector $\vec{X} = (X_1, X_2, \dots, X_k)$ with probability mass function

$$P(X_1 = n_1, X_2 = n_2, \dots, X_k = n_k) = \frac{a_{(n_1, n_2, \dots, n_k)} \theta_1^{n_1} \theta_2^{n_2} \cdots \theta_k^{n_k}}{A_{\vec{a}}(\theta_1, \theta_2, \dots, \theta_k)}$$

is called Multivariate Power Series Distribution (MPSD) with parameters $A_{\vec{a}}(\vec{x})$, $a_{(i_1, i_2, \dots, i_k)}$ and $\vec{\theta} = (\theta_1, \theta_2, \dots, \theta_k)$. Briefly $\vec{X} \sim \text{MPSD}(A_{\vec{a}}(\vec{x}), \vec{\theta})$. As follows, Θ_k denotes the set of all parameters $\vec{\theta} = (\theta_1, \theta_2, \dots, \theta_k)$ that satisfy (1.11).

This class of distributions seems to be introduced by Patil (1965) [14] and Khatri (1959) [11]. A very useful necessary and sufficient condition that characterise this family is obtained by Gerstenkorn(1981). It is well known (see e.g. Johnson *et al.* [7]) that the p.g.f. of \vec{X} is

$$(1.12) \quad \mathbb{E}z_1^{X_1} z_2^{X_2} \cdots z_k^{X_k} = \frac{A_{\vec{a}}(\theta_1 z_1, \theta_2 z_2, \dots, \theta_k z_k)}{A_{\vec{a}}(\theta_1, \theta_2, \dots, \theta_k)}, \quad (\theta_1 z_1, \theta_2 z_2, \dots, \theta_k z_k) \in \Theta_k.$$

Through the paper $k = 2, 3, \dots$, is fixed and it corresponds to the number of the coordinates. We denote by $\stackrel{d}{=}$ the coincidence in distribution, by “ \sim ” the fact that a r.v. belongs to a given class of distributions, by $G_{\xi_1, \xi_2, \dots, \xi_k}(z_1, z_2, \dots, z_k) = E(z_1^{\xi_1} \cdots z_k^{\xi_k})$, the joint p.g.f. of a random vector $(\xi_1, \xi_2, \dots, \xi_k)$ and by $FI\xi$ the index of dispersion of the r.v. ξ (i.e. the variance of ξ divided by the corresponding mean).

One can consider the different concepts for compounds. We use the following one.

Definition 1.4. Let $\vec{\xi}_i = (\xi_i^{(1)}, \xi_i^{(2)}, \dots, \xi_i^{(k)})$, $i = 1, 2, \dots$, be i.i.d. random vectors and N be a discrete r.v. independent on them. We call compound, a random vector $\vec{X}_N = (X_N^{(1)}, X_N^{(2)}, \dots, X_N^{(k)})$, defined by

$$X_N^{(j)} = I_{\{N>0\}} \sum_{i=1}^N \xi_i^{(j)} = \begin{cases} \sum_{i=1}^N \xi_i^{(j)} & \text{if } N > 0, \\ 0 & \text{otherwise,} \end{cases} \quad j = 1, 2, \dots, k.$$

The distribution of ξ is called compounding distribution.

Further on we are going to use the following properties:

1. $G_{\vec{X}_N}(z) = G_N(G_{\vec{\xi}}(z))$.
2. If $EN < \infty$ and $E\xi < \infty$, then $E\vec{X}_N = EN E\xi$ (see [17], Cor. 4.2.1).

3. If $\text{Var } N < \infty$ and coordinate-wise $\text{Var } \vec{\xi} < \infty$, $\text{Var } \vec{X}_N = \text{Var } N(E\vec{\xi})^2 + EN \text{Var } \vec{\xi}$ (see [17], Cor. 4.2.1).
4. $FI\vec{X}_N = FIN E\vec{\xi} + FI\vec{\xi}$.

Notice that properties 2. and 3. are particular cases of the well known Wald's equations.

Here we consider a multivariate distribution which coordinates are dependent compounds. In the notations of the Definition 1.4, N is PS distributed and $\vec{\xi}$ is NMn distributed. The cases when N is Poisson distributed is partially investigated in 1962, by G. Smith [20]. In Section 2, following the traditional approach about definition of distributions, first we define this distribution through its p.m.f., then we investigate its properties. We consider the case when the summands are NMn distributed. We obtain its main numerical characteristics and conditional distributions. Finally explain its relation with compounds and mixtures. We prove that the class of Compound Power Series Distributions with Negative Multinomial Summands is a particular case of Multivariate Power series distribution and find the explicit form of the parameters. We show that considered as a Mixture this distribution would be (possibly Zero-inflated) Mixed Negative Multinomial distribution with possibly scale changed Power series distributed first parameter. Using these relations we derive several properties and its main numerical characteristics. In Section 3 the risk process application is provided, together with simulations of the risk processes and estimation of ruin probabilities in a finite time interval.

2. DEFINITION AND MAIN PROPERTIES OF THE COMPOUND POWER SERIES DISTRIBUTION WITH NEGATIVE MULTINOMIAL SUMMANDS

Let us first define Compound Power series distribution with Negative multinomial summands and then to investigate its properties.

Definition 2.1. Let $\pi_j \in (0, 1)$, $j = 1, 2, \dots, k$, $\pi_0 := 1 - \pi_1 - \pi_2 - \dots - \pi_k \in (0, 1)$, $a_s \geq 0$, $s = 0, 1, \dots$ and $\theta \in \mathbb{R}$ be such that

$$(2.1) \quad 0 < g_{\vec{a}}(\theta) = \sum_{n=0}^{\infty} a_n \theta^n < \infty.$$

A random vector $\vec{X} = (X_1, X_2, \dots, X_k)$ is called Compound Power series distributed with negative multinomial summands and with parameters $g_{\vec{a}}(x)$, θ ; n , π_1, \dots, π_k , if for $i = 1, 2, \dots, k$, $m_i = 0, 1, 2, \dots$, and $(m_1, m_2, \dots, m_k) \neq (0, 0, \dots, 0)$,

$$(2.2) \quad \begin{aligned} P(X_1 = m_1, X_2 = m_2, \dots, X_k = m_k) &= \\ &= \frac{\pi_1^{m_1} \pi_2^{m_2} \dots \pi_k^{m_k}}{g_{\vec{a}}(\theta)} \sum_{j=1}^{\infty} a_j \theta^j \binom{jn + m_1 + m_2 + \dots + m_k - 1}{m_1, m_2, \dots, m_k, jn - 1} \pi_0^{nj}, \\ P(X_1 = 0, X_2 = 0, \dots, X_k = 0) &= \frac{g_{\vec{a}}(\theta \pi_0^n)}{g_{\vec{a}}(\theta)}. \end{aligned}$$

Briefly $\vec{X} \sim \text{CPSNMn}(g_{\vec{a}}(x), \theta; n, \pi_1, \pi_2, \dots, \pi_k)$ or $\vec{X} \sim \text{CPSNMn}(\vec{a}, \theta; n, \pi_1, \pi_2, \dots, \pi_k)$.¹

¹For $n = 1$ and $k = 2$ see [12].

In the next theorem we show that this distribution is a particular case of $MPSD(A(\vec{x}), \vec{\theta})$ considered in Johnson et al. [7].

Theorem 2.1. Suppose $\pi_i \in (0, 1)$, $i = 1, 2, \dots, k$, $\pi_0 := 1 - \pi_1 - \pi_2 - \dots - \pi_k \in (0, 1)$, $a_i \geq 0$, $i = 0, 1, \dots$ and $\theta \in \mathbb{R}$ are such that (2.1) is satisfied. If

$$\vec{X} \sim CPSNMn(g_{\vec{a}}(x), \theta; n, \pi_1, \pi_2, \dots, \pi_k),$$

then:

1. $\vec{X} \sim MPSD(A(\vec{x}), \vec{\theta})$, where $\vec{\theta} = (\pi_1, \pi_2, \dots, \pi_k)$, $a_{(0, \dots, 0)} = g_{\vec{a}}(\theta \pi_0^n)$.
For $(i_1, i_2, \dots, i_k) \neq (0, 0, \dots, 0)$,

$$a_{(i_1, i_2, \dots, i_k)} = \sum_{j=1}^{\infty} a_j \theta^j \binom{jn + i_1 + i_2 + \dots + i_k - 1}{i_1, i_2, \dots, i_k, jn - 1} \pi_0^{nj},$$

$$A(x_1, x_2, \dots, x_k) = g_{\vec{a}} \left\{ \theta \frac{\pi_0^n}{[1 - (x_1 + x_2 + \dots + x_k)]^n} \right\},$$

$$x_i \in (0, 1), i = 1, 2, \dots, k \text{ and } x_1 + x_2 + \dots + x_k \in (0, 1).$$

2. For $|\pi_1 z_1 + \pi_2 z_2 + \dots + \pi_k z_k| < 1$,

$$\begin{aligned} G_{X_1, X_2, \dots, X_k}(z_1, z_2, \dots, z_k) &= \frac{g_{\vec{a}} \left[\theta \left(\frac{\pi_0}{1 - (\pi_1 z_1 + \pi_2 z_2 + \dots + \pi_k z_k)} \right)^n \right]}{g_{\vec{a}}(\theta)} \\ &= \frac{g_{\vec{a}} \left[\theta \left(\frac{\pi_0}{\pi_0 + \pi_1(1 - z_1) + \pi_2(1 - z_2) + \dots + \pi_k(1 - z_k)} \right)^n \right]}{g_{\vec{a}}(\theta)}. \end{aligned}$$

3. For all $r = 2, 3, \dots, k$,

$$(X_{i_1}, X_{i_2}, \dots, X_{i_r}) \sim CPSNMn \left(g_{\vec{a}}(x), \theta; n, \frac{\pi_{i_1}}{\pi_0 + \pi_{i_1} + \pi_{i_2} + \dots + \pi_{i_r}}, \frac{\pi_{i_2}}{\pi_0 + \pi_{i_1} + \pi_{i_2} + \dots + \pi_{i_r}}, \dots, \frac{\pi_{i_r}}{\pi_0 + \pi_{i_1} + \pi_{i_2} + \dots + \pi_{i_r}} \right).$$

4. For $i = 1, 2, \dots, k$,

$$\begin{aligned} X_i &\sim CPSNBi \left(g_{\vec{a}}(x), \theta; n, \frac{\pi_i}{\pi_0 + \pi_i} \right), \\ G_{X_i}(z_i) &= \frac{g_{\vec{a}} \left[\theta \left(\frac{\pi_0}{\pi_0 + \pi_i(1 - z_i)} \right)^n \right]}{g_{\vec{a}}(\theta)}, \quad |\pi_i z_i| < \pi_0 + \pi_i, \\ EX_i &= n \theta [\log(g_{\vec{a}}(\theta))] \frac{\pi_i}{\pi_0} = n \theta \frac{\pi_i}{\pi_0} \frac{g'_{\vec{a}}(\theta)}{g_{\vec{a}}(\theta)}, \\ \text{Var } X_i &= n \frac{\pi_i \theta}{\pi_0^2} \left[n \pi_i \theta [\log(g_{\vec{a}}(\theta))]'' + [\log(g_{\vec{a}}(\theta))]' (\pi_0 + \pi_i(n + 1)) \right], \\ FIX_i &= 1 + \frac{\pi_i}{\pi_0} \left(n \theta \frac{[\log(g_{\vec{a}}(\theta))]''}{[\log(g_{\vec{a}}(\theta))]'} + n + 1 \right). \end{aligned}$$

5. For $i \neq j = 1, 2, \dots, k$,

$$(X_i, X_j) \sim \text{CPSNMn} \left(g_{\tilde{a}}(x), \theta; n, \frac{\pi_i}{\pi_0 + \pi_i + \pi_j}, \frac{\pi_j}{\pi_0 + \pi_i + \pi_j} \right),$$

$$G_{X_i, X_j}(z_i, z_j) = \frac{g_{\tilde{a}} \left[\theta \left(\frac{\pi_0}{\pi_0 + (1-z_i)\pi_i + (1-z_j)\pi_j} \right)^n \right]}{g_{\tilde{a}}(\theta)}, \quad |\pi_i z_i + \pi_j z_j| < \pi_0 + \pi_i + \pi_j,$$

$$\text{cov}(X_i, X_j) = \frac{n\pi_i\pi_j\theta}{\pi_0^2} \left\{ n\theta [\log g_{\tilde{a}}(\theta)]'' + (n+1) [\log g_{\tilde{a}}(\theta)]' \right\},$$

$$\text{cor}(X_i, X_j) = \sqrt{\frac{(FIX_i - 1)(FIX_j - 1)}{FIX_i FIX_j}}.$$

6. For $i, j = 1, 2, \dots, k, j \neq i$,

(a) For $m_j \neq 0$,

$$P(X_i = m_i | X_j = m_j) = \left(\frac{\pi_0 + \pi_j}{\pi_0 + \pi_i + \pi_j} \right)^{m_j} \frac{\pi_i^{m_i}}{m_i! (\pi_0 + \pi_i + \pi_j)^{m_i}} \cdot \frac{\sum_{s=1}^{\infty} a_s \theta^s \frac{(sn+m_i+m_j-1)!}{(sn-1)!} \left(\frac{\pi_0}{\pi_0 + \pi_i + \pi_j} \right)^{ns}}{\sum_{s=1}^{\infty} a_s \theta^s \frac{(sn+m_j-1)!}{(sn-1)!} \left(\frac{\pi_0}{\pi_0 + \pi_j} \right)^{ns}}, \quad m_i = 0, 1, \dots$$

(b) $(X_i | X_j = 0) \sim \text{CPSNMn} \left[a_s, \tilde{\theta} = \theta \left(\frac{\pi_0}{\pi_0 + \pi_j} \right)^n; n, \frac{\pi_i}{\pi_0 + \pi_i + \pi_j} \right],$

$$P(X_i = m_i | X_j = 0) = \frac{\pi_i^{m_i}}{m_i! (\pi_0 + \pi_i + \pi_j)^{m_i}} \cdot \frac{\sum_{s=1}^{\infty} a_s \frac{(sn+m_i-1)!}{(sn-1)!} \left[\theta \left(\frac{\pi_0}{\pi_0 + \pi_j} \right)^n \right]^s}{g_{\tilde{a}} \left[\theta \left(\frac{\pi_0}{\pi_0 + \pi_j} \right)^n \right]}, \quad m_i \in \mathbb{N},$$

$$P(X_i = 0 | X_j = 0) = \frac{g_{\tilde{a}} \left[\theta \left(\frac{\pi_0}{\pi_0 + \pi_i + \pi_j} \right)^n \right]}{g_{\tilde{a}} \left[\theta \left(\frac{\pi_0}{\pi_0 + \pi_j} \right)^n \right]}.$$

(c) For $i, j = 1, 2, \dots, k, j \neq i, m_j = 1, 2, \dots$,

$$E(z_i^{X_i} | X_j = m_j) = \left(\frac{\pi_0 + \pi_j}{\pi_0 + \pi_j + \pi_i - z_i \pi_i} \right)^{m_j} \cdot \frac{\sum_{s=1}^{\infty} a_s \frac{(sn+m_i-1)!}{(sn-1)!} \left[\theta \left(\frac{\pi_0}{\pi_0 + \pi_j + \pi_i - z_i \pi_i} \right)^n \right]^s}{\sum_{s=1}^{\infty} a_s \frac{(sn+m_j-1)!}{(sn-1)!} \left[\theta \left(\frac{\pi_0}{\pi_0 + \pi_j} \right)^n \right]^s},$$

$$E(z_i^{X_i} | X_j = 0) = \frac{g_{\tilde{a}} \left[\theta \left(\frac{\pi_0}{\pi_0 + \pi_j + (1-z_i)\pi_i} \right)^n \right]}{g_{\tilde{a}} \left[\theta \left(\frac{\pi_0}{\pi_0 + \pi_j} \right)^n \right]}.$$

(d) For $z_j = 1, 2, \dots$,

$$E(X_i | X_j = z_j) = \frac{\pi_i}{\pi_0 + \pi_j} \frac{\sum_{s=1}^{\infty} a_s \frac{(sn+z_j)!}{(sn-1)!} \left[\theta \left(\frac{\pi_0}{\pi_0 + \pi_j} \right)^n \right]^s}{\sum_{s=1}^{\infty} a_s \frac{(sn+z_j-1)!}{(sn-1)!} \left[\theta \left(\frac{\pi_0}{\pi_0 + \pi_j} \right)^n \right]^s},$$

$$E(X_i | X_j = 0) = \frac{n\pi_i}{\pi_0 + \pi_j} \frac{\sum_{s=1}^{\infty} s a_s \left[\theta \left(\frac{\pi_0}{\pi_0 + \pi_j} \right)^n \right]^s}{\sum_{s=0}^{\infty} a_s \left[\theta \left(\frac{\pi_0}{\pi_0 + \pi_j} \right)^n \right]^s}$$

$$= \frac{n\theta\pi_0^n\pi_i}{(\pi_0 + \pi_j)^{n+1}} \frac{g'_{\bar{a}} \left[\theta \left(\frac{\pi_0}{\pi_0 + \pi_j} \right)^n \right]}{g_{\bar{a}} \left[\theta \left(\frac{\pi_0}{\pi_0 + \pi_j} \right)^n \right]}.$$

7. $X_1 + X_2 + \dots + X_k \sim CPSNBi(g_{\bar{a}}(x), \theta; n, 1 - \pi_0)$.

8. For $i = 1, 2, \dots, k$

$$(X_i, X_1 + X_2 + \dots + X_k - X_i) \sim CPSNMn(g_{\bar{a}}(x), \theta; n, \pi_i, 1 - \pi_0 - \pi_i).$$

9. For $i = 1, 2, \dots, k, m \in \mathbb{N}$

$$(X_i | X_1 + X_2 + \dots + X_k = m) \sim Bi\left(m, \frac{\pi_i}{1 - \pi_0}\right).$$

Sketch of the proof:

1. We substitute of the considered values and function A in the necessary and sufficient condition, given in p.154, Johnson *et al.* [7], for MPSD and prove that the following two conditions are satisfied:

$$P(X_1 = 0, X_2 = 0, \dots, X_k = 0) = \frac{a(0,0,\dots,0)}{A(\theta_1, \theta_2, \dots, \theta_k)},$$

$$\frac{P(X_1 = n_1 + m_1, X_2 = n_2 + m_2, \dots, X_k = n_k + m_k)}{P(X_1 = n_1, X_2 = n_2, \dots, X_k = n_k)} =$$

$$= \frac{a(n_1+m_1, n_2+m_2, \dots, n_k+m_k)}{a(n_1, n_2, \dots, n_k)} \theta_1^{m_1} \theta_2^{m_2} \dots \theta_k^{m_k}, \quad m_i, n_i = 0, 1, \dots, \quad i = 1, 2, \dots, k.$$

2.-3. Analogously to [12], who works in case $n = 1$ and $k = 2$. Here we have used the definition of p.g.f., the definition of $g_{\bar{a}}(x)$ and the formula

$$\sum_{i_1=0}^{\infty} \sum_{i_2=0}^{\infty} \dots \sum_{i_k=0}^{\infty} \frac{(i_1 + i_2 + \dots + i_k + r - 1)!}{i_1! i_2! \dots i_k! (r - 1)!} x_1^{i_1} x_2^{i_2} \dots x_k^{i_k} = \frac{1}{(1 - x_1 - x_2 - \dots - x_k)^r}.$$

6.(a) For $m_i = 0, 1, \dots$ we substitute the proposed parameters and function in the formula of p.m.f. of PS distribution and obtain the above formula

$$\begin{aligned}
P(X_i = m_i | X_j = m_j) &= \frac{\sum_{s=1}^{\infty} a_s \theta^s \frac{(sn+m_i+m_j-1)!}{(sn-1)!m_i!} \left(\frac{\pi_0}{\pi_0+\pi_i+\pi_j}\right)^{ns} \left(\frac{\pi_i}{\pi_0+\pi_i+\pi_j}\right)^{m_i}}{\sum_{k=0}^{\infty} \sum_{s=1}^{\infty} a_s \theta^s \frac{(sn+k+m_j-1)!}{(sn-1)!k!} \left(\frac{\pi_0}{\pi_0+\pi_i+\pi_j}\right)^{ns} \left(\frac{\pi_i}{\pi_0+\pi_i+\pi_j}\right)^k} \\
&= \frac{\sum_{s=1}^{\infty} a_s \theta^s \frac{(sn+m_i+m_j-1)!}{(sn-1)!m_i!} \left(\frac{\pi_0}{\pi_0+\pi_i+\pi_j}\right)^{ns} \left(\frac{\pi_i}{\pi_0+\pi_i+\pi_j}\right)^{m_i}}{\sum_{s=1}^{\infty} \frac{a_s}{(sn-1)!} \theta^s \left(\frac{\pi_0}{\pi_0+\pi_i+\pi_j}\right)^{ns} \sum_{k=0}^{\infty} \frac{(sn+k+m_j-1)!}{k!} \left(\frac{\pi_i}{\pi_0+\pi_i+\pi_j}\right)^k} \\
&= \frac{\sum_{s=1}^{\infty} a_s \theta^s \frac{(sn+m_i+m_j-1)!}{(sn-1)!m_i!} \left(\frac{\pi_0}{\pi_0+\pi_i+\pi_j}\right)^{ns} \left(\frac{\pi_i}{\pi_0+\pi_i+\pi_j}\right)^{m_i}}{\sum_{s=1}^{\infty} \frac{a_s}{(sn-1)!} \theta^s \left(\frac{\pi_0}{\pi_0+\pi_i+\pi_j}\right)^{ns} \frac{(sn+m_j-1)!}{\left[1-\left(\frac{\pi_i}{\pi_0+\pi_i+\pi_j}\right)\right]^{sn+m_j}}} \\
&= \frac{1}{m_i!} \left(\frac{\pi_i}{\pi_0+\pi_i+\pi_j}\right)^{m_i} \left(\frac{\pi_0+\pi_j}{\pi_0+\pi_i+\pi_j}\right)^{m_j} \\
&\quad \cdot \frac{\sum_{s=1}^{\infty} a_s \theta^s \frac{(sn+m_i+m_j-1)!}{(sn-1)!} \left(\frac{\pi_0}{\pi_0+\pi_i+\pi_j}\right)^{ns}}{\sum_{s=1}^{\infty} \frac{a_s \theta^s (sn+m_j-1)!}{(sn-1)!} \left(\frac{\pi_0}{\pi_0+\pi_j}\right)^{ns}}.
\end{aligned}$$

6.(c) and **6.(d)** are analogous to [12], who work in case $n = 1$ and $k = 2$.

9. For $i = 1, 2, \dots, k$, $m \in \mathbb{N}$, we use 7., 8., the definitions about CPSNMn distribution and conditional probability, and obtain

$$\begin{aligned}
P(X_i = s | X_1 + X_2 + \dots + X_k = m) &= \frac{P(X_i = s, X_1 + X_2 + \dots + X_k = m)}{P(X_1 + X_2 + \dots + X_k = m)} \\
&= \frac{P(X_i = s, X_1 + X_2 + \dots + X_k - X_i = m - s)}{P(X_1 + X_2 + \dots + X_k = m)} \\
&= \frac{\pi_i^s (1 - \pi_0 - \pi_i)^{m-s} \sum_{j=1}^{\infty} a_j \theta^j \binom{jn+m-1}{s, m-s, jn-1} \pi_0^{nj}}{(1 - \pi_0)^m \sum_{j=1}^{\infty} a_j \theta^j \binom{jn+m-1}{m, jn-1} \pi_0^{nj}} \\
&= \frac{\pi_i^s (1 - \pi_0 - \pi_i)^{m-s} \sum_{j=1}^{\infty} a_j \theta^j \frac{(jn+m-1)!}{s!(m-s)!(jn-1)!} \pi_0^{nj}}{(1 - \pi_0)^m \sum_{j=1}^{\infty} a_j \theta^j \frac{(jn+m-1)!}{m!(jn-1)!} \pi_0^{nj}} \\
&= \binom{m}{s} \left(\frac{\pi_i}{1 - \pi_0}\right)^s \left(1 - \frac{\pi_i}{1 - \pi_0}\right)^{m-s}, \quad s = 0, 1, \dots, m.
\end{aligned}$$

We use the definition of Binomial distribution and complete the proof.

Note 2.1. The conclusion 1. in this theorem states that also in the univariate case the CPSMNn distribution is just a particular case of PSD with more complicated coefficients.

The next theorem presents this distribution as a mixture.

Theorem 2.2. Suppose $n \in \mathbb{N}$, $\pi_i \in (0, 1)$, $i = 1, 2, \dots, k$, $\pi_0 := 1 - \pi_1 - \pi_2 - \dots - \pi_k \in (0, 1)$, $a_j \geq 0$, $j = 0, 1, \dots$, $\theta \in \mathbb{R}$ are such that (2.1) is satisfied and $\vec{X} \sim \text{CPSNMn}(g_{\vec{a}}(x), \theta; n, \pi_1, \pi_2, \dots, \pi_k)$. Then there exists a probability space $(\Omega, \mathcal{A}, \mathbb{P})$, a r.v. $M \sim \text{PSD}(g_{\vec{a}}(x), \theta)$ and a random vector $\vec{Y} = (Y_1, Y_2, \dots, Y_k)$ defined on it, such that $\vec{Y}|M = m \sim \text{NMn}(nm, \pi_1, \pi_2, \dots, \pi_k)$, $m = 1, 2, \dots$,

$$P(Y_1 = 0, Y_2 = 0, \dots, Y_k = 0 | M = 0) = 1,$$

and $\vec{X} \stackrel{d}{=} \vec{Y}$. Moreover

1. For $|\pi_1 z_1 + \pi_2 z_2 + \dots + \pi_k z_k| < 1$,

$$G_{X_1, X_2, \dots, X_k}(z_1, z_2, \dots, z_k) = G_M \left[\left(\frac{\pi_0}{\pi_0 + \pi_1(1 - z_1) + \pi_2(1 - z_2) + \dots + \pi_k(1 - z_k)} \right)^n \right].$$

2. For $i = 1, 2, \dots, k$,

$$EX_i = n EM \frac{\pi_i}{\pi_0}, \quad i = 1, 2, \dots, k,$$

$$\begin{aligned} \text{Var } X_i &= \text{Var } M n^2 \frac{\pi_i^2}{\pi_0^2} + EM n \frac{\pi_i(\pi_0 + \pi_i)}{\pi_0^2} \\ &= n \frac{\pi_i}{\pi_0} EM \left[\frac{\pi_i}{\pi_0} (n FIM + 1) + 1 \right], \end{aligned}$$

$$FI X_i = 1 + \frac{\pi_i}{\pi_0} (n FIM + 1).$$

3. For $i \neq j = 1, 2, \dots, k$,

$$\text{cov}(X_i, X_j) = n \frac{\pi_i \pi_j}{\pi_0^2} \{n FIM + 1\} EM,$$

$$\text{cor}(X_i, X_j) = \sqrt{\frac{(FI Y_i - 1)(FI Y_j - 1)}{FI Y_i FI Y_j}}.$$

Note 2.2. Following analogous notations of Johnson *et al.* [7], the above two theorems state that CPSNMn distribution coincides with

$$I_{\{M > 0\}} \text{NMn}(nM, \pi_1, \pi_2, \dots, \pi_k) \bigwedge_M \text{PSD}(g_{\vec{a}}(x); \theta),$$

where $I_{M > 0}$ is a Bernoulli r.v. or indicator of the event “ $M > 0$ ”.

The following representation motivates the name of CPSNMn distribution.

Theorem 2.3. Suppose $\pi_i \in (0, 1)$, $i = 1, 2, \dots, k$, $\pi_0 = 1 - \pi_1 - \dots - \pi_k \in (0, 1)$, $a_k \geq 0$, $k = 0, 1, \dots$, and $\theta \in \mathbb{R}$ are such that (2.1) is satisfied. Let $M \sim \text{PS}(g_{\vec{a}}(x); \theta)$ and $(Y^{(1)}, \dots, Y^{(k)}) \sim \text{NMn}(n; \pi_1, \dots, \pi_k)$ be independent. Denote by $I_{\{M > 0\}}$ the Bernoulli r.v. that is an indicator of the event $\{M > 0\}$ and defined on the same probability space. Define a random vector $(T_M^{(1)}, T_M^{(2)}, \dots, T_M^{(k)})$ by

$$(2.3) \quad T_M^{(j)} = I_{\{M > 0\}} \sum_{i=1}^M Y_i^{(j)} = \begin{cases} \sum_{i=1}^M Y_i^{(j)} & \text{if } M > 0, \\ 0 & \text{otherwise,} \end{cases} \quad j = 1, 2, \dots, k.$$

Then

1. For $m \in \mathbb{N}$, $(T_M^{(1)}, T_M^{(2)}, \dots, T_M^{(k)} | M = m) \sim NMn(nm; \pi_1, \pi_2, \dots, \pi_k)$;
2. $(T_M^{(1)}, T_M^{(2)}, \dots, T_M^{(k)}) \sim CPSNMn(g_{\vec{a}}(x), \theta; n, \pi_1, \pi_2, \dots, \pi_k)$;
3. $(T_M^{(1)}, T_M^{(2)}, \dots, T_M^{(k)} | M > 0) \sim CPSNMn(g_{\vec{a}}(x), \theta; n, \pi_1, \pi_2, \dots, \pi_k)$,
where $\tilde{a}_0 = 0$, $\tilde{a}_i = a_i$, $i = 1, 2, \dots$

Sketch of the proof: We apply (1.4) and Theorem 2.2. 2. is analogous to [12], who work in case $n = 1$ and $k = 2$.

If we have no weights at coordinate planes we need to consider the following distribution.

Definition 2.2. Let $\pi_j \in (0, 1)$, $j = 1, 2, \dots, k$, $\pi_0 := 1 - \pi_1 - \pi_2 - \dots - \pi_k \in (0, 1)$, $a_s \geq 0$, $s = 0, 1, \dots$, and $\theta \in \mathbb{R}$ be such that

$$g_{\vec{a}}(\theta) = \sum_{n=0}^{\infty} a_n \theta^n < \infty.$$

A random vector $\vec{X} = (X_1, X_2, \dots, X_k)$ is called Compound Power series distributed with negative multinomial summands on \mathbb{N}^k and with parameters $g_{\vec{a}}(x)$, θ ; n , π_1, \dots, π_k , if for $i = 1, 2, \dots, k$, $m_i = 1, 2, \dots$,

$$(2.4) \quad \begin{aligned} P(X_1 = m_1, X_2 = m_2, \dots, X_k = m_k) &= \\ &= \frac{1}{\rho} \frac{\pi_1^{m_1} \pi_2^{m_2} \dots \pi_k^{m_k}}{g_{\vec{a}}(\theta \pi_0^n)} \sum_{j=1}^{\infty} a_j \theta^j \binom{jn + m_1 + m_2 + \dots + m_k - 1}{m_1, m_2, \dots, m_k, jn - 1} \pi_0^{nj}, \\ \rho &= 1 - \sum_{m=1}^k (-1)^{m+1} \sum_{1 \leq i_1 < i_2 < \dots < i_m \leq k} \frac{g_{\vec{a}}[\theta \pi_0^n (\pi_0 + \pi_{i_1} + \dots + \pi_{i_m})^{-n}]}{g_{\vec{a}}(\theta)}. \end{aligned}$$

Briefly $\vec{X} \sim CPSNMn_{\mathbb{N}^k}(g_{\vec{a}}(x), \theta; n, \pi_1, \pi_2, \dots, \pi_k)$.

The relation between $CPSNMn$ and $CPSNMn_{\mathbb{N}^k}$ distributions is given in the following theorem.

Theorem 2.4. If $\vec{X} \sim CPSNMn(g_{\vec{a}}(x), \theta; n, \pi_1, \pi_2, \dots, \pi_k)$, then

$$(X_1, X_2, \dots, X_k | X_1 \neq 0, X_2 \neq 0, \dots, X_k \neq 0) \sim CPSNMn_{\mathbb{N}^k}(g_{\vec{a}}(x), \theta; n, \pi_1, \pi_2, \dots, \pi_k).$$

3. APPLICATIONS TO RISK THEORY

In [9] we obtained the approximations of Compound Poisson risk process mixed with Pareto r.v. and provide a brief summary of previous results about risk process approximations. In this section we provide risk process application of the $CPSNMn$. Here $k, n \in \mathbb{N}$, $p_M \in (0, 1)$, $\pi_i \in (0, 1)$, $i = 1, 2, \dots, k$, $\pi_1 + \dots + \pi_k < 1$ and $\pi_0 = 1 - \pi_1 - \pi_2 - \dots - \pi_k$.

3.1. The counting process

Here we consider a discrete time counting process, satisfying the following conditions:

- C1.** The insurance company have no claims at moment $t = 0$.
- C2.** In any other moments of time $t = 1, 2, \dots$ a group of claims can arrive with probability p_M independently of others. We denote the number of groups of claims, arrived in the insurance company over an interval $[0, t]$ by $M(t)$ and by $0 < T_{G,1} < T_{G,2} < \dots$ the moments of arrivals of the corresponding group, i.e. $T_{G,k}$ is the occurrence time of the k -th group. By definition $M(0) = 0$.
- C3.** The claims can be of one of k mutually exclusive and totally exhaustive different types A_1, A_2, \dots, A_k , e.g. claims of one individual having several pension insurances.
- C4.** In any of the time points $0 < T_{G,1} < T_{G,2} < \dots$, we denote the number of claims of type $i = 1, 2, \dots, k$, arrived in the insurance company by $Y_{i,j}$, $j = 1, 2, \dots$. We assume that the random vectors $(Y_{1,j}, Y_{2,j}, \dots, Y_{k,j})$, $j = 1, 2, \dots$ are i.i.d. and

$$(Y_{1,j}, Y_{2,j}, \dots, Y_{k,j}) \sim NMn(n, \pi_1, \pi_2, \dots, \pi_k).$$

Note 3.1. Conditions C1–C2 means that the counting process of the groups of claims up to time $t > 0$ is a Binomial process. In case when the claim sizes are discrete they are considered e.g. in [5, 19, 4, 22]. The number of groups arrived up to time t is $M(t) \sim Bi(t, p_M)$ and the intervals $T_{G,1}, T_{G,2} - T_{G,1}, T_{G,3} - T_{G,2}, \dots$ between the groups arrivals are i.i.d. Geometrically distributed on $1, 2, \dots$, with parameter p_M .

C4 means that it is possible to have zero reported losses of one or of all k -types of insurance claims within one group. In that case there is a group arrived, however, the number of participants in the group is zero. This can happen e.g. when there is a claim, but it is not accepted, or it is estimated by zero value by the insurer.

Let us denote the number of claims of type $i = 1, 2, \dots, k$, arrived in the company in the interval $[0, t]$ by $N_{i,t}$. Conditions C1–C4 imply that $(N_1(0), N_2(0), \dots, N_k(0)) = (0, 0, \dots, 0)$ and, for all $t = 1, 2, \dots$,

$$N_i(t) = I\{M(t) > 0\} \sum_{j=1}^{M(t)} Y_{i,j}, \quad j = 1, 2, \dots, k.$$

Therefore

$$(N_1(t), N_2(t), \dots, N_k(t)) \sim CPSNMn\left((1+x)^t, \frac{p_M}{1-p_M}; n, \pi_1, \pi_2, \dots, \pi_k\right)$$

and

$$P(N_1(t) + N_2(t) + \dots + N_k(t) = 0) = \frac{(1-p_M)^t}{(1-p_M \pi_0^n)^t}.$$

3.2. The total claim amount process and its characteristics

Consider the total claim amount process defined as

$$(3.1) \quad S(t) = I_{\{N_1(t)>0\}} \sum_{j_1=1}^{N_1(t)} Z_{1,j_1} + I_{\{N_2(t)>0\}} \sum_{j_2=1}^{N_2(t)} Z_{2,j_2} + \cdots + I_{\{N_k(t)>0\}} \sum_{j_k=1}^{N_k(t)} Z_{k,j_k},$$

$t = 1, 2, \dots$, satisfying C1–C4.

We impose the following conditions on the claim sizes:

C5. In any of the time points $0 < T_{G,1} < T_{G,2} < \cdots$, we denote the claim sizes of the claims of type $i = 1, 2, \dots, k$ by $Z_{i,j}$, $j = 1, 2, \dots$. We assume that the random vectors $(Z_{1,j}, Z_{2,j}, \dots, Z_{k,j})$, $j = 1, 2, \dots$, are i.i.d. and the coordinates of this vector are also independent, with absolutely continuous c.d.fs. correspondingly F_i , $i = 1, 2, \dots, k$, concentrated on $(0, \infty)$.

C6. The claim arrival times and the claim sizes are assumed to be independent.

Proposition 3.1. *Consider the total claim amount process defined in (3.1) and satisfying conditions C1–C6.*

1. If $\mathbb{E}Z_{i,j} = \mu_i < \infty$, $i = 1, 2, \dots, k$, then

$$(3.2) \quad \mathbb{E}S(t) = \frac{nt p_M (1 - p_M)}{\pi_0} (\mu_1 \pi_1 + \mu_2 \pi_2 + \cdots + \mu_k \pi_k).$$

2. If additionally $\text{Var } Z_{i,j} = \sigma_i^2 < \infty$, $i = 1, 2, \dots, k$, then

$$(3.3) \quad \text{Var } S(t) = nt \frac{p_M}{\pi_0} \left\{ \sum_{i=1}^k \pi_i (\sigma_i^2 + \mu_i^2) + \frac{(1 - p_M)n + 1}{\pi_0} \left(\sum_{i=1}^k \mu_i \pi_i \right)^2 \right\},$$

$$FI S(t) = \frac{\sum_{i=1}^k \pi_i (\sigma_i^2 + \mu_i^2) + \frac{(1 - p_M)n + 1}{\pi_0} \left(\sum_{i=1}^k \mu_i \pi_i \right)^2}{(1 - p_M) (\mu_1 \pi_1 + \mu_2 \pi_2 + \cdots + \mu_k \pi_k)}.$$

Proof: [1.] is a consequence of the double expectation formula.

[2.] Using the double expectation formula, the facts that $EM(t) = tp_M$, $\text{Var } M(t) = tp_M(1 - p_M)$ and Theorem 2.2 we obtain:

$$\begin{aligned} \text{Var } S(t) &= \sum_{i=1}^k ntp_M \frac{\pi_i}{\pi_0} \sigma_i^2 + ntp_M \sum_{i=1}^k \left[(1 - p_M)n \frac{\pi_i^2}{\pi_0^2} + \frac{\pi_i(\pi_0 + \pi_i)}{\pi_0^2} \right] \mu_i^2 \\ &\quad + 2 \sum_{1 \leq i < j \leq k} \mu_i \mu_j \text{cov}(N_i(t), N_j(t)) = \end{aligned}$$

$$\begin{aligned}
&= nt \frac{p_M}{\pi_0} \left\{ \sum_{i=1}^k \pi_i \sigma_i^2 + \sum_{i=1}^k \left[(1-p_M)n \frac{\pi_i^2}{\pi_0} + \frac{\pi_i(\pi_0 + \pi_i)}{\pi_0} \right] \mu_i^2 \right\} \\
&\quad + 2 \sum_{1 \leq i < j \leq k} \mu_i \mu_j \operatorname{cov}(N_i(t), N_j(t)) \\
&= nt \frac{p_M}{\pi_0} \left\{ \sum_{i=1}^k \pi_i (\sigma_i^2 + \mu_i^2) + \sum_{i=1}^k \left[(1-p_M)n \frac{\pi_i^2}{\pi_0} + \frac{\pi_i^2}{\pi_0} \right] \mu_i^2 \right\} \\
&\quad + 2 \sum_{1 \leq i < j \leq k} \mu_i \mu_j \operatorname{cov}(N_i(t), N_j(t)) \\
&= nt \frac{p_M}{\pi_0} \left\{ \sum_{i=1}^k \pi_i (\sigma_i^2 + \mu_i^2) + \frac{(1-p_M)n+1}{\pi_0} \sum_{i=1}^k \mu_i^2 \pi_i^2 \right\} \\
&\quad + 2 \sum_{1 \leq i < j \leq k} \mu_i \mu_j \operatorname{cov}(N_i(t), N_j(t)) \\
&= nt \frac{p_M}{\pi_0} \left\{ \sum_{i=1}^k \pi_i (\sigma_i^2 + \mu_i^2) + \frac{(1-p_M)n+1}{\pi_0} \sum_{i=1}^k \mu_i^2 \pi_i^2 \right\} \\
&\quad + 2ntp_M \frac{[n(1-p_M)+1]}{\pi_0^2} \sum_{1 \leq i < j \leq k} \mu_i \mu_j \pi_i \pi_j \\
&= nt \frac{p_M}{\pi_0} \left\{ \sum_{i=1}^k \pi_i (\sigma_i^2 + \mu_i^2) + \frac{(1-p_M)n+1}{\pi_0} \left(\sum_{i=1}^k \mu_i \pi_i \right)^2 \right\}, \\
FI S(t) &= \frac{\operatorname{Var} S(t)}{\mathbb{E} S(t)} = \frac{\sum_{i=1}^k \pi_i (\sigma_i^2 + \mu_i^2) + \frac{(1-p_M)n+1}{\pi_0} \left(\sum_{i=1}^k \mu_i \pi_i \right)^2}{(1-p_M)(\mu_1 \pi_1 + \mu_2 \pi_2 + \dots + \mu_k \pi_k)}. \quad \square
\end{aligned}$$

3.3. The risk process and probabilities of ruin

Consider the following discrete time risk process

$$(3.4) \quad R_u(t) = u + ct - I_{\{N_1(t) > 0\}} \sum_{j_1=1}^{N_1(t)} Z_{1,j_1} - I_{\{N_2(t) > 0\}} \sum_{j_2=1}^{N_2(t)} Z_{2,j_2} - \dots - I_{\{N_k(t) > 0\}} \sum_{j_k=1}^{N_k(t)} Z_{k,j_k},$$

$t = 0, 1, \dots$, satisfying C1–C6. If we consider the claims in a group as one claim, we can see that it is a particular case of the Binomial risk process.²

The r.v. that describes the time of ruin with an initial capital $u \geq 0$ is defined as

$$\tau_u = \min\{t > 0: R_u(t) < 0\}.$$

The probability of ruin with infinite time and initial capital $u \geq 0$ will be denoted by $\Psi(u) = P(\tau_u < \infty)$. The corresponding probability to survive is $\Phi(u) = 1 - \Psi(u)$. Finally, $\Psi(u, t) = P(\tau_u \leq t)$ is for the probability of ruin with finite time $t = 1, 2, \dots$.

²See e.g. [5, 19, 4, 22].

If we assume that $\mathbb{E}Z_{i,j} = \mu_i < \infty$, $i = 1, 2, \dots, k$, and in a long horizon, the expected risk reserve for unit time is positive:

$$\lim_{t \rightarrow \infty} \frac{\mathbb{E}R_u(t)}{t} > 0.$$

The last is equivalent to

$$\begin{aligned} c &> \lim_{t \rightarrow \infty} \frac{\mathbb{E}S(t)}{t}, \\ c &> \frac{np_M(1-p_M)}{\pi_0} (\pi_1\mu_1 + \pi_2\mu_2 + \dots + \pi_k\mu_k). \end{aligned}$$

Note that this condition does not depend on u and it means the incomes at any $t = 1, 2, \dots$ to be bigger than the mean expenditures at that time:

$$\frac{c\pi_0}{np_M(1-p_M)(\pi_1\mu_1 + \pi_2\mu_2 + \dots + \pi_k\mu_k)} > 1.$$

Therefore, the safety loading ρ should be defined as usually as the proportion between the expected risk reserve at time t with zero initial capital, i.e. $\mathbb{E}R_0(t)$, and the expected total claim amount at same moment of time, for any fixed $t = 1, 2, \dots$:

$$\rho = \frac{c\pi_0}{np_M(1-p_M)(\pi_1\mu_1 + \pi_2\mu_2 + \dots + \pi_k\mu_k)} - 1.$$

Thus the above condition is equivalent to the safety loading condition $\rho > 0$. If this condition is not satisfied, the probability of ruin in infinite time would be 1, for any initial capital u .

The proof of the next theorem is analogous to the corresponding one in the Cramer–Lundberg model³ and in particular to those of the Polya–Aepplý risk model⁴.

Theorem 3.1. *Consider the Risk process defined in (3.4) and satisfying conditions C1–C6. Given the Laplace transforms $l_{Z_{i,1}}(s) = \mathbb{E}e^{-sZ_{i,1}}$, of $Z_{i,1}$, $i = 1, 2, \dots, k$, are finite in $-s$,*

1. *The Laplace transform of the risk process is*

$$\mathbb{E}e^{-sR_0(t)} = e^{-g(s)t}, \quad t = 0, 1, 2, \dots,$$

where

$$g(s) = sc - \log \left\{ 1 - p_M + p_M \left[\frac{\pi_0}{1 - [\pi_1 l_{Z_{1,1}}(-s) + \dots + \pi_k l_{Z_{k,1}}(-s)]} \right]^n \right\}.$$

2. *The process $R_0^*(t) = e^{-sR_0(t)+g(s)t}$, $t \geq 0$, is an $A_{R_0(\leq t)} = \sigma\{R_0(s), s \leq t\}$ -martingale.*

3.
$$\Psi(u, t) \leq e^{-su} \sup_{y \in [0, t]} e^{-yg(s)}, \quad t = 1, 2, \dots$$

4.
$$\Psi(u) \leq e^{-su} \sup_{y \geq 0} e^{-yg(s)}.$$

5. *If the Lundberg exponent ε exists, it is a strictly positive solution of the equation*

$$(3.5) \quad g(s) = 0.$$

In that case, $\Psi(u) \leq e^{-\varepsilon u}$.

³See e.g. [1] or [6], p. 10, 11.

⁴[23], Proposition 6.3.

Proof: [1.]

$$\begin{aligned}
\mathbb{E}e^{-sR_0(t)} &= \mathbb{E}e^{-s\{ct - I_{\{N_1(t)>0\}} \sum_{j_1=1}^{N_1(t)} Z_{1,j_1} - \dots - I_{\{N_k(t)>0\}} \sum_{j_k=1}^{N_k(t)} Z_{k,j_k}\}} \\
&= e^{-sct} G_{N_1(t), N_2(t), \dots, N_k(t)}(l_{Z_{1,1}}(-s), l_{Z_{2,1}}(-s), \dots, l_{Z_{k,1}}(-s)) \\
&= e^{-sct} G_{M(t)} \left[\left(\frac{\pi_0}{1 - (\pi_1 l_{Z_{1,1}}(-s) + \pi_2 l_{Z_{2,1}}(-s) + \dots + \pi_k l_{Z_{k,1}}(-s))} \right)^n \right] \\
&= e^{-sct} \left\{ 1 - p_M + p_M \left(\frac{\pi_0}{1 - (\pi_1 l_{Z_{1,1}}(-s) + \pi_2 l_{Z_{2,1}}(-s) + \dots + \pi_k l_{Z_{k,1}}(-s))} \right)^n \right\}^t \\
&= e^{-sct} e^{t \log \left\{ 1 - p_M + p_M \left(\frac{\pi_0}{1 - (\pi_1 l_{Z_{1,1}}(-s) + \pi_2 l_{Z_{2,1}}(-s) + \dots + \pi_k l_{Z_{k,1}}(-s))} \right)^n \right\}} \\
&= e^{-t \left\{ sc - \log \left[1 - p_M + p_M \left(\frac{\pi_0}{1 - (\pi_1 l_{Z_{1,1}}(-s) + \pi_2 l_{Z_{2,1}}(-s) + \dots + \pi_k l_{Z_{k,1}}(-s))} \right)^n \right] \right\}} = e^{-tg(s)}.
\end{aligned}$$

[2.] Consider $t = 0, 1, 2, \dots$ and $y \leq t$. Then, because the process $\{S(t), t = 0, 1, 2, \dots\}$ has independent and time homogeneous additive increments,

$$\begin{aligned}
\mathbb{E}(R_0^*(t) | A_{R_0(\leq y)}) &= \mathbb{E}(e^{-sct + sS(t) + g(s)t} | A_{R_0(\leq y)}) \\
&= \mathbb{E}(e^{-scy + sS(y) + g(s)y - sc(t-y) + s(S(t) - S(y)) + g(s)(t-y)} | A_{R_0(\leq y)}) \\
&= \mathbb{E}(R_0^*(y) e^{-sc(t-y) + s(S(t) - S(y)) + g(s)(t-y)} | A_{R_0(\leq y)}) \\
&= R_0^*(y) \mathbb{E}(e^{-sc(t-y) + sS(t-y) + g(s)(t-y)}) \\
&= R_0^*(y) \mathbb{E}(e^{-sR_0(t-y) + g(s)(t-y)}) \\
&= R_0^*(y) \mathbb{E}(e^{-sR_0(t-y)}) e^{g(s)(t-y)} \\
&= R_0^*(y) e^{-g(s)(t-y)} e^{g(s)(t-y)} = R_0^*(y).
\end{aligned}$$

[3.] Following the traditional approach we start with the definition of R_0^* and use that for $R_0^*(0) = 1$. Because τ_u is a random stopping time, by Doob's martingale stopping theorem, the stopped process $R_0^*(\min(\tau_u, t))$, is again a martingale. Therefore, for any $0 \leq t < \infty$, by the double expectations formula,

$$\begin{aligned}
1 &= R_0^*(0) = \mathbb{E}R_0^*(0) = \mathbb{E}R_0^*(\min(\tau_u, t)) \\
&= \mathbb{E}(R_0^*(\min(\tau_u, t)) | \tau_u \leq t) P(\tau_u \leq t) + \mathbb{E}(R_0^*(\min(\tau_u, t)) | \tau_u > t) P(\tau_u > t) \\
&\geq \mathbb{E}(R_0^*(\min(\tau_u, t)) | \tau_u \leq t) P(\tau_u \leq t) \\
&= \mathbb{E}(e^{-sR_0(\min(\tau_u, t)) + g(s) \min(\tau_u, t)} | \tau_u \leq t) P(\tau_u \leq t) \\
&= e^{su} \mathbb{E}(e^{g(s) \min(\tau_u, t)} | \tau_u \leq t) P(\tau_u \leq t) \\
&\geq e^{su} \mathbb{E}(e^{g(s) \tau_u} | \tau_u \leq t) P(\tau_u \leq t) \\
&= e^{su} \mathbb{E}(e^{g(s) \tau_u} | \tau_u \leq t) \Psi(u, t) \geq e^{su} \inf_{y \in [0, t]} e^{g(s)y} \Psi(u, t).
\end{aligned}$$

[4.] This is an immediate consequence of 3., when $t \rightarrow \infty$.

[5.] This is an immediate consequence of the inequality

$$1 \geq \mathbb{E}(e^{-\varepsilon R_0(\min(\tau_u, t)) + g(\varepsilon) \min(\tau_u, t)} | \tau_u \leq t) P(\tau_u \leq t),$$

applied for $t \rightarrow \infty$ and the fact that $R_0(s) = R_u(s) - u$. □

Remark 3.1. In general, to compute solution of the equation (3.5) is a difficult task and it can be done only numerically, since e.g. for exponential claims it involves roots of algebraic equations of high order. These solutions can be however also negative or/and complex conjugates. To illustrate the complexity of this setup, let us consider $k = 1$ and respective equation (for special choice of parameters) $s = \log(1 - 0.5 + 0.5((1-p)/(1-p/(1-s)))^n)$. Then, real solutions are plotted at Figure 1 for $n = 1, \dots, 10$, $p \in (0, 1)$, therein we can see the complexity of such computations.

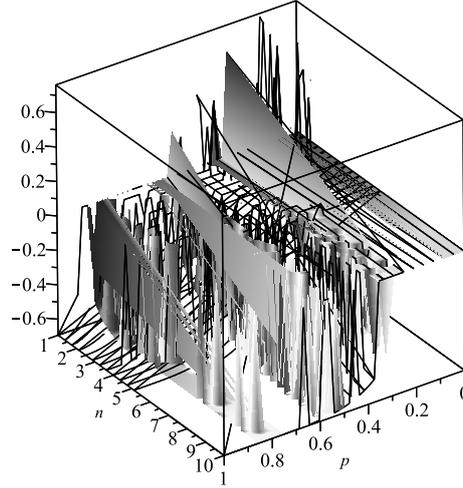


Figure 1: Real solutions of equation (3.5) for $\pi_M = \pi_0 = \pi_1 = 0.5$.

Theorem 3.2. Consider the Risk process defined in (3.4) and satisfying conditions C1–C6. Suppose that it satisfies the net profit condition. Denote by

$$\sigma_S = \sqrt{\text{Var } S(1)} = \sqrt{n \frac{p_M}{\pi_0} \left\{ \sum_{i=1}^k \pi_i (\sigma_i^2 + \mu_i^2) + \frac{(1-p_M)n+1}{\pi_0} \left(\sum_{i=1}^k \mu_i \pi_i \right)^2 \right\}}.$$

Define

$$R_m(t) = \frac{u_m + cmt - S(mt)}{\sigma_S \sqrt{m}},$$

where

$$\frac{u_m}{\sigma_S} \sim u_0 \sqrt{m}, \quad m \rightarrow \infty,$$

and

$$\frac{\rho_m \mu_S}{\sigma_S} \sim \frac{\rho_0}{\sqrt{m}}, \quad m \rightarrow \infty,$$

$$\mu_S = \mathbb{E}S(1) = \frac{np_M(1-p_M)}{\pi_0} (\mu_1\pi_1 + \mu_2\pi_2 + \dots + \mu_k\pi_k),$$

then

$$R_m(t) \Rightarrow u_0 + \rho_0 t + W(p_M t), \quad m \rightarrow \infty.$$

3.4. Simulations of the risk processes and estimation of the probabilities of ruin

In this subsection we provide a brief simulation study on probabilities of ruin in a finite time in the model (3.4). For any of them 10 000 sample paths were created and the relative frequencies of those which goes at least once below zero was determined. The number of groups is $k = 20$. The parameters of the NMn distribution are $n = 40$ and $p = (0.002, 0.004, 0.006, 0.008, 0.01, 0.012, 0.014, 0.016, 0.018, 0.02, 0.022, 0.024, 0.026, 0.028, 0.03, 0.032, 0.034, 0.036, 0.038, 0.04)$. Different parameters on different coordinates allow higher flexibility of the model. The probability of arrival of a group in a fixed time point is $p_M = 0.4$, and premium income rate is $c = 0.1$.

Example 3.1. Exponential claim sizes. For computations of ruin probabilities under exponential claims we consider parameter vector $\lambda = (10, 20, 30, 40, 50, 60, 70, 80, 90, 100, 110, 120, 130, 140, 150, 160, 170, 180, 190, 200)$. The i -th coordinate describe the parameter of the Exponential distribution of the claim sizes within the i -th group. The resulting probabilities for ruin for different initial capitals u and time intervals $[0, t]$ are presented in the Table 1. The corresponding 10 000 sample paths of the risk process are depicted on Figure 2.

Table 1: Probabilities of ruin for exponential claims.

u	t					
	2	5	10	20	50	100
0	0.3224	0.5581	0.7037	0.8053	0.9016	0.9515
1	0.0001	0.0036	0.0370	0.1546	0.4234	0.6586
2	0	0	0.0005	0.0072	0.1131	0.3386
3	0	0	0	0.0001	0.0159	0.1366
4	0	0	0	0	0.0012	0.0364
5	0	0	0	0	0	0.0081

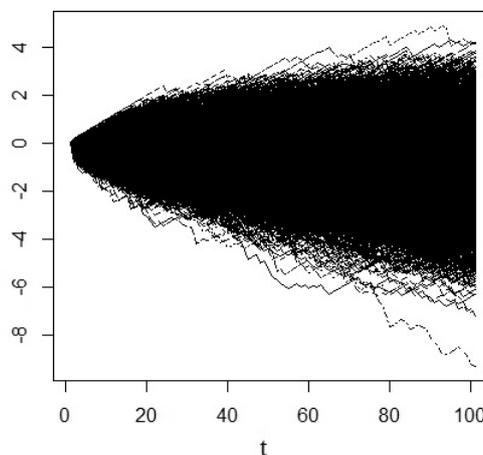


Figure 2: 10 000 sample paths of the risk process (3.4) for Exponential individual claim sizes, $t = 100$.

Example 3.2. Gamma claim sizes. Table 2 presents the probabilities for ruin in case when the claim sizes are Gamma distributed with parameters $\alpha = seq$ ($from = 0.001$, $to = 0.001 + (k - 1) * 0.005$, $by = 0.005$) and $\beta = seq$ ($from = 1$, $to = 1 + (k - 1) * 0.2$, $by = 0.2$), where seq is the function for creating a sequence in R software, see [16]. The corresponding 10 000 sample paths of the risk process are depicted on Figure 3.

Table 2: Probabilities of ruin for gamma claims.

u	t					
	2	5	10	20	50	100
0	0.164	0.294	0.442	0.529	0.706	0.787
1	0.042	0.085	0.183	0.273	0.490	0.578
2	0.017	0.046	0.086	0.160	0.317	0.458
3	0.010	0.024	0.051	0.098	0.202	0.358
4	0.003	0.008	0.019	0.057	0.139	0.248
5	0.000	0.002	0.013	0.027	0.071	0.171

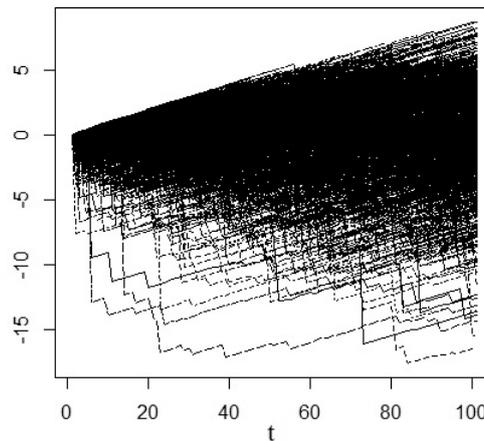


Figure 3: 10 000 sample paths of the risk process (3.4) for Gamma individual claim sizes, $t = 100$.

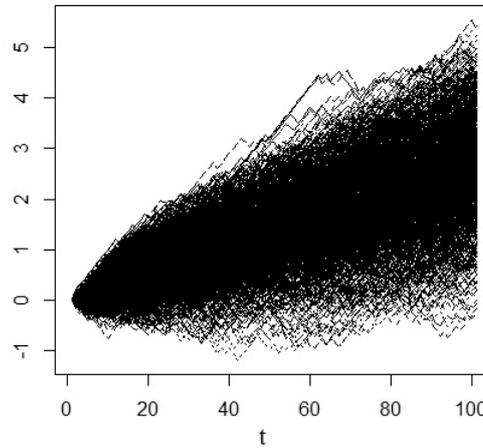
Example 3.3. Uniform claim sizes. The ruin probabilities presented in the Table 3 are calculated under assumption for uniform claim sizes with left and right bounds of the intervals, presented correspondingly via parameter vectors $Umin = seq$ ($from = 0.0001$, $to = 0.0001 + (k - 1) * 0.0001$, $by = 0.0001$) and $Umax = Umin + 0.01$.

The corresponding 10 000 sample paths of the risk process are depicted on Figure 4.

Analogously, the probabilities for ruin in a finite time interval, for different claim sizes with finite variance, and related with the risk process (3.4) can be estimated. The corresponding confidence intervals can be calculated using the Central Limit Theorem, applied to relative frequencies.

Table 3: Probabilities of ruin for uniform claims.

u	t					
	2	5	10	20	50	100
0	0.133	0.253	0.361	0.402	0.392	0.412
1	0.000	0.000	0.000	0.001	0.001	0.000
2	0.000	0.000	0.000	0.000	0.000	0.000
3	0.000	0.000	0.000	0.000	0.000	0.000
4	0.000	0.000	0.000	0.000	0.000	0.000
5	0.000	0.000	0.000	0.000	0.000	0.000

**Figure 4:** 10 000 sample paths of the risk process (3.4) for Uniform individual claim sizes, $t = 100$.

4. CONCLUSIONS

The paper shows that CPSMNn distribution is easy to work with, and it can be very useful for modelling of the number of claims in Risk theory. Recently [26] and [2] have published another important application of multivariate negative binomial distribution in actuarial risk theory. Both models show that they are suitable for capturing the overdispersion phenomena. These distributions provide a flexible modelling of the number of claims that have appeared up to time t . The number of summands of the random sum reflects the number of groups of claims that have occurred up to this moment. The negative multinomial summands and their dependence structure describe types of claims within a group which are different from those given by [26] and [2]. From mathematical point of view our paper describes completely novel presentations of the CPSMNn distributions. Thus we can conclude by following conclusions:

- These distributions are a particular case of Multivariate PSD.
- Considered as a mixture, CPSMNn would be called (possibly Zero-inflated) Mixed NMn with scale changed PSD first parameter. More precisely,

$$I_{\{M>0\}}NMn(nM, \pi_1, \pi_2, \dots, \pi_k) \bigwedge_M PSD(g_{\bar{\alpha}}(x); \theta),$$

where $I_{M>0}$ is a Bernoulli r.v. or indicator of the event " $M > 0$ ".

- CPSMNn is particular case of compounds or random sums $(T_M^{(1)}, T_M^{(2)}, \dots, T_M^{(k)})$, where

$$T_M^{(j)} = I_{\{M>0\}} \sum_{i=1}^M Y_i^{(j)} = \begin{cases} \sum_{i=1}^M Y_i^{(j)} & \text{if } M > 0, \\ 0 & \text{otherwise,} \end{cases} \quad j = 1, 2, \dots, k.$$

These observations allow us to make the first complete characterization of Compound power series distribution with negative multinomial summands and to give an example of their application in modelling the main process in the Insurance risk theory.

ACKNOWLEDGMENTS

The work was supported by project Fondecyt Proyecto Regular No. 1151441, Project LIT-2016-1-SEE-023, and by the bilateral projects Bulgaria–Austria, 2016–2019, “Feasible statistical modeling for extremes in ecology and finance”, BNSF, Contract number 01/8, 23/08/2017 and WTZ Project No. BG 09/2017, <https://pavlinakj.wordpress.com/>.

The authors are very grateful to the Editor, Associated Editor and the Reviewers for their valuable comments.

REFERENCES

- [1] ASMUSSEN, S. (2000). *Ruin Probabilities*, World scientific.
- [2] BADESCU, A.L.; LAN, G.; LIN, X.S. and TANG, D. (2015). Modeling Correlated Frequencies with Application in Operational Risk Management, *Journal of Operational Risk*, **10**(1), 1–43.
- [3] BATES, G.E. and NEYMAN, J. (1952). Contributions to the Theory of Accident Proneness. 1. An Optimistic Model of the Correlation Between Light and Severe Accidents, *California Univ. Berkeley*, **132**, 215–253.
- [4] DICKSON, DAVID C.M. (1994). Some comments on the compound binomial model, *ASTIN Bulletin: The Journal of the IAA*, **24**(1), 33–45.
- [5] GERBER, HANS U. (1988). Mathematical fun with the compound binomial process, *ASTIN Bulletin: The Journal of the IAA*, **18**(2), 161–168.
- [6] GRANDSELL, JAN (1991). *Aspects of Risk Theory*, Springer-Verlag, New York, Berlin, Heidelberg.
- [7] JOHNSON, N.; KOTZ, S. and BALAKRISHNAN, N. (1997). *Discrete Multivariate Distributions*, Wiley, New York.
- [8] JOHNSON, N.; KOTZ, S. and KEMP, A.W. (2005). *Univariate Discrete Distributions*, John Wiley and Sons, New York.
- [9] JORDANOVA, P. and STEHLÍK, M. (2016). Mixed Poisson process with Pareto mixing variable and its risk applications, *Lithuanian Mathematical Journal*, **56**(2), 189–206.
- [10] JOSE, K.K. and JACOB, S. (2016). Type II Bivariate Generalized Power Series Poisson Distribution and its Applications in Risk Analysis, [doi:10.20944/preprints201608.0209.v1](https://doi.org/10.20944/preprints201608.0209.v1).

- [11] KHATRI, C.G. (1959). On certain properties of power-series distributions, *JSTOR*, **46**(3/4), 486–490.
- [12] KOSTADINOVA, K. and MINKOVA, L.D. (2016). Type II family of Bivariate Inflated-parameter Generalized Power Series Distributions, *Serdica Mathematical Journal*, **42**(1), 27–42.
- [13] NOACK, A. (1950). A class of random variables with discrete distributions, *The Annals of Mathematical Statistics*, **21**(1), 127–132.
- [14] PATIL, G.P. (1966). On multivariate generalized power series distribution and its applications to the multinomial and negative multinomial, *Sankhya: The Indian Journal of Statistics, Series A*, **28**(2/3), 225–238.
- [15] PHATAK, A.G. and SREEHARI, M. (1981). Some characterizations of a bivariate geometric distribution, *Journal of Indian Statistical Association*, **19**, 141–146.
- [16] R DEVELOPMENT CORE TEAM (2005). *R: a language and environment for statistical computing*, R Foundation for Statistical Computing, Vienna, Austria, www.R-project.org.
- [17] ROLSKI, T.; SCHMIDT, H. and TEUGELS, J. (1999). *Stochastic Processes for Insurance and Finance*, John Wiley and Sons, New York.
- [18] SIBUYA, M.; YOSHIMURA, I. and SHIMIZU, R. (1964). Negative multinomial distribution, *Annals of the Institute of Statistical Mathematics*, **16**(1), 409–426.
- [19] SHIU, ELIAS S.W. (1989). The probability of eventual ruin in the compound binomial model, *ASTIN Bulletin*, **19**(2), 179–190.
- [20] SMITH, G.E.J. (1965). *Parameter estimation in some multivariate compound distributions*, PhD Thesis, The University of British Columbia, Canada.
- [21] SRIVASTAVA, R.C. and BAGCHI, S.N. (1985). On some characterizations of the univariate and multivariate geometric distributions, *J. Indian Statist. Assoc.*, **23**(1), 27–33.
- [22] LI, SHUANMING and SENDOVA, KRISTINA P. (2013). The finite-time ruin probability under the compound binomial risk model, *European Actuarial Journal*, **3**(1), 249–271.
- [23] MINKOVA, L.D. (2012). *Distributions in Insurance Risk Models*, Doctor of Science Thesis, FMI, Sofia University, Bulgaria.
- [24] TWEEDIE, M.C.K. (1952). The estimation of parameters from sequentially sampled data on a discrete distribution, *Journal of the Royal Statistical Society, Series B (Methodological)*, **14**(2), 238–245.
- [25] WISHART, J. (1949). Cumulants of multivariate multinomial distribution, *Biometrika*, **36**(2), 47–58.
- [26] TANG, D.; BADESCU, A.L.; LIN, X.S. and VALDEZ, E.A. (2015). Multivariate Pascal Mixture Regression Models for Correlated Claim Frequencies, <https://ssrn.com/abstract=2618265>.

CHARACTERIZATION OF THE MAXIMUM PROBABILITY FIXED MARGINALS $r \times c$ CONTINGENCY TABLES

Author: FRANCISCO REQUENA
– Department of Statistics and O.R., University of Granada, Granada, Spain
fcoreque@ugr.es

Received: January 2016

Revised: July 2017

Accepted: October 2017

Abstract:

- In this paper operators $i[j]$ and $[j]k$ are defined, whose effects on an $r \times c$ contingency table X are to subtract 1 from x_{ij} and to add 1 to x_{kj} , respectively, so that the composition $i[j]k$ of the two operators changes the j -th column of the contingency table without altering its total. Also a *loop* is defined as a composition of such operators that leaves unchanged both row and column totals. This is used to characterize the $r \times c$ contingency tables of maximum probability over the fixed marginals reference set (under the hypothesis of row and column independence). Another characterization of such maximum probability tables is given using the concept of associated U tables, a $U = \{u_{ij}\}$ table being defined as a table such that $u_{ij} > 0$, $1 \leq i \leq r$ and $1 \leq j \leq c$, and for a given set of values r_h , $1 \leq h < r$, $u_{h+1,j} = r_h u_{hj}$ for all j . Finally, a necessary and sufficient condition for the uniqueness of a maximum probability table in the fixed marginals reference set is provided.

Key-Words:

- $r \times c$ contingency table; maximum probability $r \times c$ contingency table; network algorithm; Fisher's exact test.

AMS Subject Classification:

- 62H05, 62H17.

1. INTRODUCTION

Let $X = \{x_{ij}\}$ denote an $r \times c$ contingency table, with $x_{ij} \in \mathbb{N}$ the entry in row i and column j , and let R_1, \dots, R_r be the sums of rows, C_1, \dots, C_c the sums of columns and $N = \sum_i R_i = \sum_j C_j$. Given the marginal sums R_i and C_j , $i = 1, \dots, r$, $j = 1, \dots, c$, let

$$\mathcal{F} = \left\{ X \mid \sum_{j=1}^c x_{ij} = R_i, \sum_{i=1}^r x_{ij} = C_j \right\}$$

be the reference set of all possible $r \times c$ tables with the aforementioned marginal sums. Then, under the hypothesis of row and column independence, it is well known that for $X \in \mathcal{F}$,

$$(1.1) \quad P(X) = \frac{\prod_i R_i! \prod_j C_j!}{N! \prod_{ij} x_{ij}!}.$$

A problem that is of interest is that of obtaining a table $X \in \mathcal{F}$ which maximizes (1.1), i.e. a maximum probability fixed marginals $r \times c$ table (MPT). This problem arises, for example, as part of the best known and most efficient algorithm for calculating the p -value of Fisher's exact test in unordered $r \times c$ contingency tables: the network algorithm of Mehta and Patel [2]. The application of this algorithm to an observed $r \times c$ table requires, for many of the nodes in the network, the calculation of the longest subpath from each node to the terminal node, and this involves (many) repeated applications of the calculation of maximum probability $r \times c'$ tables ($c' \leq c$) for given fixed marginal sums.

Methods for obtaining these MPTs have been proposed by Mehta and Patel [2] and by Joe [1]. The most general is that of Joe, which is based on a necessary condition for the MPTs, and generally involves the (recursive) construction of a subset of \mathcal{F} in which the MPTs are contained, and obtaining these by inspecting the probabilities of the tables of this subset. However, the computation time for the Joe method grows exponentially when r or c increase, and it is practically unviable for relatively large values of r and c .

In the particular case of $2 \times c$ tables, Requena and Martín [3] present a necessary and sufficient condition for the MPTs. Based on this characterization, Requena and Martín [4] propose a general and very efficient method for obtaining the MPTs, and Requena and Martín [5] present some modifications in the network algorithm of Mehta and Patel for $2 \times c$ tables, which produce a drastic reduction in computation time.

In order to obtain general and more efficient methods for obtaining the MPTs, in the general case of $r \times c$ tables, it is important that these methods are based on necessary and sufficient conditions for the MPTs. In this sense, in this paper, two necessary and sufficient conditions are presented in order to characterize the MPTs. However, this characterization is not a generalization of the one previously shown in Requena and Martín [3]; it is completely different, although logically in the particular case of $2 \times c$ tables, the characterization presented in this paper is equivalent to that of Requena and Martín [3].

In Section 2 of this paper, we define and study the concepts of sequence and loop which we will use in the characterization of the MPTs, which is presented in Sections 3 and 5.

In Section 3 we present the characterization as a more theoretical result, while in Section 5, with a more applied purpose, the characterization is presented in terms of a particular type of tables (U tables), which we define and study in Section 4. Finally, in Section 6 we provide a necessary and sufficient condition of the uniqueness of the MPT.

2. SEQUENCES AND LOOPS

The characterization of the MPTs which we set out in the following sections is based on the concepts of *sequence* and *loop*. In order to define these concepts, we will start by defining some operators, which are applied to an $r \times c$ table $X = \{x_{ij}\}$.

We define the operator $i[j]$ whose effect on X is to subtract 1 from x_{ij} leaving all the other entries unchanged, and the operator $[j]k$ whose effect on X is to add 1 to x_{kj} leaving all the other entries unchanged. Based on these operators, we define the operator $i[j]k$ as the composition of $i[j]$ with $[j]k$ ($i[j] \circ [j]k = i[j]k = [j]k \circ i[j]$). It is clear that $i[j]k$ changes the j -th column of the table without altering its sum. Also, as $i[j]i$ is the identity operator, $i[j]$ and $[j]i$ are *inverse* of each other.

Definition 2.1. Given an $r \times c$ table $X = \{x_{ij}\}$, and given the rows i_0, i_1, \dots, i_k (with $i_{h-1} \neq i_h$) and columns j_1, \dots, j_k (not all equal), a *sequence* is the composition of $i_0[j_1]i_1$ with $i_1[j_2]i_2, \dots$ with $i_{k-1}[j_k]i_k$, which for simplicity we denote by $i_0[j_1]i_1[j_2]i_2 \cdots i_{k-1}[j_k]i_k$ ($1 \leq i_h \leq r$ and $1 \leq j_h \leq c$).

Definition 2.2. Given an $r \times c$ table $X = \{x_{ij}\}$, a *loop* is a sequence in which $i_k = i_0$, i.e. $i_0[j_1]i_1[j_2]i_2 \cdots i_{k-1}[j_k]i_0$.

From this point onward in the text, when we write a sequence as

$$i_0[\cdot]i_1 \cdots i_{k-1}[\cdot]i_k$$

it will be understood that it is a sequence for an unspecified set of columns j_1, \dots, j_k .

In terms of the effect of applying a sequence or a loop to a table X , we can understand a sequence or a loop as a succession of operators $i_{h-1}[j_h]$ and $[j_h]i_h$ (or as a succession of operators $i_{h-1}[j_h]i_h$), $h = 1, 2, \dots, k$, applied in a successive manner: each operator is applied to the table obtained by applying the previous one (the first one is applied to X). For example, applying the sequence $1[2]2[3]4$ to a 4×4 table has the effect of adding 1 in x_{43} , subtracting 1 in x_{23} , adding 1 in x_{22} and subtracting 1 in x_{12} . A sequence applied to X does not alter the column sums, but it alters the i_0 -th and the i_k -th row sum. However a loop does not alter neither the column sums nor the row sums.

Logically, if one removes pairs of inverse operators from a sequence (or loop) in an appropriate way, one would obtain a new and more reduced sequence (or loop), but one which would have the same effect on X as the previous one. In this sense, we give the following definition:

Definition 2.3. Given an $r \times c$ table $X = \{x_{ij}\}$, two sequences (or two loops) are *equivalent* when they have the same effect on X .

Thus, we have classes of equivalent sequences (or loops). Within a same class, the difference between two sequences (or two loops) is a set of pairs of inverse operators.

In the same way, we will define the equivalence between a sequence (or loop) and a group of several sequences (or loops), based on the understanding that the sequences (or loops) which compose the group are applied in a successive manner: each sequence (or loop) is applied to the table obtained by applying the previous one.

Because the effect of an operator $[j]i$ on X is to add 1 to x_{ij} , and the effect of an operator $i[j]$ is to subtract 1 from x_{ij} , and denoting the number of operators $[j]i$ and $i[j]$ in the loop by n_{ij} and n'_{ij} , respectively, any loop can be represented by means of a table $D = \{d_{ij}\}$, defined as $d_{ij} = n_{ij} - n'_{ij}$, $i = 1, \dots, r$ and $j = 1, \dots, c$. It is easy to see that a table D defined thus has all its marginal sums equal to 0. Reciprocally, any table $D = \{d_{ij}\}$, with d_{ij} being integer numbers and marginal sums equal to 0, will represent a loop or a group of loops. Moreover, applying a loop to a table X is equivalent to adding the corresponding table D to it, thereby obtaining a new table X' with entries $x'_{ij} = x_{ij} + d_{ij}$, and with the same marginal sums as X . But X' is not necessarily an $r \times c$ table, because some of the entries x'_{ij} could be negative. If this happens (although we can consider such a loop) we would not consider that table X' . This is taken into account in Section 3.

Example 2.1. Let us consider the loop $2[1]3[4]1[1]3[2]4[3]2$. Applying this loop to a 4×4 table X , we will obtain a new 4×4 table X' . Let us see it for some i 's and j 's. For $i = 2$ and $j = 1$, because there is only one operator $2[1]$ ($n'_{21} = 1$) and no operator $[1]2$ ($n_{21} = 0$), $d_{21} = 0 - 1 = -1$ and we have to subtract 1 from x_{21} ($x'_{21} = x_{21} - 1$). Likewise, for $i = 3$ and $j = 1$ there are two operators $[1]3$ ($n_{31} = 2$) and no operator $3[1]$ ($n'_{31} = 0$), therefore $d_{31} = 2 - 0 = 2$ and we have to add 2 to x_{31} ($x'_{31} = x_{31} + 2$). In a similar way for the other i 's and j 's. The complete table D that represent this loop is

-1	0	0	1
-1	0	1	0
2	-1	0	-1
1	1	-1	0

and adding this table D to the table X we obtain the table X' .

If in a sequence (or loop) B we invert the order of the i_h 's and also of the j_h 's (each operator would be substituted by its inverse one), we will obtain a new sequence (or loop): we will call it *inverse* of B . For example, the sequence inverse of $1[2]2[3]4$ is $4[3]2[2]1$. Furthermore, if $D = \{d_{ij}\}$ represents a loop B , then $-D = \{-d_{ij}\}$ will represent the inverse of B .

Let us now define a particular type of loop which we will use in the characterization of the MPTs.

Definition 2.4. For $1 < k \leq \min(r, c)$, an *order k simple loop*, is a loop $i_0[j_1]i_1[j_2] \cdots i_{k-1}[j_k]i_0$ in which all the k rows i_0, i_1, \dots, i_{k-1} are different and all the k columns j_1, \dots, j_k are different. We will call these the k rows and the k columns of the loop.

Observe that such simple loop leaves $r-k$ rows and $c-k$ columns of X unchanged.

In order to distinguish them from the general case, we will write the tables D which represent order k simple loops as $E = \{e_{ij}\}$. In an order k simple loop, since all of its rows i_h (and all of its columns j_h) are different, both n_{ij} and n'_{ij} can only be equal to 1 or 0, and $n_{ij} + n'_{ij} \leq 1$. Therefore, the corresponding table E will have all of its entries e_{ij} equal to 0, except for a 1 and a -1 in each of the k rows and in each of the k columns of the loop. Moreover, any table E of this type will represent an order k simple loop. For example, the table

$$\begin{array}{|c|c|c|c|c|} \hline 0 & -1 & 0 & 1 & 0 \\ \hline -1 & 1 & 0 & 0 & 0 \\ \hline 1 & 0 & 0 & -1 & 0 \\ \hline 0 & 0 & 0 & 0 & 0 \\ \hline \end{array}$$

represents the order 3 simple loop $2[1]3[4]1[2]2$.

It is obvious that if one subtracts from a table D (different to any table E) a table E whose $e_{ij} \neq 0$ have the same sign as the corresponding d_{ij} in D , this will result in another type D table (or type E). Therefore, it is easy to deduce that any table D is the sum of several type E tables, i.e. any loop (represented by D) can be broken down into a group of simple loops, which together are equivalent to D . For example:

$$\begin{array}{|c|c|c|c|} \hline 0 & 0 & 1 & -1 \\ \hline -1 & 3 & -2 & 0 \\ \hline 0 & -1 & 0 & 1 \\ \hline 1 & -2 & 1 & 0 \\ \hline \end{array} = \begin{array}{|c|c|c|c|} \hline 0 & 0 & 0 & 0 \\ \hline -1 & 1 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline 1 & -1 & 0 & 0 \\ \hline \end{array} + \begin{array}{|c|c|c|c|} \hline 0 & 0 & 0 & 0 \\ \hline 0 & 1 & -1 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline 0 & -1 & 1 & 0 \\ \hline \end{array} + \begin{array}{|c|c|c|c|} \hline 0 & 0 & 1 & -1 \\ \hline 0 & 1 & -1 & 0 \\ \hline 0 & -1 & 0 & 1 \\ \hline 0 & 0 & 0 & 0 \\ \hline \end{array}$$

The loop represented by the table D on the left-hand side is broken down into the (equivalent) group of three simple loops on the right-hand side (the first and the second are order 2 and the third order 3).

Finally, for any $r \times c$ table $X = \{x_{ij}\}$ and for any loop, from this point onward we will use expressions of the type

$$(2.1) \quad Q = \prod_{j \in J} \frac{x_{bj} + 1}{x_{aj}},$$

where J represents the set of columns j 's (which are not necessarily all different) corresponding to the $[j]$'s of the operators $a[j]b$ in the loop, and a and b are the rows of these operators. In this type of expression, a one-to-one relation between the terms of the product and the set of operators $a[j]b$ of the loop is established. For example, for the loop $2[1]4[5]1[3]2$

$$Q = \frac{x_{41} + 1}{x_{21}} \cdot \frac{x_{15} + 1}{x_{45}} \cdot \frac{x_{23} + 1}{x_{13}}.$$

3. CHARACTERIZATION OF THE MAXIMUM PROBABILITY $r \times c$ TABLES

The simple loops defined in the previous section are used in the following result to characterize the MPTs.

Theorem 3.1. *The necessary and sufficient condition for $X = \{x_{ij}\} \in \mathcal{F}$ to be an MPT is that*

$$(3.1) \quad \prod_{j \in J} \frac{x_{bj} + 1}{x_{aj}} \geq 1$$

for every order k simple loop $E = \{e_{ij}\}$ and every k , $1 < k \leq \min(r, c)$, where J is the set of the k columns of the loop, and for each $j \in J$, b and a are the rows such that $e_{bj} = 1$ and $e_{aj} = -1$.

Proof: Let X be an MPT, and let us consider $X' = X + E$, for any order k simple loop E , $1 < k \leq \min(r, c)$. All of the elements of X and X' will be identical, except $x'_{bj} = x_{bj} + 1$ and $x'_{aj} = x_{aj} - 1$ for $j \in J$, and a , b and J previously defined. Firstly, if E is an order k simple loop such that x'_{aj} is a negative integer (for some a and j), that is, X' is not an $r \times c$ table, then $x_{aj} = 0$ and, hence, the condition (3.1) is fulfilled for that E . Secondly, if (on the contrary) E is such that X' is an $r \times c$ table ($X' \in \mathcal{F}$), then from expression (1.1), and because $P(X) \geq P(X')$, we obtain

$$\frac{P(X)}{P(X')} = \prod_{j \in J} \frac{(x_{bj} + 1)! (x_{aj} - 1)!}{x_{bj}! x_{aj}!} = \prod_{j \in J} \frac{x_{bj} + 1}{x_{aj}} \geq 1.$$

Therefore, for an MPT, (3.1) is fulfilled for all E of order k .

In order to prove the sufficient condition one must note, in the first place, that if an $r \times c$ table fulfils (3.1) for every simple loop, it will fulfil said expression for an order k simple loop E , and also for the inverse loop $-E$ (which is also an order k simple loop). For this reason

$$(3.2) \quad \prod_{j \in J} \frac{x_{aj} + 1}{x_{bj}} \geq 1$$

will also be fulfilled with a , b and J defined for E as in the formulation of the theorem. Thus, for each E , (3.1) and (3.2) will be fulfilled. The proof of the sufficient condition in the case that there is only one $r \times c$ table of \mathcal{F} fulfilling (3.1) is trivial. Therefore, we will assume that there is more than one. Let $X \in \mathcal{F}$ be an MPT which will obviously fulfil (3.1). It will be necessary to prove that for any $X' \in \mathcal{F}$ satisfying (3.1) for all order k simple loops, $P(X') = P(X)$ must be fulfilled.

It is clear that X' can always be written as $X' = X + D$, when D is a table representing a loop (or group of loops), which can be broken down into a group of tables E 's (simple loops). According to this type of decomposition (as we have seen in the previous section), for any of these E 's, considering $j \in J$, a and b defined as before, the signs of e_{aj} and e_{bj} should be

the same as those of their corresponding d_{aj} and d_{bj} in table D . Moreover, because X' fulfils (3.1) and (3.2) for any of these E 's, we can write

$$(3.3) \quad \prod_{j \in J} \frac{x'_{bj}}{x'_{aj} + 1} = \prod_{j \in J} \frac{x_{bj} + d_{bj}}{x_{aj} + d_{aj} + 1} \leq 1$$

for each E , and because X also fulfils (3.1), we have that

$$(3.4) \quad \prod_{j \in J} \frac{x_{bj} + 1}{x_{aj}} \geq 1$$

for each E . From this expression, and because the d_{bj} 's are positive and the d_{aj} 's are negative,

$$(3.5) \quad Q' = \prod_{j \in J} \frac{x_{bj} + d_{bj}}{x_{aj} + d_{aj} + 1} \geq 1$$

will be fulfilled. Moreover, if any $d_{bj} > 1$ or any $|d_{aj}| > 1$ we will obtain $Q' > 1$, which would contradict (3.3). Hence $d_{bj} = 1$ and $d_{aj} = -1$, and from (3.3) and (3.5) we obtain

$$(3.6) \quad \prod_{j \in J} \frac{x_{bj} + 1}{x_{aj}} = 1.$$

Since the above is valid for any of the E 's in which D is broken down, on the one hand we will obtain $|d_{hl}| \leq 1$ for all h and l , hence for every $d_{hl} \neq 0$ there will be one and only one of the loops E 's such that $e_{hl} = d_{hl}$. On the other hand, considering the expression (3.6) for all the E 's in which D has been broken down, we will obtain

$$(3.7) \quad \frac{\prod_{hl \in D+} (x_{hl} + 1)}{\prod_{hl \in D-} x_{hl}} = 1$$

where $D+$ and $D-$ are the sets of subindices hl such that $d_{hl} = 1$ and $d_{hl} = -1$, respectively. Finally, from (1.1)

$$\frac{P(X)}{P(X')} = \frac{\prod_{hl \in D+} (x_{hl} + 1)! \prod_{hl \in D-} (x_{hl} - 1)!}{\prod_{hl \in D+} x_{hl}! \prod_{hl \in D-} x_{hl}!} = \frac{\prod_{hl \in D+} (x_{hl} + 1)}{\prod_{hl \in D-} x_{hl}}$$

is obtained, and from (3.7) we will obtain $P(X') = P(X)$. □

From this theorem, and from what has been said in the proof, the two following results are easily deduced:

Theorem 3.2. *If X is an MPT and E an order k simple loop for which*

$$(3.8) \quad \prod_{j \in J} \frac{x_{bj} + 1}{x_{aj}} = 1$$

holds, where J , a and b are defined as in Theorem 3.1, then $X' = X + E$ is also an MPT.

Theorem 3.3. *If two tables, X and X' , belonging to \mathcal{F} are MPTs, then the difference between both tables is one or several simple loops, such that (3.8) holds for X and for each of these simple loops. Moreover the following always holds*

$$|x'_{hl} - x_{hl}| \leq 1, \quad \forall h, l.$$

Finally, the following result extends expression (3.1) to any loop.

Theorem 3.4. *Given the expression Q defined in (2.1), if X is an MPT, for any loop the following always holds*

$$(3.9) \quad Q = \prod_{j \in J} \frac{x_{bj} + 1}{x_{aj}} \geq 1$$

where J is the set of columns j 's corresponding to the $[j]$'s of the loop, and $a[j]b$ are the operators that compose the loop.

Proof: From Theorem 3.1, for simple loops it is obvious that (3.9) is fulfilled. In the case of non-simple loops, if the loop is represented by a table D , it can be decomposed into a set of n simple loops. Representing the expression (2.1) for the simple loop h ($1 \leq h \leq n$) by Q_h , we will obtain that

$$Q = \prod_{h=1}^n Q_h$$

and because $Q_h \geq 1$ for all h (from Theorem 3.1), we obtain $Q \geq 1$. Finally, if the loop is not represented explicitly by any type D table, there will always be an equivalent loop represented by a table D , and the difference between both loops will only be a set of pairs of inverse operators. Without loss of generality, let us suppose that the difference is the pair $a[b]$, $[b]a$. Then, according to what was said when defining expression (2.1), and decomposing (as before) the loop D into n simple loops, we will obtain

$$Q = \frac{x_{ab} + 1}{x_{ab}} \prod_{h=1}^n Q_h > 1.$$

Therefore, to sum up, (3.9) is fulfilled for every loop. □

4. A PARTICULAR TYPE OF TABLES: THE U TABLES

We will now proceed to define and study a type of tables (U tables) that is particularly important in a new characterization of the MPTs.

Definition 4.1. A U table is a table $\{u_{ij}\}$ with r rows and c columns ($1 \leq i \leq r$ and $1 \leq j \leq c$), in which u_{ij} are strictly positive real values ($u_{ij} > 0$) and such that, for a given set of values r_h , $1 \leq h < r$, $u_{h+1,j} = r_h u_{hj}$ for all j .

Given this definition, from this point onward $r_h = u_{h+1,j}/u_{hj}$ will represent the ratio between the consecutive rows h and $h+1$ of the U table. On the other hand, it is obvious that for any two rows h and i , the ratio $r_{hi} = u_{ij}/u_{hj}$ is constant for all j , and $r_{ih} = 1/r_{hi}$. In particular, $r_{h,h+1} = r_h$. Moreover, for $h < i$, r_{hi} coincides with the product of the ratios between consecutive rows from row h to row i , i.e., $r_{hi} = r_h r_{h+1} \cdots r_{i-1}$. So we will also denote this product by r_{hi} . For example, $r_{14} = r_1 r_2 r_3$. Furthermore, it will always be understood that $r_{hh} = 1$.

Let us consider some properties of this type of table, the proofs for which are very straightforward.

Property 4.1. If any row or column of a U table is multiplied by a constant, or the rows (or columns) of a U table are interchanged, another U table is obtained.

Property 4.2. For any two rows h and i of a U table, $r_{hi} = 1/r_{ih}$ is always fulfilled. Moreover, given the rows h , s and i ($h \leq s \leq i$) of a U table, $r_{hi} = r_{hs} r_{si}$ will always hold.

Property 4.3. In a U table $\{u_{ij}\}$ the following always holds:

$$\prod_{j \in J} \frac{u_{bj}}{u_{aj}} = 1$$

for every order k simple loop $E = \{e_{ij}\}$, and every k , $1 < k \leq \min(r, c)$, where J is the set of the k columns of the loop, and for each $j \in J$, b and a are the rows such that $e_{bj} = 1$ and $e_{aj} = -1$.

The following are two examples of U tables.

Example 4.1. A table $\{u_{ij}\}$ in which all the elements in each row are equal (that is, $u_{ij} = A_i > 0$ for all j) is a U table.

Example 4.2. Given an $r \times c$ table, with marginal sums $\{R_i\}$ and $\{C_j\}$, the table of expected frequencies $\{E_{ij}\}$, defined as $E_{ij} = R_i C_j / N$, is a U table. In this case, the ratios between the rows are $r_{hi} = R_i / R_h$. It would also be a U table if R_i and C_j were strictly positive real values.

The following definition establishes a link between the U tables and the $r \times c$ tables.

Definition 4.2. We say that a U table $\{u_{ij}\}$ is associated with an $r \times c$ table $X = \{x_{ij}\}$ if the following holds

$$(4.1) \quad 0 \leq u_{ij} - x_{ij} \leq 1, \quad \forall i, j, \quad 1 \leq i \leq r, \quad 1 \leq j \leq c.$$

From this definition and from the definition of U tables, it is easy to deduce that the U table associated with an $r \times c$ table, if it exists, is not necessarily unique (and generally it is not so).

Given an $r \times c$ table $X = \{x_{ij}\}$, is there always a U table $\{u_{ij}\}$ associated with it? In order for such a U table to exist, $u_{ij} = x_{ij} + \varepsilon_{ij}$ would have to be fulfilled for all i and j , with $0 \leq \varepsilon_{ij} \leq 1$. Now, because the ratios between the rows in the U table would be

$$r_{hi} = u_{ij}/u_{hj} = (x_{ij} + \varepsilon_{ij})/(x_{hj} + \varepsilon_{hj}), \quad \forall j,$$

and so, for a given j , the minimum and the maximum value for r_{hi} would be $x_{ij}/(x_{hj} + 1)$ and $(x_{ij} + 1)/x_{hj}$, respectively, then, r_{hi} should fulfil

$$m_{hi}^o \leq r_{hi} \leq M_{hi}^o,$$

where

$$m_{hi}^o = \max_j \{x_{ij}/(x_{hj} + 1)\} \quad \text{and} \quad M_{hi}^o = \min_j \{(x_{ij} + 1)/x_{hj}\}.$$

In the particular case of consecutive rows (i.e., $i = h + 1$), the limits for the ratios r_h would be

$$(4.2) \quad m_{h,h+1}^o \leq r_h \leq M_{h,h+1}^o, \quad 1 \leq h < r.$$

Moreover, because $r_{hi} = r_h r_{h+1} \cdots r_{i-1}$, the limits for the products of ratios r_{hi} should likewise be

$$(4.3) \quad m_{hi}^o \leq r_{hi} \leq M_{hi}^o, \quad 1 \leq h < i - 1 < r.$$

Therefore, in principle, in order for the said U table to exist, there must be a set of ratios r_h that fulfil (4.2) and whose products r_{hi} fulfil (4.3).

Remark 4.1. If $X = \{x_{ij}\}$ is an MPT, applying expression (3.1) to all order 2 simple loops, we obtain $x_{ij'}/(x_{hj'} + 1) \leq (x_{ij} + 1)/x_{hj}$ for all h, i, j and j' . Hence the following will always be fulfilled

$$m_{hi}^o \leq M_{hi}^o, \quad \forall h, i, \quad 1 \leq h < i \leq r.$$

Remark 4.2. For any two rows h and i , and from the definition of the limits m_{hi}^o and M_{hi}^o , we easily obtain that

$$M_{hi}^o = 1/m_{ih}^o.$$

We can use this expression to obtain m_{pq}^o and M_{pq}^o for $p > q$.

We will call these limits $m_{h,h+1}^o$, $M_{h,h+1}^o$, m_{hi}^o and M_{hi}^o (for each ratio r_h and each product r_{hi}) *initial limits* and, in general, we will refer to them (without specifying the subindices) as limits m^o 's and limits M^o 's.

Example 4.3. In order for there to be a U table associated with the 3×3 table

16	10	6
11	7	5
5	2	2

there must be a set of ratios, r_1 and r_2 , that fulfils (4.2) and whose product $r_{13} = r_1 r_2$ fulfils (4.3). In this case, the initial limits are: $0.714 \leq r_1 \leq 0.750$, $0.417 \leq r_2 \leq 0.429$ and

$0.294 \leq r_{13} \leq 0.300$. In principle, we can take appropriate values of r_1 and r_2 in order to construct an associated U table.

If there are appropriate values of r_h such that (4.2) and (4.3) are fulfilled, and considering the initial limits m^o 's and M^o 's as the current limits for r_h and r_{hi} , they can be redefined (in the sense that we will show below) given the limits of the other products and ratios, thus obtaining new and more accurate limits for r_h and r_{hi} . In general we will denote these new limits as m_{hi} and M_{hi} .

For example, in $3 \times c$ tables, for the product r_{13} , because $r_{13} = r_1 r_2$, the restriction $m_{12}^o m_{23}^o \leq r_{13} \leq M_{12}^o M_{23}^o$ must also be fulfilled, which means r_{13} should fulfil $m_{13} \leq r_{13} \leq M_{13}$, and the new limits will be

$$m_{13} = \max\{m_{13}^o, m_{12}^o m_{23}^o\} \quad \text{and} \quad M_{13} = \min\{M_{13}^o, M_{12}^o M_{23}^o\}.$$

Likewise, for the ratio r_1 , because $r_1 = r_{13}/r_2$, the restriction $m_{13}^o m_{32}^o \leq r_1 \leq M_{13}^o M_{32}^o$ must also be fulfilled, and the new limits for r_1 will be

$$m_{12} = \max\{m_{12}^o, m_{13}^o m_{32}^o\} \quad \text{and} \quad M_{12} = \min\{M_{12}^o, M_{13}^o M_{32}^o\}.$$

In a similar way, the new limits for r_2 are

$$m_{23} = \max\{m_{23}^o, m_{21}^o m_{13}^o\} \quad \text{and} \quad M_{23} = \min\{M_{23}^o, M_{21}^o M_{13}^o\}.$$

Example 4.3 revisited. Starting from the previously calculated initial limits in Example 4.3, we calculate the new limits at the second stage as indicated in the previous paragraph, and we obtain

$$0.298 \leq r_{13} \leq 0.300, \quad 0.714 \leq r_1 \leq 0.720 \quad \text{and} \quad 0.417 \leq r_2 \leq 0.420.$$

In $4 \times c$ tables, for the product r_{13} , because from Property 4.2 $r_{13} = r_1 r_2$ and $r_{13} = r_{14}/r_3$, the restrictions $m_{12}^o m_{23}^o \leq r_{13} \leq M_{12}^o M_{23}^o$ and $m_{14}^o m_{43}^o \leq r_{13} \leq M_{14}^o M_{43}^o$ must also be fulfilled, which means r_{13} has to fulfil $m_{13} \leq r_{13} \leq M_{13}$, and the new limits will be:

$$m_{13} = \max\{m_{13}^o, m_{12}^o m_{23}^o, m_{14}^o m_{43}^o\}$$

and

$$M_{13} = \min\{M_{13}^o, M_{12}^o M_{23}^o, M_{14}^o M_{43}^o\}.$$

Likewise, for the ratio r_2 , because $r_2 = r_{13}/r_1$, $r_2 = r_{24}/r_3$ and $r_2 = r_{14}/(r_1 r_3)$, the following restrictions must be fulfilled:

$$\begin{aligned} m_{21}^o m_{13}^o &\leq r_2 \leq M_{21}^o M_{13}^o, \\ m_{24}^o m_{43}^o &\leq r_2 \leq M_{24}^o M_{43}^o, \\ m_{21}^o m_{14}^o m_{43}^o &\leq r_2 \leq M_{21}^o M_{14}^o M_{43}^o. \end{aligned}$$

Thus r_2 must fulfil that $m_{23} \leq r_2 \leq M_{23}$, and the new limits will be:

$$\begin{aligned} m_{23} &= \max\{m_{23}^o, m_{21}^o m_{13}^o, m_{24}^o m_{43}^o, m_{21}^o m_{14}^o m_{43}^o\}, \\ M_{23} &= \min\{M_{23}^o, M_{21}^o M_{13}^o, M_{24}^o M_{43}^o, M_{21}^o M_{14}^o M_{43}^o\}. \end{aligned}$$

In a similar way for r_1 , r_3 , r_{14} and r_{24} .

Remark 4.3. It is evident that the new limits will fulfil $m_{hi}^o \leq m_{hi}$ and $M_{hi} \leq M_{hi}^o$, and if $m_{hi} \leq M_{hi}$, the new intervals (m_{hi}, M_{hi}) will be contained in the corresponding initial (current) intervals (m_{hi}^o, M_{hi}^o) , both for the ratios r_h and for the products r_{hi} .

Remark 4.4. For any two rows h and i , from Remark 4.2 and from the definition of the new limits, we easily obtain that $M_{hi} = 1/m_{ih}$.

Now, taking the limits m_{hi} and M_{hi} as the current limits for the ratios and products, we can recalculate the limits in the same sense as before, obtaining new limits (for the ratios and products) which we will also denote as m_{hi} and M_{hi} . Thus we will have a recursive process, where, at each stage, the newly calculated limits will have the same property as the current limits. At each stage, we always obtain the new limits m_{hi} and M_{hi} for $h < i$, and we can use Remark 4.4 for $h > i$. In general, and at any stage of the process, we will refer to these limits (without specifying the subindices) as limits m 's and limits M 's.

In this process, because from Property 4.2, $r_{hi} = r_{hs}r_{si}$, $1 \leq h < s < i \leq r$, and $r_{hi} = r_{h'i'}/(r_{h'h}r_{i'i'})$, $1 \leq h' \leq h < i \leq i' \leq r$, it is easy to see that the general expressions of the new limits, m_{hi} and M_{hi} , for r_{hi} ($1 \leq h < i \leq r$) in terms of the current limits can be written as:

$$(4.4) \quad m_{hi} = \max_{i',h',s} \left\{ m_{hs}m_{si}, h < s < i; m_{hh'}m_{h'i'}m_{i'i}, 1 \leq h' \leq h < i \leq i' \leq r \right\},$$

$$(4.5) \quad M_{hi} = \min_{i',h',s} \left\{ M_{hs}M_{si}, h < s < i; M_{hh'}M_{h'i'}M_{i'i}, 1 \leq h' \leq h < i \leq i' \leq r \right\},$$

where the terms on the right-hand side of the expressions correspond to the current limits of the ratios and products (these will coincide with the initial limits m^o 's and M^o 's in the first stage of the process), and where we understand that $m_{qq} = M_{qq} = 1$.

In particular, taking $i = h + 1$ in (4.4) and (4.5) we will obtain the limits for the ratios r_h :

$$(4.6) \quad m_{h,h+1} = \max_{i',h'} \left\{ m_{hh'}m_{h'i'}m_{i',h+1}, 1 \leq h' \leq h < i' \leq r \right\},$$

$$(4.7) \quad M_{h,h+1} = \min_{i',h'} \left\{ M_{hh'}M_{h'i'}M_{i',h+1}, 1 \leq h' \leq h < i' \leq r \right\}.$$

If all the intervals (m_{hi}, M_{hi}) are not empty ($m_{hi} \leq M_{hi}$) (this we will see in Section 5), and because the new intervals (m_{hi}, M_{hi}) are contained in the corresponding current intervals, the process will converge and we will be able to obtain *final limits* for the ratios and products, and we will continue to represent these by m_{hi} and M_{hi} .

Example 4.3 revisited. For the 3×3 table of Example 4.3, given the second stage limits, the new limits obtained at the third stage are the same limits as at the second stage. Therefore, the final limits are

$$0.298 \leq r_{13} \leq 0.300, \quad 0.714 \leq r_1 \leq 0.720 \quad \text{and} \quad 0.417 \leq r_2 \leq 0.420.$$

Once the final limits have been obtained, we can answer the question posed previously more precisely. Given an $r \times c$ table X , in order for there to be a U table associated with it, there

must be a set of ratios r_h that fulfils (4.2) and whose products r_{hi} fulfil (4.3), but taking (in these expressions) the final limits m_{hi} and M_{hi} instead of the initial ones. In greater detail, and taking r_h successively, there must be: first, a value r_1 such that $m_{12} \leq r_1 \leq M_{12}$; second, a value r_2 such that $m_{23} \leq r_2 \leq M_{23}$ and with the product $r_1 r_2 = r_{13}$ such that $m_{13} \leq r_1 r_2 \leq M_{13}$, i.e. a value r_2 such that

$$\max\{m_{23}, m_{13}/r_1\} \leq r_2 \leq \min\{M_{23}, M_{13}/r_1\},$$

and so on. Moreover, the associated U table $\{u_{ij}\}$ would be of the form: $u_{1j} = x_{1j} + \varepsilon_{1j}$ ($0 \leq \varepsilon_{1j} \leq 1$) and $u_{ij} = u_{1j} r_{1i}$, $1 < i \leq r$, $1 \leq j \leq c$.

We can express all this in general form by saying that, given an $r \times c$ table X , in order for there to be a U table associated with X , it must be possible to take successively a set of ratios r_h , $h = 1, 2, \dots, r - 1$, such that

$$(4.8) \quad \max_{1 \leq s \leq h} \{m_{s,h+1}/r_{sh}\} \leq r_h \leq \min_{1 \leq s \leq h} \{M_{s,h+1}/r_{sh}\}$$

(in which $r_{sh} = r_s r_{s+1} \cdots r_{h-1}$ and we understand that $r_{hh} = 1$) and a set of ε_{1j} ($1 \leq j \leq c$) (for the first row of the U table) such that (4.1) is fulfilled.

Further on in this paper, we will see that an associated U table exists for the MPTs, and only for these.

Example 4.3 revisited. Given the final limits we have calculated in this example, in order to obtain a U table associated with the 3×3 table, we can take $r_1 = 0.716$ (for example). In this case, from (4.8) we have to take a value r_2 such that $0.417 \leq r_2 \leq 0.419$: it may be $r_2 = 0.418$. With these ratios, and taking appropriate values for ε_{1j} , for example, $\varepsilon_{11} = 0.74$, $\varepsilon_{12} = 0.02$ and $\varepsilon_{13} = 0.99$ (we will see how to take these values in Section 5) we obtain the associated U table

$16 + 0.74$	$10 + 0.02$	$6 + 0.99$	=	16.74	10.02	6.99
$16.74 \cdot 0.716$	$10.02 \cdot 0.716$	$6.99 \cdot 0.716$		11.98	1.174	5.005
$16.74 \cdot 0.716 \cdot 0.418$	$10.02 \cdot 0.716 \cdot 0.418$	$6.99 \cdot 0.716 \cdot 0.418$		5.010	2.999	2.092

5. CHARACTERIZATION OF THE MPTs IN TERMS OF THE U TABLES

In order to characterize the MPTs in terms of the U tables we will use products of limits M 's, m 's, M^o 's and m^o 's (which we will denote by ΠM , Πm , ΠM^o and Πm^o , respectively), the subindices of which are chained in the sense that we are going to define.

Definition 5.1. We will say that ΠM is a product whose subindices are *chained* if it can be written as $M_{i_0 i_1} M_{i_1 i_2} \cdots M_{i_{h-1} i_h}$. When $i_h = i_0$ we will say that the subindices of the product are *circularly chained*. We will say the same for products Πm , ΠM^o and Πm^o .

From this definition we see that the chained subindices of a product

$$M_{i_0 i_1} M_{i_1 i_2} \cdots M_{i_{h-1} i_h}$$

form a sequence $i_0[\cdot]i_1[\cdot]i_2 \cdots i_{h-1}[\cdot]i_h$ (without specifying the columns), and if $i_h = i_0$ they would form a loop. Let us give some examples. The subindices of the product $M_{23}M_{35}M_{54}$ are chained, and they form the sequence $2[\cdot]3[\cdot]5[\cdot]4$. The subindices of the product $m_{13}m_{34}m_{42}m_{21}$ are circularly chained, and they form the loop $1[\cdot]3[\cdot]4[\cdot]2[\cdot]1$. Logically there will be some products whose subindices are not chained, e.g. $M_{12}M_{34}$.

We will now provide a result about products of terms M 's, which we will use in the characterization of the MPTs in terms of the U tables.

Theorem 5.1. *Given an MPT $X = \{x_{ij}\}$, for a product of terms M 's (ΠM) whose subindices are circularly chained, the following will always hold*

$$(5.1) \quad \Pi M \geq 1.$$

Proof: As we have seen previously, the terms M 's and m 's are obtained by a recursive process from the M 's and m 's of the previous step. Specifically, from (4.5), M_{hi} is either a product $M_{hs}M_{si}$ or a product $M_{hh'}M_{h'i'}M_{i'i}$ of terms of the previous step, whose subindices (in both cases) are chained, and they form a sequence beginning in row h and ending in row i . Now, by going back one step in the recursive process, the same can be applied to each of these M_{hs} , M_{si} , $M_{hh'}$, \dots . In this way, by going back to the initial step in the process, we will obtain that M_{hi} can always be written as a product of terms M^o 's whose subindices are chained, and they form a sequence beginning in row h and ending in row i .

Thus, and in accordance with what has just been said, the product on the left-hand side of (5.1), whose subindices are circularly chained, can always be expressed as a product ΠM^o whose subindices are circularly chained. Therefore, in order to demonstrate the theorem we will have to prove that $\Pi M^o \geq 1$ whenever the subindices of the product are circularly chained, and they form a loop. Let the product be

$$\Pi M^o = M_{i_0 i_1}^o M_{i_1 i_2}^o \cdots M_{i_{s-1} i_0}^o$$

which, according to the definition of the terms M^o 's, can be written as

$$\Pi M^o = \prod_{h=1}^s \frac{x_{i_h j_h} + 1}{x_{i_{h-1} j_h}}$$

where j_1, j_2, \dots, j_s are columns that correspond to the terms M^o 's of the product, and where $i_s = i_0$. Then, if $j_1 = j_2 = \cdots = j_s$ it is evident that $\Pi M^o > 1$. Otherwise, we can consider the loop $i_0[j_1]i_1[j_2]i_2 \cdots i_{s-1}[j_s]i_0$, which is determined by the subindices of the product, and from the expression (3.9) of the Theorem 3.4 one will obtain $\Pi M^o \geq 1$. \square

The following result characterizes the MPTs in terms of the U tables.

Theorem 5.2. *An $r \times c$ table $X = \{x_{ij}\}$ is an MPT if, and only if, a U table $\{u_{ij}\}$ exists associated with it.*

Proof: Let $\{u_{ij}\}$ be a U table associated with X . From the Property 4.3,

$$\prod_{j \in J} \frac{u_{bj}}{u_{aj}} = 1$$

for every order k simple loop, $1 < k \leq \min(r, c)$, and J , a and b defined as in the said property. In addition, from (4.1) we will obtain

$$x_{aj} \leq u_{aj}, \quad x_{bj} \leq u_{bj} \quad \text{and} \quad x_{bj} + 1 \geq u_{bj}.$$

Hence, for every order k simple loop, $1 < k \leq \min(r, c)$,

$$1 = \prod_{j \in J} \frac{u_{bj}}{u_{aj}} \leq \prod_{j \in J} \frac{x_{bj} + 1}{x_{aj}}$$

and, therefore, X fulfils the condition of Theorem 3.1 and will be an MPT.

It remains to be demonstrated that if X is an MPT, there will always be a U table associated with it. For this purpose, and in accordance with what was said previously in Section 4, on the one hand we must demonstrate that there will always be a set of ratios r_h , $1 \leq h < r$, which fulfil (4.8).

Firstly, in order to demonstrate that we can always take at least one value for each of the ratios r_h , $1 \leq h < r$, within the respective intervals $(m_{h,h+1}, M_{h,h+1})$, which would be true if $m_{hi} \leq M_{hi}$ for $h < i$, it will be sufficient to prove that any of the expressions which appear on the right-hand side of (4.4) is less than or equal to any of those on the right-hand side of (4.5), i.e.

$$\begin{aligned} M_{hs'}M_{s'i}/(m_{hs}m_{si}) &= M_{hs'}M_{s'i}M_{is}M_{sh} \geq 1, \\ M_{hh'}M_{h'i'}M_{i'i}/(m_{hs}m_{si}) &= M_{hh'}M_{h'i'}M_{i'i}M_{is}M_{sh} \geq 1, \\ M_{hs}M_{si}/(m_{hh'}m_{h'i'}m_{i'i}) &= M_{hs}M_{si}M_{i'i'}M_{i'h'}M_{h'h} \geq 1, \\ M_{hh'}M_{h'i'}M_{i'i}/(m_{hh''}m_{h''i''}m_{i''i}) &= M_{hh'}M_{h'i'}M_{i'i}M_{i''i''}M_{i''h''}M_{h''h} \geq 1, \end{aligned}$$

for all s , h' and i' within the established limits in (4.4) and (4.5), and all s' , h'' and i'' with the same limits of s , h' and i' , respectively. But all these inequalities are true from Theorem 5.1, because the subindices of each one of the products are circularly chained.

Secondly, we will demonstrate by induction that, given the final limits, there will always be at least one set of these ratios (taken successively, r_1, r_2, \dots, r_{r-1}) that fulfil (4.8). We can always take one value for the first ratio r_1 from inside (m_{12}, M_{12}) , and it is obvious that this r_1 fulfils (4.8).

Now we have to prove that if we take a subset of ratios r_1, r_2, \dots, r_{h-1} ($1 < h < r$) such that they fulfil (4.8), we can always take an r_h which also fulfils (4.8). It is easy to see that if r_1, r_2, \dots, r_{h-1} fulfil (4.8), we will have

$$(5.2) \quad m_{ss'} \leq r_{ss'} \leq M_{ss'}, \quad 1 \leq s < s' \leq h.$$

Now, for an r_h that fulfils (4.8) to exist it will be enough to prove that

$$(5.3) \quad m_{s,h+1}/r_{sh} \leq M_{s',h+1}/r_{s'h}, \quad \forall s, s', \quad 1 \leq s \leq h, \quad 1 \leq s' \leq h.$$

For $s = s'$ it is obvious that this inequality is fulfilled, because we have already proved that $m_{hi} \leq M_{hi}$. For $s < s'$, and taking into account Remark 4.4, expression (5.3) is reduced to

$$m_{s,h+1} m_{h+1,s'} \leq r_{ss'}, \quad 1 \leq s < s' \leq h,$$

which is true, because from (4.4) and (5.2) we have

$$m_{s,h+1} m_{h+1,s'} \leq m_{ss'} \leq r_{ss'}, \quad 1 \leq s < s' \leq h.$$

This is proved in a similar way for $s > s'$.

Finally, given a set of ratios r_1, r_2, \dots, r_{r-1} that fulfil (4.8), we must demonstrate that there will always be a U table $\{u_{ij}\}$, with $u_{1j} = x_{1j} + \varepsilon_{1j}$ ($0 \leq \varepsilon_{1j} \leq 1$) and $u_{ij} = u_{1j} r_{1i}$, $1 < i \leq r$, $1 \leq j \leq c$, which is associated with X . In other words, we have to prove that there will always be values ε_{1j} , $1 \leq j \leq c$, such that (4.1) is fulfilled, i.e., such that

$$0 \leq (x_{1j} + \varepsilon_{1j}) r_{1i} - x_{ij} \leq 1, \quad 1 \leq i \leq r,$$

from which it follows that the ε_{1j} , $1 \leq j \leq c$, should satisfy

$$(5.4) \quad \max_{1 \leq i \leq r} \left\{ \frac{x_{ij}}{r_{1i}} - x_{1j} \right\} \leq \varepsilon_{1j} \leq \min_{1 \leq i \leq r} \left\{ \frac{x_{ij} + 1}{r_{1i}} - x_{1j} \right\}.$$

Let us see that for every j there is a value ε_{1j} which satisfies (5.4). For this purpose, it is enough to prove that for any j and any i and i' the following holds:

$$x_{ij}/r_{1i} \leq (x_{i'j} + 1)/r_{1i'}.$$

For $i = i'$ it is trivial that this is true. For $i < i'$ it is also true, because $r_{1i'}/r_{1i} = r_{ii'}$, and because from (4.8) (taking $s = i$ and $h + 1 = i'$) and from the definition of the limits M 's and the limits M^o 's we can obtain

$$r_{ii'} \leq M_{ii'} \leq M_{ii'}^o \leq (x_{i'j} + 1)/x_{ij}, \quad \forall j.$$

This is proved in a similar way for $i > i'$. □

The last part of the proof of the Theorem 5.2 shows us how we can easily obtain a U table associated with an MPT. This is summarized in the next result.

Corollary 5.1. *Given an MPT $X = \{x_{ij}\}$ and a set of ratios r_h , $1 \leq h < r$, fulfilling (4.8), a table $\{u_{ij}\}$ with*

$$u_{ij} = \begin{cases} x_{1j} + \varepsilon_{1j}, & i = 1, \quad 1 \leq j \leq c, \\ u_{1j} r_{1i}, & 1 < i \leq r, \quad 1 \leq j \leq c, \end{cases}$$

where $r_{1i} = r_1 r_2 \cdots r_{i-1}$ and ε_{1j} satisfies (5.4) for all j , is a U table associated with X .

6. ON THE UNIQUENESS OF A MAXIMUM PROBABILITY $r \times c$ TABLE

The following result, based on the previous results, is a necessary and sufficient condition which characterizes the uniqueness of an MPT.

Theorem 6.1. *An MPT $X = \{x_{ij}\}$ is unique if and only if*

$$(6.1) \quad \prod_{j \in J} \frac{x_{bj} + 1}{x_{aj}} > 1$$

for every order k simple loop $E = \{e_{ij}\}$ and every k , $1 < k \leq \min(r, c)$, where J is the set of the k columns of the loop, and for each $j \in J$, b and a are the rows such that $e_{bj} = 1$ and $e_{aj} = -1$.

Proof: Let X be the unique MPT, which obviously will fulfil (3.1). If (3.8) is fulfilled for a simple order k loop E , then, from Theorem 3.2, $X' = X + E$ would be an MPT, which would contradict the initial hypothesis and, thus, (6.1) is fulfilled. Reciprocally, let X be an MPT fulfilling (6.1), and let us suppose that X' is also an MPT. Then, from Theorem 3.3, the difference between both will be one or several simple loops, such that (3.8) will be fulfilled for X and for each of these simple loops, which would contradict (6.1). Hence, X is the only MPT. \square

7. CONCLUSIONS

The most efficient algorithm (network algorithm) to calculate the p -value of the Fisher's exact test in an $r \times c$ table requires us to calculate many times maximum probability $r \times c'$ ($c' \leq c$) contingency tables, and to perform a great amount of comparisons in which the probabilities of these tables are involved. At present, the general method to obtain maximum probability fixed marginals contingency tables is based on a necessary condition for these tables, which makes that method insufficiently efficient, especially for a relatively large r or c . In this paper, we present two necessary and sufficient conditions for these maximum probability tables. This characterization, especially that which is expressed based on U tables, will allow us to construct a general algorithm for obtaining the aforementioned maximum probability contingency tables.

ACKNOWLEDGMENTS

This research was supported by the Spanish Ministry of Economy and Competitiveness, Grant Number MTM2016-76938-P. The author would like to thank the referee for many helpful comments which improved the presentation of this paper.

REFERENCES

- [1] JOE, H. (1988). Extreme Probabilities for Contingency Tables under Row and Column Independence with Application to Fisher's Exact Test, *Comm. Statist. Theory Meth.*, **17**, 3677–3685.
- [2] MEHTA, C.R. and PATEL, R.N. (1983). A Network Algorithm for Performing Fisher's Exact Test in $r \times c$ Contingency Tables, *J. Amer. Statist. Assoc.*, **78**, 427–434.
- [3] REQUENA, F. and MARTÍN, N. (2000). Characterization of Maximum Probability Points in the Multivariate Hypergeometric Distribution, *Statist. Probab. Lett.*, **50**, 39–47.
- [4] REQUENA, F. and MARTÍN, N. (2003). The Maximum Probability $2 \times c$ Contingency Tables and the Maximum Probability Points of the Multivariate Hypergeometric Distribution, *Comm. Statist. Theory Meth.*, **32**, 1737–1752.
- [5] REQUENA, F. and MARTÍN, N. (2006). A Major Improvement to the Network Algorithm for Fisher's Exact Test in $2 \times c$ Contingency Tables, *Comput. Statist. Data Anal.*, **51**, 490–498.

ON THE OCCURRENCE OF BOUNDARY SOLUTIONS IN TWO-WAY INCOMPLETE TABLES

Authors: SAYAN GHOSH

– Department of Mathematics, Indian Institute of Technology Bombay,
Powai, Mumbai 400076, India
sayang@math.iitb.ac.in

PALANIAPPAN VELLAISAMY

– Department of Mathematics, Indian Institute of Technology Bombay,
Powai, Mumbai 400076, India
pv@math.iitb.ac.in

Received: August 2017

Revised: October 2017

Accepted: October 2017

Abstract:

- The analysis of incomplete contingency tables is an important problem, which is also of practical interest. In this paper, we consider boundary solutions under nonignorable nonresponse models in two-way incomplete tables with data on both variables missing. We establish a result similar to [9] on sufficient conditions for the occurrence of boundary solutions. We also provide a new result, which connects the forms of boundary solutions under various parameterizations of the missing data models. This result helps us to obtain the exact form of boundary solutions in the above tables, which improves a claim made in [2] and avoids computational burden. A counterexample is provided to show that the sufficient conditions for the occurrence of boundary solutions are not necessary, thereby disproving a conjecture of [7]. Finally, we establish new necessary conditions for the occurrence of boundary solutions under nonignorable nonresponse models in square two-way incomplete tables, and show that they are not sufficient. These conditions are simple and easy to check as they depend only on the observed cell counts. They are useful and important for model selection also. Some real life data sets are analyzed to illustrate the results.

Key-Words:

- *incomplete tables; boundary solutions; Log-linear models; NMAR models.*

AMS Subject Classification:

- 62H17.

1. INTRODUCTION

Contingency tables with fully observed counts and partially classified margins (nonresponses) are called incomplete tables. The following three types of missing data mechanisms have been proposed in the literature ([8]): missing completely at random (MCAR), missing at random (MAR) and not missing at random (NMAR). The missing mechanism is said to be (a) MCAR when missingness is independent of both observed and unobserved data, (b) MAR when missingness depends only on observed data, and (c) NMAR if missingness depends on unobserved data. Nonresponses are called ignorable when the missing data mechanism is MAR or MCAR, and the parameters governing the missing data mechanism are distinct from those to be estimated. They are nonignorable when the missing data mechanism is NMAR.

Log-linear models have generally been used to study missing data mechanisms in incomplete tables (see [9] and references therein). However, under nonignorable models, a boundary solution occurs when the cell probabilities of non-respondents are estimated to be zeros for certain levels of the missing variables. That is, the maximum likelihood estimators (MLE's) of the parameters lie on the boundary of the parameter space. Note that the problem of boundary solutions is an important one as it has serious consequences for statistical inference. For example, the observed counts cannot be reproduced by a perfect fit model (a model for which the estimated expected counts are equal to the observed counts) if boundary solutions occur. This implies that the fit is inadequate and the parameter estimates are imprecise. The log likelihood function is flat and, therefore, convergence of the EM algorithm to the boundary MLE's requires a lot of iterations. Also, the eigenvalues of the covariance matrix are inappropriate (either around zero or negative), which implies some parameter estimates have large estimated standard errors and wide confidence intervals. Hence, it is useful to study various forms of boundary solutions and explore conditions for their occurrence in incomplete tables.

Consider two categorical variables with I and J levels. Then an $I \times J \times 2$ table and an $I \times J \times 2 \times 2$ table represent two-way incomplete tables with data on one of the variables and data on both the variables missing respectively. The problem of boundary solutions was first considered by [1] who proposed a sufficient condition for their occurrence in a $2 \times 2 \times 2$ incomplete table. [2] studied the problem for an $I \times J \times 2 \times 2$ incomplete table, which has non-monotone missing value patterns. For an $I \times J \times 2$ incomplete table with simple monotone missing value patterns, [10] and [3] described the problem geometrically, while [4] discussed properties of MLE's in case of boundary solutions. [9] proposed sufficient conditions for the occurrence of boundary solutions under various NMAR models in an $I \times I \times 2 \times 2$ incomplete table. Recently, [5] provided forms of boundary solutions in arbitrary three-way and n -dimensional incomplete tables with one or more variables missing, and also established sufficient conditions for their occurrence under various NMAR models. In this paper, we consider the above and other related issues for an $I \times J \times 2 \times 2$ table. Note that a lower dimensional incomplete table is not a special case of a higher dimensional one and hence any result for the former cannot be obtained directly from that for the latter.

The purpose of this paper is to provide a comprehensive treatment of the problem of boundary solutions in two-way incomplete tables with both variables missing. To this effect,

we first introduce some notations and consider various identifiable NMAR log-linear models (Models [M1]–[M5]) for an $I \times J \times 2 \times 2$ incomplete table. The problem of boundary solutions, along with their forms under the above models, is discussed in Section 3. We formally define boundary solutions for an $I \times J \times 2 \times 2$ incomplete table by extending the definition of [1], which are unavailable in the literature. A novel result (Proposition 3.1) is provided, which gives the relationship among forms of boundary solutions according to various parameterizations for the missing data models. This helps us to theoretically justify and deduce the exact boundary solutions in those models directly without having to obtain them empirically (see pp. 39–40 of [9]) using the EM algorithm. In Section 4, we illustrate this result using some data analysis examples from [2], thereby improving a claim made by them on the forms of boundary solutions in $I \times J \times 2 \times 2$ tables, which also eliminates computations.

In Section 5, we provide a result (Theorem 5.1) on sufficient conditions for the occurrence of boundary solutions in the above tables, which is similar to Theorem 1 of [9] but proved using direct arguments instead of contrapositive ones used in [9]. While [9] consider only Model [M5] in Theorem 1, we consider Models [M1]–[M5] in Theorem 5.1. A counterexample is provided to show that the sufficient conditions for the occurrence of boundary solutions are not necessary, which refutes a conjecture due to [7].

Finally, we propose new necessary conditions in Theorem 5.2 for the occurrence of boundary solutions under Models [M1]–[M5] in square two-way incomplete tables, and later show that they are not sufficient through a counterexample. Such conditions do not exist in the literature. Note that these conditions help us to identify the non-occurrence of boundary solutions, which is very useful for fitting appropriate models to the incomplete data (model selection). Also, these conditions involve only the observed cell counts and their sums in the tables, and hence can be easily verified. Section 6 provides some concluding remarks.

2. NMAR LOG-LINEAR MODELS

Suppose Y_1 and Y_2 are two categorical variables having I and J levels respectively. For $i = 1, 2$, let R_i denote the missing indicator for Y_i so that $R_i = 1$ or 2 if Y_i is observed or unobserved. Then we have an $I \times J \times 2 \times 2$ incomplete table, corresponding to Y_1 , Y_2 , R_1 and R_2 , with cell counts $\mathbf{y} = \{y_{ijkl}\}$ where $1 \leq i \leq I$, $1 \leq j \leq J$ and $1 \leq k, l \leq 2$. The vector of observed counts is $\mathbf{y}_{\text{obs}} = (\{y_{ij11}\}, \{y_{i+12}\}, \{y_{+j21}\}, y_{++22})$, where $\{y_{ij11}\}$ are the fully observed counts and $\{y_{i+12}\}, \{y_{+j21}\}, y_{++22}$ are the partially classified counts also known as the supplementary margins. All cell counts are assumed to be positive. The fully observed counts are those for which data on both Y_1 and Y_2 is available, while data on at most Y_1 or Y_2 is available for the supplementary margins. Note that ‘+’ denotes summation over levels of the corresponding variable. For example, y_{+j21} denotes the number of observations corresponding to $Y_2 = j$ for which data on Y_2 is observed but data on Y_1 is missing. Let $\pi = \{\pi_{ijkl}\}$ be the vector of cell probabilities, $\mu = \{\mu_{ijkl}\}$ be the vector of expected counts and $N = \sum_{i,j,k,l} y_{ijkl}$ the total number of cell counts. For $I = J = 2$, we have the $2 \times 2 \times 2 \times 2$ incomplete table given by Table 1.

Table 1: $2 \times 2 \times 2 \times 2$ Incomplete Table.

		$R_2 = 1$		$R_2 = 2$
		$Y_2 = 1$	$Y_2 = 2$	Y_2 missing
$R_1 = 1$	$Y_1 = 1$	y_{1111}	y_{1211}	y_{1+12}
	$Y_1 = 2$	y_{2111}	y_{2211}	y_{2+12}
$R_1 = 2$	Y_1 missing	y_{+121}	y_{+221}	y_{++22}

We consider Poisson sampling for convenience, that is, $Y_{ijkl} \sim P(\mu_{ijkl})$. The likelihood function of μ is

$$(2.1) \quad L(\mu; \mathbf{y}_{\text{obs}}) = \frac{e^{-\sum_{i,j,k,l} \mu_{ijkl}} \prod_{i,j} \mu_{ij11}^{y_{ij11}} \prod_i \mu_{i+12}^{y_{i+12}} \prod_j \mu_{+j21}^{y_{+j21}} \mu_{++22}^{y_{++22}}}{\prod_{i,j,k,l} y_{ijkl}!}$$

so that the log-likelihood function of μ is

$$(2.2) \quad \begin{aligned} l(\mu; \mathbf{y}_{\text{obs}}) &= \sum_{i,j} y_{ij11} \log \mu_{ij11} + \sum_i y_{i+12} \log \mu_{i+12} + \sum_j y_{+j21} \log \mu_{+j21} \\ &\quad + y_{++22} \log \mu_{++22} - \mu_{++++} + \Delta, \end{aligned}$$

where Δ is independent of μ_{ijkl} 's. For an $I \times J \times 2 \times 2$ incomplete table, [2] proposed the following log-linear model (with no three-way or four-way interactions):

$$(2.3) \quad \begin{aligned} \log \mu_{ijkl} &= \lambda + \lambda_{Y_1}(i) + \lambda_{Y_2}(j) + \lambda_{R_1}(k) + \lambda_{R_2}(l) + \lambda_{Y_1 Y_2}(i, j) \\ &\quad + \lambda_{Y_1 R_1}(i, k) + \lambda_{Y_2 R_1}(j, k) + \lambda_{Y_1 R_2}(i, l) + \lambda_{Y_2 R_2}(j, l) + \lambda_{R_1 R_2}(k, l), \end{aligned}$$

where the sum over any argument of a log-linear parameter is zero, for example, $\sum_i \lambda_{Y_1 Y_2}(i, j) = \sum_j \lambda_{Y_1 Y_2}(i, j) = 0$. To study the various missing mechanisms of Y_1 and Y_2 , [2] introduced the following notations:

$$\begin{aligned} a_{ij} &= \frac{P(R_1 = 2, R_2 = 1 | Y_1 = i, Y_2 = j)}{P(R_1 = 1, R_2 = 1 | Y_1 = i, Y_2 = j)} = \frac{\pi_{ij21}}{\pi_{ij11}} = \frac{\mu_{ij21}}{\mu_{ij11}}, \\ b_{ij} &= \frac{P(R_1 = 1, R_2 = 2 | Y_1 = i, Y_2 = j)}{P(R_1 = 1, R_2 = 1 | Y_1 = i, Y_2 = j)} = \frac{\pi_{ij12}}{\pi_{ij11}} = \frac{\mu_{ij12}}{\mu_{ij11}}, \quad m_{ij11} = N\pi_{ij11}, \\ g &= \frac{P(R_1 = 1, R_2 = 1 | Y_1 = i, Y_2 = j) P(R_1 = 2, R_2 = 2 | Y_1 = i, Y_2 = j)}{P(R_1 = 1, R_2 = 2 | Y_1 = i, Y_2 = j) P(R_1 = 2, R_2 = 1 | Y_1 = i, Y_2 = j)}. \end{aligned}$$

Remark 2.1. Under (2.3), it can be shown that $a_{ij} = \exp[-2\{\lambda_{R_1}(1) + \lambda_{Y_1 R_1}(i, 1) + \lambda_{Y_2 R_1}(j, 1) + \lambda_{R_1 R_2}(1, 1)\}]$ and $b_{ij} = \exp[-2\{\lambda_{R_2}(1) + \lambda_{Y_1 R_2}(i, 1) + \lambda_{Y_2 R_2}(j, 1) + \lambda_{R_1 R_2}(1, 1)\}]$. Also, we have $g = \frac{\pi_{ij11}\pi_{ij22}}{\pi_{ij12}\pi_{ij21}} = \frac{\mu_{ij11}\mu_{ij22}}{\mu_{ij12}\mu_{ij21}}$. Hence

$$\begin{aligned} \log g &= \log \mu_{ij11} + \log \mu_{ij22} - \log \mu_{ij12} - \log \mu_{ij21} \\ \implies \log g &= \lambda_{R_1 R_2}(1, 1) + \lambda_{R_1 R_2}(2, 2) - \lambda_{R_1 R_2}(1, 2) - \lambda_{R_1 R_2}(2, 1) \quad (\text{from (2.3)}) \\ \implies \log g &= 4\lambda_{R_1 R_2}(1, 1) \quad (\because \lambda_{R_1 R_2}(1, 2) = -\lambda_{R_1 R_2}(2, 2) = \lambda_{R_1 R_2}(2, 1) = -\lambda_{R_1 R_2}(1, 1)) \\ \implies g &= \exp[4\lambda_{R_1 R_2}(1, 1)], \end{aligned}$$

which is independent of i and j .

Note that $m_{ij11} = \mu_{ij11}$ and g denotes the odds ratio between the missing indicators of Y_1 and Y_2 . Also, $\mu_{ij21} = m_{ij11}a_{ij}$, $\mu_{ij12} = m_{ij11}b_{ij}$ and $\mu_{ij22} = m_{ij11}a_{ij}b_{ij}g$. Note that a_{ij} is the conditional odds of Y_1 being missing given Y_2 is observed, while b_{ij} is the conditional odds of Y_2 being missing given Y_1 is observed. Here, a_{ij} and b_{ij} describe the missing mechanisms of Y_1 and Y_2 , respectively. Denote a_{ij} (b_{ij}) by α_i (β_i) or α_j (β_j) or $\alpha_{..}$ ($\beta_{..}$) if it depends only on i or j or none, respectively. Then we have the following definition.

Definition 2.1. The missing mechanism of Y_1 under (2.3) is NMAR if $a_{ij} = \alpha_i$, MAR if $a_{ij} = \alpha_j$ and MCAR if $a_{ij} = \alpha_{..}$. Similarly, the missing mechanism of Y_2 is NMAR if $b_{ij} = \beta_j$, MAR if $b_{ij} = \beta_i$ and MCAR if $b_{ij} = \beta_{..}$.

Using Definition 2.1 and the above notations, there are nine possible identifiable models (see pp. 647–648 of [2]) based on different missing mechanisms for Y_1 and Y_2 . The equivalent log-linear models can be obtained as submodels of (2.3). As an example, consider the model (α_i, β_i) , for which the missing mechanism is NMAR for Y_1 and MAR for Y_2 . Using the expressions of a_{ij} and b_{ij} in Remark 2.1, the corresponding log-linear model is obtained from (2.3) by substituting $\lambda_{Y_2R_1}(j, k) = \lambda_{Y_2R_2}(j, l) = 0$. The following are the five models when the missing mechanism is NMAR for Y_1 or Y_2 .

1. Model M1 (NMAR for Y_1 , MCAR for Y_2):

$$\begin{aligned} \log \mu_{ijkl} &= \lambda + \lambda_{Y_1}(i) + \lambda_{Y_2}(j) + \lambda_{R_1}(k) + \lambda_{R_2}(l) + \lambda_{Y_1Y_2}(i, j) \\ &\quad + \lambda_{Y_1R_1}(i, k) + \lambda_{R_1R_2}(k, l). \end{aligned}$$

2. Model M2 (NMAR for Y_2 , MCAR for Y_1):

$$\begin{aligned} \log \mu_{ijkl} &= \lambda + \lambda_{Y_1}(i) + \lambda_{Y_2}(j) + \lambda_{R_1}(k) + \lambda_{R_2}(l) + \lambda_{Y_1Y_2}(i, j) \\ &\quad + \lambda_{Y_2R_2}(j, l) + \lambda_{R_1R_2}(k, l). \end{aligned}$$

3. Model M3 (NMAR for Y_1 , MAR for Y_2):

$$\begin{aligned} \log \mu_{ijkl} &= \lambda + \lambda_{Y_1}(i) + \lambda_{Y_2}(j) + \lambda_{R_1}(k) + \lambda_{R_2}(l) + \lambda_{Y_1Y_2}(i, j) + \lambda_{Y_1R_1}(i, k) \\ &\quad + \lambda_{Y_1R_2}(i, l) + \lambda_{R_1R_2}(k, l). \end{aligned}$$

4. Model M4 (NMAR for Y_2 , MAR for Y_1):

$$\begin{aligned} \log \mu_{ijkl} &= \lambda + \lambda_{Y_1}(i) + \lambda_{Y_2}(j) + \lambda_{R_1}(k) + \lambda_{R_2}(l) + \lambda_{Y_1Y_2}(i, j) + \lambda_{Y_2R_1}(j, k) \\ &\quad + \lambda_{Y_2R_2}(j, l) + \lambda_{R_1R_2}(k, l). \end{aligned}$$

5. Model M5 (NMAR for both Y_1 and Y_2):

$$\begin{aligned} \log \mu_{ijkl} &= \lambda + \lambda_{Y_1}(i) + \lambda_{Y_2}(j) + \lambda_{R_1}(k) + \lambda_{R_2}(l) + \lambda_{Y_1Y_2}(i, j) + \lambda_{Y_1R_1}(i, k) \\ &\quad + \lambda_{Y_2R_2}(j, l) + \lambda_{R_1R_2}(k, l). \end{aligned}$$

Note that for Models [M1]–[M5], there is an association term between a variable and its missing indicator if the missing mechanism is NMAR for that variable (for example, the term $\lambda_{Y_1R_1}(i, k)$ in Model [M1]), between a variable and the other missing indicator if the missing mechanism is MAR for that variable (for example, the term $\lambda_{Y_2R_1}(j, k)$ in Model [M4]) and none if the missing mechanism is MCAR for a variable (for example, $\lambda_{Y_1R_1}(i, k)$ and $\lambda_{Y_2R_1}(j, k)$ are absent in Model [M2]).

3. BOUNDARY SOLUTIONS IN NMAR MODELS

In this section, we consider boundary solutions under non-ignorable nonresponse (NMAR) models for an $I \times J \times 2 \times 2$ incomplete table. We first define boundary solutions under the above models and then present a result relating the forms of boundary solutions in terms of various parameterizations of the models.

For an incomplete table, boundary solutions in NMAR models occur when the MLE's of nonresponse cell probabilities are all zeros for certain levels of the missing variables. For an $I \times J \times 2$ incomplete table, where data on only Y_2 is missing, [1] defined boundary solutions in the NMAR model for Y_2 as $\hat{\pi}_{ij2} = 0$ for at least one pair (i, j) . For the same model, [4] showed that boundary solutions are given by $\hat{\pi}_{+j2} = 0$ for at least one and at most $(J-1)$ values of Y_2 . [1] defined a nonresponse boundary solution under NMAR models in general to be a stationary point that lies on a boundary of the space of parameters modeling the nonignorable nonresponse. Using this, we may extend their definition to an $I \times J \times 2 \times 2$ table as follows.

Definition 3.1. Consider an $I \times J \times 2 \times 2$ incomplete table, and let $1 \leq i \leq I$, $1 \leq j \leq J$ and $k, l = 1, 2$. Then we have the following:

1. A nonresponse boundary solution under the NMAR models for Y_1 only, that is, Models [M1] and [M3], is an MLE given by $\hat{\pi}_{ij2l} = 0$ for at least one combination (i, j, l) .
2. A nonresponse boundary solution under the NMAR models for Y_2 only, that is, Models [M2] and [M4], is an MLE given by $\hat{\pi}_{ijk2} = 0$ for at least one combination (i, j, k) .
3. A nonresponse boundary solution under the NMAR model for both Y_1 and Y_2 , that is, Model [M5], is an MLE given by $\hat{\pi}_{ij2l} = 0$ for at least one combination (i, j, l) or $\hat{\pi}_{ijk2} = 0$ for at least one combination (i, j, k) .

Note that in the literature, boundary solutions have usually been defined in terms of cell probabilities because the cell probabilities are in some sense natural to the model for the incomplete table, whereas the loglinear parameters are not. The next proposition explores the relationships among boundary solutions under Models [M1]–[M5] in terms of MLE's of nonresponse cell probabilities, some specific log-linear parameters and α_i or β_j for two-way incomplete tables with both variables missing.

Proposition 3.1. For an $I \times J \times 2 \times 2$ incomplete table, we have the following:

1. For Models [M1] and [M3], if boundary solutions occur, then they are given by $\hat{\lambda}_{Y_1 R_1}(i, 2) = -\infty \Leftrightarrow \hat{\pi}_{i+2+} = 0 \Leftrightarrow \hat{\alpha}_i = 0$ for at least one and at most $(I-1)$ values of Y_1 .
2. For Models [M2] and [M4], if boundary solutions occur, then they are given by $\hat{\lambda}_{Y_2 R_2}(j, 2) = -\infty \Leftrightarrow \hat{\pi}_{+j+2} = 0 \Leftrightarrow \hat{\beta}_j = 0$ for at least one and at most $(J-1)$ values of Y_2 .
3. For Model [M5], if boundary solutions occur, then they are given by $\hat{\lambda}_{Y_1 R_1}(i, 2) = -\infty$ or $\hat{\lambda}_{Y_2 R_2}(j, 2) = -\infty \Leftrightarrow \hat{\pi}_{i+2+} = 0$ or $\hat{\pi}_{+j+2} = 0 \Leftrightarrow \hat{\alpha}_i = 0$ for at least one and at most $(I-1)$ values of Y_1 or $\hat{\beta}_j = 0$ for at least one and at most $(J-1)$ values of Y_2 .

Proof: See Appendix A.1. □

From the proof of Proposition 3.1 in Appendix A.1, note that the one-to-one relation between the cell probabilities and the log-linear parameters cannot be used to derive the connection between the different forms of boundary solutions. This is because it is not obvious which specific log-linear parameters have infinite MLE's just by noting the zero MLE's of the nonresponse cell probabilities when boundary solutions occur.

4. SOME EXAMPLES OF BOUNDARY SOLUTIONS IN NMAR MODELS

In this section, we reanalyze some examples in [2], illustrating the result in Section 3. We use Proposition 3.1 to investigate a claim made by [2] regarding forms and occurrence of boundary solutions in an $I \times J \times 2 \times 2$ incomplete table. This improvement is useful as it avoids computation and provides the exact boundary solutions under a NMAR model by simply noting the level(s) of the variable(s) for which the MLE's of the parameters are negative or infinite.

First, we present the correct expression of the likelihood ratio statistic for missing data models in such a table. Consider testing the goodness of fit of a null model (here one of the Models [M1]–[M5]) against the alternative model (perfect fit model). Let $\{\hat{\mu}_{ijkl}\}$ and $\{\tilde{\mu}_{ijkl}\}$ denote the MLE's of the expected counts under a null model and a perfect fit model respectively. Also, let L_0 and L_1 denote the log-likelihoods for the null and the alternative models, respectively. Then the likelihood ratio statistic is given by

$$\begin{aligned}
 G^2 &= -2(L_0 - L_1) \\
 &= -2 \left[\sum_{i,j} y_{ij11} \ln \left(\frac{\hat{\mu}_{ij11}}{\tilde{\mu}_{ij11}} \right) + \sum_i y_{i+12} \ln \left(\frac{\hat{\mu}_{i+12}}{\tilde{\mu}_{i+12}} \right) + \sum_j y_{+j21} \ln \left(\frac{\hat{\mu}_{+j21}}{\tilde{\mu}_{+j21}} \right) \right. \\
 &\quad \left. + y_{++22} \ln \left(\frac{\hat{\mu}_{++22}}{\tilde{\mu}_{++22}} \right) - \hat{\mu}_{++++} + \tilde{\mu}_{++++} \right] \\
 (4.1) \quad &= -2 \left[\sum_{i,j} y_{ij11} \ln \left(\frac{\hat{m}_{ij11}}{y_{ij11}} \right) + \sum_i y_{i+12} \ln \left(\frac{\sum_j \hat{m}_{ij11} \hat{b}_{ij}}{y_{i+12}} \right) \right. \\
 &\quad \left. + \sum_j y_{+j21} \ln \left(\frac{\sum_i \hat{m}_{ij11} \hat{a}_{ij}}{y_{+j21}} \right) + y_{++22} \ln \left(\frac{\sum_{i,j} \hat{m}_{ij11} \hat{a}_{ij} \hat{b}_{ij} \hat{g}}{y_{++22}} \right) \right. \\
 &\quad \left. - \sum_{i,j} \hat{m}_{ij11} (1 + \hat{a}_{ij} + \hat{b}_{ij} + \hat{a}_{ij} \hat{b}_{ij} \hat{g}) + N \right].
 \end{aligned}$$

Note that the last two terms of (4.1) are missing in the expression of G^2 in [2] (see p. 646). Observe that in general, $\sum_{i,j} \hat{m}_{ij11} (1 + \hat{a}_{ij} + \hat{b}_{ij} + \hat{a}_{ij} \hat{b}_{ij} \hat{g}) \neq N$, unless the hypothetical (null) model is a perfect fit model for example, in which case $G^2 = 0$.

Using Definition 2.1 and the notations in Section 2, Models [M1]–[M5] can be represented as follows — Model [M1]: $(\alpha_{i.}, \beta_{.})$, Model [M2]: $(\alpha_{.}, \beta_{.j})$, Model [M3]: $(\alpha_{i.}, \beta_{i.})$, Model [M4]: $(\alpha_{.j}, \beta_{.j})$ and Model [M5]: $(\alpha_{i.}, \beta_{.j})$. Accordingly, the expression of G^2 in (4.1) for each of the

above models may be obtained by making suitable substitutions and using the MLE's in [2] (see pp. 647–648). For example, the MLE's under the model $(\alpha_i, \beta_{..})$ are

$$\hat{m}_{ij11} = \frac{y_{ij11} y_{i+1+} y_{++11}}{y_{i+11} y_{++1+}}, \quad \sum_i \hat{m}_{ij11} \hat{\alpha}_i = y_{+j21}, \quad \hat{\beta}_{..} = \frac{y_{++12}}{y_{++11}}, \quad \hat{g} = \frac{y_{++11} y_{++22}}{y_{++12} y_{++21}}.$$

Hence, from (4.1), the likelihood ratio statistic is

$$G^2 = -2 \left[\sum_{i,j} y_{ij11} \ln \left(\frac{y_{i+1+} y_{++11}}{y_{i+11} y_{++1+}} \right) + \sum_i y_{i+12} \ln \left(\frac{y_{i+1+} y_{++12}}{y_{i+12} y_{++1+}} \right) \right].$$

[2] mentioned that if any solution $\hat{\alpha}_i$ or $\hat{\beta}_{.j}$ to the systems of equations $\sum_i N \hat{\pi}_{ij11} \hat{\alpha}_i = y_{+j21}$ and $\sum_j N \hat{\pi}_{ij11} \hat{\beta}_{.j} = y_{i+12}$ respectively is negative, then boundary solutions occur, that is, the MLE lies on the boundary of the parameter space. Closed-form boundary MLE's under Models [M1]–[M5] may then be obtained (see p. 649 of [2]) by setting certain parameter estimates ($\hat{\alpha}_i$ or $\hat{\beta}_{.j}$) to 0 in the likelihood equations obtained from (2.2) for the models. They claimed that counterintuitively, the parameter estimate set to 0 need not be the estimate with a negative value as the solution to the above systems of equations. In particular, for a $2 \times 2 \times 2 \times 2$ incomplete table, they suggested examining both boundaries $\hat{\alpha}_1 = 0$ and $\hat{\alpha}_2 = 0$; similarly $\hat{\beta}_{.1} = 0$ and $\hat{\beta}_{.2} = 0$ to determine the minimum value of G^2 , which corresponds to the MLE. We improve this claim and thereby obviate computations by showing that the MLE indeed always occurs on the specific boundary (level(s) of the variable(s)) for which $\hat{\alpha}_i$ or $\hat{\beta}_{.j}$ is negative. In the next three examples, we use Proposition 3.1 to illustrate this point for Models [M1]–[M5].

Example 4.1. Consider the data in Table 2 discussed in [2], which cross-classifies mother's self-reported smoking status (Y_1) ($Y_1 = 1(2)$ for smoker (non-smoker)) with newborn's weight (Y_2) ($Y_2 = 1(2)$ if weight < 2500 grams (≥ 2500 grams)). The supplementary margins contain data on only smoking status, data on only newborn's weight and missing data on both variables.

Table 2: Birth weight and smoking: observed counts.

		$R_2 = 1$		$R_2 = 2$
		$Y_2 = 1$	$Y_2 = 2$	Y_2 missing
$R_1 = 1$	$Y_1 = 1$	4512	21009	1049
	$Y_1 = 2$	3394	24132	1135
$R_1 = 2$	Y_1 missing	142	464	1224

[2] mentioned that $\hat{\alpha}_2 < 0$ is obtained on fitting models [M1], [M3] and [M5] to the data in Table 2. Also, the value of G^2 corresponding to $\hat{\alpha}_2 = 0$ is larger than that corresponding to $\hat{\alpha}_1 = 0$ for all the above models, which is incorrect as shown below. When we fit the same models to the data in Table 2 using the 'MASS' package in R software, we obtain $\hat{\alpha}_1 = 0.0493$ and $\hat{\alpha}_2 = -0.0237$ under Models [M1], [M3] and [M5], that is, boundary solutions occur in each of the models.

Also, $G^2 = 55.2198$ (12.4682) under Model [M1], $G^2 = 55.2168$ (12.4638) under Model [M3] and $G^2 = 55.214$ (12.464) under Model [M5] when $\hat{\alpha}_1 = 0$ ($\hat{\alpha}_2 = 0$). The G^2 values for $\hat{\alpha}_2 = 0$ upon rounding off in each of the models match those given in Table V of [2]. Hence, G^2 is minimum for $\hat{\alpha}_2 = 0$ in each case, which implies that boundary solutions are given by $\hat{\alpha}_2 = 0$ or equivalently $\hat{\pi}_{2+2+} = 0$. This result is consistent with points 1 and 3 of Proposition 3.1. Further, it is the exact form of boundary solutions that we obtain on fitting Models [M1], [M3] and [M5] to the data in Table 2 using the EM algorithm (see the ‘ecm.cat’ function of ‘cat’ package in R software).

Example 4.2. Consider the example given in the last paragraph of p. 646 in [2]. The model [M1] was fitted to the following data: $y_{1111} = 100$, $y_{1211} = 40$, $y_{2111} = 50$, $y_{2211} = 1000$, $y_{1+12} = 0$, $y_{2+12} = 0$, $y_{+121} = 100$, $y_{+221} = 10$ and $y_{++22} = 0$. They mentioned that though $\hat{\alpha}_1 < 0$, G^2 is minimum for $\hat{\alpha}_2 = 0$ implying that the MLE is on the boundary $\hat{\alpha}_2 = 0$. However, we obtain $\hat{\alpha}_1 = 1.0153$ (> 0) and $\hat{\alpha}_2 = -0.0306$ on fitting Model [M1] to the above data. Also, note that $\hat{g} = \frac{y_{++11}y_{++22}}{y_{++12}y_{++21}}$ (see p. 649 of [2]) is undefined since $y_{++12} = 0$. Hence, we introduce the following changes: $y_{1+12} = 1$, $y_{2+12} = 1$ and $y_{++22} = 2$ as shown in Table 3.

Table 3: Modified $2 \times 2 \times 2 \times 2$ table.

		$R_2 = 1$		$R_2 = 2$
		$Y_2 = 1$	$Y_2 = 2$	Y_2 missing
$R_1 = 1$	$Y_1 = 1$	100	40	1
	$Y_1 = 2$	50	1000	1
$R_1 = 2$	Y_1 missing	100	10	2

On fitting models [M1], [M3] and [M5] to the data in Table 3, we obtain $\hat{\alpha}_1 = 1.0098$ under [M1], and $\hat{\alpha}_1 = 1.0153$ under [M3] and [M5], along with $\hat{\alpha}_2 = -0.0306$ under all the above models, which implies boundary solutions occur in each case. Also, $G^2 = 426.1604$ (17.4704) under Model [M1], $G^2 = 424.3288$ (15.669) under Model [M3] and $G^2 = 424.3188$ (15.664) under Model [M5] when $\hat{\alpha}_1 = 0$ ($\hat{\alpha}_2 = 0$). Hence, G^2 is minimum for $\hat{\alpha}_2 = 0$ in each model, which implies that boundary solutions are given by $\hat{\pi}_{2+2+} = 0$. This result is consistent with points 1 and 3 of Proposition 3.1. Further, it is the exact form of boundary solutions that we obtain on fitting Models [M1], [M3] and [M5] to the data in Table 3 using the EM algorithm.

Example 4.3. Consider the data in Table 2 discussed in Example 4.1. We introduce the following changes corresponding to supplementary margins in Table 2: $464 \rightarrow 700$ and $1135 \rightarrow 750$. The modified table is shown in Table 4.

Table 4: Birth weight and smoking: observed counts (modified).

		$R_2 = 1$		$R_2 = 2$
		$Y_2 = 1$	$Y_2 = 2$	Y_2 missing
$R_1 = 1$	$Y_1 = 1$	4512	21009	1049
	$Y_1 = 2$	3394	24132	750
$R_1 = 2$	Y_1 missing	142	700	1224

When we fit the models [M2], [M4] and [M5] to the data in Table 4, we obtain $\hat{\beta}_{.1} = 0.2538$ under [M2], and $\hat{\beta}_{.1} = 0.2543$ under [M4] and [M5] along with $\hat{\beta}_{.2} = -0.0047$ under all the above models, that is, boundary solutions occur in each of the models. Also, $G^2 = 98.5962$ (3.3548) under Model [M2], $G^2 = 96.1622$ (0.922) under Model [M4] and $G^2 = 96.162$ (0.9276) under Model [M5] when $\hat{\beta}_{.1} = 0$ ($\hat{\beta}_{.2} = 0$). The G^2 values in brackets above match those obtained using the EM algorithm. Hence, G^2 is minimum for $\hat{\beta}_{.2} = 0$ in each case, which implies that boundary solutions are given by $\hat{\beta}_{.2} = 0$ or equivalently $\hat{\pi}_{+2+2} = 0$. This result is consistent with points 2 and 3 of Proposition 3.1. Further, it is the exact form of boundary solutions that we obtain on fitting Models [M2], [M4] and [M5] to the data in Table 4 using the EM algorithm.

5. CONDITIONS FOR THE OCCURRENCE OF BOUNDARY SOLUTIONS

In this section, we discuss sufficient conditions and also propose necessary conditions for the occurrence of boundary solutions in two-way incomplete tables with both variables missing. We show that the sufficient conditions are not necessary, which disproves a conjecture made by [7]. Further, we prove that the proposed necessary conditions are not sufficient. Both sets of conditions are simple to verify since they involve only the observed cell counts in the tables. The sufficient conditions and the necessary conditions are of practical utility in identifying the occurrence and non-occurrence, respectively of boundary solutions in such tables.

5.1. Sufficient conditions for the occurrence of boundary solutions

Following [9], define the four odds based on the observed (joint/marginal) cell counts for any pair (j, j') of Y_2 :

$$(5.1) \quad \begin{aligned} \nu_i(j, j') &= \frac{\hat{\pi}_{ij11}}{\hat{\pi}_{ij'11}}, & \nu_n(j, j') &= \min_i \{\nu_i(j, j')\}, & \nu_m(j, j') &= \max_i \{\nu_i(j, j')\}, \\ \nu(j, j') &= \frac{y_{+j21}}{y_{+j'21}}. \end{aligned}$$

Similarly, for a given pair (i, i') of Y_1 , define the four odds using the observed cell counts:

$$(5.2) \quad \begin{aligned} \omega_j(i, i') &= \frac{\hat{\pi}_{ij11}}{\hat{\pi}_{i'j11}}, & \omega_n(i, i') &= \min_j \{\omega_j(i, i')\}, & \omega_m(i, i') &= \max_j \{\omega_j(i, i')\}, \\ \omega(i, i') &= \frac{y_{i+12}}{y_{i'+12}}. \end{aligned}$$

Note that $\nu_i(j, j')$ and $\omega_j(i, i')$ are called the response odds, while $\nu(j, j')$ and $\omega(i, i')$ are called the nonresponse odds. Using the MLE's of $\{\pi_{ij11}\}$ under Models [M1]–[M5] (see pp. 647–648 of [2]), we deduce that $\nu_i(j, j') = \frac{y_{ij11}}{y_{ij'11}}$ and $\omega_j(i, i') = \frac{y_{ij11}}{y_{i'j11}}$, which involve only the fully observed counts.

Theorem 1 of [9] deals with sufficient conditions for the occurrence of boundary solutions only under Model [M5]. However, in the next result, we provide such conditions for the occurrence of boundary solutions under Models [M1]–[M5]. Also, we provide a proof which is similar to that of Theorem 1 of [9], but we give direct arguments, which are different from the contrapositive ones used by [9].

Theorem 5.1. *Consider the following conditions for an $I \times I \times 2 \times 2$ contingency table:*

1. $\nu(j, j') \notin (\nu_n(j, j'), \nu_m(j, j'))$ for at least one pair (j, j') of Y_2 ,
2. $\omega(i, i') \notin (\omega_n(i, i'), \omega_m(i, i'))$ for at least one pair (i, i') of Y_1 .

Then we have the following:

- (a) *Boundary solutions in NMAR models for only Y_1 (Models [M1] and [M3]) occur if Condition 1 holds.*
- (b) *Boundary solutions in NMAR models for only Y_2 (Models [M2] and [M4]) occur if Condition 2 holds.*
- (c) *Boundary solutions in the NMAR model for both Y_1 and Y_2 (Model [M5]) occur if Condition 1 or Condition 2 holds.*

Proof: See Appendix A.2. □

5.2. The sufficient conditions are not necessary

The next example shows that the sufficient conditions for the occurrence of boundary solutions mentioned in Theorem 5.1 are not necessary. This result has not been discussed in the literature earlier. In fact, [7] proved that the above conditions are both necessary and sufficient for a $2 \times 2 \times 2 \times 2$ incomplete table. They conjectured that a similar result would hold for general two-way incomplete tables as well.

Example 5.1. Consider Table 5 discussed in [9], which cross-classifies data on bone mineral density (Y_1) and family income (Y_2) in a $3 \times 3 \times 2 \times 2$ incomplete table. Both variables Y_1 and Y_2 have three levels. The total count is 2998 out of which data on Y_1 and Y_2 are available for 1844 persons, data on Y_1 only for 231 persons, data on Y_2 only for 878 persons, and data on neither of them for 45 persons.

Table 5: Bone mineral density (Y_1) and family income (Y_2).

		$R_2 = 1$			$R_2 = 2$
		$Y_2 = 1$	$Y_2 = 2$	$Y_2 = 3$	Missing
$R_1 = 1$	$Y_1 = 1$	621	290	284	135
	$Y_1 = 2$	260	131	117	69
	$Y_1 = 3$	93	30	18	27
$R_1 = 2$	Missing	456	156	266	45

Now, we introduce the following changes corresponding to supplementary margins in Table 5: $266 \rightarrow 125$, $69 \rightarrow 60$ and $27 \rightarrow 20$. The modified table is shown in Table 6.

Table 6: Modified Table 5.

		$R_2 = 1$			$R_2 = 2$
		$Y_2 = 1$	$Y_2 = 2$	$Y_2 = 3$	Missing
$R_1 = 1$	$Y_1 = 1$	621	290	284	135
	$Y_1 = 2$	260	131	117	60
	$Y_1 = 3$	93	30	18	20
$R_1 = 2$	Missing	456	156	125	45

From Table 6, $\nu(1, 2) = 456/156 = 2.92$, $\nu(1, 3) = 456/125 = 3.65$, $\nu(2, 3) = 156/125 = 1.25$, $\omega(1, 2) = 135/60 = 2.25$, $\omega(1, 3) = 135/20 = 6.75$ and $\omega(2, 3) = 60/20 = 3.00$. Let $I_\nu(j, j') = (\nu_n(j, j'), \nu_m(j, j'))$ and $I_\omega(i, i') = (\omega_n(i, i'), \omega_m(i, i'))$. Then from Table 6, it can be shown that $\nu(1, 2) \in I_\nu(1, 2) = (260/131, 93/30)$, $\nu(1, 3) \in I_\nu(1, 3) = (621/284, 93/18)$, $\nu(2, 3) \in I_\nu(2, 3) = (290/284, 30/18)$, $\omega(1, 2) \in I_\omega(1, 2) = (290/131, 284/117)$, $\omega(1, 3) \in I_\omega(1, 3) = (621/93, 284/18)$ and $\omega(2, 3) \in I_\omega(2, 3) = (260/93, 117/18)$ so that the sufficient conditions for the occurrence of boundary solutions in Theorem 5.1 are not satisfied. The MLE's of the parameters obtained on fitting Models [M1]–[M5] in various subtables of Table 6 are shown in Table 7.

Table 7: MLE's of parameters in subtables of Table 6.

Subtable	NMAR model	MLE's	Boundary solutions
Y_1	[M1]	$\hat{\alpha}_1 = 0.6556$, $\hat{\alpha}_2 = -1.0537$, $\hat{\alpha}_3 = 3.4109$	$\hat{\pi}_{2+2+} = 0$
Y_2	[M2]	$\hat{\beta}_{.1} = 0.1355$, $\hat{\beta}_{.2} = 0.3420$, $\hat{\beta}_{.3} = -0.1846$	$\hat{\pi}_{+3+2} = 0$
$Y_1 Y_2$	[M1]	$\hat{\alpha}_1 = 0.6556$, $\hat{\alpha}_2 = -1.0537$, $\hat{\alpha}_3 = 3.4109$	$\hat{\pi}_{2+2+} = 0$
	[M3]	$\hat{\alpha}_1 = 0.6534$, $\hat{\alpha}_2 = -1.0551$, $\hat{\alpha}_3 = 3.4874$	$\hat{\pi}_{2+2+} = 0$
$Y_1 Y_2$	[M2]	$\hat{\beta}_{.1} = 0.1355$, $\hat{\beta}_{.2} = 0.3420$, $\hat{\beta}_{.3} = -0.1846$	$\hat{\pi}_{+3+2} = 0$
	[M4]	$\hat{\beta}_{.1} = 0.1421$, $\hat{\beta}_{.2} = 0.3289$, $\hat{\beta}_{.3} = -0.1712$	$\hat{\pi}_{+3+2} = 0$
$Y_1 Y_2$	[M5]	$\hat{\alpha}_1 = 0.6534$, $\hat{\alpha}_2 = -1.0551$, $\hat{\alpha}_3 = 3.4874$	$\hat{\pi}_{2+2+} = 0$
		$\hat{\beta}_{.1} = 0.1421$, $\hat{\beta}_{.2} = 0.3289$, $\hat{\beta}_{.3} = -0.1712$	$\hat{\pi}_{+3+2} = 0$

From Table 7, note that in each subtable, at least one of $\hat{\alpha}_i$ and $\hat{\beta}_{.j}$ is negative, which imply that boundary solutions occur. The forms of boundary solutions under the Models [M1]–[M5] are also the same as described in Section 3. This shows that for an $I \times J \times 2 \times 2$ incomplete table, where $I, J \geq 3$, the sufficient conditions for the occurrence of boundary solutions under Models [M1]–[M5] in Theorem 5.1 are not necessary.

5.3. Necessary conditions for the occurrence of boundary solutions

We next state below a result due to [6], which will be used later to obtain a result on the occurrence of boundary solutions.

Lemma 5.1. Suppose $A = (a_{ij})$ is a matrix with $a_{ij} \geq 0$ for $i \neq j = 1, 2, \dots, n$ and $a_{ii} > 0$. Also, let $\mathbf{b} = (b_j)$, where $b_j > 0$ for $1 \leq j \leq n$. If

$$(5.3) \quad b_i > \sum_{j \neq i=1}^n a_{ij} \frac{b_j}{a_{jj}}, \quad \forall 1 \leq i \leq n,$$

then A is invertible and $A^{-1}\mathbf{b} > \mathbf{0}$.

Using Lemma 5.1, the next result provides necessary conditions for the occurrence of boundary solutions under Models [M1]–[M5] in square two-way incomplete tables.

Theorem 5.2. For an $I \times I \times 2 \times 2$ incomplete table, consider the following conditions:

1. $y_{+j21} \leq \sum_{i \neq j=1}^I \hat{\mu}_{ji11} \frac{y_{+i21}}{\hat{\mu}_{ii11}}$ for at least one $j = 1, 2, \dots, I$,
2. $y_{i+12} \leq \sum_{j \neq i=1}^I \hat{\mu}_{ij11} \frac{y_{j+12}}{\hat{\mu}_{jj11}}$ for at least one $i = 1, 2, \dots, I$,

where $\hat{\mu}_{ij11}$ is the MLE of μ_{ij11} . Also, let $\{\hat{\mu}_{ij11}\} > 0$, $\{y_{i+12}\} > 0$ and $\{y_{+j21}\} > 0$. Then we have the following:

- (a) If boundary solutions under Models [M1] and [M3] occur, then only Condition 1 holds.
- (b) If boundary solutions under Models [M2] and [M4] occur, then only Condition 2 holds.
- (c) If boundary solutions under the Model [M5] occur, then Condition 1 or Condition 2 holds.

Proof: See Appendix A.3. □

Henceforth, we denote $A = (a_{ij}) = (\hat{\mu}_{ij11})$, $\mathbf{b} = (b_j) = (y_{+j21})$ and $\mathbf{b}^* = (b_i^*) = (y_{i+12})$ for $1 \leq i \leq I$, $1 \leq j \leq I$. The example below is an application of Theorem 5.2.

Example 5.2. From Table 6 in Example 5.1, we have the following:

$$A = \begin{pmatrix} 621 & 290 & 284 \\ 260 & 131 & 117 \\ 93 & 30 & 18 \end{pmatrix}, \quad \mathbf{b} = (456, 156, 125), \quad \mathbf{b}^* = (135, 60, 20).$$

The MLE's $\hat{\boldsymbol{\alpha}} = (\hat{\alpha}_i)$ and $\hat{\boldsymbol{\beta}} = (\hat{\beta}_j)$ under Model [M5] satisfy respectively the systems $A^T \boldsymbol{\alpha} = \mathbf{b}$ from (A.11) and $A \boldsymbol{\beta} = \mathbf{b}^*$ from (A.12) for $i, j = 1, 2, 3$. From Table 7, we observe that if Model [M5] is fitted to the data in Table 6, then we obtain $\hat{\alpha}_2 < 0$ and $\hat{\beta}_3 < 0$, that is, boundary solutions occur. Now we need to verify if both Conditions 1 and 2 of Theorem 5.2 hold. For the matrix A^T and the vector \mathbf{b} , we have

$$\begin{aligned} 456 &< a_{12} \times \frac{b_2}{a_{22}} + a_{13} \times \frac{b_3}{a_{33}} = 260 \times \frac{156}{131} + 93 \times \frac{125}{18} = 955.4516, \\ 156 &< a_{21} \times \frac{b_1}{a_{11}} + a_{23} \times \frac{b_3}{a_{33}} = 290 \times \frac{456}{621} + 30 \times \frac{125}{18} = 421.2802, \\ 125 &< a_{31} \times \frac{b_1}{a_{11}} + a_{32} \times \frac{b_2}{a_{22}} = 284 \times \frac{456}{621} + 117 \times \frac{156}{131} = 347.8693, \end{aligned}$$

so that Condition 1 in Theorem 5.2 is satisfied. Also, for the matrix A and the vector \mathbf{b}^* , we have

$$\begin{aligned} 135 &< a_{12} \times \frac{b_2^*}{a_{22}} + a_{13} \times \frac{b_3^*}{a_{33}} = 290 \times \frac{60}{131} + 284 \times \frac{20}{18} = 448.38, \\ 60 &< a_{21} \times \frac{b_1^*}{a_{11}} + a_{23} \times \frac{b_3^*}{a_{33}} = 260 \times \frac{135}{621} + 117 \times \frac{20}{18} = 186.5217, \\ 20 &< a_{31} \times \frac{b_1^*}{a_{11}} + a_{32} \times \frac{b_2^*}{a_{22}} = 93 \times \frac{135}{621} + 30 \times \frac{60}{131} = 33.9578, \end{aligned}$$

so that Condition 2 in Theorem 5.2 is satisfied. Further, from Table 7, we observe that boundary solutions also occur if Models [M1]–[M4] are fitted to data in Table 6. Then only Condition 1 is satisfied if boundary solutions under [M1] and [M3] occur, while only Condition 2 is satisfied if boundary solutions under [M2] and [M4] occur. This is because the MLE $\hat{\boldsymbol{\alpha}} = (\hat{\alpha}_i)$ under Models [M1] and [M3] satisfies the system $A^T \boldsymbol{\alpha} = \mathbf{b}$, while the MLE $\hat{\boldsymbol{\beta}} = (\hat{\beta}_j)$ under Models [M2] and [M4] satisfies the system $A\boldsymbol{\beta} = \mathbf{b}^*$.

5.4. The necessary conditions are not sufficient

The next example shows that the necessary conditions for the occurrence of boundary solutions in Theorem 5.2 are not sufficient.

Example 5.3. In Example 5.2, replace 456 by 366 in \mathbf{b} and 20 by 15 in \mathbf{b}^* so that $\mathbf{b} = (366, 156, 125)$ and $\mathbf{b}^* = (135, 60, 15)$ now. For the matrix A^T and the vector \mathbf{b} , we have

$$\begin{aligned} 366 &< a_{12} \times \frac{b_2}{a_{22}} + a_{13} \times \frac{b_3}{a_{33}} = 260 \times \frac{156}{131} + 93 \times \frac{125}{18} = 955.4516, \\ 156 &< a_{21} \times \frac{b_1}{a_{11}} + a_{23} \times \frac{b_3}{a_{33}} = 290 \times \frac{366}{621} + 30 \times \frac{125}{18} = 379.2512, \\ 125 &< a_{31} \times \frac{b_1}{a_{11}} + a_{32} \times \frac{b_2}{a_{22}} = 284 \times \frac{366}{621} + 117 \times \frac{156}{131} = 306.7099, \end{aligned}$$

so that Condition 1 in Theorem 5.2 is satisfied. Also, for the matrix A and the vector \mathbf{b}^* , we have

$$\begin{aligned} 135 &< a_{12} \times \frac{b_2^*}{a_{22}} + a_{13} \times \frac{b_3^*}{a_{33}} = 290 \times \frac{60}{131} + 284 \times \frac{15}{18} = 369.4911, \\ 60 &< a_{21} \times \frac{b_1^*}{a_{11}} + a_{23} \times \frac{b_3^*}{a_{33}} = 260 \times \frac{135}{621} + 117 \times \frac{15}{18} = 154.0217, \\ 15 &< a_{31} \times \frac{b_1^*}{a_{11}} + a_{32} \times \frac{b_2^*}{a_{22}} = 93 \times \frac{135}{621} + 30 \times \frac{60}{131} = 33.9578, \end{aligned}$$

so that Condition 2 in Theorem 5.2 is satisfied. Now, when we solve the system $A^T \boldsymbol{\alpha} = \mathbf{b}$, then we obtain the MLE's $\hat{\alpha}_1 = 0.0133$, $\hat{\alpha}_2 = 0.7796$ and $\hat{\alpha}_3 = 1.6671$. So, there are no boundary solutions under Model [M3]. Similarly, the system $A\boldsymbol{\beta} = \mathbf{b}^*$ yields the MLE's $\hat{\beta}_1 = 0.041$, $\hat{\beta}_2 = 0.3655$ and $\hat{\beta}_3 = 0.0126$, that is, there are no boundary solutions under Model [M4]. Since the MLE's in Model [M5] satisfy both the systems $A^T \boldsymbol{\alpha} = \mathbf{b}$ and $A\boldsymbol{\beta} = \mathbf{b}^*$, there are no boundary solutions under [M5] as well. Similar results hold for Models [M1] and [M2]. Hence, the conditions in Theorem 5.2 are not sufficient for the occurrence of boundary solutions under Models [M1]–[M5].

5.5. Importance of the necessary conditions

Here, we discuss additional details about Theorem 5.2 and discuss its simplicity and effectiveness.

From Theorem 5.2, note that if $\{y_{i+12}\}$, $\{y_{+j21}\}$, and/or $\{\hat{\mu}_{ii11}\}$ are large, then Conditions 1 and 2 may not hold. Indeed, if the inequalities in Conditions 1 and 2 are reversed for all $1 \leq i \leq I$ and $1 \leq j \leq I$, then from statements (a), (b) and (c) of Theorem 5.2, boundary solutions do not occur on fitting Models [M1]–[M5] in an $I \times I \times 2 \times 2$ incomplete table.

It is known that when boundary solutions occur, perfect fit models (here Models [M3], [M4] and [M5]) cannot reproduce the observed counts, indicating poor fit and imprecision of the parameter estimates. The MLE's of the parameters under NMAR models lie on the boundary of the parameter space and the log likelihood function tends to be flat, which makes derivation of the MLE's computationally intensive. Also, the corresponding covariance matrix has unreasonable eigenvalues (close to either zero or negative), which implies the estimated standard errors for some parameter estimates are large. Hence, for model selection, we prefer NMAR models which don't yield boundary solutions upon fitting them to the given data.

Theorem 5.1 provides conditions, which help us identify the occurrence of boundary solutions. However, boundary solutions may occur under some NMAR models if any of the sufficient conditions in Theorem 5.1 does not hold. This implies that Theorem 5.1 cannot always provide us the set of plausible NMAR models for model selection. However, note that Theorem 5.2 is very useful in this regard since it gives us an insight into verifying the non-occurrence of boundary solutions under each of the NMAR models [M1]–[M5]. That is, if any of the necessary conditions in Theorem 5.2 does not hold, then we know for sure that boundary solutions do not occur. This always helps us to obtain the list of candidate NMAR models suitable for fitting the given data. Hence, Theorem 5.2 is more reliable than Theorem 5.1 for the purpose of model selection in square two-way incomplete tables.

The non-boundary MLE's of μ_{ij11} are $\hat{\mu}_{ij11} = \frac{y_{ij11}y_{i+1+}y_{++11}}{y_{i+11}y_{++1+}}$ under Model [M1], $\hat{\mu}_{ij11} = \frac{y_{ij11}y_{+j+1}y_{++11}}{y_{+j11}y_{++1+}}$ under Model [M2], and $\hat{\mu}_{ij11} = y_{ij11}$ under Models [M3], [M4] and [M5] (see pp. 647–648 of [2]), which involve only the observed cell counts and their sums. Hence, from Theorem 5.2, there is no need to solve any system of likelihood equations, use the EM algorithm or compute odds (based on the observed (joint/marginal) cell counts) to check for the non-occurrence of boundary solutions in an $I \times I \times 2 \times 2$ incomplete table.

Remark 5.1. If $A_D = \text{diag}(a_{11}, \dots, a_{ii})$, then from [6], the solutions $\boldsymbol{\alpha} = (\alpha_i)$ of the system $A^T \boldsymbol{\alpha} = \mathbf{b}$ may be obtained iteratively as follows:

$$(5.4) \quad \begin{aligned} \alpha^{(0)} &= A_D^{-1} \mathbf{b} \\ \alpha^{(n+1)} &= \alpha^{(n)} + A_D^{-1} (\mathbf{b} - A^T \alpha^{(n)}), \quad n = 0, 1, 2, \dots \end{aligned}$$

Similarly, the solutions $\boldsymbol{\beta} = (\beta_j)$ of the system $A\boldsymbol{\beta} = \mathbf{b}^*$ may be obtained iteratively as follows:

$$(5.5) \quad \begin{aligned} \beta^{(0)} &= A_D^{-1} \mathbf{b}^* \\ \beta^{(n+1)} &= \beta^{(n)} + A_D^{-1} (\mathbf{b}^* - A\beta^{(n)}), \quad n = 0, 1, 2, \dots \end{aligned}$$

Both the sequences (5.4) and (5.5) converge to the solutions of the respective systems.

6. CONCLUSIONS

In this paper, we have discussed the problem of boundary solutions that occur under various NMAR models for an $I \times J \times 2 \times 2$ table. We formally define boundary solutions for such a table and provide a result (Proposition 3.1) that theoretically connects and justifies various forms of these solutions under alternative parametrizations of the missing data models. This eliminates the need of using the EM algorithm (see pp. 39–40 of [9]) to empirically obtain the forms of the solutions in two-way incomplete tables. The above result is then used to improve a claim in [2] regarding the occurrence of boundary solutions. We give the precise forms of such solutions by just noting the corresponding level(s) of the variable(s) in the table, which reduces computational burden.

As discussed earlier, boundary solutions pose a lot of problems for estimation and inference under NMAR models in incomplete tables. Hence, it is important to investigate sufficient and necessary conditions for their occurrence in such tables. We have provided a result (Theorem 5.1) on sufficient conditions for the occurrence of boundary solutions in an $I \times J \times 2 \times 2$ table. While [9] consider only Model [M5], we consider Models [M1]–[M5] in Theorem 5.1. We use a similar approach but give direct arguments instead of contrapositive ones used in Theorem 1 of [9] for proving Theorem 5.1. [7] conjectured that these conditions would also be necessary for general two-way incomplete tables. However, we show by a counterexample that this is not the case for $I, J \geq 3$, thereby disproving the conjecture.

We have also established necessary conditions in Theorem 5.2 for the occurrence of boundary solutions in an $I \times J \times 2 \times 2$ table, which have not been discussed in the literature so far. As discussed in Section 5.5, these conditions are of practical utility to identify the non-occurrence of boundary solutions and hence for model selection. However, we show by a counterexample that these conditions are not sufficient. Note that a major advantage of the proposed sufficient conditions and necessary conditions is that they depend only on the observed cell counts in the table or their sums. As mentioned in [9], this makes the verification process much easier, and avoids using the EM algorithm or solving likelihood equations. Finally, all the above results are illustrated using six data analysis examples. It would be helpful to obtain a set of conditions involving only the observed cell counts, which are sufficient as well as necessary for the occurrence of boundary solutions in two-way incomplete tables with both variables missing.

APPENDIX

A.1. Proof of Proposition 3.1

From Definition 3.1, it follows that if boundary solutions occur under the Models [M1]–[M5], then the MLE's of the cell probabilities except some of the nonresponse ones are all non-zero. On substituting $k = l = 1$ (for response cell probabilities) in the above models and using the parameter constraints, we can then deduce that the MLE's of the constant, the main effects and the association terms between Y_i 's, between R_i 's, and between Y_i and R_j for $i \neq j$ are all finite. This is because non-zero terms (response cell probabilities) on the LHS of the log-linear models imply that the log-linear parameters on the RHS are finite.

Consider part 1 first. For the Models [M1] and [M3], the log-linear parameters modelling the non-ignorable nonresponse (NMAR) mechanism of Y_1 are $\lambda_{R_1}(k)$ and $\lambda_{Y_1 R_1}(i, k)$. If boundary solutions occur, then they are of the form $\hat{\pi}_{ij2l} = 0$ (see point 1 of Definition 3.1), which implies $\hat{\lambda}_{Y_1 R_1}(i, 2) = -\infty$ for at least one i since the other parameters are finite as mentioned above. Then under Model [M1], we have

$$\begin{aligned} \hat{\pi}_{i+2+} &= \sum_{j,l} \hat{\pi}_{ij2l} \\ &= \frac{1}{N} \sum_{j,l} \exp\left\{ \hat{\lambda} + \hat{\lambda}_{Y_1}(i) + \hat{\lambda}_{Y_2}(j) + \hat{\lambda}_{R_1}(2) + \hat{\lambda}_{R_2}(l) + \hat{\lambda}_{Y_1 R_1}(i, 2) \right. \\ &\quad \left. + \hat{\lambda}_{Y_1 Y_2}(i, j) + \hat{\lambda}_{R_1 R_2}(2, l) \right\} \\ &= 0 \end{aligned}$$

for at least one i . Conversely, we have

$$\begin{aligned} \hat{\pi}_{i+2+} = 0 \text{ (for at least one } i) &\implies \\ &\implies \sum_{j,l} \exp\left\{ \hat{\lambda} + \hat{\lambda}_{Y_1}(i) + \hat{\lambda}_{Y_2}(j) + \hat{\lambda}_{R_1}(2) + \hat{\lambda}_{R_2}(l) + \hat{\lambda}_{Y_1 R_1}(i, 2) \right. \\ &\quad \left. + \hat{\lambda}_{Y_1 Y_2}(i, j) + \hat{\lambda}_{R_1 R_2}(2, l) \right\} = 0 \\ &\implies \hat{\lambda}_{Y_1 R_1}(i, 2) = -\infty \text{ for at least one } i, \end{aligned}$$

so that $\hat{\lambda}_{Y_1 R_1}(i, 2) = -\infty \Leftrightarrow \hat{\pi}_{i+2+} = 0$ for at least one i under Model [M1]. The same can be shown for Model [M3]. Under Models [M1] and [M3], $a_{ij} = \exp[2\{\lambda_{R_1}(2) + \lambda_{Y_1 R_1}(i, 2) + \lambda_{R_1 R_2}(2, 1)\}]$. Since a_{ij} depends only on i , we have $a_{ij} = \alpha_{i.}$. It is clear that $\hat{\alpha}_{i.} = 0 \Leftrightarrow \hat{\lambda}_{Y_1 R_1}(i, 2) = -\infty$. Also, note that by definition of a_{ij} , if $\hat{\alpha}_{i.} = 0$ for all $1 \leq i \leq I$, then $y_{+j21} = 0$ for all $1 \leq j \leq J$, which is a contradiction since supplementary margins are assumed to be positive. Hence, under Models [M1] and [M3], boundary solutions are given by $\hat{\lambda}_{Y_1 R_1}(i, 2) = -\infty \Leftrightarrow \hat{\pi}_{i+2+} = 0 \Leftrightarrow \hat{\alpha}_{i.} = 0$ for at least one and at most $(I - 1)$ values of Y_1 .

Consider part 2 now. Under Models [M2] and [M4], the log-linear parameters modelling the NMAR mechanism of Y_2 are $\lambda_{R_2}(l)$ and $\lambda_{Y_2 R_2}(j, l)$. Also, $b_{ij} = \exp[2\{\lambda_{R_2}(2) + \lambda_{Y_2 R_2}(j, 2) + \lambda_{R_1 R_2}(1, 2)\}]$. Since b_{ij} depends only on j , we have $b_{ij} = \beta_{.j}$. Then it can be shown similarly as above that boundary solutions in this case are given by $\hat{\lambda}_{Y_2 R_2}(j, 2) = -\infty \Leftrightarrow \hat{\pi}_{+j+2} = 0 \Leftrightarrow \hat{\beta}_{.j} = 0$ for at least one and at most $(J - 1)$ values of Y_2 .

Finally, consider part 3. Under Model [M5], the log-linear parameters modelling the NMAR mechanisms of Y_1 and Y_2 are $\lambda_{R_1}(k)$, $\lambda_{R_2}(l)$, $\lambda_{Y_1 R_1}(i, k)$ and $\lambda_{Y_2 R_2}(j, l)$. The proof for the form of boundary solutions under Model [M5] follows on similar lines as for Models [M1]–[M4] shown above.

A.2. Proof of Theorem 5.1

From [2], the MLE's $\hat{\alpha}_i$ under the NMAR model for only Y_1 (Models [M1] and [M3]) satisfy

$$(A.1) \quad \sum_i N \hat{\pi}_{ij11} \hat{\alpha}_i = y_{+j21}, \quad \forall 1 \leq j \leq I,$$

while the MLE's $\hat{\beta}_j$ under the NMAR model for only Y_2 (Models [M2] and [M4]) satisfy

$$(A.2) \quad \sum_j N \hat{\pi}_{ij11} \hat{\beta}_j = y_{i+12}, \quad \forall 1 \leq i \leq I.$$

The MLE's $\hat{\alpha}_i$ and $\hat{\beta}_j$ under the NMAR model for both Y_1 and Y_2 (Model [M5]) satisfy both (A.1) and (A.2). Note that boundary solutions in Models [M1] and [M3] occur if $\hat{\alpha}_i \leq 0$ for at least one and at most $(I - 1)$ values of Y_1 , while boundary solutions in Models [M2] and [M4] occur if $\hat{\beta}_j \leq 0$ for at least one and at most $(I - 1)$ values of Y_2 . Also note that boundary solutions under [M5] occur if at least one of the following holds:

- (i) $\hat{\alpha}_i \leq 0$ for at least one and at most $(I - 1)$ values of Y_1 ,
- (ii) $\hat{\beta}_j \leq 0$ for at least one and at most $(I - 1)$ values of Y_2 .

From (5.1) and (A.1), we have

$$\nu(j, j') = \frac{y_{+j21}}{y_{+j'21}} = \frac{\sum_i \hat{\pi}_{ij11} \hat{\alpha}_i}{\sum_i \hat{\pi}_{ij'11} \hat{\alpha}_i},$$

$$(A.3) \quad \nu_m(j, j') - \nu(j, j') = \frac{\sum_{i \neq m_1} (\hat{\pi}_{m_1 j 11} \hat{\pi}_{ij'11} - \hat{\pi}_{m_1 j'11} \hat{\pi}_{ij11}) \hat{\alpha}_i}{\hat{\pi}_{m_1 j'11} \sum_i \hat{\pi}_{ij'11} \hat{\alpha}_i},$$

$$(A.4) \quad \nu(j, j') - \nu_n(j, j') = \frac{\sum_{i \neq n_1} (\hat{\pi}_{n_1 j'11} \hat{\pi}_{ij11} - \hat{\pi}_{n_1 j11} \hat{\pi}_{ij'11}) \hat{\alpha}_i}{\hat{\pi}_{n_1 j'11} \sum_i \hat{\pi}_{ij'11} \hat{\alpha}_i},$$

where m_1 and n_1 are the levels of Y_1 corresponding to $\nu_m(j, j')$ and $\nu_n(j, j')$ respectively. From (5.1), we get

$$(A.5) \quad \nu_n(j, j') = \frac{\hat{\pi}_{n_1 j 11}}{\hat{\pi}_{n_1 j'11}} < \nu_i(j, j') = \frac{\hat{\pi}_{ij11}}{\hat{\pi}_{ij'11}} < \nu_m(j, j') = \frac{\hat{\pi}_{m_1 j 11}}{\hat{\pi}_{m_1 j'11}}.$$

From (A.5), we have the following inequalities:

$$(A.6) \quad \hat{\pi}_{m_1 j 11} \hat{\pi}_{ij'11} > \hat{\pi}_{m_1 j'11} \hat{\pi}_{ij11}, \quad \hat{\pi}_{n_1 j'11} \hat{\pi}_{ij11} > \hat{\pi}_{n_1 j11} \hat{\pi}_{ij'11} \quad \text{for } i \neq m_1, n_1.$$

Consider part (a). Suppose Condition 1 holds, which implies that (A.3) and (A.4) are of opposite signs. Using this fact and (A.6), we observe that $\hat{\alpha}_i < 0$ for at least one and at most $(I - 1)$ values of Y_1 , that is, boundary solutions of the form $\hat{\pi}_{i+2+} = 0$ occur.

Again from (5.2) and (A.2), we have

$$\omega(i, i') = \frac{y_{i+12}}{y'_{i'+12}} = \frac{\sum_j \hat{\pi}_{ij11} \hat{\beta}_{.j}}{\sum_j \hat{\pi}_{i'j11} \hat{\beta}_{.j}},$$

$$(A.7) \quad \omega_m(i, i') - \omega(i, i') = \frac{\sum_{j \neq m_2} (\hat{\pi}_{im_211} \hat{\pi}_{i'j11} - \hat{\pi}_{i'm_211} \hat{\pi}_{ij11}) \hat{\beta}_{.j}}{\hat{\pi}_{i'm_211} \sum_i \hat{\pi}_{i'j11} \hat{\beta}_{.j}},$$

$$(A.8) \quad \omega(i, i') - \omega_n(i, i') = \frac{\sum_{j \neq n_2} (\hat{\pi}_{i'n_211} \hat{\pi}_{ij11} - \hat{\pi}_{in_211} \hat{\pi}_{i'j11}) \hat{\beta}_{.j}}{\hat{\pi}_{i'n_211} \sum_i \hat{\pi}_{i'j11} \hat{\beta}_{.j}},$$

where m_2 and n_2 are the levels of Y_2 corresponding to $\omega_m(i, i')$ and $\omega_n(i, i')$ respectively. From (5.2), we get

$$(A.9) \quad \omega_n(i, i') = \frac{\hat{\pi}_{in_211}}{\hat{\pi}_{i'n_211}} < \omega_j(i, i') = \frac{\hat{\pi}_{ij11}}{\hat{\pi}_{i'j11}} < \omega_m(i, i') = \frac{\hat{\pi}_{im_211}}{\hat{\pi}_{i'm_211}}.$$

From (A.9), we have the following inequalities:

$$(A.10) \quad \hat{\pi}_{m_2j11} \hat{\pi}_{i'j11} > \hat{\pi}_{m_2j'11} \hat{\pi}_{ij11}, \quad \hat{\pi}_{n_2j'11} \hat{\pi}_{ij11} > \hat{\pi}_{n_2j11} \hat{\pi}_{i'j11} \quad \text{for } j \neq m_2, n_2.$$

Now consider part (b). Assume Condition 2 holds, which implies that (A.7) and (A.8) are of opposite signs. Using this fact and (A.10), we observe that $\hat{\beta}_{.j} < 0$ for at least one and at most $(I - 1)$ values of Y_2 , that is, boundary solutions of the form $\hat{\pi}_{+j+2} = 0$ occur.

Finally consider part (c). Assume at least one of Conditions 1 and 2 holds. The cases when only Condition 1 holds or only Condition 2 holds follow from the proofs of part (a) and part (b), respectively. So it is sufficient here to assume both Conditions 1 and 2 hold. This implies, from part (a), $\hat{\alpha}_i < 0$ for at least one and at most $(I - 1)$ values of Y_1 , that is, boundary solutions of the form $\hat{\pi}_{i+2+} = 0$ occur. Also from part (b), we have $\hat{\beta}_{.j} < 0$ for at least one and at most $(I - 1)$ values of Y_2 , that is, boundary solutions of the form $\hat{\pi}_{+j+2} = 0$ occur. This completes the proof.

A.3. Proof of Theorem 5.2

From Theorem 5.1, the MLE's $\hat{\alpha}_i$ and $\hat{\beta}_{.j}$ under Model [M5] satisfy

$$(A.11) \quad \sum_i \hat{\mu}_{ij11} \hat{\alpha}_i = y_{+j21} \quad \text{for } j = 1, \dots, I,$$

$$(A.12) \quad \sum_j \hat{\mu}_{ij11} \hat{\beta}_{.j} = y_{i+12} \quad \text{for } i = 1, \dots, I.$$

Also, the MLE $\hat{\alpha}_i$ under Models [M1] and [M3] satisfy (A.11) only, while the MLE $\hat{\beta}_{.j}$ under Models [M2] and [M4] satisfy (A.12) only. Note that boundary solutions under [M5] occur if at least one of the following conditions hold:

- (i) $\hat{\alpha}_i \leq 0$ for at least one and at most $(I - 1)$ values of Y_1 ,
- (ii) $\hat{\beta}_{.j} \leq 0$ for at least one and at most $(I - 1)$ values of Y_2 .

Also, boundary solutions in Models [M1] and [M3] are given by only Condition (i), while boundary solutions in Models [M2] and [M4] are given by only Condition (ii). In Lemma 5.1, take $A = (\hat{\mu}_{ij11})$, $\mathbf{b} = (b_j) = (y_{+j21})$ and $\mathbf{b}^* = (b_i^*) = (y_{i+12})$ for $1 \leq i \leq I$, $1 \leq j \leq I$.

Then (A.11) may be written as $A^T\boldsymbol{\alpha} = \mathbf{b}$, while (A.12) may be written as $A\boldsymbol{\beta} = \mathbf{b}^*$, where $\boldsymbol{\alpha} = (\alpha_i)$ and $\boldsymbol{\beta} = (\beta_j)$. We prove Theorem 5.2 by contrapositive.

Consider part (a) first. Suppose Condition 1 in Theorem 5.2 does not hold. Then by Lemma 5.1, $\boldsymbol{\alpha} = (A^T)^{-1}\mathbf{b} > \mathbf{0}$. In other words, $\hat{\alpha}_i > 0$ for all $1 \leq i \leq I$, that is, boundary solutions under Models [M1] and [M3] do not occur.

Consider part (b) now. Assume Condition 2 in Theorem 5.2 does not hold. Then by Lemma 5.1, $\boldsymbol{\beta} = A^{-1}\mathbf{b}^* > \mathbf{0}$. In other words, $\hat{\beta}_j > 0$ for all $1 \leq j \leq I$, that is, boundary solutions under Models [M2] and [M4] do not occur.

Finally consider part (c). Assume both Conditions 1 and 2 in Theorem 5.2 do not hold. Then by Lemma 5.1, both $\hat{\alpha}_i > 0$ and $\hat{\beta}_j > 0$ for all $1 \leq i \leq I$, $1 \leq j \leq I$, that is, boundary solutions under Model [M5] do not occur.

Hence, the result follows.

ACKNOWLEDGMENTS

The research of S. Ghosh was supported by UGC, Govt. of India grant F.2-2/98 (SA-I). The authors are also grateful to the referees for carefully reading the manuscript and suggesting numerous improvements.

REFERENCES

- [1] BAKER, S.G. and LAIRD, N.M. (1988). Regression analysis for categorical variables with outcome subject to nonignorable nonresponse, *Journal of the American Statistical Association*, **83**, 62–69.
- [2] BAKER, S.G.; ROSENBERGER, W.F. and DERSIMONIAN, R. (1992). Closed-form estimates for missing counts in two-way contingency tables, *Statistics in Medicine*, **11**, 643–657.
- [3] CLARKE, P.S. (2002). On boundary solutions and identifiability in categorical regression with non-ignorable non-response, *Biometrical Journal*, **44**, 701–717.
- [4] CLARKE, P.S. and SMITH, P.W.F. (2005). On maximum likelihood estimation for log-linear models with non-ignorable non-responses, *Statistics and Probability Letters*, **73**, 441–448.
- [5] GHOSH, S. and VELLAISAMY, P. (2016). On the occurrence of boundary solutions in multi-dimensional incomplete tables, *Statistics and Probability Letters*, **119**, 63–75.
- [6] KAYKOBAD, M. (1985). Positive Solutions of Positive Linear Systems, *Linear Algebra and its Applications*, **64**, 133–140.
- [7] KIM, S. and PARK, Y. (2014). Power-linear models for incomplete contingency tables with nonignorable non-responses, *Statistics*, DOI: 10.1080/02331888.2013.869595.
- [8] LITTLE, J.A. and RUBIN, D.B. (2002). *Statistical Analysis with Missing Data*, second ed., Wiley, New York.
- [9] PARK, Y.; KIM, D. and KIM, S. (2014). Identification of the occurrence of boundary solutions in a contingency table with nonignorable nonresponse, *Statistics and Probability Letters*, **93**, 34–40.
- [10] SMITH, P.W.F.; SKINNER, C.J. and CLARKE, P.S. (1999). Allowing for non-ignorable nonresponse in the analysis of voting intention data, *Journal of the Royal Statistical Society: Series C*, **48**, 563–577.

DEPTH-BASED SIGNED-RANK TESTS FOR BIVARIATE CENTRAL SYMMETRY

Authors: SAKINEH DEGHAN

– Department of Statistics, Shahid Beheshti University,
Tehran, Iran
sa_dehghan@sbu.ac.ir

MOHAMMAD REZA FARIDROHANI

– Department of Statistics, Faculty of Mathematical Sciences,
Shahid Beheshti University, Tehran, Iran
m_faridrohani@sbu.ac.ir

Received: February 2017

Revised: October 2017

Accepted: October 2017

Abstract:

- In this paper, distribution-free, affine invariant, signed-rank test statistics are proposed for the hypothesis that a bivariate distribution is centrally symmetric about an arbitrary specified point. The proposed tests are based on the concept of data depth. However, our tests are inherently orthogonal invariant, an affine invariant version of them is provided by using Tyler's estimator of scatter. The limiting null distribution of proposed tests is derived and the performance of the proposed tests is evaluated through a Monte Carlo study. This study demonstrates that the tests always detect asymmetry and they are convenient to determine small departures from the null hypothesis with high power. Also it shows that the tests perform well comparing other procedures in the literature.

Key-Words:

- *affine invariance; central symmetry; depth function; distribution-free; Tyler's estimator of scatter.*

AMS Subject Classification:

- 49A05, 78B26.

1. INTRODUCTION

Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ denote independent copies of the bivariate random vector $\mathbf{X} = (X_1, X_2)^\top$ from a continuous bivariate population. One problem which has been considered in the literature is to test whether the distribution is symmetric about an unknown center against the alternative that the symmetry is lost (Heathcote *et al.* [16], Koltchinskii and Li [22], Neuhaus and Zhu [33], Manzotti *et al.* [30] and Henze *et al.* [17]). Moreover, in the univariate case, we can mention to Cassart *et al.* [6]. Unlike the univariate case, there are several concepts of multivariate symmetry including spherical, elliptical, central and angular symmetry. It is worth noting that the mentioned arrangement of the multivariate symmetry concepts are ordered in increasing generality. To read more about different types of multivariate symmetry see Serfling [38].

A different problem is the testing of the hypothesis that the bivariate distribution is symmetric about a known center $\boldsymbol{\mu}_0$ against the alternative that the distribution is symmetric about $\boldsymbol{\mu} \neq \boldsymbol{\mu}_0$. There is a substantial literature for this problem. Under the multivariate normality assumption, it is common to use Hotelling's T^2 test [20]. A multivariate affine-invariant sign test based on counts called interdirections has been presented by Randles [35]. In the sequence, Peter and Randles [41] based on the notion of interdirection, provided affine invariant signed rank test and signed sum test, respectively. Optimal affine invariant tests based on interdirections and pseudo-Mahalanobis ranks have been developed by Hallin and Paindaveine [12]. Hallin and Paindaveine [11] also presented an alternative version of these procedures in which interdirections are replaced by angles between the observations standardized via Tyler's estimator of scatter [40]. Mottonen and Oja [32] developed the tests based on spatial signs and ranks. Hettmansperger *et al.* [18] and Hettmansperger *et al.* [19] extended the bivariate tests of Brown and Hettmansperger [5] to the multivariate case. An affine invariant sign test by applying the Tyler's transformation on data points has been presented by Randles [36]. The affine invariant signed rank test, modified from sign test of Randles [36], was suggested by Mahfoud and Randles [29]. The tests described in preceding paragraph can serve as important preliminaries before applying these corresponding location tests. Moreover, there are several tests for testing of the hypothesis that the bivariate distribution is symmetric against not only location parameter but also regression and serial dependence alternatives e.g. Hallin and Paindaveine [13], [14] and [15].

Another problem that has received attention is to test whether the distribution is symmetric about known center $\boldsymbol{\mu}_0$ against the alternative that either the symmetry is lost or the location parameter is changed. Our paper deals with the latter problem. Indeed, the purpose of this paper is to develop affine invariant tests for testing the central symmetry of the bivariate distribution about a known center $\boldsymbol{\mu}_0$. Baringhaus [4] introduced the rotation invariant tests, for testing the spherical symmetry of the multivariate distribution about known center. For central symmetry that it is a weaker assumption than spherical and elliptical symmetry, the tests have been developed employing the empirical characteristic functions by Ghosh and Ruymgaart [10]. Aki [1] proposed a rotation invariant test based on the empirical distribution function. An extension of McWilliams' univariate run test (McWilliams [31]) into a test of bivariate central symmetry based on the depth function have been presented by Dyckerhoff *et al.* [8]. Although, this test is affine invariant, it suffers from low power in distinguishing most of the alternative hypotheses to central symmetry. Recently Einmahl and Gan [9] proposed two versions of a rotation invariant test based on empirical measures of opposite regions.

In this paper, we aim to propose test statistics for central symmetry in such a way that they would be affine invariant, distribution-free and have good power against alternatives to the null hypothesis. The test statistics are created based on sum of the signed-ranks where the sign and rank functions are determined through the depth function. Based on a given depth function, this procedure results in an orthogonal invariant test statistic. An affine invariant version of this test is provided by applying Tyler's transformation (Tyler [40]) on data points. The affine invariance property ensures that the performance of the test does not depend on the underlying coordinate system.

The word of depth has been used for the first time by Tukey [39] to introduce the halfspace depth function. In the sequence, different depth functions have been introduced and the multivariate data have been ordered as center-outward based on them. This center-outward ranking has been widely applied in multivariate nonparametric inference. Liu and Singh [27] presented a quality index and provided some multivariate rank tests for difference between two independent distributions based on it. In the following, a distribution-free test was presented based on both the depth function and the principal components by Rousseeuw [37] for the multivariate two-sample location-scale model. Based on DD plots (depth vs. depth plots) introduced by Liu *et al.* [26], two tests have been provided by Li and Liu [24] for location difference between two multivariate distributions. In addition, Liu and Singh [28] introduced some rank tests for multivariate scale difference between two or more independent populations. Depth-based run tests for bivariate central symmetry is introduced by Dyckerhoff *et al.* [8].

The remainder of this paper is organized as follows. In Section 2, we review briefly the concept of depth function and ranking based on it. The proposed test statistics will be described in Section 3 and the asymptotic properties of those are also investigated. Finally, in Section 4, a Monte Carlo study evaluates the finite sample performance of the proposed test statistics in accordance with other tests. All technical proofs are deferred to the Appendix.

2. DEPTH FUNCTION

Let \mathbf{X} be a p -dimensional random vector defined on a probability space (Ω, \mathcal{F}, P) . We denote F as a distribution function corresponding to P . A depth function associated with a distribution function F on \mathbb{R}^p is defined to provide a center-outward ordering of points of \mathbb{R}^p relative to F . Based on depth function, a corresponding notion of center or multidimensional median could be defined. The higher depth values refer to the points near to the center, whereas the lower values refer to the outer points of the center. A formal definition of "statistical depth function" is presented by Zuo and Serfling [42] as a function $D(\cdot, F): \mathbb{R}^p \rightarrow \mathbb{R}$ satisfying the following properties:

- P1.** Affine invariance: for any nonsingular $p \times p$ matrix \mathbf{A} and p -vector \mathbf{b} , $D(\mathbf{A}\mathbf{x} + \mathbf{b}, F_{\mathbf{A}\mathbf{x}+\mathbf{b}}) = D(\mathbf{x}, F)$.
- P2.** Maximality at center: if F is symmetric about $\boldsymbol{\theta}$ in some sense, then $D(\boldsymbol{\theta}, F) = \sup_{\mathbf{x} \in \mathbb{R}^p} D(\mathbf{x}, F)$.
- P3.** Monotonicity relative to deepest point: if $D(\boldsymbol{\theta}, F) \geq D(\mathbf{x}, F)$ for any $\mathbf{x} \in \mathbb{R}^p$ then $D(\boldsymbol{\theta} + \alpha(\mathbf{x} - \boldsymbol{\theta}), F) \geq D(\mathbf{x}, F)$ for each $\alpha \in [0, 1]$ and $\mathbf{x} \in \mathbb{R}^p$.
- P4.** Vanishing at infinity: as $\|\mathbf{x}\| \rightarrow \infty$, $D(\mathbf{x}, F) \rightarrow 0$.

Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be a random sample from p -dimensional distribution function F . The sample version of the depth function $D(\cdot, F)$ will be obtained by replacing F with the sample distribution F_n .

Remark 2.1. If the sample depth function $D(\cdot, F_n)$ satisfies property P1, then it will also be invariant under data-dependent nonsingular transformations.

Different depth functions have been proposed by some authors, which the definition of some of them that we deal with in this paper are given as follows.

Definition 2.1 (Tukey [39]). The halfspace depth of $\mathbf{x} \in \mathbb{R}^p$ with respect to F is defined as

$$HD(\mathbf{x}, F) = \inf_H \left\{ P(H) : H \text{ is a closed halfspace in } \mathbb{R}^p \text{ and } \mathbf{x} \in H \right\}$$

and the sample halfspace depth function is

$$HD(\mathbf{x}, F_n) = \frac{\min_{\|\mathbf{u}\|=1} \#\{i : \mathbf{u}^\top \mathbf{X}_i \leq \mathbf{u}^\top \mathbf{x}, i = 1, \dots, n\}}{n}.$$

Definition 2.2 (Liu [25]). The simplicial depth of \mathbf{x} with respect to F is defined as

$$SD(\mathbf{x}, F) = P_F(\mathbf{x} \in S[\mathbf{X}_1, \dots, \mathbf{X}_{p+1}]),$$

where $S[\mathbf{X}_1, \dots, \mathbf{X}_{p+1}]$ is a closed simplex with $\mathbf{X}_1, \dots, \mathbf{X}_{p+1}$ vertices. The sample version of $SD(\mathbf{x}, F)$ is given by the fraction of the sample random simplices containing the point \mathbf{x} .

Definition 2.3 (Liu [27]). The Mahalanobis depth of \mathbf{x} with respect to F is given by

$$MD(\mathbf{x}, F) = \frac{1}{1 + (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})},$$

where $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ are the mean vector and dispersion matrix of F distribution, respectively. The sample version of Mahalanobis depth is provided by replacing $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ with their sample estimates.

Additionally, some other depth functions have been introduced such as Oja depth (Oja [34]) and zonoid depth (Koshevov and Mosler [23]). A more recent proposal for data depth is the Monge–Kantorovich depth (Chernozhukov *et al.* [7]) based on the Monge–Kantorovich theory of measure transportation.

Now, we present the definition of center-outward ranking of data points.

Definition 2.4. Assume that $\mathbf{X}_1, \dots, \mathbf{X}_n$ is a random sample from distribution function F in \mathbb{R}^p . The center-outward rank \mathbf{X}_i within the sample $\mathbf{X}_1, \dots, \mathbf{X}_n$ is

$$\#\left\{ \mathbf{X}_j \in \{\mathbf{X}_1, \dots, \mathbf{X}_n\} : D(\mathbf{X}_j, F_n) \geq D(\mathbf{X}_i, F_n) \right\},$$

where F_n is the sample distribution function.

Thus, the center-outward ranking is defined in such a way that a larger rank is assigned to a more outlying point w.r.t. $\mathbf{X}_1, \dots, \mathbf{X}_n$. If there are no ties, rank 1 and rank n are assigned to the deepest point and the most outlying point, respectively.

3. THE PROPOSED TESTS

Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be independently and identically distributed as $\mathbf{X} = (X_1, X_2)^\top$, where \mathbf{X} has an arbitrary bivariate continuous distribution F . The null hypothesis of interest is that, the random vector \mathbf{X} has a distribution centrally symmetric about the known point $\boldsymbol{\mu}_0$. The random vector \mathbf{X} is centrally symmetric around $\boldsymbol{\mu}_0$ provided $\mathbf{X} - \boldsymbol{\mu}_0$ and $\boldsymbol{\mu}_0 - \mathbf{X}$ have the same distribution. Since it is assumed that the symmetry point is known, it is possible to take $\boldsymbol{\mu}_0 = 0$, without loss of generality. So, the hypothesis that the probability distribution is centrally symmetric about $\boldsymbol{\mu}_0$, reduces to the hypothesis $H_0: \mathbf{X} \stackrel{d}{=} -\mathbf{X}$, where $\stackrel{d}{=}$ denotes “equal in distribution”. We now describe the procedure for defining affine invariant tests. Let us look at tests that they are only invariant with respect to orthogonal transformations of the data in Subsection 3.1, and then proceed to provide our main affine-invariant tests in Subsection 3.2.

3.1. The orthogonal invariant tests

Let $D(\cdot, F)$ be a depth function on \mathbb{R}^2 associated with a distribution function F . Now, under the given depth function $D(\cdot, F)$, we derive a test statistic using depth-based ranks and signs of $\mathbf{X}_1, \dots, \mathbf{X}_n$. To define the proposed test statistic, we need to order the points $\mathbf{X}_1, \dots, \mathbf{X}_n$ in terms of the evidence they provide against the null hypothesis. To this end, we order the points $\mathbf{X}_1, \dots, \mathbf{X}_n$ as center-outward, such that the larger ranks correspond to the closer points to the null symmetry center and the smaller ranks correspond to the outer ones. Let F_n and F_n^s denote the sample distribution function of random sample $\mathbf{X}_1, \dots, \mathbf{X}_n$ and the symmetrized sample $(\pm \mathbf{X}_1, \dots, \pm \mathbf{X}_n)$, respectively. Employing property P2 of the depth function, to obtain center-outward rank of points relative to the null symmetry center instead of the median of $\mathbf{X}_1, \dots, \mathbf{X}_n$, the points are ordered based on $D(\cdot, F_n^s)$ rather than $D(\cdot, F_n)$. More precisely, define

$$(3.1) \quad R_i = \#\left\{ \mathbf{X}_j \in \{\mathbf{X}_1, \dots, \mathbf{X}_n\} : D(\mathbf{X}_j, F_n^s) \geq D(\mathbf{X}_i, F_n^s) \right\}, \quad i = 1, \dots, n.$$

If ties occur in this ranking, the ranks within each ties-class have been assigned based on increasing values at the corresponding index set of that. This assignment is allocated to induce invariance property on proposed test statistic.

The test statistic is sum of the signed-ranks of points. The sign of each bivariate point can be determined as the sign of the first or second component of it. Specifically, the sign of a bivariate point is equal to 1 if its first (or second) component is nonnegative and otherwise is equal to -1 . This definition of sign, leads to a test statistic which is not only noninvariant, but also it is not able to detect all different types of departures from the null hypothesis. Moreover, the sign of \mathbf{X}_i , $i = 1, \dots, n$, could be defined as the spatial sign vector $\mathbf{X}_i / \|\mathbf{X}_i\|$ with $\|\cdot\|$ denoting the Euclidean norm in \mathbb{R}^2 . By this definition of sign, the resulted test statistic is not strictly distribution-free. To overcome these limitations, we will determine sign of points based on a data-dependent line passing through the origin instead of the horizontal or vertical axis of the coordinate plane. In what follows, we will describe how to obtain this line.

Let sample median \mathbf{M}_n be a point among $\mathbf{X}_1, \dots, \mathbf{X}_n$ with maximum sample depth $D(\cdot, F_n^s)$. If there is more than one sample point with the highest depth value $D(\cdot, F_n^s)$, \mathbf{M}_n will be the point with minimum index among those data points. Let

$$\theta_{\mathbf{M}_n} = -\arctan\left(\frac{M_{n1}}{M_{n2}}\right) + \frac{\pi}{2}$$

be the angle between the bivariate vector $\mathbf{M}_n = (M_{n1}, M_{n2})^\top$ and the horizontal-axis. Note that $\theta_{\mathbf{M}_n} \in [0, \pi)$. Related point $\mathbf{Z}_{ni} = (Z_{ni1}, Z_{ni2})^\top$ is given by rotating \mathbf{X}_i counter-clockwise by angle $\frac{\pi}{2} - \theta_{\mathbf{M}_n}$, for all $i = 1, \dots, n$. Based on the sample depth function $D(\cdot, F_n)$, the proposed test statistic is defined as

$$(3.2) \quad T_{n,D} = \frac{6}{n(n+1)(2n+1)} \left(\sum_{i=1}^n \delta_{ni} R_i \right)^2,$$

where R_i is expressed in (3.1) and the random variable δ_{ni} is defined as

$$(3.3) \quad \delta_{ni} = \begin{cases} 1, & Z_{ni2} \geq 0, \\ -1, & Z_{ni2} < 0, \end{cases}$$

for all $i = 1, \dots, n$. The large values of the test statistic $T_{n,D}$ reject H_0 in favor of alternative hypothesis.

Note that the sign of bivariate points is determined based on a data-dependent line passing through the origin that is perpendicular to depth based median. Indeed, the reason for restricting to dimension two is that this procedure is employed to divide plane R^2 into two unique halfspaces based on two points (the origin and the depth based median), whereas by this procedure dividing hyperplane R^p ($p > 2$) into two unique halfspaces would not be possible.

In what follows, we present the desirable property of orthogonal invariance of $T_{n,D}$ and asymptotic distribution of $T_{n,D}$ under the null hypothesis is developed. The proofs are provided in the Appendix.

Theorem 3.1. *If the sample depth function $D(\cdot, F_n)$ satisfies property P1, then the test statistic $T_{n,D}$ will be invariant under orthogonal transformations; that is,*

$$T_{n,D}(\mathbf{X}_1, \dots, \mathbf{X}_n) = T_{n,D}(\mathbf{A}\mathbf{X}_1, \dots, \mathbf{A}\mathbf{X}_n)$$

for any 2×2 orthogonal matrix \mathbf{A} .

Theorem 3.2. *If the sample depth function satisfies property P1, then under the null hypothesis of centrally symmetric about 0, $T_{n,D}$ converges in distribution to a chi-square random variable with 1 degree of freedom.*

By applying this theorem, the null hypothesis will be rejected at level α when

$$T_{n,D} \geq \chi_{1,1-\alpha}^2,$$

where $\chi_{1,1-\alpha}^2$ denotes the $1-\alpha$ quantile of the chi-square distribution with 1 degree of freedom.

As mentioned in Theorem 3.2, the asymptotic null distribution of the test statistics presented here is chi-square with one degree of freedom. One would expect, for location alternatives, a chi-square with two degrees of freedom (the dimension of the information matrix for location). It should be remembered that the main object of this paper is proposing several test statistics for testing that the distribution is symmetric about a specified value against the alternative that either the symmetry is lost or the location parameter is changed. Indeed, this alternative is different from location alternatives.

In our proof of Theorem 3.2 we show that R_i 's, $i = 1, \dots, n$, are identically and uniformly distributed on the set $\{1, 2, \dots, n\}$ and δ_{ni} 's, $i = 1, \dots, n$, are i.i.d. random variables as distributed independently of R_i and taking the values 1 and -1 each with probability $1/2$. These traits immediately imply that under the null hypothesis and the conditions of Theorem 3.2, our test statistic $T_{n,D}$ is strictly distribution-free.

3.2. The affine invariant tests

As shown, Theorem 3.1 indicates that $T_{n,D}$ is orthogonal invariant. In this subsection, we would extend $T_{n,D}$ to be affine invariant, preserving the asymptotic behavior of $T_{n,D}$. To achieve the affine invariant version of the proposed test statistics, we can apply the Tyler's auxiliary transformation (Tyler [40]) on data points. Tyler [40] proposed the data-dependent $p \times p$ scatter matrix \mathbf{V}_n , that is a positive definite and symmetric matrix, satisfying $\text{trace}(\mathbf{V}_n) = p$ and

$$(3.4) \quad \frac{1}{n} \sum_{i=1}^n \left(\frac{\mathbf{\Gamma}_n \mathbf{X}_i}{\|\mathbf{\Gamma}_n \mathbf{X}_i\|} \right) \left(\frac{\mathbf{\Gamma}_n \mathbf{X}_i}{\|\mathbf{\Gamma}_n \mathbf{X}_i\|} \right)^\top = \frac{1}{p} \mathbf{I}_p,$$

where \mathbf{X}_i , $i = 1, \dots, n$, is a random vector in \mathbb{R}^p , $\mathbf{\Gamma}_n^\top \mathbf{\Gamma}_n = \mathbf{V}_n^{-1}$ such that $\mathbf{\Gamma}_n$ is an upper triangular nonsingular matrix with 1 on the first element on the diagonal and \mathbf{I}_p is the p -dimensional identity matrix. This scatter matrix is unique up to multiplication by a positive constant if the sample comes from a continuous p -dimensional distribution and $n > p(p-1)$ (Tyler [40]). An iterative computation scheme has been developed to compute this matrix by Randles [36].

We define $\mathbf{W}_{ni} = \mathbf{\Gamma}_n \mathbf{X}_i$, $i = 1, \dots, n$, and $F_{\mathbf{w}_n}^s$ as the sample distribution of the symmetrized sample $(\pm \mathbf{W}_{n1}, \dots, \pm \mathbf{W}_{nn})$. Let

$$(3.5) \quad R_{ni} = \#\left\{ \mathbf{W}_{nj} \in \{\mathbf{W}_{n1}, \dots, \mathbf{W}_{nn}\} : D(\mathbf{W}_{nj}, F_{\mathbf{w}_n}^s) \geq D(\mathbf{W}_{ni}, F_{\mathbf{w}_n}^s) \right\}, \quad i = 1, \dots, n,$$

and

$$\theta_{\mathbf{M}_{\mathbf{w}_n}} = -\arctan\left(\frac{M_{\mathbf{w}_n1}}{M_{\mathbf{w}_n2}}\right) + \frac{\pi}{2},$$

where $\mathbf{M}_{\mathbf{w}_n} = (M_{\mathbf{w}_n1}, M_{\mathbf{w}_n2})^\top$ refers to the sample median among $\mathbf{W}_{n1}, \dots, \mathbf{W}_{nn}$ based on $D(\cdot, F_{\mathbf{w}_n}^s)$. In the following, points $\mathbf{W}_{n1}, \dots, \mathbf{W}_{nn}$ are rotated counter-clockwise by angle $\frac{\pi}{2} - \theta_{\mathbf{M}_{\mathbf{w}_n}}$, which we call them as $\mathbf{V}_{n1}, \dots, \mathbf{V}_{nn}$.

Now, based on $D(\cdot, F_n)$, the affine invariant test statistic is defined as

$$(3.6) \quad T_{n,D}^* = \frac{6}{n(n+1)(2n+1)} \left(\sum_{i=1}^n \gamma_{ni} R_{ni} \right)^2,$$

where γ_{ni} is specified in the same way as δ_{ni} , through \mathbf{V}_{ni} instead of \mathbf{Z}_{ni} , for all $i = 1, \dots, n$.

It is worth to note that, the test statistic $T_{n,D}^*$ is also distribution-free. The affine invariance property and asymptotic null distribution of $T_{n,D}^*$ are presented in the following Theorems.

Theorem 3.3. *If the sample depth function $D(\cdot, F_n)$ satisfies property P1 and $n > 2$, the test statistic $T_{n,D}^*$ will be affine invariant; that is,*

$$T_{n,D}^*(\mathbf{X}_1, \dots, \mathbf{X}_n) = T_{n,D}^*(\mathbf{A}\mathbf{X}_1, \dots, \mathbf{A}\mathbf{X}_n)$$

for any 2×2 nonsingular matrix \mathbf{A} .

Theorem 3.4. *If the sample depth function satisfies property P1, then under the null hypothesis of centrally symmetric about 0, $T_{n,D}^*$ converges in distribution to a chi-square random variable with 1 degree of freedom.*

4. SIMULATION STUDY

In this section, an extensive simulation study is conducted to evaluate the finite sample behavior of the proposed test procedure. Two characteristics of interest are the empirical level and power of the proposed testing procedure. To assess the effects of different depth rankings on the performance of our test statistic, we determined three versions of $T_{n,D}$, derived from the simplicial, halfspace, and Mahalanobis depth functions as $T_{n,SD}$, $T_{n,HD}$ and $T_{n,MD}$, respectively. In the same way, $T_{n,SD}^*$, $T_{n,HD}^*$, $T_{n,MD}^*$ will be defined corresponding to $T_{n,D}^*$. The performance of our test statistics is compared with the affine invariant run test based on the simplicial depth function that we refer to $R_{n,SD}$ hereafter (Dyckerhoff *et al.* [8]) and the two rotation invariant tests Q_n^1 and Q_n^2 proposed by Einmahl and Gan [9]. Q_n^1 refers to their main test, and Q_n^2 is given by Q_n^1 adding a weight function to it (we avoid presenting the details of these test statistics).

To illustrate the effect of the sample size on the finite sample behavior of our proposed test statistics, we set the sample sizes as $n = 100$ and 200 . Moreover, the nominal level was set at 0.05 throughout. In each setting, 2000 independent random samples were generated to calculate the proportion of replications for which the null hypothesis is rejected. To examine the finite sample behavior of test statistics under the null and alternative hypotheses, we have simulated samples from several bivariate distribution families, including Azzalini's skew-normal distribution (Azzalini and Dalla Valle [3]), Azzalini's skew- t distribution (Azzalini and Capitanio [2]), perturbed symmetric beta distribution (Azzalini and Capitanio [2]) and sinh-arsinh distribution (Jones and Pewsey [21]). Indeed, we consider different types of skewness over very light-tailed distributions to very heavy-tailed ones. In what follows, we provide an overview of these families.

- *Bivariate skew-normal distribution:* Let \mathbf{X} be defined as

$$\mathbf{X} = \begin{cases} \mathbf{Y}, & \text{if } Z > \Delta^\top \mathbf{Y}, \\ -\mathbf{Y}, & \text{if } Z \leq \Delta^\top \mathbf{Y}, \end{cases}$$

where $\mathbf{Y} \sim N_2(0, \Sigma)$, $\Delta = (\Delta_1, \Delta_2)^\top$ is the shape parameter, and Z is distributed independently of \mathbf{Y} according to $N(0, 1)$. The random vector \mathbf{X} is known as bivariate skew-normal random vector and it may be written as $\mathbf{X} \sim SN_2(0, \Sigma, \Delta)$.

- *Bivariate skew-t distribution:* Let $\mathbf{T} = V^{-\frac{1}{2}}\mathbf{X}$, where the random vector \mathbf{X} follows the distribution $SN_2(0, \mathbf{\Sigma}, \mathbf{\Delta})$ and νV is distributed independently of \mathbf{X} according to a chi-squared distribution with ν degrees of freedom. We will say that \mathbf{T} has a bivariate skew- t distribution and write $\mathbf{T} \sim ST_2(0, \mathbf{\Sigma}, \mathbf{\Delta}, \nu)$.
- *Bivariate perturbed symmetric beta distribution:* Let $\mathbf{Y} = (2B_1 - 1, 2B_2 - 1)^\top$, where B_1 and B_2 have beta distributions $B(a, a)$ and $B(b, b)$, respectively. The random vector \mathbf{Y} can be treated as a central and non-elliptical symmetric random vector. Define the random vector \mathbf{X} as

$$\mathbf{X} = \begin{cases} \mathbf{Y}, & \text{if } Z < w(\mathbf{Y}), \\ -\mathbf{Y}, & \text{if } Z > w(\mathbf{Y}), \end{cases}$$

where Z (independently of \mathbf{Y}) has distribution function $G(\cdot)$. The distribution function $G(\cdot)$ and function $w(\cdot)$ are given as

$$G(z) = \frac{e^z}{1 + e^z} \quad \text{and} \quad w(\mathbf{y}) = \frac{\sin(p_1 y_1 + p_2 y_2)}{1 + \cos(q_1 y_1 + q_2 y_2)},$$

where p_1, p_2, q_1 and q_2 are additional parameters. Then, we will say that \mathbf{X} has a perturbed symmetric beta distribution.

- *Bivariate sinh-arcsinh distribution:* This family is generated by sinh-arcsinh transformation on a primary symmetric distribution. We consider the bivariate normal distribution as the primary distribution. The desirable property of this transformation is to induce skewness on the primary distribution and distributions with heavier/lighter tails than the primary one. Suppose random vector $\mathbf{Z} = (Z_1, Z_2)^\top$ follows $N_2(0, \mathbf{\Sigma})$. Define the bivariate vector $\mathbf{X} = (X_1, X_2)^\top$ as

$$(4.1) \quad X_j = \sinh \left[\frac{1}{\delta_j} (\sinh^{-1}(Z_j) + \Delta_j) \right], \quad j = 1, 2,$$

where Δ_j and δ_j denote the measure of skewness and tail weight in direction of j -th component of \mathbf{Z} , respectively. Amount of skewness increases with increasing positive Δ_j or decreasing negative Δ_j . Additionally, distributions with heavier and lighter tails than the bivariate normal distribution are generated by taking $0 < \delta_j < 1$ and $\delta_j > 1$, respectively.

In this study, we generate samples from the aforementioned distribution families with $\mathbf{\Sigma} = (1 - \rho)\mathbf{I}_2 + \rho\mathbf{J}_2$ with $\rho = -0.5, 0$ and 0.5 , and \mathbf{J}_2 denoting the 2×2 matrix with all entries equal 1 and $\Delta_i = k\eta, i = 1, 2$, with $\eta = (0.15, 0.15)^\top$ and $k = 0, 1, 2$ and 3 . We consider $\nu = 1, 3, 6, 10$ and 20 for bivariate skew- t distribution and $\delta_i = 0.5, 0.75, 1, 2$ and $5, i = 1, 2$, for bivariate sinh-arcsinh distribution.

Table 1 and Figures 1 and 2 provide the empirical rejection probabilities for sample size $n = 100$ and for bivariate skew-normal, skew- t and sinh-arcsinh distribution, respectively. Inspection of the table and figures confirms that the performance of our test statistics is not affected by different depth ranking. In all of them, the empirical rejection probabilities corresponding to $k = 0$ represents the proportion of rejection under the null hypothesis. These results demonstrate that all the tests would be accurate in estimating the nominal level, except $R_{n,SD}$ which it has been underestimated in some cases. Since the performance of test

statistics, even affine invariant test statistics are affected by correlation structure of primary distribution, we provide three possibilities for ρ as $-0.5, 0$ and 0.5 . From the represented results in Table 1 and Figure 1, it is obvious that all empirical powers will be increased by increasing the value of ρ for bivariate skew-normal and skew- t distributions. This situation is reversed for bivariate sinh-arsinh distribution in Figure 2 except for $T_{n,D}$.

Table 1: Empirical rejection probabilities (out of 2000 replications) for bivariate skew-normal distribution with $n = 100$, $\rho = -0.5, 0$ and 0.5 , and $\Delta_i = k\eta$, $i = 1, 2$, with $\eta = (0.15, 0.15)^T$ and $k = 0, 1, 2, 3$.

Test	$\rho = -0.5$				$\rho = 0$				$\rho = 0.5$			
	$k=0$	$k=1$	$k=2$	$k=3$	$k=0$	$k=1$	$k=2$	$k=3$	$k=0$	$k=1$	$k=2$	$k=3$
$T_{n,SD}^*$	0.046	0.105	0.294	0.520	0.046	0.175	0.488	0.694	0.046	0.245	0.607	0.763
$T_{n,HD}^*$	0.047	0.107	0.295	0.518	0.047	0.169	0.484	0.692	0.047	0.244	0.602	0.766
$T_{n,MD}^*$	0.048	0.103	0.291	0.511	0.048	0.177	0.486	0.688	0.048	0.241	0.604	0.760
$T_{n,SD}$	0.047	0.072	0.175	0.282	0.043	0.171	0.469	0.658	0.047	0.336	0.779	0.877
$T_{n,HD}$	0.044	0.074	0.169	0.280	0.044	0.171	0.462	0.644	0.044	0.340	0.779	0.881
$T_{n,MD}$	0.049	0.076	0.172	0.282	0.048	0.172	0.466	0.648	0.049	0.337	0.774	0.876
$R_{n,SD}$	0.043	0.048	0.081	0.159	0.043	0.057	0.135	0.320	0.043	0.070	0.190	0.460
Q_n^1	0.049	0.098	0.269	0.544	0.050	0.155	0.439	0.783	0.049	0.182	0.572	0.883
Q_n^2	0.054	0.080	0.173	0.355	0.048	0.098	0.250	0.507	0.053	0.127	0.342	0.657

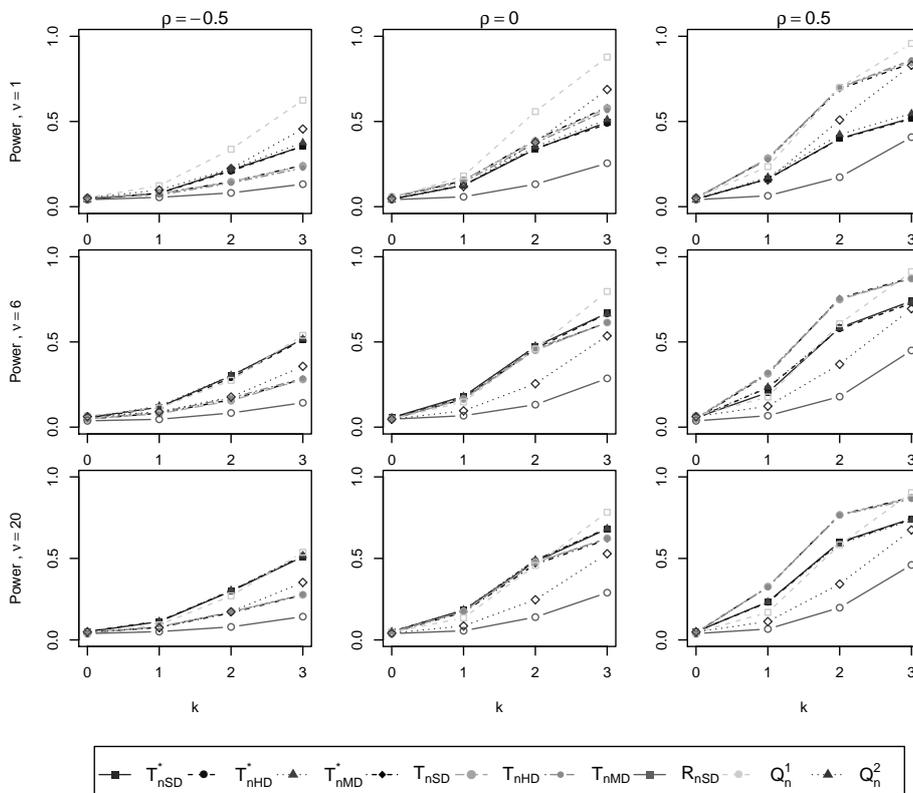


Figure 1: Empirical rejection probabilities (out of 2000 replications) for bivariate skew- t distribution with $n = 100$, $\rho = -0.5, 0$ and 0.5 , $\nu = 1, 6, 20$ and $\Delta_i = k\eta$, $i = 1, 2$, with $\eta = (0.15, 0.15)^T$ and $k = 0, 1, 2, 3$.

Table 1 shows that $T_{n,D}^*$ outperforms $R_{n,SD}$ and Q_n^2 in all cases, performs virtually as well as Q_n^1 for $k = 1$ and 2 and has slightly lower power than Q_n^1 for $k = 3$. Moreover, $T_{n,D}$ outperforms $R_{n,SD}$ in all cases, Q_n^2 when $\rho = 0$ and 0.5 and Q_n^1 when $\rho = 0.5$. Figure 1 indicates that $T_{n,D}^*$ and $T_{n,D}$ outperform $R_{n,SD}$ for all values of ν and ρ . In addition $T_{n,D}^*$ has higher power than Q_n^2 except when $\nu = 1$, and $T_{n,D}$ overcomes Q_n^2 except when $\nu = 1$ and $\rho = -0.5$. In comparison on Q_n^1 , $T_{n,D}^*$ performs better when $k = 1, 2, v = 6, 20$ and all values of ρ , and $T_{n,D}$ performs better when $k = 1$ and $2, v = 6, 20$ and $\rho = 0$ and 0.5 . Indeed, the empirical power of our tests increases as degrees of freedom increases. In Figure 2, superiority of our affine invariant tests is clear in most cases especially for $\rho = 0.5$ and $k = 1$ and 2 .

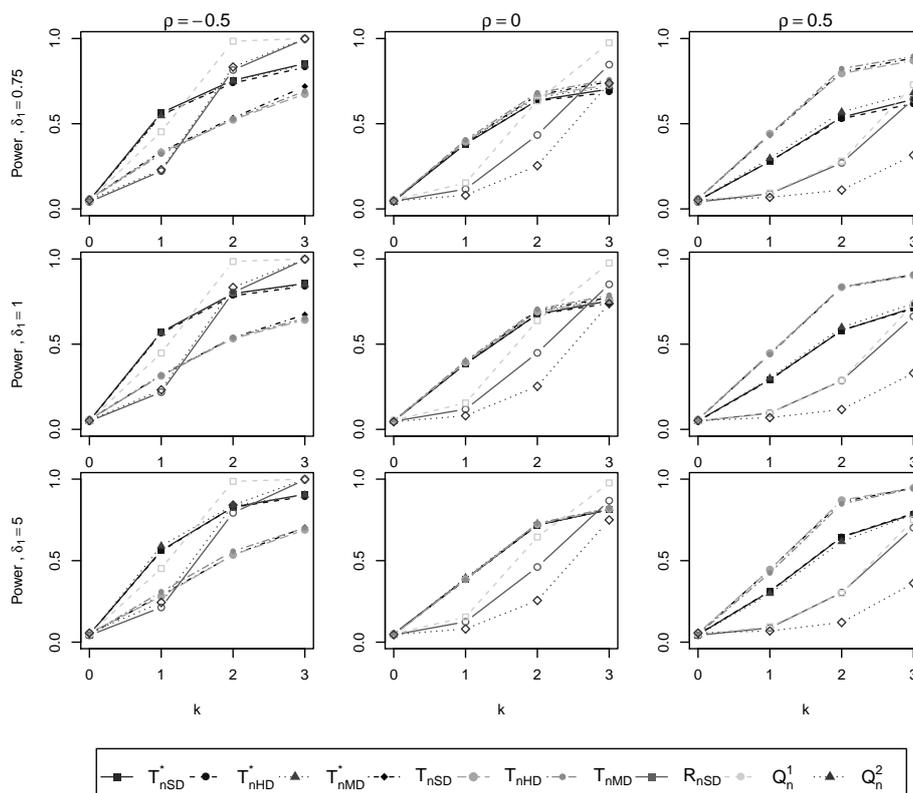


Figure 2: Empirical rejection probabilities (out of 2000 replications) for bivariate sinh–arcsinh distribution with $n = 100$, $\rho = -0.5, 0$ and 0.5 , $\delta_i = 0.75, 1, 5$ and $\Delta_i = k\eta$, $i = 1, 2$, with $\eta = (0.15, 0.15)^T$ and $k = 0, 1, 2, 3$.

Table 2 provide the empirical rejection probabilities for sample size $n = 200$ and for bivariate skew-normal distribution. In Figures 3 and 4, we plot the empirical rejection probabilities against k corresponding to some values of parameters of the same populations and tests with Figures 1 and 2 respectively, for sample size $n = 200$. Note that, as expected, the empirical powers increase with the sample size. These simulations lead to almost the same conclusions as in $n = 100$.

These simulations demonstrate that our tests are more powerful for small and moderate departures from the null hypothesis and for light-tailed distributions. As expected, the performance of affine invariant tests $T_{n,D}^*$ is less affected by changing the value of ρ rather than the orthogonal invariant tests $T_{n,D}$. The results show that, compared to $T_{n,D}^*$ test, $T_{n,D}$ performs better when $\rho = 0.5$, is comparable when $\rho = 0$ and performs worse when $\rho = -0.5$.

Table 2: Empirical rejection probabilities (out of 2000 replications) for bivariate skew-normal distribution with $n = 200$, $\rho = -0.5, 0$ and 0.5 , and $\Delta_i = k\eta$, $i = 1, 2$, with $\eta = (0.15, 0.15)^\top$ and $k = 0, 1, 2, 3$.

Test	$\rho = -0.5$				$\rho = 0$				$\rho = 0.5$			
	$k=0$	$k=1$	$k=2$	$k=3$	$k=0$	$k=1$	$k=2$	$k=3$	$k=0$	$k=1$	$k=2$	$k=3$
$T_{n,SD}^*$	0.052	0.204	0.486	0.690	0.052	0.317	0.681	0.813	0.052	0.414	0.747	0.834
$T_{n,HD}^*$	0.047	0.200	0.486	0.692	0.047	0.313	0.678	0.808	0.047	0.419	0.740	0.828
$T_{n,MD}^*$	0.052	0.206	0.488	0.696	0.052	0.319	0.680	0.804	0.052	0.420	0.748	0.832
$T_{n,SD}$	0.050	0.122	0.276	0.387	0.053	0.311	0.652	0.776	0.050	0.554	0.888	0.921
$T_{n,HD}$	0.048	0.121	0.283	0.391	0.050	0.305	0.645	0.770	0.048	0.560	0.879	0.918
$T_{n,MD}$	0.052	0.125	0.279	0.391	0.051	0.313	0.645	0.769	0.052	0.559	0.886	0.920
$R_{n,SD}$	0.049	0.061	0.117	0.250	0.049	0.076	0.218	0.519	0.049	0.093	0.347	0.754
Q_n^1	0.054	0.162	0.506	0.850	0.053	0.257	0.766	0.984	0.054	0.318	0.875	0.998
Q_n^2	0.056	0.130	0.342	0.665	0.053	0.156	0.500	0.856	0.056	0.212	0.647	0.942

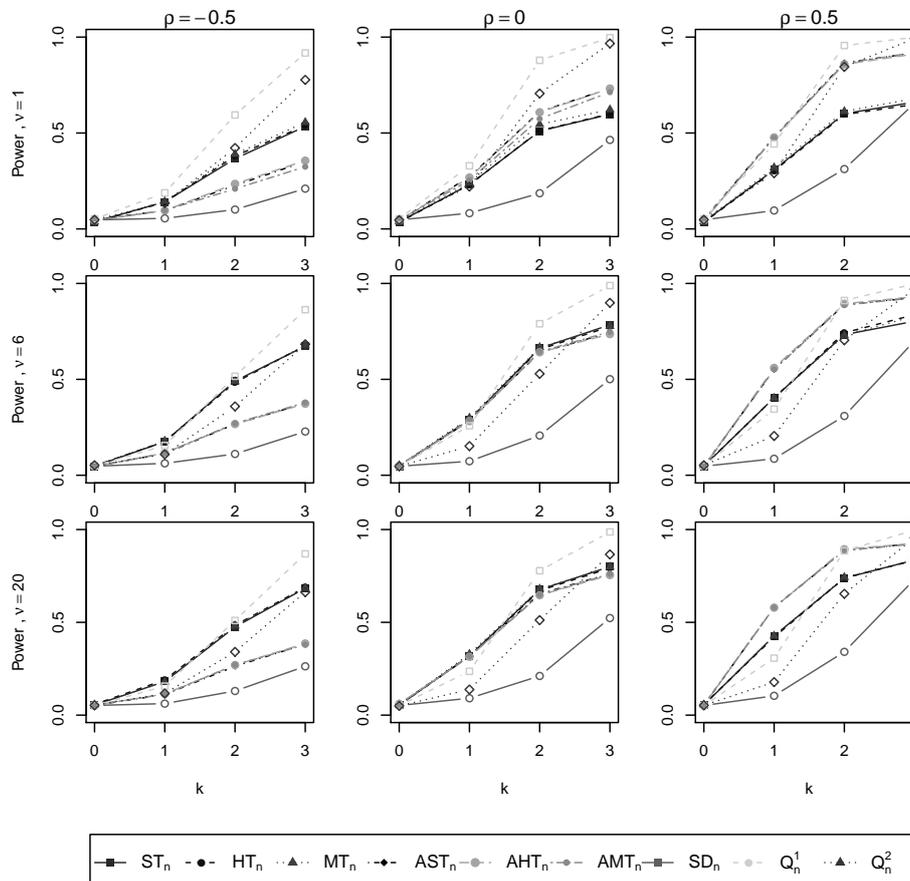


Figure 3: Empirical rejection probabilities (out of 2000 replications) for bivariate skew- t distribution with $n = 200$, $\rho = -0.5, 0$ and 0.5 , $\nu = 1, 6, 20$ and $\Delta_i = k\eta$, $i = 1, 2$, with $\eta = (0.15, 0.15)^\top$ and $k = 0, 1, 2, 3$.

Finally, to complete our simulations, for sample sizes $n = 100$ and 200 , we generate samples from bivariate perturbed symmetric beta distribution with several choices of the parameters such that different situations of asymmetry can be considered. A thorough investigation of Table 3 and 4 indicated that our tests overcome $R_{n,SD}$ and Q_n^2 in all cases and Q_n^1 in some cases.

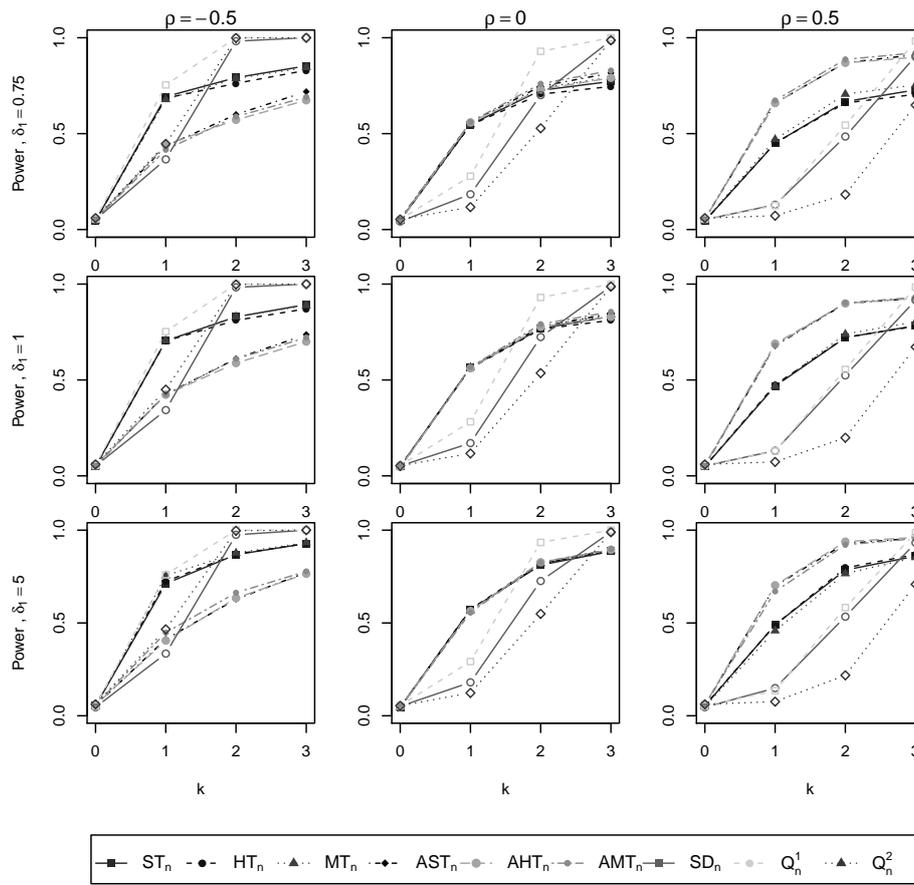


Figure 4: Empirical rejection probabilities (out of 2000 replications) for bivariate sinh–arcsinh distribution with $n = 200$, $\rho = -0.5, 0$ and 0.5 , $\delta_i = 0.75, 1, 5$ and $\Delta_i = k\eta$, $i = 1, 2$, with $\eta = (0.15, 0.15)^T$ and $k = 0, 1, 2, 3$.

Table 3: Empirical rejection probabilities (out of 2000 replications) for bivariate perturbed symmetric beta distribution with $n = 100$, $a, b = 0.5, 1$ and 3 , $p_1 = q_1 = 1$ and p_2 and $q_2 = 0.5, 1$ and 2 .

a, b	p_2	q_2	Test								
			$T_{n,SD}^*$	$T_{n,HD}^*$	$T_{n,MD}^*$	$T_{n,SD}$	$T_{n,HD}$	$T_{n,MD}$	$R_{n,SD}$	Q_n^1	Q_n^2
3, 3	2	0.5	0.194	0.194	0.194	0.175	0.184	0.175	0.060	0.170	0.112
	1	1	0.147	0.149	0.147	0.134	0.145	0.137	0.048	0.110	0.078
	0.5	2	0.165	0.168	0.164	0.147	0.158	0.154	0.051	0.098	0.073
3, 0.5	2	0.5	0.238	0.237	0.205	0.355	0.349	0.332	0.069	0.319	0.232
	1	1	0.284	0.290	0.270	0.364	0.351	0.340	0.077	0.265	0.177
	0.5	2	0.598	0.604	0.619	0.636	0.632	0.618	0.163	0.415	0.216
1, 1	2	0.5	0.241	0.244	0.244	0.221	0.215	0.222	0.076	0.320	0.223
	1	1	0.305	0.315	0.323	0.300	0.297	0.306	0.077	0.277	0.165
	0.5	2	0.581	0.593	0.606	0.558	0.573	0.585	0.165	0.397	0.201
0.5, 0.5	2	0.5	0.222	0.214	0.218	0.200	0.192	0.191	0.069	0.381	0.301
	1	1	0.470	0.474	0.486	0.445	0.453	0.464	0.113	0.505	0.324
	0.5	2	0.815	0.838	0.869	0.779	0.807	0.833	0.353	0.835	0.563

Table 4: Empirical rejection probabilities (out of 2000 replications) for bivariate perturbed symmetric beta distribution with $n = 200$, $a, b = 0.5, 1$ and 3 , $p_1 = q_1 = 1$ and p_2 and $q_2 = 0.5, 1$ and 2 .

a, b	p_2	q_2	Test								
			$T_{n,SD}^*$	$T_{n,HD}^*$	$T_{n,MD}^*$	$T_{n,SD}$	$T_{n,HD}$	$T_{n,MD}$	$R_{n,SD}$	Q_n^1	Q_n^2
3, 3	2	0.5	0.310	0.303	0.310	0.300	0.299	0.300	0.087	0.329	0.204
	1	1	0.252	0.250	0.250	0.235	0.237	0.233	0.083	0.194	0.120
	0.5	2	0.279	0.278	0.279	0.268	0.265	0.266	0.089	0.173	0.098
3, 0.5	2	0.5	0.399	0.399	0.338	0.624	0.619	0.595	0.111	0.611	0.471
	1	1	0.469	0.472	0.444	0.631	0.627	0.608	0.112	0.491	0.310
	0.5	2	0.791	0.810	0.808	0.896	0.893	0.875	0.283	0.729	0.400
1, 1	2	0.5	0.370	0.369	0.380	0.364	0.363	0.365	0.105	0.578	0.411
	1	1	0.509	0.520	0.531	0.499	0.499	0.505	0.117	0.512	0.315
	0.5	2	0.793	0.816	0.819	0.788	0.801	0.812	0.262	0.694	0.378
0.5, 0.5	2	0.5	0.352	0.350	0.352	0.325	0.336	0.339	0.101	0.656	0.559
	1	1	0.700	0.738	0.747	0.686	0.725	0.733	0.175	0.804	0.592
	0.5	2	0.916	0.935	0.931	0.913	0.934	0.943	0.594	0.993	0.891

A second Monte Carlo study is provided in order to evaluate the performance of our tests for pure location alternatives. In this study the performance of tests considered in first Monte Carlo study compared with the Hotelling's T^2 and the tests due to Hallin and Paindaveine [12] computed with the sign score function, van der Waerden score function and Wilcoxon score function and denoted by HS_n , HN_n and HR_n , respectively.

We set the sample size as $n = 50$. In each setting, 2000 independent random samples were generated to calculate the proportion of replications for which the null hypothesis is rejected. For each replication, the all tests were performed at the significance level $\alpha = 0.05$. To examine the finite sample behavior of test statistics under the null and alternative hypotheses, we have simulated samples from the t family of distributions and the exponential power family of distributions. In what follows, we provide an overview of these families.

A p -dimensional random vector \mathbf{X} has a multivariate t -distribution with ν degree of freedom if its density function has the form

$$f_{\boldsymbol{\mu}, \boldsymbol{\Sigma}}(\mathbf{x}) = \frac{\Gamma((p+\nu)/2)}{\Gamma(\nu/2) (\pi\nu)^{p/2}} |\boldsymbol{\Sigma}|^{-1/2} \left[1 + \frac{1}{\nu} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]^{-(p+\nu)/2},$$

where $\boldsymbol{\mu} = (\mu_1, \dots, \mu_p)^\top \in \mathbb{R}^p$ and $\boldsymbol{\Sigma}$ is a symmetric $p \times p$ positive definite matrix.

The density function of a p -dimensional random vector \mathbf{X} from the exponential power family of distributions is

$$f_{\boldsymbol{\mu}, \boldsymbol{\Sigma}}(\mathbf{x}) = \frac{\nu \Gamma(p/2)}{\Gamma(p+2\nu) (\pi c_0)^{p/2}} |\boldsymbol{\Sigma}|^{-1/2} \exp \left\{ -\frac{(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})}{c_0} \right\}^\nu,$$

where

$$c_0 = \frac{p \Gamma(p/2\nu)}{\Gamma((p+2)/2\nu)}$$

and $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ are defined as above.

We generate samples from the aforementioned distribution families with $\Sigma = \mathbf{I}$ and $\mu = k\Delta$ with $\Delta = (0.2, 0.2)^\top$ and $\Delta = (0.1, 0.1)^\top$ for the t family of distributions and the exponential power family of distributions, respectively and $k = 0, 1, 2, 3$. We consider $\nu = 1, 6$ and 10 for t -distribution family and $\nu = 0.5, 1$ and 2 for the exponential power family of distributions.

Inspection of Tables 5 and 6 demonstrated that the performance of our tests is comparable to the other tests. The proposed tests overcome $R_{n,SD}$ and Q_n^2 in most cases and Q_n^1 in some cases. It worth to note that all tests which are defined in the similar way of our proposed test e.g. $R_{n,SD}$, Q_n^1 and Q_n^2 are not expected to perform as well as T^2 , HS_n , HN_n and HR_n . In other hand, the results confirm that the performance of our test statistics is not affected by different depth ranking.

Table 5: Empirical rejection probabilities (out of 2000 replications) for bivariate t -distribution with $n = 50$, $\Sigma = \mathbf{I}$, $\nu = 1, 6, 20$ and $\mu = k\Delta$ with $\Delta = (0.2, 0.2)^\top$ and $k = 0, 1, 2, 3$.

Test	$\nu = 1$				$\nu = 6$				$\nu = 20$			
	$k=0$	$k=1$	$k=2$	$k=3$	$k=0$	$k=1$	$k=2$	$k=3$	$k=0$	$k=1$	$k=2$	$k=3$
T^2	0.015	0.033	0.084	0.164	0.058	0.287	0.809	0.987	0.056	0.374	0.911	1
HS_n	0.050	0.203	0.634	0.902	0.053	0.285	0.805	0.993	0.054	0.322	0.851	1
HN_n	0.040	0.122	0.351	0.622	0.048	0.270	0.801	0.988	0.049	0.345	0.894	1
HR_n	0.044	0.120	0.343	0.599	0.049	0.279	0.804	0.988	0.057	0.354	0.901	1
$T_{n,SD}^*$	0.042	0.085	0.150	0.431	0.050	0.212	0.542	0.757	0.056	0.236	0.601	0.7715
$T_{n,HD}^*$	0.040	0.083	0.146	0.425	0.053	0.204	0.541	0.750	0.056	0.223	0.592	0.768
$T_{n,MD}^*$	0.039	0.093	0.156	0.432	0.051	0.205	0.526	0.724	0.060	0.234	0.593	0.762
$T_{n,SD}$	0.052	0.100	0.171	0.429	0.046	0.188	0.523	0.749	0.047	0.219	0.570	0.762
$T_{n,HD}$	0.050	0.104	0.173	0.435	0.047	0.182	0.520	0.746	0.048	0.216	0.575	0.761
$T_{n,MD}$	0.051	0.090	0.155	0.389	0.045	0.183	0.510	0.722	0.048	0.223	0.567	0.749
$R_{n,SD}$	0.038	0.073	0.105	0.342	0.040	0.085	0.229	0.569	0.040	0.090	0.287	0.650
Q_n^1	0.057	0.223	0.450	0.937	0.066	0.203	0.671	0.945	0.058	0.197	0.661	0.946
Q_n^2	0.039	0.160	0.338	0.862	0.049	0.121	0.412	0.791	0.038	0.106	0.386	0.787

Table 6: Empirical rejection probabilities (out of 2000 replications) for bivariate power family of distributions with $n = 50$, $\Sigma = \mathbf{I}$, $\nu = 0.5, 1, 2$ and $\mu = k\Delta$ with $\Delta = (0.1, 0.1)^\top$ and $k = 0, 1, 2, 3$.

Test	$\nu = 0.5$				$\nu = 1$				$\nu = 2$			
	$k=0$	$k=1$	$k=2$	$k=3$	$k=0$	$k=1$	$k=2$	$k=3$	$k=0$	$k=1$	$k=2$	$k=3$
T^2	0.044	0.161	0.515	0.867	0.048	0.392	0.940	1	0.043	0.468	0.981	1
HS_n	0.043	0.203	0.633	0.930	0.044	0.311	0.873	1	0.047	0.317	0.868	0.999
HN_n	0.036	0.158	0.528	0.864	0.038	0.348	0.917	1	0.036	0.447	0.978	1
HR_n	0.040	0.155	0.519	0.849	0.044	0.361	0.922	1	0.038	0.467	0.982	1
$T_{n,SD}^*$	0.053	0.127	0.355	0.601	0.054	0.241	0.619	0.782	0.048	0.282	0.669	0.793
$T_{n,HD}^*$	0.056	0.133	0.363	0.605	0.057	0.240	0.612	0.775	0.051	0.282	0.671	0.779
$T_{n,MD}^*$	0.054	0.130	0.349	0.570	0.055	0.248	0.618	0.761	0.044	0.289	0.685	0.793
$T_{n,SD}$	0.049	0.118	0.344	0.583	0.047	0.214	0.605	0.777	0.036	0.273	0.646	0.784
$T_{n,HD}$	0.046	0.120	0.343	0.596	0.046	0.210	0.604	0.772	0.041	0.268	0.648	0.772
$T_{n,MD}$	0.048	0.123	0.335	0.565	0.048	0.222	0.600	0.758	0.035	0.273	0.660	0.786
$R_{n,SD}$	0.044	0.057	0.133	0.322	0.038	0.095	0.288	0.676	0.042	0.090	0.362	0.763
Q_n^1	0.049	0.179	0.524	0.854	0.045	0.185	0.640	0.949	0.059	0.141	0.533	0.935
Q_n^2	0.035	0.118	0.366	0.670	0.037	0.099	0.350	0.763	0.041	0.070	0.233	0.683

5. CONCLUSION

This paper concerns with the problem of detecting central symmetry of a bivariate distribution. To this end, based on depth function, we introduced a family of signed-rank test which is orthogonal invariant and distribution-free. Affine invariant tests were obtained by applying our proposed test to the standardized data with Tyler's matrix. The proposed orthogonal and affine invariant tests have the same asymptotic properties. In simulation study, the finite sample behavior of the proposed test procedure was evaluated over distributions family from very light to very heavy-tailed distributions with different kinds of skewness. The simulations confirmed that our affine invariant tests successfully can distinguish different asymmetries and shifting the location parameter. Moreover, we observed that they performed as good as their competitors and actually in many cases they even outperform them.

A. APPENDIX

Proof of Theorem 3.1: According to the construction of \mathbf{Z}_{ni} , it is clear that $\mathbf{Z}_{ni} = \mathbf{B}_{\mathbf{X}_n} \mathbf{X}_i$, $i = 1, \dots, n$, where

$$(A.1) \quad \mathbf{B}_{\mathbf{X}_n} = \begin{bmatrix} \cos\left(\frac{\pi}{2} - \theta_{\mathbf{M}_n}\right) & -\sin\left(\frac{\pi}{2} - \theta_{\mathbf{M}_n}\right) \\ \sin\left(\frac{\pi}{2} - \theta_{\mathbf{M}_n}\right) & \cos\left(\frac{\pi}{2} - \theta_{\mathbf{M}_n}\right) \end{bmatrix}.$$

Let \mathbf{A} be an arbitrary 2×2 orthogonal matrix. Define $\tilde{\mathbf{Z}}_{ni} = \mathbf{B}_{\mathbf{A}\mathbf{X}_n} \mathbf{A}\mathbf{X}_i$ for all $i = 1, \dots, n$, where

$$\mathbf{B}_{\mathbf{A}\mathbf{X}_n} = \begin{bmatrix} \cos\left(\frac{\pi}{2} - \theta_{\tilde{\mathbf{M}}_n}\right) & -\sin\left(\frac{\pi}{2} - \theta_{\tilde{\mathbf{M}}_n}\right) \\ \sin\left(\frac{\pi}{2} - \theta_{\tilde{\mathbf{M}}_n}\right) & \cos\left(\frac{\pi}{2} - \theta_{\tilde{\mathbf{M}}_n}\right) \end{bmatrix},$$

with $\theta_{\tilde{\mathbf{M}}_n} \in [0, \pi)$, as the angle between horizontal-axis and the sample median $\tilde{\mathbf{M}}_n$ that is obtained in the same way as \mathbf{M}_n , through $\mathbf{A}\mathbf{X}_i$'s instead of \mathbf{X}_i 's, $i = 1, \dots, n$. The orthogonality of matrix \mathbf{A} implies that there exists an angle $\alpha \in [0, 2\pi)$ such that

$$(A.2) \quad \mathbf{A} = \begin{bmatrix} \cos(\alpha) & -\sin(\alpha) \\ \sin(\alpha) & \cos(\alpha) \end{bmatrix} \quad \text{or} \quad \mathbf{A} = \begin{bmatrix} \cos(\alpha) & \sin(\alpha) \\ \sin(\alpha) & -\cos(\alpha) \end{bmatrix}.$$

Property P1 of the sample depth function shows that

$$(A.3) \quad \tilde{\mathbf{M}}_n = \mathbf{A}\mathbf{M}_n.$$

Let matrix \mathbf{A} be defined as the left side of (A.2), then (A.3) results in

$$\theta_{\tilde{\mathbf{M}}_n} = \begin{cases} \alpha + \theta_{\mathbf{M}_n}, & 0 \leq \alpha + \theta_{\mathbf{M}_n} < \pi, \\ \alpha + \theta_{\mathbf{M}_n} - \pi, & \pi \leq \alpha + \theta_{\mathbf{M}_n} < 2\pi, \\ \alpha + \theta_{\mathbf{M}_n} - 2\pi, & 2\pi \leq \alpha + \theta_{\mathbf{M}_n} < 3\pi. \end{cases}$$

Using the trigonometric relationships, it is straightforward to verify that $\mathbf{B}_{\mathbf{A}\mathbf{X}_n} \mathbf{A} = \mathbf{B}_{\mathbf{X}_n}$, or $\mathbf{B}_{\mathbf{A}\mathbf{X}_n} \mathbf{A} = -\mathbf{B}_{\mathbf{X}_n}$. Thus

$$(A.4) \quad \tilde{\mathbf{Z}}_{ni} = \mathbf{Z}_{ni} \quad \text{or} \quad \tilde{\mathbf{Z}}_{ni} = -\mathbf{Z}_{ni}, \quad i = 1, \dots, n.$$

Now, let matrix \mathbf{A} be according to the right side of (A.2), similarly we have

$$\theta_{\tilde{\mathbf{M}}_n} = \begin{cases} \alpha - \theta_{\mathbf{M}_n} + \pi, & -\pi < \alpha - \theta_{\mathbf{M}_n} < 0, \\ \alpha - \theta_{\mathbf{M}_n}, & 0 \leq \alpha - \theta_{\mathbf{M}_n} < \pi, \\ \alpha - \theta_{\mathbf{M}_n} - \pi, & \pi \leq \alpha - \theta_{\mathbf{M}_n} < 2\pi, \end{cases}$$

and

$$\mathbf{B}_{\mathbf{A}\mathbf{X}_n}\mathbf{A} = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} \mathbf{B}_{\mathbf{X}_n} \quad \text{or} \quad \mathbf{B}_{\mathbf{A}\mathbf{X}_n}\mathbf{A} = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} \mathbf{B}_{\mathbf{X}_n}.$$

Hence

$$(A.5) \quad \tilde{\mathbf{Z}}_{ni} = (-\mathbf{Z}_{ni1}, \mathbf{Z}_{ni2})^\top \quad \text{or} \quad \tilde{\mathbf{Z}}_{ni} = (\mathbf{Z}_{ni1}, -\mathbf{Z}_{ni2})^\top, \quad i = 1, \dots, n.$$

The proof of affine invariance of $T_{n,D}$ will be completed by using (A.4), (A.5) and property P1 of the sample depth function. \square

Proof of Theorem 3.2: Under the null hypothesis, $\mathbf{X}_1, \dots, \mathbf{X}_n$ are i.i.d. from F , where $F(\cdot)$ is centrally symmetric distribution about the origin. Hence, we have

$$(A.6) \quad (\mathbf{X}_1, \dots, \mathbf{X}_n) \stackrel{d}{=} (\eta_1 \mathbf{X}_1, \dots, \eta_n \mathbf{X}_n),$$

where η_i 's, $i = 1, \dots, n$, are i.i.d. random variables taking the values 1 and -1 each with probability $1/2$. It is clear that

$$(A.7) \quad (\pm \mathbf{X}_1, \dots, \pm \mathbf{X}_n) = (\pm \eta_1 \mathbf{X}_1, \dots, \pm \eta_n \mathbf{X}_n).$$

Additionally, $\mathbf{M}_n \equiv \mathbf{M}(\mathbf{X}_1, \dots, \mathbf{X}_n)$ is considered as a point with maximum sample depth with respect to the symmetrized sample $(\pm \mathbf{X}_1, \dots, \pm \mathbf{X}_n)$ (if there is more than one sample point with the highest depth value, \mathbf{M}_n will be defined as the point with minimum index among those data points). By this definition of \mathbf{M}_n , there exists $i \in \{1, \dots, n\}$ such that

$$(A.8) \quad \mathbf{M}(\mathbf{X}_1, \dots, \mathbf{X}_n) = \mathbf{X}_i.$$

From property P1 and equation (A.7),

$$(A.9) \quad \mathbf{M}(\eta_1 \mathbf{X}_1, \dots, \eta_n \mathbf{X}_n) = \eta_i \mathbf{X}_i.$$

Hence from (A.8) and (A.9), we have

$$(A.10) \quad \mathbf{M}(\mathbf{X}_1, \dots, \mathbf{X}_n) = \eta \mathbf{M}(\eta_1 \mathbf{X}_1, \dots, \eta_n \mathbf{X}_n),$$

where $\eta = 1$ or -1 . Thus $\mathbf{B}_{\mathbf{X}_n}$ where is defined as (A.1) will be same whether it is obtained from either $\mathbf{X}_1, \dots, \mathbf{X}_n$ or $\eta_1 \mathbf{X}_1, \dots, \eta_n \mathbf{X}_n$. Hence (A.6) implies that

$$(A.11) \quad (\mathbf{B}_{\mathbf{X}_n} \mathbf{X}_1, \dots, \mathbf{B}_{\mathbf{X}_n} \mathbf{X}_n) \stackrel{d}{=} (\eta_1 \mathbf{B}_{\mathbf{X}_n} \mathbf{X}_1, \dots, \eta_n \mathbf{B}_{\mathbf{X}_n} \mathbf{X}_n).$$

This yields that δ_{ni} 's, $i = 1, \dots, n$, are independent and identically distributed random variables that take the values 1 and -1 with probability $1/2$. Let $\mathbf{Z}_{ni} = \delta_{ni} \mathbf{Y}_{ni}$, where $\mathbf{Y}_{ni} = (Y_{ni1}, Y_{ni2})^\top$

for all $i = 1, \dots, n$. (A.11) denotes that $\mathbf{Z}_{n1}, \dots, \mathbf{Z}_{nm}$ distributed as centrally symmetric random vectors about origin. Thus, for $\mathbf{y} = (y_1, y_2)^\top \in \mathbb{R}^2$ and $i = 1, \dots, n$,

$$\begin{aligned}
P_{H_0}(Y_{ni1} \leq y_1, Y_{ni2} \leq y_2, \delta_{ni} = 1) &= P_{H_0}(\delta_{ni}Y_{ni1} \leq y_1, \delta_{ni}Y_{ni2} \leq y_2, \delta_{ni} = 1) \\
&= P_{H_0}(Z_{ni1} \leq y_1, Z_{ni2} \leq y_2, \delta_{ni} = 1) \\
&= P_{H_0}(Z_{ni1} \leq y_1, Z_{ni2} \leq y_2, Z_{ni2} > 0) \\
&= P_{H_0}(Z_{ni1} \leq y_1, 0 < Z_{ni2} \leq y_2) \\
&= P_{H_0}(-Z_{ni1} \leq y_1, 0 < -Z_{ni2} \leq y_2) \\
&= P_{H_0}(Z_{ni1} \geq -y_1, -y_2 \leq Z_{ni2} < 0) \\
&= P_{H_0}(Z_{ni1} \geq -y_1, Z_{ni2} \geq -y_2, Z_{ni2} < 0) \\
&= P_{H_0}(-Z_{ni1} \leq y_1, -Z_{ni2} \leq y_2, \delta_{ni} = -1) \\
&= P_{H_0}(\delta_{ni}Z_{ni1} \leq y_1, \delta_{ni}Z_{ni2} \leq y_2, \delta_{ni} = -1) \\
&= P_{H_0}(Y_{ni1} \leq y_1, Y_{ni2} \leq y_2, \delta_{ni} = -1)
\end{aligned}$$

and, for $j \neq i$,

$$\begin{aligned}
P_{H_0}(Y_{ni1} \leq y_1, Y_{ni2} \leq y_2, \delta_{ni} = 1, \delta_{nj} = 1) &= P_{H_0}(Y_{ni1} \leq y_1, Y_{ni2} \leq y_2, \delta_{ni} = 1, \delta_{nj} = 1) \\
&\quad + P_{H_0}(Y_{ni1} \leq y_1, Y_{ni2} \leq y_2, \delta_{ni} = -1, \delta_{nj} = 1) \\
&= P_{H_0}(Z_{ni1} \leq y_1, Z_{ni2} \leq y_2, Z_{ni2} > 0, Z_{nj2} > 0) \\
&\quad + P_{H_0}(Z_{ni1} \geq -y_1, Z_{ni2} \geq -y_2, Z_{ni2} < 0, Z_{nj2} > 0) \\
&= P_{H_0}(Z_{ni1} \leq y_1, Z_{ni2} \leq y_2, Z_{ni2} > 0, Z_{nj2} < 0) \\
&\quad + P_{H_0}(Z_{ni1} \geq -y_1, Z_{ni2} \geq -y_2, Z_{ni2} < 0, Z_{nj2} < 0) \\
&= P_{H_0}(Y_{ni1} \leq y_1, Y_{ni2} \leq y_2, \delta_{ni} = 1, \delta_{nj} = -1) \\
&\quad + P_{H_0}(Y_{ni1} \leq y_1, Y_{ni2} \leq y_2, \delta_{ni} = -1, \delta_{nj} = -1) \\
&= P_{H_0}(Y_{ni1} \leq y_1, Y_{ni2} \leq y_2, \delta_{nj} = -1).
\end{aligned}$$

Hence these imply that δ_{ni} , for $i = 1, \dots, n$, is independent of $\mathbf{Y}_{n1}, \dots, \mathbf{Y}_{nn}$. Now, suppose that $F_{\mathbf{Z}_n}^s$ and $F_{\mathbf{Y}_n}^s$ be the sample distribution functions of $\{\pm \mathbf{Z}_{n1}, \dots, \pm \mathbf{Z}_{nn}\}$ and $\{\pm \mathbf{Y}_{n1}, \dots, \pm \mathbf{Y}_{nn}\}$, respectively. Since $\{\pm \mathbf{Z}_{n1}, \dots, \pm \mathbf{Z}_{nn}\} = \{\pm \mathbf{Y}_{n1}, \dots, \pm \mathbf{Y}_{nn}\}$, it is clear that $F_{\mathbf{Z}_n}^s = F_{\mathbf{Y}_n}^s$. This equality, along with $D(\mathbf{Z}_{ni}, F_{\mathbf{Z}_n}^s) = D(-\mathbf{Z}_{ni}, F_{\mathbf{Z}_n}^s)$ (resulted from property P1 by considering $\mathbf{A} = -I_2$ and $\mathbf{b} = \mathbf{0}$) conclude that $D(\mathbf{Z}_{ni}, F_{\mathbf{Z}_n}^s) = D(\mathbf{Y}_{ni}, F_{\mathbf{Y}_n}^s)$, for all $i = 1, \dots, n$. Additionally, from property P1 of the sample depth function and Remark 2.1, we see that $D(\mathbf{X}_i, F_n^s) = D(\mathbf{Z}_{ni}, F_{\mathbf{Z}_n}^s)$. Hence $D(\mathbf{X}_i, F_n^s) = D(\mathbf{Y}_{ni}, F_{\mathbf{Y}_n}^s)$. This shows that R_i is a function of $\mathbf{Y}_{n1}, \dots, \mathbf{Y}_{nn}$ and thus is independent of δ_{ni} , $i = 1, 2, \dots, n$. Under null hypothesis, R_1, \dots, R_n have the discrete uniform distribution on $\{1, \dots, n\}$. Then the expectation and variance of $T_{n,D}^{1/2}$ are given as

$$E(T_{n,D}^{1/2}) = \sqrt{\frac{6}{n(n+1)(2n+1)}} E\left(\sum_{i=1}^n \delta_{in} R_i\right) = 0$$

and

$$\text{Var}(T_{n,D}^{1/2}) = \frac{6}{n(n+1)(2n+1)} \sum_{i=1}^n E(R_i^2) + \frac{6}{n(n+1)(2n+1)} \sum_{i=1}^n \sum_{\substack{j=1 \\ i \neq j}}^n E(\delta_{ni} \delta_{nj} R_i R_j) = 1,$$

respectively. Because of the dependency between summands in $T_{n,D}^{1/2}$, the central limit theory

is not applied. In the other hand, $T_{n,D}^{1/2}$ is equal in distribution to

$$K_n = \sqrt{\frac{6}{n(n+1)(2n+1)}} \sum_{i=1}^n \delta_i i,$$

where δ_i 's, $i=1, \dots, n$, are independent random variables with probability 1/2 of being 1 or -1 . Since K_n is sum of independent random variables and the Lyapunov's condition

$$\lim_{n \rightarrow \infty} \left(\sqrt{\frac{6}{n(n+1)(2n+1)}} \right)^{(2+\delta)} \sum_{i=1}^n (E|\delta_i i|)^{(2+\delta)} = \lim_{n \rightarrow \infty} \left(\frac{n^3}{3} \right)^{-(2+\delta)/2} \sum_{i=1}^n i^{2+\delta} = 0$$

is satisfied for $\delta = 1$, then the asymptotic null distribution is obtained by Lyapunov's central limit theorem. \square

Proof of Theorem 3.3: It is clear that $V_{ni} = B_{W_n} W_{ni}$, $i=1, \dots, n$, where

$$(A.12) \quad B_{W_n} = \begin{bmatrix} \cos\left(\frac{\pi}{2} - \theta_{M_{W_n}}\right) & -\sin\left(\frac{\pi}{2} - \theta_{M_{W_n}}\right) \\ \sin\left(\frac{\pi}{2} - \theta_{M_{W_n}}\right) & \cos\left(\frac{\pi}{2} - \theta_{M_{W_n}}\right) \end{bmatrix}.$$

Let A be an arbitrary 2×2 nonsingular matrix and define $\tilde{V}_{ni} = B_{A W_n} \tilde{W}_{ni}$, where

$$B_{A W_n} = \begin{bmatrix} \cos\left(\frac{\pi}{2} - \theta_{\tilde{M}_{W_n}}\right) & -\sin\left(\frac{\pi}{2} - \theta_{\tilde{M}_{W_n}}\right) \\ \sin\left(\frac{\pi}{2} - \theta_{\tilde{M}_{W_n}}\right) & \cos\left(\frac{\pi}{2} - \theta_{\tilde{M}_{W_n}}\right) \end{bmatrix},$$

with $\theta_{\tilde{M}_{W_n}} \in [0, \pi)$, as the angle between horizontal-axis and the sample median \tilde{M}_{W_n} that is obtained in the same way as M_{W_n} , through \tilde{W}_{ni} 's instead of W_{ni} 's, $i=1, \dots, n$. Moreover, $\tilde{W}_{ni} = \Gamma_{A X_n} A X_i$, where $\Gamma_{A X_n}$ is Tyler's matrix defined in terms of the transformed data points $A X_i$, for all $i=1, \dots, n$.

If $n > 2$, Randles [36] indicated that Γ_n satisfies the condition

$$(A.13) \quad A^T \Gamma_{A X_n}^T \Gamma_{A X_n} A = k \Gamma_n^T \Gamma_n,$$

where k is a positive scalar that may depends on A and the data. This equation clearly shows that there exists an orthogonal matrix $H = k^{-1/2} \Gamma_{A X_n} A \Gamma_n^{-1}$ such that

$$(A.14) \quad \sqrt{k} H \Gamma_n = \Gamma_{A X_n} A.$$

It follows easily that

$$(A.15) \quad \tilde{W}_{ni} = \Gamma_{A X_n} A X_i = \sqrt{k} H \Gamma_n X_i = \sqrt{k} H W_{ni}.$$

Additionally, property P1 of the sample depth function along with Remark 2.1 and equation (A.15) show that $\tilde{M}_{W_n} = \sqrt{k} H M_{W_n}$. Thus, the result follows from Theorem 3.1. \square

Proof of Theorem 3.4: The Tyler's matrix $\Gamma_n \equiv \Gamma(X_1, \dots, X_n)$ is invariant under sign changes among the X_i 's (Randles [36]), that is

$$(A.16) \quad \Gamma(X_1, \dots, X_n) = \Gamma(\eta_1 X_1, \dots, \eta_n X_n).$$

Hence, by (A.6) we have

$$(A.17) \quad (B_{W_n} \Gamma_n X_1, \dots, B_{W_n} \Gamma_n X_n) \stackrel{d}{=} (\eta_1 B_{W_n} \Gamma_n X_1, \dots, \eta_n B_{W_n} \Gamma_n X_n).$$

where B_{W_n} is defined as (A.12). Additionally, from property P1 of the sample depth function and Remark 2.1, it is straightforward to verify that $R_{ni} = R_i$ for all $i=1, \dots, n$. The rest of the proof proceeds as in Theorem 3.2. \square

REFERENCES

- [1] AKI, S. (1993). On nonparametric tests for symmetry in R^m , *Annals of the Institute of Statistical Mathematics*, **45**(4), 787–800.
- [2] AZZALINI, A. and CAPITANIO, A. (2003). Distributions generated by perturbation of symmetry with emphasis on a multivariate skew t -distribution, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **65**(2), 367–389.
- [3] AZZALINI, A. and DALLA VALLE, A. (1996). The multivariate skew-normal distribution, *Biometrika*, **83**(4), 715–726.
- [4] BARINGHAUS, L. (1991). Testing for spherical symmetry of a multivariate distribution, *Annals of Statistics*, **19**(2), 899–917.
- [5] BROWN, B.M. and HETTMANSPERGER, T.P. (1987). Affine invariant rank methods in the bivariate location model, *Journal of the Royal Statistical Society. Series B (Methodological)*, 301–310.
- [6] CASSART, D.; HALLIN, M. and PAINDAVEINE, D. (2011). A class of optimal tests for symmetry based on local Edgeworth approximations, *Bernoulli*, **17**(3), 1063–1094.
- [7] CHERNOZHUKOV, V.; GALICHON, A.; HALLIN, M. and HENRY, M. (2017). Monge–Kantorovich depth, quantiles, ranks and signs, *The Annals of Statistics*, **45**(1), 223–256.
- [8] DYCKERHOFF, R.; LEY, C. and PAINDAVEINE, D. (2015). Depth-based runs tests for bivariate central symmetry, *Annals of the Institute of Statistical Mathematics*, **67**(5), 917–941.
- [9] EINMAHL, J.H. and GAN, Z. (2016). Testing for central symmetry, *Journal of Statistical Planning and Inference*, **169**, 27–33.
- [10] GHOSH, S. and RUYMGAART, F.H. (1992). Applications of empirical characteristic functions in some multivariate problems, *Canadian Journal of Statistics*, **20**(4), 429–440.
- [11] HALLIN, M. and PAINDAVEINE, D. (2002a). Multivariate Signed Ranks: Randles’ Interdirections or Tyler’s Angles, *Statistical Data Analysis Based on the L_1 -Norm and Related Methods*, 271–282.
- [12] HALLIN, M. and PAINDAVEINE, D. (2002b). Optimal tests for multivariate location based on interdirections and pseudo-mahalanobis ranks, *The Annals of Statistics*, **30**(4), 1103–1133.
- [13] HALLIN, M. and PAINDAVEINE, D. (2002). Optimal procedures based on interdirections and pseudo-Mahalanobis ranks for testing multivariate elliptic white noise against ARMA dependence, *Bernoulli*, **8**(6), 787–815.
- [14] HALLIN, M. and PAINDAVEINE, D. (2004). Rank-based optimal tests of the adequacy of an elliptic VARMA model, *Annals of Statistics*, **32**, 2642–2678.
- [15] HALLIN, M. and PAINDAVEINE, D. (2005). Affine-invariant aligned rank tests for the multivariate general linear model with VARMA errors, *Journal of Multivariate Analysis*, **93**(1), 122–163.
- [16] HEATHCOTE, C.R.; RACHEV, S.T. and CHENG, B. (1995). Testing multivariate symmetry, *Journal of Multivariate Analysis*, **54**(1), 91–112.
- [17] HENZE, N.; KLAR, B. and MEINTANIS, S.G. (2003). Invariant tests for symmetry about an unspecified point based on the empirical characteristic function, *Journal of Multivariate Analysis*, **87**(2), 275–297.
- [18] HETTMANSPERGER, T.P.; NYBLUM, J. and OJA, H. (1994). Affine invariant multivariate one-sample sign tests, *Journal of the Royal Statistical Society. Series B (Methodological)*, 221–234.
- [19] HETTMANSPERGER, T.P.; MOTTONEN, J. and OJA, H. (1997). Affine-invariant multivariate one-sample signed-rank tests, *Journal of the American Statistical Association*, **92**(440), 1591–1600.

- [20] HOTELLING, H. (1931). The generalization of Student's ratio, *Annals of Mathematical Statistics*, **2**(3), 360–378.
- [21] JONES, M.C. and PEWSEY, A. (2009). Sinh–arcsinh distributions, *Biometrika*, **96**(4), 761–780.
- [22] KOLTCHINSKII, V.I. and LI, L. (1998). Testing for spherical symmetry of a multivariate distribution, *Journal of Multivariate Analysis*, **65**(2), 228–244.
- [23] KOSHEVEY, G. and MOSLER, K. (1997). Zonoid trimming for multivariate distributions, *The Annals of Statistics*, 1998–2017.
- [24] LI, J. and LIU, R.V. (2004). New nonparametric tests of multivariate locations and scales using data depth, *Statistical Science*, **19**(4), 686–696.
- [25] LIU, R.Y. (1988). On a notion of simplicial depth, *Proceedings of the National Academy of Sciences*, **85**(6), 1732–1734.
- [26] LIU, R.Y.; PARELIUS, J.M. and SINGH, K. (1999). Multivariate analysis by data depth: descriptive statistics, graphics and inference (with discussion and a rejoinder by Liu and Singh), *Annals of Statistics*, **27**(3), 783–858.
- [27] LIU, R.Y. and SINGH, K. (1993). A quality index based on data depth and multivariate rank tests, *Journal of the American Statistical Association*, **88**(421), 252–260.
- [28] LIU, R.Y. and SINGH, K. (2006). Rank tests for multivariate scale difference based on data depth, *DIAMCS Series in Discrete Mathematics and Theoretical Computer Science*, **72**, 17–35.
- [29] MAHFOUD, Z.R. and RANDLES, R.H. (2005). On multivariate signed-rank tests, *J. Nonparametric Statistics*, **17**, 201–216.
- [30] MANZOTTI, A.; PEREZ, F.J. and QUIROZ, A.J. (2002). A statistic for testing the null hypothesis of elliptical symmetry, *Journal of Multivariate Analysis*, **81**(2), 274–285.
- [31] MCWILLIAMS, T.P. (1990). A distribution-free test for symmetry based on a runs statistic, *Journal of the American Statistical Association*, **85**(412), 1130–1133.
- [32] MOTTONEN, J. and OJA, H. (1995). Multivariate Spatial Sign and Rank Methods, *Journal of Nonparametric Statistics*, **5**, 201–213.
- [33] NEUHAUS, G. and ZHU, L.X. (1998). Permutation tests for reflected symmetry, *Journal of Multivariate Analysis*, **67**(2), 129–153.
- [34] OJA, H. (1983). Descriptive statistics for multivariate distributions, *Statistics & Probability Letters*, **1**(6), 327–332.
- [35] RANDLES, R.H. (1989). A distribution-free multivariate sign test based on interdirections, *J. Amer. Statist. Assoc.*, **84**, 1045–1050.
- [36] RANDLES, R.H. (2000). A simpler, affine-invariant, multivariate, distribution-free sign test, *Journal of the American Statistical Association*, **95**(452), 1263–1268.
- [37] ROUSSON, V. (2002). On distribution-free tests for the multivariate two-sample location-scale model, *Journal of Multivariate Analysis*, **80**(1), 43–57.
- [38] SERFLING, R.J. (2006). Multivariate symmetry and asymmetry, *Encyclopedia of Statistical Sciences*.
- [39] TUKEY, J.W. (1975). Mathematics and the picturing of data. In *Proceedings of the International Congress of Mathematicians*, **2**, 523–531.
- [40] TYLER, D.E. (1987). A distribution-free M -estimator of multivariate scatter, *Annals of Statistics*, **15**(1), 234–251.
- [41] PETERS, D. and RANDLES, R.H. (1990). A multivariate signed-rank test for the one-sample location problem, *J. Amer. Statist. Assoc.*, **85**, 552–557.
- [42] ZUO, Y. and SERFLING, R. (2000). General notions of statistical depth function, *Annals of Statistics*, **28**(2), 461–482.

PREDICTION INTERVALS FOR TIME SERIES AND THEIR APPLICATIONS TO PORTFOLIO SELECTION

Authors: SHIH-FENG HUANG

– Department of Applied Mathematics, National University of Kaohsiung, Taiwan
huangsf@nuk.edu.tw

HSIANG-LING HSU

– Institute of Statistics, National University of Kaohsiung, Taiwan
hsuhl@nuk.edu.tw

Received: March 2017

Revised: November 2017

Accepted: December 2017

Abstract:

- This study considers prediction intervals for time series and applies the results to portfolio selection. The dynamics of the high and low underlying returns are depicted by time series models, which lead to a prediction interval of future returns. We propose an innovative criterion for portfolio selection based on the prediction interval. A new concept of coherent risk measures for the interval of returns is introduced. An empirical study is conducted with the stocks of the Dow Jones Industrial Average Index. A self-financing trading strategy is established by daily reallocating the holding positions via the proposed portfolio selection criterion. The numerical results indicate that the proposed prediction interval has promising coverage, efficiency and accuracy for prediction. The proposed portfolio selection criterion constructed from the prediction intervals is capable of suggesting an optimal portfolio according to the economic conditions.

Key-Words:

- *coherent risk measure; portfolio selection; prediction interval.*

AMS Subject Classification:

- 37M10, 91G10.

1. INTRODUCTION

We propose to obtain prediction intervals of a time series by constructing interval-valued time series (ITS) models. The proposed method is used to integrate the information of the daily high, low and closing prices of a stock and is applied to the problem of portfolio selection. Optimal portfolio selection has been extensively discussed in the fields of financial investment and risk management. Markowitz [22, 23] introduced a mean-variance portfolio optimization procedure by using the standard deviation of a portfolio as the measure of risk and assuming that the returns of the underlying assets are independent and identically distributed (i.i.d.). During the past decade, risk measures other than the standard deviation have been considered for selecting investment portfolios. For example, the value-at-risk (VaR), conditional VaR (CVaR) and spectral risk measure (SRM) are commonly used risk measures by market practitioners and analysts in the recent literature on portfolio selection (Rockafellar and Uryasev [25, 26], Acerbi [1], Krokmal *et al.* [20] and Adam *et al.* [2]). However, many empirical findings indicate that the return processes of the underlying assets in financial markets usually exhibit autocorrelation, negative skewness, kurtosis, conditional heteroscedasticity and tail dependence (Tsay [30]). To reflect these features, time series models are used to depict the dynamics of the underlying asset returns for portfolio selection (Harris and Mazibas [16]). However, the development of the above portfolio selection issue uses only information about the closing prices of the underlying assets. The daily high and low prices of a stock are public information and can be observed in the market. The main purpose of this study is to apply daily high and low price information to portfolio selection by ITS models.

One of the main techniques for analyzing ITS is to fit univariate time series models to the interval bounds (Teles and Brito [28]). Maia *et al.* [21] proposed fitting univariate ARIMA models to the midpoints and ranges of the observed interval process and used these models to forecast the interval bounds. Recently, many more complicated ITS models have been proposed and applied to solve problems in various fields. For example, He and Hu [17] used the interval computing approach to forecast the annual and quarterly variability of the stock market. Arroyo *et al.* [3, 4] discussed financial applications based on forecasting with ITS data. García-Ascanio and Maté [13] used vector autoregressive (VAR) models to forecast electric power demand. Yang *et al.* [32] proposed autoregressive conditional interval-valued models with exogenous explanatory interval variables to forecast crude oil prices. Rodrigues and Salish [27] used threshold models to analyze and forecast ITS and applied their model to a weekly sample of S&P500 index returns. Fischer *et al.* [12] predicted stock return volatility using regression models for return intervals. The results of these studies showed that the interval forecasts obtained by ITS perform better than those obtained by the classic approach based on fitting a single time series model to closing prices.

Following Markowitz's [22, 23] approach, the basic idea of various portfolio selection criteria is to determine asset allocations by maximizing the expected investment returns subject to a risk limit of the investment. In addition to daily high and low prices, we also consider the closing prices of a stock. Subsequently, the daily high (low) log returns should be defined as the differences between the logarithms of the daily high (low) price and the last closing price. Therefore, we propose fitting time series models to the daily high and low log returns rather than fitting ITS models directly to the interval bounds of stock prices.

Furthermore, an innovative criterion for portfolio selection is proposed based on the predicted interval of the log returns. Specifically, we maximize the expected high log returns of a portfolio subject to a limitation on the predicted low log returns. We also introduce the concept of a coherent risk measure for the interval of returns, which extends the axioms of the coherent risk measure proposed by Artzner *et al.* [6] for classic financial risk management. In the empirical investigation, we employ the stocks of the companies on the Dow Jones Industrial Average Index (DJIA Index) during the financial crisis period (from July 2, 2007 to June 24, 2009) and under improved market conditions (from July 1, 2014 to June 23, 2016). For each time period, the first 250 daily data are used to fit a time series model to determine the initial trading strategy. A self-financing trading strategy is constructed by daily reallocating the holding weights of the optimal portfolio via the proposed scheme, where a rolling scheme is employed and the time series model is updated with the previous 250 daily historical data. The numerical results indicate that the proposed interval estimation has promising coverage, efficiency and accuracy for predicting high and low prices. Moreover, the proposed portfolio suggests conservative investments during 2008–2009 but aggressive investments during 2015–2016.

The rest of this paper is organized as follows. Section 2 introduces the model assumptions and the prediction interval for ITS. The proposed criterion for portfolio selection using the prediction intervals is introduced in Section 3. Section 4 presents a study to compare the coverage, efficiency and accuracy of the proposed interval estimation for ITS data with those of various approaches in the literature. An empirical study to assess the performance of the self-financing trading strategy constructed by the proposed criterion of portfolio selection is presented in Section 5. Conclusions are given in Section 6.

2. THE PROPOSED INTERVAL TIME SERIES MODEL

Let $P_{m,t}^C$ be the daily closing price of the m -th underlying stock price at time t , and let $P_{m,t}^H$ and $P_{m,t}^L$ be the intraday high and low stock prices, respectively, $m = 1, \dots, p$. Denote the set of information up to time t by \mathcal{F}_t . To obtain a one-step-ahead prediction interval of the price of the m -th underlying stock for a given \mathcal{F}_t , a classic approach is to fit a time series model for the historical closing prices, $P_{m,s}^C$, $s = 1, \dots, t$, and then derive a 95% prediction interval, for example, for $P_{m,t+1}^C$, from the fitted model. Recently, many studies have proposed fitting ITS models for interval observations $[P_{m,s}^L, P_{m,s}^H]$, $s = 1, \dots, t$, and then obtaining an interval estimation of $[P_{m,t+1}^L, P_{m,t+1}^H]$ from the fitted ITS model (see Arroyo *et al.* [3, 4], Teles and Brito [29] and the references therein).

We propose an alternative approach to obtain an estimate of $[P_{m,t+1}^L, P_{m,t+1}^H]$ conditional on \mathcal{F}_t based on the following daily low and high log returns at time t :

$$(2.1) \quad X_{m,t}^{(CL)} = \log(P_{m,t}^L/P_{m,t-1}^C) \quad \text{and} \quad X_{m,t}^{(CH)} = \log(P_{m,t}^H/P_{m,t-1}^C).$$

The definitions of $X_{m,t}^{(CL)}$ and $X_{m,t}^{(CH)}$ are similar to the classic daily log returns, $X_{m,t} = \log(P_{m,t}^C/P_{m,t-1}^C)$ discussed widely in the literature of finance and statistics. $X_{m,t}^{(CL)}$ and $X_{m,t}^{(CH)}$ are capable of depicting realistic investment characteristics. Suppose that an investor buys a given stock on the previous day with closing price $P_{m,t-1}^C$ and sells it on day t .

Then, the investor's return belongs to the interval $[X_{m,t}^{(CL)}, X_{m,t}^{(CH)}]$ depending on when he/she sells the stock during day t . According to the definitions of $X_{m,t}^{(CL)}$ and $X_{m,t}^{(CH)}$ in (2.1), we have the following inequality

$$(2.2) \quad X_{m,t}^{(CL)} \leq_{st} X_{m,t}^{(CH)}$$

since $P_{m,t}^L \leq_{st} P_{m,t}^H$, for all $t = 0, 1, \dots$, and $m = 1, \dots, p$, where the notation $A \leq_{st} B$ means that random variable A is stochastically less than or equal to random variable B . Hence, $\mathbf{X}_{m,t}^{(CI)} = [X_{m,t}^{(CL)}, X_{m,t}^{(CH)}]$, $t = 1, 2, \dots$, also form an ITS, and the prediction interval of $[P_{m,t+1}^L, P_{m,t+1}^H]$ can be obtained. For example, let $[\hat{P}_{m,t+1}^L, \hat{P}_{m,t+1}^H]$ denote the prediction of $[P_{m,t+1}^L, P_{m,t+1}^H]$ conditional on \mathcal{F}_t . By using (2.1), our proposed scheme is to model the interval observations, $[X_{m,s}^{(CL)}, X_{m,s}^{(CH)}]$, $s = 1, \dots, t$, and then estimate $[\hat{P}_{m,t+1}^L, \hat{P}_{m,t+1}^H]$ by

$$\left[P_{m,t}^C \exp\{\hat{X}_{m,t+1}^{(CL)}\}, P_{m,t}^C \exp\{\hat{X}_{m,t+1}^{(CH)}\} \right],$$

where $\hat{X}_{m,t+1}^{(CL)}$ and $\hat{X}_{m,t+1}^{(CH)}$ are the predictions of $X_{m,t+1}^{(CL)}$ and $X_{m,t+1}^{(CH)}$, respectively, which can be obtained from the time series models defined below. Traditionally, ITS data are formed by only the high and low prices (Arroyo *et al.* [3, 4] and Maia *et al.* [21]). This study includes the closing prices in the model and investigates whether this additional information can improve the interval prediction.

To jointly model $X_{m,t}^{(h)}$, $h = CL, CH$, we need to capture the features inherent in the data. For example, $X_{m,t}^{(h)}$, $h = CL, CH$ could be conditionally heteroscedastic and auto- and cross-correlated. To characterize these features, a two-stage procedure is proposed to model the dynamics of $X_{m,t}^{(h)}$, $h = CL, CH$. The first stage is to adjust the conditional heteroscedasticity of $X_{m,t}^{(h)}$ marginally for $h = CL, CH$. The second stage is to simultaneously model the auto- and cross-correlation of the adjusted time series.

In the first stage, we propose to de-GARCH $X_{m,t}^{(h)}$ to obtain volatility-adjusted returns. De-GARCHing is a widely used technique for modeling multivariate time series. For example, Engle [10, 11] proposed a dynamic conditional correlation (DCC) model to capture time-varying correlations. The first step of their scheme is to de-GARCH the data. Härdle *et al.* [15] also used de-GARCHing with a GARCH(1,1) model to analyze the multi-dimensional dependencies of time series data with a hidden Markov model for hierarchical Archimedean copulae. Grigoryeva *et al.* [14] proposed a method based on various state space models to extract global stochastic (GST) financial trends from non-synchronous financial data. They mentioned that de-GARCHing is commonly used for GST. In this study, we propose to fit $X_{m,t}^{(h)}$ with a univariate ARMA-GARCH model and let

$$(2.3) \quad \tilde{X}_{m,t}^{(h)} = (X_{m,t}^{(h)} - \mu_m^{(h)}) / \sigma_{m,t}^{(h)}$$

be the de-GARCHed process of $X_{m,t}^{(h)}$, $h = CL, CH$, where $\mu_m^{(h)}$ is the stationary (unconditional) mean of $X_{m,t}^{(h)}$ and $\sigma_{m,t}^{(h)}$ is the conditional standard deviation of $X_{m,t}^{(h)}$, which is estimated from the univariate GARCH-type model

$$(2.4) \quad \sigma_{m,t}^{(h)} = g_{m,t-1}^{(h)}(X_{m,s}^{(h)}, \sigma_{m,s}^{(h)}, s < t),$$

which is \mathcal{F}_{t-1} -measurable. This type of model (2.4) is capable of describing many features of financial data, for example, conditional heteroscedasticity, volatility clustering and asymmetry. It also includes various univariate financial time series models that are widely used

by practitioners in economics, statistics and finance (see Engle [9], Bollerslev [8], Nelson [24], Tsay [30] and the references therein). In particular, we employ the stationary mean (not the conditional mean) to define the proposed de-GARCHed process in (2.3). The main reason for this design is to retain the autocorrelation in $\tilde{X}_{m,t}^{(h)}$, $h = CL, CH$ and to model the auto- and cross-correlation of $\tilde{X}_{m,t}^{(h)}$, $h = CL, CH$ simultaneously in the second stage of the proposed procedure.

In the second stage, we employ the following vector autoregressive-moving-average model of orders p and q , denoted by VARMA(p, q), to depict the dynamics of the two de-GARCHed processes, $\{\tilde{X}_{m,t}^{(h)}, t = 1, \dots, T\}$, $h = CL, CH$,

$$(2.5) \quad \begin{pmatrix} \tilde{X}_{m,t}^{(CL)} \\ \tilde{X}_{m,t}^{(CH)} \end{pmatrix} = \sum_{i=1}^p \begin{pmatrix} \phi_{m,i}^{LL} & \phi_{m,i}^{LH} \\ \phi_{m,i}^{HL} & \phi_{m,i}^{HH} \end{pmatrix} \begin{pmatrix} \tilde{X}_{m,t-i}^{(CL)} \\ \tilde{X}_{m,t-i}^{(CH)} \end{pmatrix} + \begin{pmatrix} \varepsilon_{m,t}^{(CL)} \\ \varepsilon_{m,t}^{(CH)} \end{pmatrix} + \sum_{j=1}^q \begin{pmatrix} \theta_{m,j}^{LL} & \theta_{m,j}^{LH} \\ \theta_{m,j}^{HL} & \theta_{m,j}^{HH} \end{pmatrix} \begin{pmatrix} \varepsilon_{m,t-j}^{(CL)} \\ \varepsilon_{m,t-j}^{(CH)} \end{pmatrix},$$

for $m = 1, \dots, p$, where $(\varepsilon_{m,t}^{(CL)}, \varepsilon_{m,t}^{(CH)})^\top$, $t = 1, \dots, T$, are uncorrelated random vectors of a bivariate normal distribution with mean zero and covariance matrix Σ . In addition, $(\varepsilon_{m,t}^{(CL)}, \varepsilon_{m,t}^{(CH)})^\top$, $t = 1, \dots, T$, are assumed to be independent of $(\tilde{X}_{m,s}^{(CL)}, \tilde{X}_{m,s}^{(CH)})^\top$, $s < t$.

Denote the 1-step-ahead predictions of $X_{m,t+1}^{(h)}$ conditional on \mathcal{F}_t by $\hat{X}_{m,t}^{(h)}(1) = E_t(X_{m,t+1}^{(h)})$, $h = CL, CH$, where $E_t(X)$ denotes the conditional expectation of X given \mathcal{F}_t . From (2.3)–(2.5), we have

$$(2.6) \quad \begin{aligned} \hat{X}_{m,t}^{(CL)}(1) &= E_t(X_{m,t+1}^{(CL)}) = \mu_m^{(CL)} + \sigma_{m,t+1}^{(CL)} E_t(\tilde{X}_{m,t+1}^{(CL)}) \\ &= \mu_m^{(CL)} + \sigma_{m,t+1}^{(CL)} \left\{ \sum_{i=1}^p \left(\phi_{m,i}^{LL} \tilde{X}_{m,t+1-i}^{(CL)} + \phi_{m,i}^{LH} \tilde{X}_{m,t+1-i}^{(CH)} \right) \right. \\ &\quad \left. + \sum_{j=1}^q \left(\theta_{m,j}^{LL} \varepsilon_{m,t+1-j}^{(CL)} + \theta_{m,j}^{LH} \varepsilon_{m,t+1-j}^{(CH)} \right) \right\} \end{aligned}$$

and

$$(2.7) \quad \begin{aligned} \hat{X}_{m,t}^{(CH)}(1) &= E_t(X_{m,t+1}^{(CH)}) = \mu_m^{(CH)} + \sigma_{m,t+1}^{(CH)} E_t(\tilde{X}_{m,t+1}^{(CH)}) \\ &= \mu_m^{(CH)} + \sigma_{m,t+1}^{(CH)} \left\{ \sum_{i=1}^p \left(\phi_{m,i}^{HL} \tilde{X}_{m,t+1-i}^{(CL)} + \phi_{m,i}^{HH} \tilde{X}_{m,t+1-i}^{(CH)} \right) \right. \\ &\quad \left. + \sum_{j=1}^q \left(\theta_{m,j}^{HL} \varepsilon_{m,t+1-j}^{(CL)} + \theta_{m,j}^{HH} \varepsilon_{m,t+1-j}^{(CH)} \right) \right\}. \end{aligned}$$

To guarantee the mathematical coherence $\hat{X}_{m,t+1}^{(CL)} \leq_{st} \hat{X}_{m,t+1}^{(CH)}$ in their predictions, let

$$\hat{X}_{m,t+1}^{(CL)} = \min \left\{ \hat{X}_{m,t}^{(CL)}(1), \hat{X}_{m,t}^{(CH)}(1) \right\}$$

and

$$\hat{X}_{m,t+1}^{(CH)} = \max \left\{ \hat{X}_{m,t}^{(CL)}(1), \hat{X}_{m,t}^{(CH)}(1) \right\},$$

and $[\widehat{X}_{m,t+1}^{(CL)}, \widehat{X}_{m,t+1}^{(CH)}]$ forms a prediction interval of X_{t+1} conditional on \mathcal{F}_t . In our empirical study, there are 250(days) \times 30(companies) \times 2(time periods) = 15,000 prediction intervals, and the situation of $\widehat{X}_{m,t}^{(CL)}(1) > \widehat{X}_{m,t}^{(CH)}(1)$ occurs only 8 times. The numerical results indicate that the proposed scheme is capable of guaranteeing $\widehat{X}_{m,t+1}^{(CL)} \leq_{st} \widehat{X}_{m,t+1}^{(CH)}$ in most cases.

3. APPLICATION OF THE ITS PREDICTION TO PORTFOLIO SELECTION

In this section, we propose an innovative portfolio selection scheme on the basis of the ITS prediction with models (2.6) and (2.7). The literature contains many different models from models (2.6) and (2.7) for analyzing ITS. Nevertheless, the proposed portfolio selection scheme is not restricted to our considered model.

The classic portfolio optimization problem is represented as follows:

$$(3.1) \quad \max_{\mathbf{c}_t} E_t \left(\sum_{m=1}^p c_{m,t} X_{m,t+1} \right) \quad \text{subject to} \quad \mathbf{c}_t \geq 0, \quad \sum_{m=1}^p c_{m,t} \leq 1 \quad \text{and} \quad \rho_t \leq L,$$

where $\mathbf{c}_t = (c_{1,t}, \dots, c_{p,t})^\top$, $c_{m,t}$ denotes the holding position of $X_{m,t}$ at time t , $\mathbf{c}_t \geq 0$ is the no short-selling constraint, $\sum_{m=1}^p c_{m,t} \leq 1$ is the budget constraint, ρ_t is the value of a pre-determined risk measure at time t , and L is a pre-specified upper bound of the investment risk. The main objective is to select the holding positions \mathbf{c}_t at time t . In the portfolio selection literature, when $X_{m,t}$, $t = 1, 2, \dots$, are assumed to be i.i.d. for each $m = 1, \dots, p$, Markowitz [22, 23] used the standard deviation of a portfolio, Rockafellar and Uryasev [25, 26] and Krokmal *et al.* [20] employed the CVaR, and Adam *et al.* [2] considered the SRM as the risk measure to determine \mathbf{c}_t . Recently, Harris and Mazibas [16] and Huang *et al.* [18] further considered fitting time series models for the underlying asset returns, $X_{m,t}$, $m = 1, \dots, p$, $t = 1, 2, \dots$, with the CVaR and SRM to solve (3.1).

In this study, we determine the allocations of the underlying assets with the following criterion:

$$(3.2) \quad \max_{\mathbf{c}_t} E_t \left(\sum_{m=1}^p c_{m,t} X_{m,t+1}^{(CH)} \right) \quad \text{subject to} \quad \mathbf{c}_t \geq 0, \quad \sum_{m=1}^p c_{m,t} \leq 1$$

$$\quad \text{and} \quad - \sum_{m=1}^p c_{m,t} E_t \left(X_{m,t+1}^{(CL)} \mid X_{m,t+1}^{(CL)} \leq q_{\alpha,m,t+1} \right) \leq L,$$

where $X_{m,t+1}^{(h)} = \mu_m^{(h)} + \sigma_{m,t+1}^{(h)} \widetilde{X}_{m,t+1}^{(h)}$, $h = CH, CL$, follows models (2.3)–(2.5), and $q_{\alpha,m,t+1}$ is the α -th quantile of $X_{m,t+1}^{(CL)}$ conditional on \mathcal{F}_t . In practice, since the expected values of daily stock returns are usually very close to 0, one can select a sufficiently small α such that $q_{\alpha,m,t+1} < 0$. The main concept behind (3.2) is to maximize the potential high portfolio returns subject to a predetermined limitation, L , on the corresponding potential low and nonpositive returns. In contrast to (3.1), we use $E_t(\sum_{m=1}^p c_{m,t} X_{m,t+1}^{(CH)})$ to replace $E_t(\sum_{m=1}^p c_{m,t} X_{m,t+1})$ and use

$$(3.3) \quad - \sum_{m=1}^p c_{m,t} E_t \left(X_{m,t+1}^{(CL)} \mid X_{m,t+1}^{(CL)} \leq q_{\alpha,m,t+1} \right)$$

as the risk measure ρ_t in (3.1). In addition, the values of $E_t(X_{m,t+1}^{(CH)})$ and $E_t(X_{m,t+1}^{(CL)})$, $m = 1, \dots, p$, are estimated by the models defined in (2.5). Moreover, the optimal allocations $c_{m,t}$, $m = 1, \dots, p$, are in linear forms in the objective function and constraints in (3.2). Consequently, the optimal allocations in (3.2) can be obtained by linear programming, which is a popular technique for various portfolio selection criteria (Markowitz [22, 23], Rockafellar and Uryasev [25, 26], Adam *et al.* [2] and Huang *et al.* [18]).

In the following, we introduce the concept of a coherent risk measure for the intervals of returns, which provides economic and financial reasons to use (3.3) as a risk constraint in (3.2). In financial risk management, Artzner *et al.* [6] introduced the following concept of the coherent risk measure for classic portfolio selection. Let \mathcal{G} be the set of random portfolio returns, ρ be a risk measure, which is a mapping from \mathcal{G} into \mathbb{R} , and X denote the return of an asset. A risk measure is called coherent if it satisfies the following properties:

- (A1) Translation invariance: If A is a deterministic portfolio with guaranteed return α , then for all $X \in \mathcal{G}$, we have $\rho(X + A) = \rho(X) - \alpha$.
- (A2) Subadditivity: For all X and $Y \in \mathcal{G}$, $\rho(X + Y) \leq \rho(X) + \rho(Y)$.
- (A3) Positive homogeneity: For all $\lambda \geq 0$ and all $X \in \mathcal{G}$, $\rho(\lambda X) = \lambda \rho(X)$.
- (A4) Monotonicity: For all X and $Y \in \mathcal{G}$ with $X \leq Y$, we have $\rho(Y) \leq \rho(X)$.

The economic explanations of these four properties are as follows. Translation invariance implies that the addition of a definite amount of capital reduces the risk by the same amount. Subadditivity implies that diversification is beneficial. Positive homogeneity implies that the risk of a position is proportional to its size. Monotonicity implies that a portfolio with greater future returns has less risk.

In this study, we consider an interval of returns denoted by $\mathbf{X}^I = [X^L, X^H]$, where X^L and X^H are the low and high returns of an asset, respectively. To extend the concepts of (A1)–(A4) from random variables to random intervals, we propose the following properties for a risk measure of the interval of returns. Let \mathcal{G}_I be the set of random intervals of portfolio returns and $\rho_I: \mathcal{G}_I \rightarrow \mathbb{R}$ be a corresponding risk measure.

- (A1') Translation invariance for the interval of returns: If A is a deterministic portfolio with guaranteed return α , then for all $\mathbf{X}^I \in \mathcal{G}_I$, we have $\rho_I(\mathbf{X}^I + A) = \rho_I(\mathbf{X}^I) - \alpha$, where we use $\mathbf{X}^I + A$ to denote $[X^L + A, X^H + A]$.
- (A2') Subadditivity for the interval of returns: For all \mathbf{X}^I and $\mathbf{Y}^I \in \mathcal{G}_I$, $\rho_I(\mathbf{X}^I + \mathbf{Y}^I) \leq \rho_I(\mathbf{X}^I) + \rho_I(\mathbf{Y}^I)$, where $\mathbf{X}^I + \mathbf{Y}^I = [X^L + Y^L, X^H + Y^H]$. In addition, one can also use the Cartesian join of \mathbf{X}^I and \mathbf{Y}^I , denoted by $\mathbf{X}^I \oplus \mathbf{Y}^I = [\min(X^L, Y^L), \max(X^H, Y^H)]$, to define the subadditivity, that is, $\rho_I(\mathbf{X}^I \oplus \mathbf{Y}^I) \leq \rho_I(\mathbf{X}^I) + \rho_I(\mathbf{Y}^I)$.
- (A3') Positive homogeneity for the interval of returns: For all $\lambda \geq 0$ and all $\mathbf{X}^I \in \mathcal{G}_I$, $\rho_I(\lambda \mathbf{X}^I) = \lambda \rho_I(\mathbf{X}^I)$.
- (A4') Monotonicity for the interval of returns: For all \mathbf{X}^I and $\mathbf{Y}^I \in \mathcal{G}_I$ with $\mathbf{X}^I \leq \mathbf{Y}^I$, where $\mathbf{X}^I \leq \mathbf{Y}^I$ if and only if $X^L \leq Y^L$ and $X^H \leq Y^H$, we have $\rho_I(\mathbf{Y}^I) \leq \rho_I(\mathbf{X}^I)$.

The economic explanations of (A1')–(A4') are similar to those of (A1)–(A4). Specifically, the monotonicity for the interval of returns (A4') implies only that a portfolio with greater future interval of returns has less risk. For the case of $\mathbf{X}^I \subset \mathbf{Y}^I$, the relationship between $\rho_I(\mathbf{Y}^I)$ and $\rho_I(\mathbf{X}^I)$ is not clear. If a risk measure for the interval of returns satisfies (A1')–(A4'), we call it a coherent risk measure for the interval of returns. In the following proposition, a coherent risk measure for the interval of returns is proposed.

Proposition 3.1. *Let $\mathbf{X}^I = [X^L, X^H]$ be an interval of returns, and let*

$$\rho_I(\mathbf{X}^I) = -E(X^{(L)} | X^{(L)} \leq q_\alpha),$$

where q_α is the α -th quantile of $X^{(L)}$. Then, $\rho_I(\cdot)$ is a coherent risk measure for the interval of returns.

Proof: Note that for a random variable X^L , $-E(X^L | X^L \leq q_\alpha)$ is the so-called expected shortfall, which is a coherent risk measure. Therefore, it is straightforward to obtain that $\rho_I(\mathbf{X}^I) = -E(X^L | X^L \leq q_\alpha)$ satisfies (A1'), (A2'), (A3'), and (A4'), where the interval addition in (A2') is defined by the usual way $\mathbf{X}^I + \mathbf{Y}^I = [X^L + Y^L, X^H + Y^H]$. In the following, we prove that $\rho_I(\mathbf{X}^I)$ also satisfies $\rho_I(\mathbf{X}^I \oplus \mathbf{Y}^I) \leq \rho_I(\mathbf{X}^I) + \rho_I(\mathbf{Y}^I)$.

Let $q_{\alpha,0}$, $q_{\alpha,X}$ and $q_{\alpha,Y}$ be the α -th quantile of $\min(X^L, Y^L)$, X^L and Y^L , respectively. Apparently, $q_{\alpha,0} \leq \min(q_{\alpha,X}, q_{\alpha,Y})$ for any $\alpha \in (0, 1)$. Let α be small enough such that $\max(q_{\alpha,X}, q_{\alpha,Y}) < 0$. Consequently, for all \mathbf{X}^I and $\mathbf{Y}^I \in \mathcal{G}_1$, we have

$$\begin{aligned} \rho_I(\mathbf{X}^I \oplus \mathbf{Y}^I) &= -E\left[\min(X^L, Y^L) \mid \min(X^L, Y^L) \leq q_{\alpha,0}\right] \\ &= -\frac{1}{\alpha} E\left[\min(X^L, Y^L) I(\min(X^L, Y^L) \leq q_{\alpha,0})\right] \\ &\leq -\frac{1}{\alpha} \left\{E[X^L I(X^L \leq q_{\alpha,X})] + E[Y^L I(Y^L \leq q_{\alpha,Y})]\right\} \\ &= -\left\{E(X^L | X^L \leq q_{\alpha,X}) + E(Y^L | Y^L \leq q_{\alpha,Y})\right\} \\ &= \rho_I(\mathbf{X}^I) + \rho_I(\mathbf{Y}^I), \end{aligned}$$

where $I(\cdot)$ is an indicator function and the inequality holds by using the facts that

$$-\min(X^L, Y^L) I(\min(X^L, Y^L) \leq q_{\alpha,0}) \leq -X^L I(X^L \leq q_{\alpha,X}) - Y^L I(Y^L \leq q_{\alpha,Y}),$$

almost surely, for $q_{\alpha,0} \leq \min(q_{\alpha,X}, q_{\alpha,Y}) \leq \max(q_{\alpha,X}, q_{\alpha,Y}) < 0$. Thus, (A2') with the Cartesian join also holds and the proof is complete. \square

By Proposition 3.1, the measurement defined in (3.3) can be rewritten as

$$\sum_{m=1}^p c_{m,t} \rho_I(\mathbf{X}_{m,t+1}^{(CI)} | \mathcal{F}_t),$$

which is a linear combination of coherent risk measures for the interval of returns, where

$$(3.4) \quad \rho_I(\mathbf{X}_{m,t+1}^{(CI)} | \mathcal{F}_t) = -E_t\left(X_{m,t+1}^{(CL)} \mid X_{m,t+1}^{(CL)} \leq q_{\alpha,m,t+1}\right).$$

Due to the convexity of the coherent risk measure, we have

$$(3.5) \quad \rho_I \left(\sum_{m=1}^p c_{m,t} \mathbf{X}_{m,t+1}^{(CI)} \mid \mathcal{F}_t \right) \leq \sum_{m=1}^p c_{m,t} \rho_I (\mathbf{X}_{m,t+1}^{(CI)} \mid \mathcal{F}_t).$$

For a portfolio with allocations $c_{m,t}$, $m = 1, \dots, p$, set up at time t , the left side of (3.5) represents the risk of the worst case occurring at time $t + 1$ since each underlying return reaches the bottom of the corresponding prediction interval. However, if a limitation is set on $\rho_I(\sum_{m=1}^p c_{m,t} \mathbf{X}_{m,t+1}^{(CI)} \mid \mathcal{F}_t)$ in the portfolio selection criterion (3.2), the optimal allocations $c_{m,t}$, $m = 1, \dots, p$, are difficult to obtain directly using linear programming since $\rho_I(\cdot \mid \mathcal{F}_t)$ is a nonlinear function of $c_{m,t}$, $m = 1, \dots, p$. A similar situation is encountered in the classic portfolio selection problem shown in (3.1) when using the expected shortfall as the risk measure. Rockafellar and Uryasev [25, 26] proposed a method to overcome this difficulty by considering more latent variables, but the computational cost also increased. Therefore, we set a limitation on the right side of (3.5), and the optimal allocations can be obtained directly using linear programming.

In the following sections, we consider several scenarios to investigate the coverage, efficiency and accuracy of the proposed interval estimation and the performance of the proposed criterion for portfolio selection.

4. EVALUATION OF THE PROPOSED INTERVAL ESTIMATION METHOD

Let $Y_t = [P_t^L, P_t^H]$ denote the realized ITS of the stock prices and \hat{Y}_t be an estimation of Y_t , $t = 1, \dots, T$. In this section, we use the four measures to evaluate the performance of the proposed interval estimation (He and Hu [17], Rodrigues and Salish [27] and Xiong *et al.* [31]). The first measure is the coverage rate

$$R_C = \frac{1}{T} \sum_{t=1}^T \frac{w(Y_t \cap \hat{Y}_t)}{w(Y_t)},$$

where $w(\cdot)$ denotes the width of the interval, R_C indicates what part of the realized ITS of the stock prices is covered by its forecast.

The second measure is the efficiency rate

$$R_E = \frac{1}{T} \sum_{t=1}^T \frac{w(Y_t \cap \hat{Y}_t)}{w(\hat{Y}_t)},$$

which provides information about what part of the forecast covers the realized ITS. It should be noted that R_C and R_E must be considered simultaneously; otherwise, incorrect conclusions may be drawn. For example, if Y_t is a subinterval of \hat{Y}_t , then R_C will be 1, but R_E might be much less than 1, which indicates that the predicted interval is much wider than the realized ITS. Therefore, we only conclude that the forecast is satisfactory when R_C and R_E are reasonably high and the difference between them is small.

The third measure is the accuracy ratio

$$R_A = \frac{1}{T} \sum_{t=1}^T \frac{w(Y_t \cap \hat{Y}_t)}{w(Y_t \cup \hat{Y}_t)}.$$

A prediction with a larger R_A performs better than a prediction with a smaller one.

The fourth measure is the U_I criterion

$$U_I = \sqrt{\frac{\sum_{t=1}^T (P_t^H - \widehat{P}_t^H)^2 + \sum_{t=1}^T (P_t^L - \widehat{P}_t^L)^2}{\sum_{t=1}^T (P_t^H - P_{t-1}^H)^2 + \sum_{t=1}^T (P_t^L - P_{t-1}^L)^2}},$$

which is derived from Theil's U statistic and compares the performance of an estimated method with a naïve estimate $[P_{t-1}^L, P_{t-1}^H]$ of $[P_t^L, P_t^H]$. The U_I statistic is less than one if the predictor performs better than the naïve predictor.

In addition to the proposed interval estimation $\widehat{Y}_t^{(p)}$, three commonly used interval predictors are considered in our comparison studies. One is fitting time series models to the log return process $X_t = \log(P_t^C/P_{t-1}^C)$ and then deriving the corresponding 95% confidence interval of P_{t+1}^C . We denote this estimation of Y_t by $\widehat{Y}_t^{(1)}$.

The second estimation of Y_t is the popular center-range prediction interval, which is obtained by separately fitting time series models to the processes of the center, $P_t^M = (P_t^H + P_t^L)/2$, and the range, $P_t^R = (P_t^H - P_t^L)/2$, of the price intervals and then deriving an interval estimation of $[P_{t+1}^L, P_{t+1}^H]$ conditional on \mathcal{F}_t . We denote the second estimation by $\widehat{Y}_t^{(2)}$.

The third alternative estimation of Y_t is derived from a linear interval-data model motivated from Fischer *et al.* [12]. The center-range-representation of interval data can also be expressed as the following regression model

$$(4.1) \quad Y_t = \beta_0 + \beta_1 Y_{t-1}^C + \beta_2 Y_{t-1}^R + \delta_t,$$

where $Y_t^C = [P_t^M, P_t^M]$, $Y_t^R = [-P_t^R, P_t^R]$, δ_t is an interval-valued random error, and $\beta_0 = [a, a]$ and (a, β_1, β_2) are unknown parameters. Blanco-Fernández *et al.* [7] derived the estimation procedures for (4.1), and the obtained predictor is denoted as $\widehat{Y}_t^{(3)}$.

We conduct the comparison study using the stock prices of the 30 companies of the DJIA Index during the financial crisis period (from July 2, 2007 to June 24, 2009) and under improved market conditions (from July 1, 2014 to June 23, 2016). The 1-step-ahead prediction intervals during the two time periods (from June 27, 2008 to June 24, 2009 and from June 29, 2015 to June 23, 2016) are obtained with the previous 250 daily historical high and low returns. We adopt an ARMA(p, q)-GARCH(p_0, q_0) model, where $p, q \in \{0, 1, 2, 3, 4, 5\}$ and $p_0, q_0 \in \{0, 1\}$, to obtain the de-GARCHed process defined in (2.3) for $h = CL$ and CH , separately. The multivariate portmanteau test (Tsay [30], Chapter 8) is used for testing the auto- and cross-correlation in $\{(\widetilde{X}_{m,t}^{(CL)}, \widetilde{X}_{m,t}^{(CH)})\}$, $t = 1, \dots, T$. If the de-GARCHed processes have significant auto- and cross-correlation, we model the vector time series $(\widetilde{X}_{m,s}^{(CL)}, \widetilde{X}_{m,t}^{(CH)})^\top$ with VARMA(p_1, q_1) defined in (2.5), where (p_1, q_1) are selected from $\{(1, 0), (0, 1) \text{ and } (u, v), u, v = 1, 2, 3\}$ based on the Bayesian information criterion (BIC). Table 1 summarizes the p -values of the multivariate portmanteau test for the de-GARCHed processes and the residual processes $\{(\varepsilon_{m,t}^{(CL)}, \varepsilon_{m,t}^{(CH)})\}$, $t = 1, \dots, T$. In Table 1, all the de-GARCHed processes have significant auto- and cross-correlation during 2008–2009, and most (around 96.2%) of the de-GARCHed processes have significant auto- and cross-correlation during 2015–2016. More than 99.4% of the p -values of the fitted residual processes during the two time periods are greater than 0.01, which indicates that the above scheme is capable of removing most of the auto- and cross-correlation of the de-GARCHed processes.

Table 1: The proportions of the p -values of the multivariate portmanteau test for testing auto- and cross-correlation in the de-GARCHed processes $\{(\tilde{X}_{m,t}^{(CL)}, \tilde{X}_{m,t}^{(CH)}), t = 1, \dots, T\}$ (shown in rows) and the residual processes $\{(\varepsilon_{m,t}^{(CL)}, \varepsilon_{m,t}^{(CH)}), t = 1, \dots, T\}$ (shown in columns).

(a) 2008–2009

de-GARCHed \ residual	residual	p -value < 0.01	p -value \geq 0.01
p -value < 0.01		0.001	0.999
p -value \geq 0.01		0.000	0.000

(b) 2015–2016

de-GARCHed \ residual	residual	p -value < 0.01	p -value \geq 0.01
p -value < 0.01		0.011	0.951
p -value \geq 0.01		0.000	0.038

Figure 1 summarizes the proportions of selected orders (p_1, q_1) in the two time periods, where the 3.8% de-GARCHed processes without significant auto- and cross-correlation during 2015–2016 are denoted by VARMA(0,0). VARMA(1,1) is the most commonly selected model during the financial crisis period, whereas VARMA(1,0) and VARMA(1,1) are frequently selected under improved market conditions.

Table 2: The average values of R_C , R_E , R_A and U_I of $\hat{Y}_t^{(1)}$, $\hat{Y}_t^{(2)}$, $\hat{Y}_t^{(3)}$ and $\hat{Y}_t^{(p)}$ in June 27, 2008 – June 24, 2009 and June 29, 2015 – June 23, 2016, in the top panel. The bottom panel presents the improvement of $\hat{Y}_t^{(p)}$ for each $\hat{Y}_t^{(i)}$, $i = 1, 2, 3$, by calculating $(\hat{Y}_t^{(p)} - \hat{Y}_t^{(i)})/\hat{Y}_t^{(i)}$ for R_C , R_E , and R_A , and $(\hat{Y}_t^{(i)} - \hat{Y}_t^{(p)})/\hat{Y}_t^{(i)}$ for U_I .

Average values								
	2008–2009				2015–2016			
	$\hat{Y}_t^{(1)}$	$\hat{Y}_t^{(2)}$	$\hat{Y}_t^{(3)}$	$\hat{Y}_t^{(p)}$	$\hat{Y}_t^{(1)}$	$\hat{Y}_t^{(2)}$	$\hat{Y}_t^{(3)}$	$\hat{Y}_t^{(p)}$
R_C	0.57	0.96	0.61	0.63	0.53	0.96	0.55	0.60
R_E	0.55	0.34	0.53	0.60	0.51	0.30	0.51	0.56
R_A	0.42	0.33	0.41	0.46	0.39	0.30	0.39	0.44
U_I	0.98	1.70	0.99	0.87	0.99	1.82	0.97	0.88

Improvement of $\hat{Y}_t^{(p)}$ for each $\hat{Y}_t^{(i)}$, $i = 1, 2, 3$						
	2008–2009			2015–2016		
	$\hat{Y}_t^{(1)}$	$\hat{Y}_t^{(2)}$	$\hat{Y}_t^{(3)}$	$\hat{Y}_t^{(1)}$	$\hat{Y}_t^{(2)}$	$\hat{Y}_t^{(3)}$
R_C	10.4%	–35.0%	1.7%	14.6%	–37.3%	9.4%
R_E	8.5%	75.5%	13.0%	9.0%	85.6%	10.5%
R_A	9.8%	38.5%	12.3%	12.4%	47.3%	10.9%
U_I	9.8%	47.9%	10.99%	10.5%	51.4%	9.2%

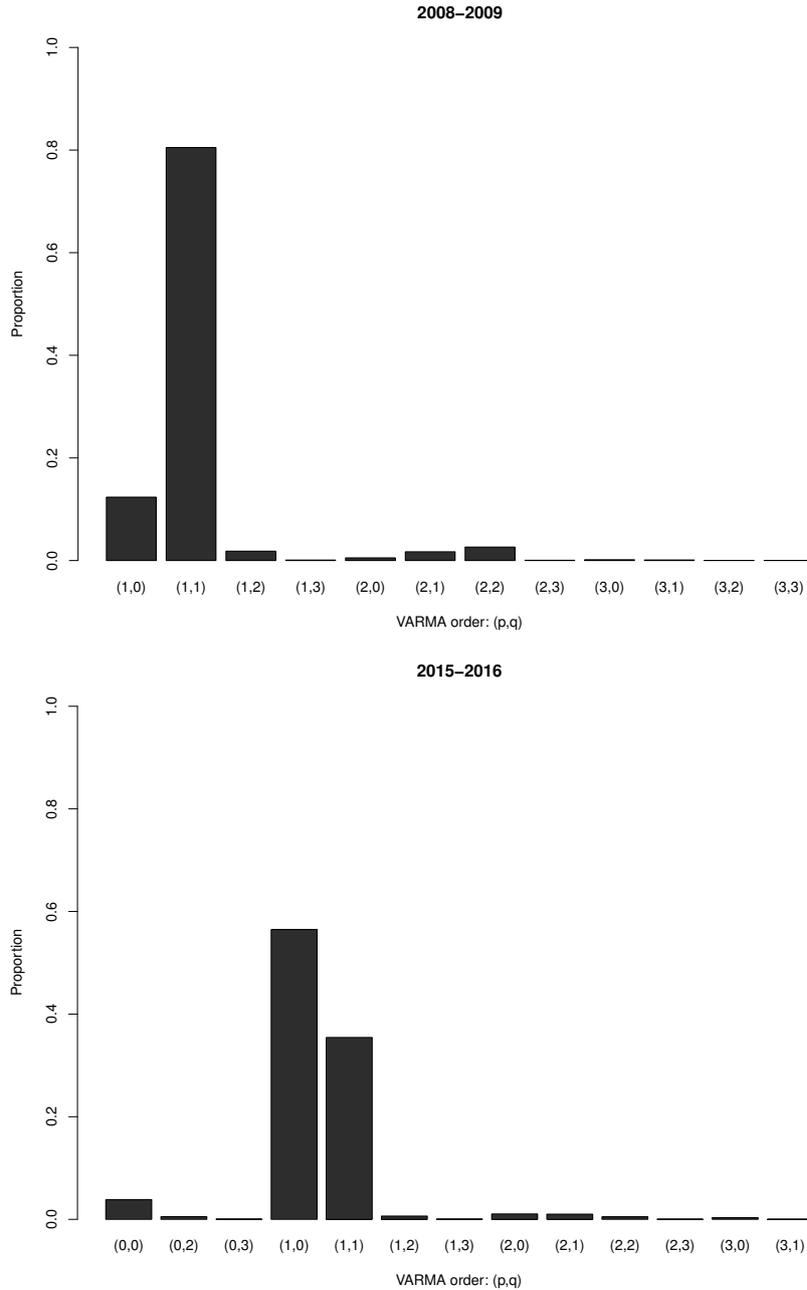


Figure 1: Summaries of the selected orders of VARMA for 15,000 prediction intervals in the two time periods (June 27, 2008 to June 24, 2009 and June 29, 2015 to June 23, 2016).

Table 2 presents the average values of R_C , R_E , R_A and U_I of $\hat{Y}_t^{(p)}$ and $\hat{Y}_t^{(i)}$, $i = 1, 2, 3$, in the top panel. In the bottom panel, we present the improvement of $\hat{Y}_t^{(p)}$ for each $\hat{Y}_t^{(i)}$, $i = 1, 2, 3$, by calculating $(\hat{Y}_t^{(p)} - \hat{Y}_t^{(i)})/\hat{Y}_t^{(i)}$ for R_C , R_E , and R_A and $(\hat{Y}_t^{(i)} - \hat{Y}_t^{(p)})/\hat{Y}_t^{(i)}$ for U_I . The numerical results indicate that $\hat{Y}_t^{(p)}$ performs better than $\hat{Y}_t^{(i)}$, $i = 1, 2, 3$, in terms of R_E , R_A and U_I . Although $\hat{Y}_t^{(2)}$ has a larger R_C than $\hat{Y}_t^{(p)}$, the improvement of $\hat{Y}_t^{(p)}$ in R_E is much greater than the loss of $\hat{Y}_t^{(p)}$ in R_C . In particular, the popular center-range prediction interval $\hat{Y}_t^{(2)}$ has U_I greater than 1, which indicates that the proposed prediction interval $\hat{Y}_t^{(p)}$ is more reliable than $\hat{Y}_t^{(2)}$. By contrast, $\hat{Y}_t^{(p)}$ outperforms $\hat{Y}_t^{(1)}$ and $\hat{Y}_t^{(3)}$, especially in 2015-2016, with an improvement in the 4 measures of at least 9.0%.

Figure 2 presents the average values of R_C, R_E, R_A and U_I for the four prediction intervals for the 30 companies of the DJIA Index from June 27, 2008 to June 24, 2009. The results of the time period from June 29, 2015 to June 23, 2016 are given in Figure 3. These figures reveal similar findings as those in Table 2. The proposed prediction interval $\widehat{Y}_t^{(p)}$ has the best performance with respect to R_A and U_I and performs robustly in R_C and R_E , especially in 2015–2016.

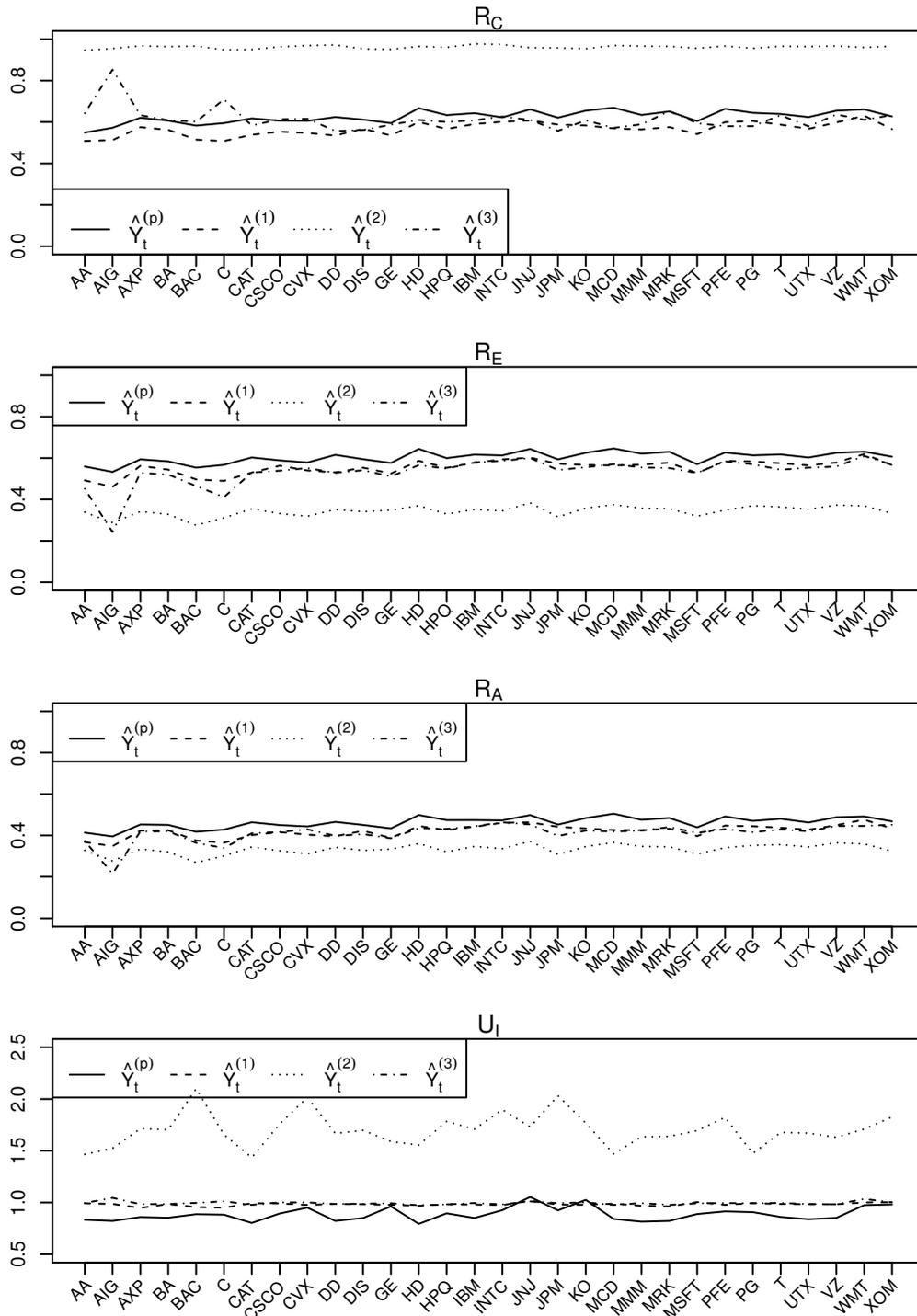


Figure 2: The average values of R_C, R_E, R_A and U_I of $\widehat{Y}_t^{(p)}$ (solid line), $\widehat{Y}_t^{(1)}$ (dashed line), $\widehat{Y}_t^{(2)}$ (dotted line) and $\widehat{Y}_t^{(3)}$ (dash-dotted line) for 30 different time series from June 27, 2008 to June 24, 2009.

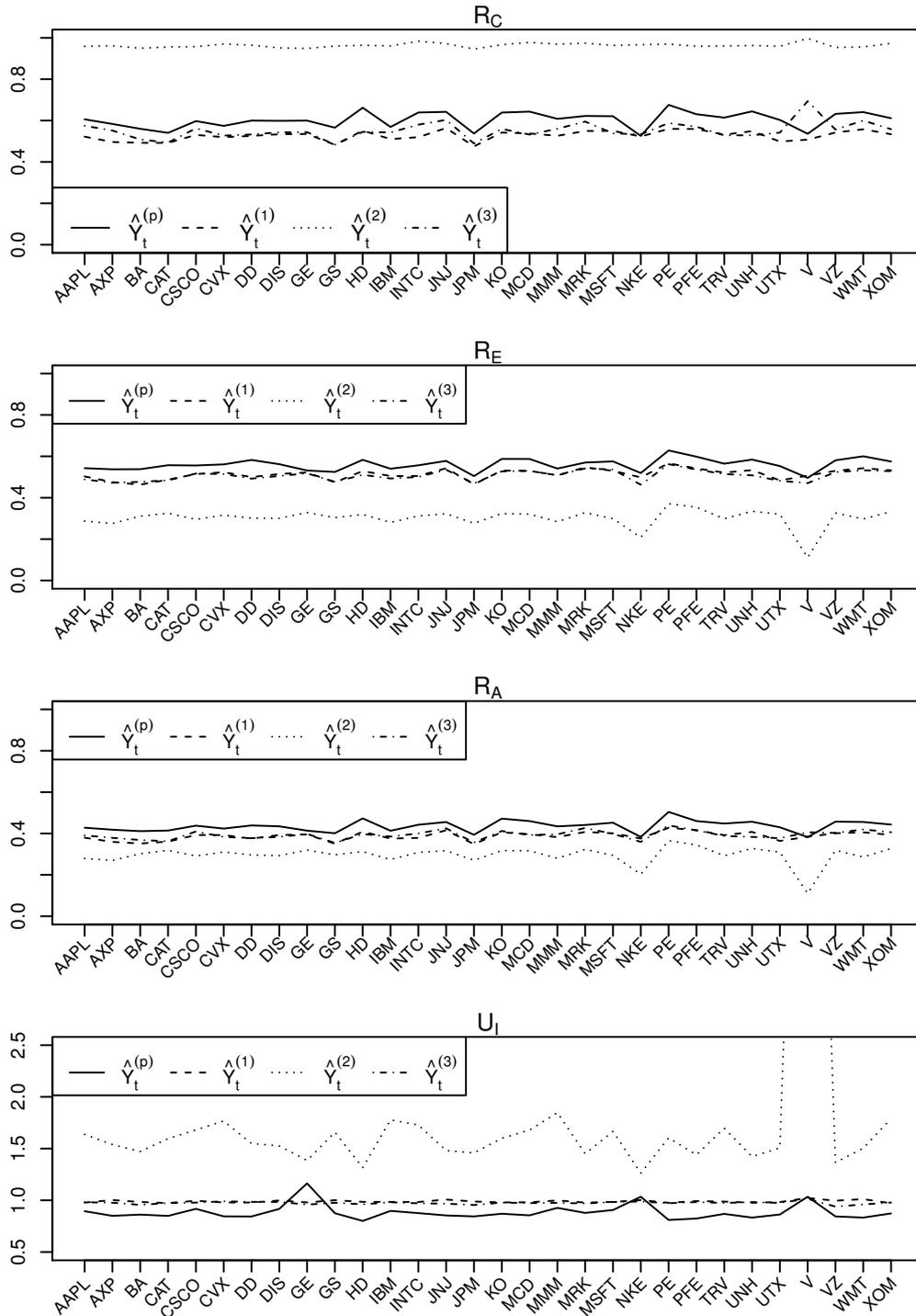


Figure 3: The average values of R_C , R_E , R_A and U_I of $\hat{Y}_t^{(p)}$ (solid line), $\hat{Y}_t^{(1)}$ (dashed line), $\hat{Y}_t^{(2)}$ (dotted line) and $\hat{Y}_t^{(3)}$ (dash-dotted line) for 30 different time series from June 29, 2015 to June 23, 2016.

The main reason for the good performance of $\hat{Y}_t^{(p)}$ is that $\hat{Y}_t^{(p)}$ uses more information than the other predictors. All the other predictors involve (traditional) ITS, that is, they are formed exclusively with the high and low returns and do not consider past closing prices. Therefore, the predictors $\hat{Y}_t^{(i)}$, $i = 1, 2, 3$, have a clear disadvantage relative to $\hat{Y}_t^{(p)}$; consequently, the latter should show much better performance.

5. EMPIRICAL STUDY

In this section, an empirical study is designed to investigate the performance of the proposed criterion for selecting the optimal portfolio using the stock prices of the companies of the DJIA Index. The DJIA Index was launched on October 29, 2002. This Index covers the top 30 companies by total market capitalization and is reviewed quarterly in January, April, July and October every year. Suppose that a self-financing trading strategy, which daily reallocates the holding weights of the portfolio, is employed from the beginning of each period. The proposed criterion is used to reallocate the optimal portfolios daily during the financial crisis period and under improved market conditions by fitting the time series models defined in (2.5) with the previous 250 daily historical high and low returns for each underlying asset. Then, the corresponding 250 one-day-ahead returns of the optimal portfolios are computed and compared with the DJIA Index. In the following, we illustrate the details of the construction of the self-financing trading strategy during the financial crisis period:

1. Let DJ_t be the value of the DJIA Index at time t , where $t = 0$ stands for the date of June 27, 2008.
2. Let V_t denote the value of the self-financing portfolio at time t . Further, let V_0 be the value of the DJIA Index on June 27, 2008. The initial allocations of the underlying assets, $c_{m,0}$, are obtained by solving (3.2), where the high and low return processes of each underlying asset are fitted by model (2.5) based on $X_{m,t}^{(CH)}$ and $X_{m,t}^{(CL)}$ for $t = -250, \dots, -1$, $m = 1, \dots, p$ and $p = 30$. Moreover, since $\sum_{m=1}^p c_{m,0}$ can be less than 1, the amount $V_0 \sum_{m=1}^p c_{m,0}$ is invested in risky assets and the rest of the portfolio value, denoted by $C_0 = V_0(1 - \sum_{m=1}^p c_{m,0})$, is invested in the risk-free market.
3. At time $t = 1$, the value of the portfolio is

$$V_{1-} = b^{(0)} \sum_{m=1}^p c_{m,0} P_{m,1}^C + e^{r_d} C_0,$$

prior to the adjustment of the holding portfolio, where

$$b^{(0)} = \frac{V_0 \sum_{m=1}^p c_{m,0}}{\sum_{m=1}^p c_{m,0} P_{m,0}^C}$$

and r_d is the daily risk-free interest rate. We reestimate the dynamic models of each return process using the data $P_{m,t}$, $t = -249, \dots, 0$, and compute the updated optimal allocations, which are proportional to $c_{m,1}$ obtained by solving (3.2), where the value of the updated portfolio, denoted by V_1 , is the same as V_{1-} to satisfy self-financing. That is,

$$V_1 = V_{1-} = b^{(1)} \sum_{m=1}^p c_{m,1} P_{m,1}^C + C_1,$$

where

$$b^{(1)} = \frac{V_{1-} \sum_{m=1}^p c_{m,1}}{\sum_{m=1}^p c_{m,1} P_{m,1}^C}$$

and $C_1 = V_{1-}(1 - \sum_{m=1}^p c_{m,1})$ denotes the amount invested in the risk-free market after reallocation.

4. Repeat Step 3 until June 24, 2009.

In addition to adjusting the allocations of the above self-financing trading strategy daily, we proposed dynamic adjustment of the risk limitation L in (3.2) by considering

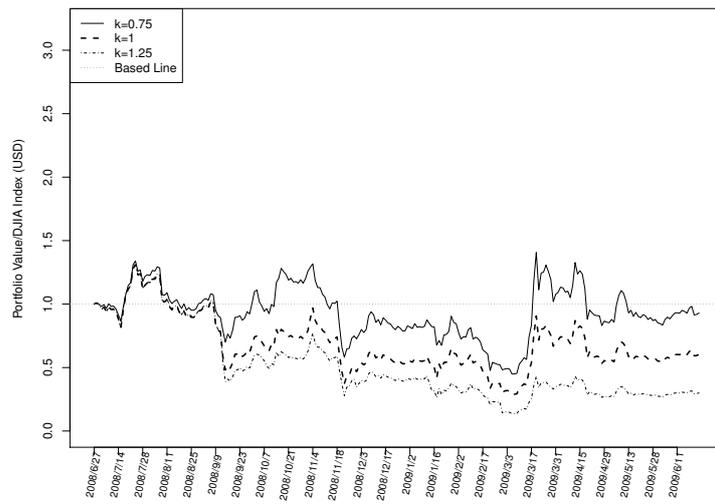
$$(5.1) \quad L = \frac{k}{p} \sum_{m=1}^p \rho_I(\mathbf{X}_{m,t+1}^{(CI)} | \mathcal{F}_t)$$

at time t , where k is a positive constant and $\rho_I(\mathbf{X}_{m,t+1}^{(CI)} | \mathcal{F}_t)$ is defined in (3.4). The L defined in (5.1) is a special case of (3.3) with $c_{m,t} = 1/p$, for $m = 1, \dots, p$, multiplied by k . In other words, we set the limitation of the investment risk in (3.2) by considering the trading strategy of an equally weighted portfolio. Moreover, conditional on \mathcal{F}_t and by (2.3)–(2.5), $X_{m,t+1}^{(CL)} = \mu_m^{(CL)} + \sigma_{m,t+1}^{(CL)} \tilde{X}_{m,t+1}^{(CL)}$ is normally distributed with conditional mean $\hat{X}_{m,t}^{(CL)}(1)$ defined in (2.6) and conditional standard deviation $\sigma_{m,t+1}^{(CL)}$. Consequently, (3.4) yields

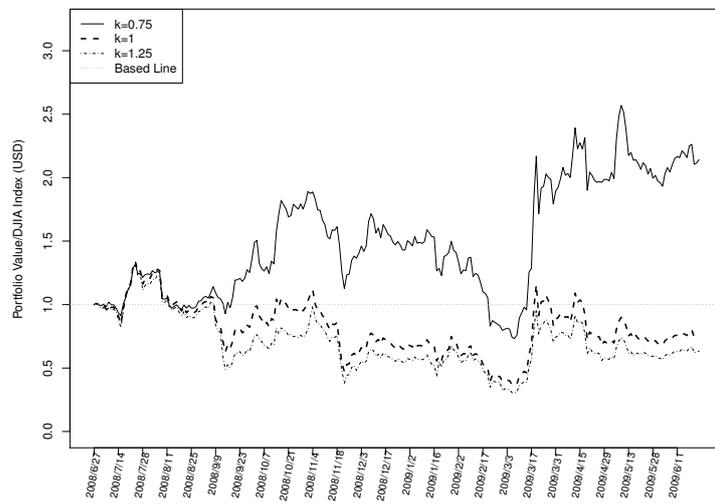
$$\begin{aligned} \rho_I(\mathbf{X}_{m,t+1}^{(CI)} | \mathcal{F}_t) &= -E_t \left(X_{m,t+1}^{(CL)} \mid X_{m,t+1}^{(CL)} \leq q_{\alpha,m,t+1} \right) \\ &= -\hat{X}_{m,t}^{(CL)}(1) + \sigma_{m,t+1}^{(CL)} \phi \left((q_{\alpha,m,t+1} - \hat{X}_{m,t}^{(CL)}(1)) / \sigma_{m,t+1}^{(CL)} \right) / \alpha, \end{aligned}$$

where $\phi(\cdot)$ is the density function of the standard normal distribution.

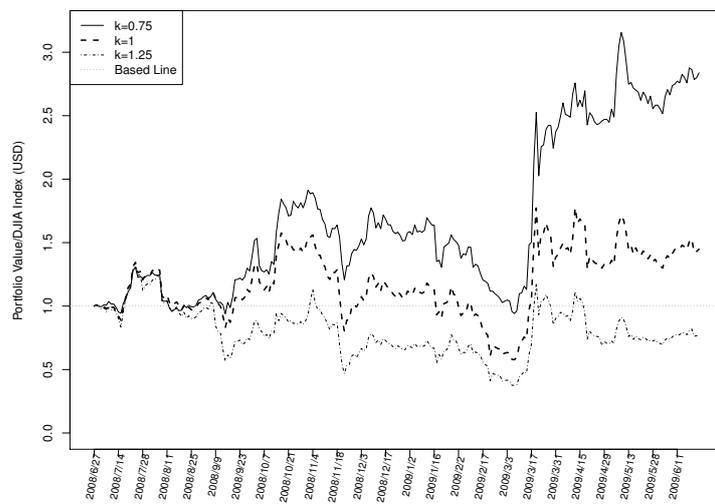
The numerical results are presented in Figures 4 and 5 with $\alpha = 0.05, 0.20$ and 0.35 and $r_d = 0$. Figure 6 presents the values of L in 2008–2009 and 2015–2016 with different settings of α and $k = 1$. Figure 6 shows that a portfolio constructed by (3.2) with a large α is more conservative than a portfolio constructed with a small α since the values of L with $\alpha = 0.35$ are smaller than their counterparts. In Figure 4, the solid, dashed and dash-dotted lines denote the ratios of the capitals of the proposed trading strategy with $k = 0.75, 1$ and 1.25 , respectively, to the DJIA Index in 2008–2009, and the results for 2015–2016 are presented in Figure 5. For a fixed α , a portfolio with a small k is more conservative than one with a large k . In Figure 4, the proposed portfolio selection criterion (3.2) with L defined in (5.1) suggests a conservative portfolio during the financial crisis in 2008–2009 since the case with $k = 0.75$ performs better than the others for each α . In particular, the portfolio with $(\alpha, k) = (0.35, 0.75)$ has the best performance among all scenarios. For 2015–2016, compared with the portfolios selected in 2008–2009, the results presented in Figure 5 indicate that (3.2) suggests aggressive portfolios, decreasing α from 0.35 to 0.05 or 0.20 with $k = 0.75$ or increasing k from 0.75 to 1.00 with $\alpha = 0.35$. In view of the results in Figures 4 and 5, the proposed portfolio selection criterion (3.2) is capable of adjusting its suggestions according to the economic conditions.



(a) $\alpha = 0.05$

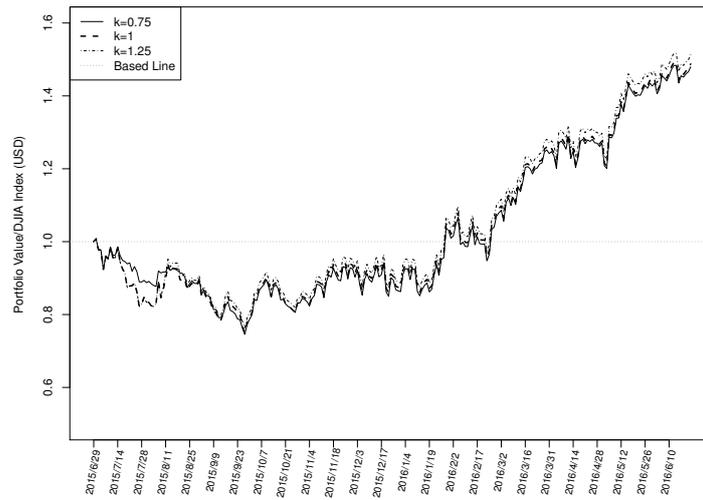


(b) $\alpha = 0.20$

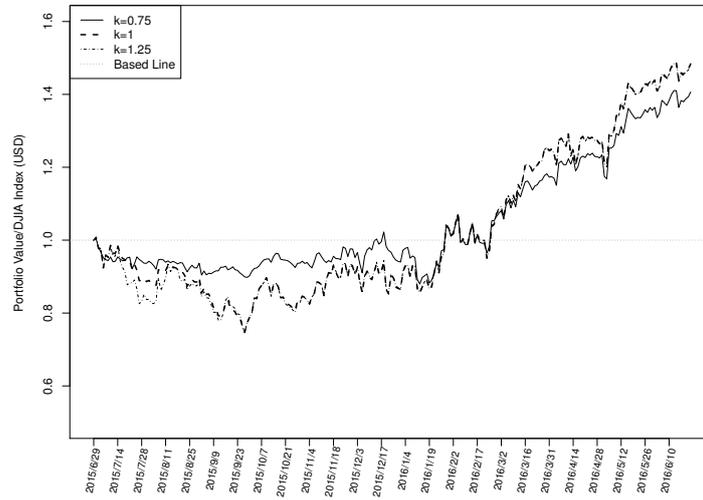


(c) $\alpha = 0.35$

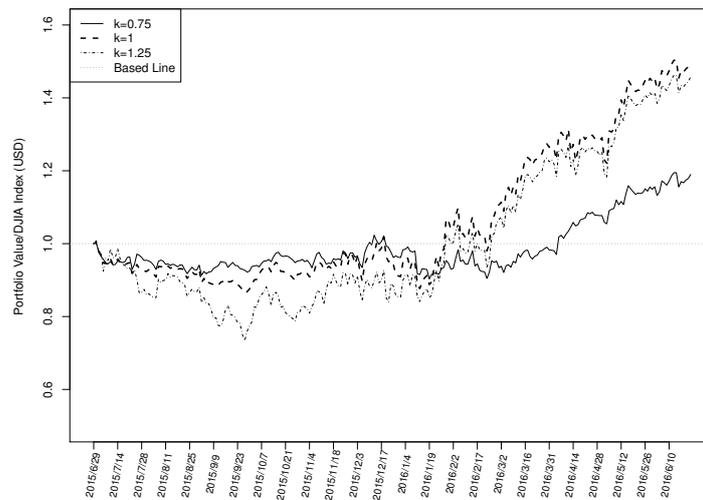
Figure 4: The ratios of the capitals of different trading strategies to the Dow Jones Industrial Average Index in 2008–2009.



(a) $\alpha = 0.05$

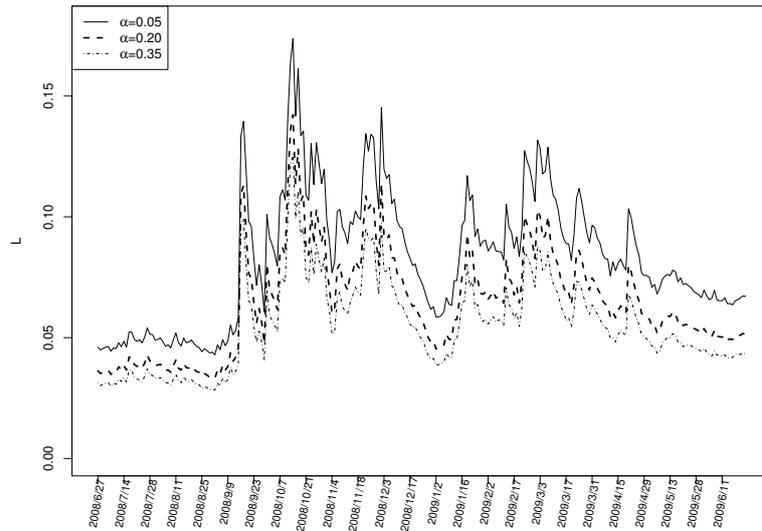


(b) $\alpha = 0.20$

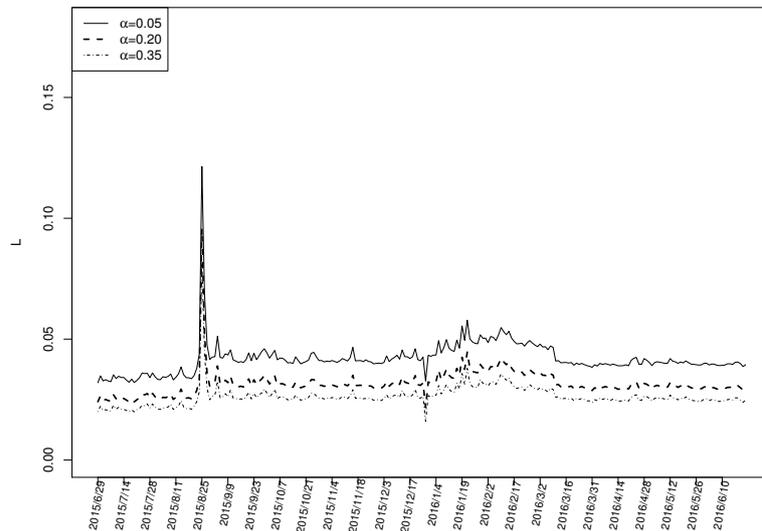


(c) $\alpha = 0.35$

Figure 5: The ratios of the capitals of different trading strategies to the Dow Jones Industrial Average Index in 2015–2016.



(a) L in 2008–2009



(b) L in 2015–2016

Figure 6: The values of L in 2008–2009 and 2015–2016.

6. CONCLUSION

In this study, we propose a prediction interval for future stock prices by fitting time series models to the high and low return processes. The proposed interval estimator is shown to have promising coverage, efficiency and accuracy. In particular, the numerical results of the U_I index indicate that the proposed interval estimator reduces the prediction error of the naïve interval predictor more remarkably than three popular interval estimators discussed in the literature. Consequently, an innovative criterion for portfolio selection is proposed on the basis of our interval estimator. The allocations of the underlying assets in the proposed optimal criterion are determined by maximizing the potential high portfolio returns subject to a predetermined limitation on the corresponding potential low and nonpositive returns.

An empirical study is conducted to investigate the investment returns of the proposed optimal portfolio. A dynamic self-financing trading strategy is established by investing in the stocks of the 30 companies of the DJIA Index and adjusting the asset allocations by the proposed method daily during the financial crisis period and a period with improved market conditions. The numerical results indicate that the proposed portfolio selection criterion constructed from the prediction intervals is capable of suggesting an optimal portfolio according to the economic conditions.

This study demonstrates that ITS data, including daily closing, high, and low prices, are capable of improving the performance of investment decisions and risk management by means of the proposed scheme. Additionally, better prediction performance is expected if intra-daily ITS data are available. This is an interesting direction for future studies.

ACKNOWLEDGMENTS

Huang's research was supported by the grant MOST 104-2118-M-390-003-MY2 from the Ministry of Science and Technology of Taiwan. Hsu's research was supported by the Ministry of Science and Technology of Taiwan, Grant No. MOST 105-2118-M-390-004.

REFERENCES

- [1] ACERBI, C. (2002). Spectral measures of risk: A coherent representation of subjective risk aversion, *Journal of Banking & Finance*, **26**, 1505–1518.
- [2] ADAM, A.; HOUKARI, M. and LAURENT, J.-P. (2008). Spectral risk measures and portfolio selection, *Journal of Banking & Finance*, **32**, 1870–1882.
- [3] ARROYO, J.; ESPÍNOLA, R. and MATÉ, C. (2011). Different approaches to forecast interval time series: A comparison in finance, *Computational Economics*, **37**, 169–191.
- [4] ARROYO, J.; GONZÁLEZ-RIVERA, G. and MATÉ, C. (2010). *Forecasting with interval and histogram data: Some financial applications*. In “Handbook of empirical economics and finance” (A. Ullah, D. Giles, N. Balakrishnan, W. Schucany and E. Schilling, Eds.), Chapman and Hall/CRC, New York, 247–280.
- [5] ARROYO, J. and MATÉ, C. (2006). *Introducing interval time series: Accuracy measures*. In “COMPSTAT 2006 – Proceedings in Computational Statistic” (A. Rizzi and M. Vichi, Eds.), Springer, Heidelberg, 1139–1146.
- [6] ARTZNER, P.; DELBAEN, F.; EBER, J.-M. and HEATH, D. (1999). Coherent measures of risk, *Mathematical Finance*, **9**, 203–228.
- [7] BLANCO-FERNÁNDEZ, A.; CORRAL, N. and GONZÁLEZ-RODRÍGUEZ, G. (2011). Estimation of a flexible simple linear model for interval data based on set arithmetic, *Computational Statistics & Data Analysis*, **55**, 2568–2578.
- [8] BOLLERSLEV, T. (1986). Generalized autoregressive conditional heteroskedasticity, *Journal of Econometrics*, **31**, 307–327.
- [9] ENGLE, R. (1982). Autoregressive conditional heteroscedasticity with estimates of the variance of United Kingdom inflation, *Econometrica*, **50**, 987–1007.

- [10] ENGLE, R. (2002). Dynamic conditional correlation – a simple class of multivariate GARCH models, *Journal of Business & Economic Statistics*, **20**, 339–350.
- [11] ENGLE, R. (2009). *Anticipating correlations – A New Paradigm for Risk Management*, Princeton University Press, New Jersey.
- [12] FISCHER, H.; BLANCO-FERNÁNDEZ, Á. and WINKER, P. (2016). Predicting stock return volatility: Can we benefit from regression models for return intervals? *Journal of Forecasting*, **35**, 113–146.
- [13] GARCÍA-ASCANIO, C. and MATÉ, C. (2010). Electric power demand forecasting using interval time series: A comparison between var and iMLP, *Energy Policy*, **38**, 715–725.
- [14] GRIGORYEVA, L.; ORTEGA, J.-P. and PERESETSKY, A. (2017). Volatility forecasting using global stochastic financial trends extracted from non-synchronous data, *Econometrics and Statistics*, in Press.
- [15] HÄRDLE, W.; OKHRIN, O. and WANG, W. (2015). Hidden Markov structures for dynamic copulae, *Econometric Theory*, **31**, 981–1015.
- [16] HARRIS, R.D.F. and MAZIBAS, M. (2013). Dynamic hedge fund portfolio construction: A semi-parametric approach, *Journal of Banking & Finance*, **37**, 139–149.
- [17] HE, L.T. and HU, C. (2009). Impacts of interval computing on stock market variability forecasting, *Computational Economics*, **33**, 263–276.
- [18] HUANG, S.F.; LIN, C.H. and LIN, T.Y. (2017). *Portfolio selection with spectral risk measures*. In “Applied Quantitative Finance”, 3rd Edition (W. Härdle, C. Chen and L. Overbeck, Eds.), Springer, GmbH Germany, 39–56.
- [19] ICHINO, M. and YAGUCHI, H. (1994). Generalized Minkowski metrics for mixed feature-type data analysis, *IEEE Trans. on Systems, Man and Cybernetics*, **24**, 698–708.
- [20] KROKHMAL, P.; PALMQUIST, J. and URYASEV, S. (2002). Portfolio optimization with conditional value-at-risk objective and constraints, *Journal of Risk*, **4**, 43–68.
- [21] MAIA, A.L.S.; DE CARVALHO, F.A.T. and LUDERMIR, T.B. (2008). Forecasting models for interval-valued time series, *Neurocomputing*, **71**, 3344–3352.
- [22] MARKOWITZ, H. (1952). Portfolio selection, *The Journal of Finance*, **7**, 77–91.
- [23] MARKOWITZ, H. (1959). *Portfolio Selection: Efficient Diversification of Investments*, Wiley, New York.
- [24] NELSON, D.B. (1990). Stationarity and persistence in the garch(1,1) model, *Econometric Theory*, **6**, 318–334.
- [25] ROCKAFELLAR, R.T. and URYASEV, S. (2000). Optimization of conditional value-at-risk, *Journal of Risk*, **2**, 21–41.
- [26] ROCKAFELLAR, R.T. and URYASEV, S. (2002). Conditional value-at-risk for general loss distributions, *Journal of Banking & Finance*, **26**, 1443–1471.
- [27] RODRIGUES, P.M.M. and SALISH, N. (2015). Modeling and forecasting interval time series with threshold models, *Advances in Data Analysis and Classification*, **9**, 41–57.
- [28] TELES, P. and BRITO, P. (2005). *Modelling interval time series data*. In “Proceedings of the 3rd IASC World Conference on Computational Statistics and Data Analysis”, Limassol, Cyprus.
- [29] TELES, P. and BRITO, P. (2015). Modeling interval time series with space-time processes, *Communications in Statistics – Theory and Methods*, **44**, 3599–3627.
- [30] TSAY, R.S. (2010). *Analysis of Financial Time Series* (3rd ed.), Wiley, New Jersey.
- [31] XIONG, T.; LI, C.; BAO, Y.; HU, Z. and ZHANG, L. (2015). A combination method for interval forecasting of agricultural commodity futures prices, *Knowledge-Based Systems*, **77**, 92–102.
- [32] YANG, W.; HAN, A.; CAI, K. and WANG, S. (2012). Acix model with interval dummy variables and its application in forecasting interval-valued crude oil prices, *Procedia Computer Science*, **9**, 1273–1282.

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.

Editorial Board

Editor-in-Chief

Isabel Fraga Alves, University of Lisbon, Portugal

Co-Editor

Giovani L. Silva, University of Lisbon, Portugal

Associate Editors

Marília Antunes, University of Lisbon, Portugal

Barry Arnold, University of California, USA

Narayanaswamy Balakrishnan, McMaster University, Canada

Jan Beirlant, Katholieke Universiteit Leuven, Belgium

Graciela Boente (2019-2020), University of Buenos Aires, Argentina

Paula Brito, University of Porto, Portugal

Vanda Inácio de Carvalho, University of Edinburgh, UK

Arthur Charpentier, Université du Québec à Montréal, Canada

Valérie Chavez-Demoulin, University of Lausanne, Switzerland

David Conesa, University of Valencia, Spain

Charmaine Dean, University of Waterloo, Canada

Jorge Milhazes Freitas, University of Porto, Portugal

Alan Gelfand, Duke University, USA

Stéphane Girard, Inria Grenoble Rhône-Alpes, France

Wenceslao Gonzalez-Manteiga, University of Santiago de Compostela, Spain

Marie Kratz, ESSEC Business School, France

Victor Leiva, Pontificia Universidad Católica de Valparaíso, Chile

Maria Nazaré Mendes-Lopes, University of Coimbra, Portugal

Fernando Moura, Federal University of Rio de Janeiro, Brazil

John Nolan, American University, USA

Paulo Eduardo Oliveira, University of Coimbra, Portugal

Pedro Oliveira, University of Porto, Portugal

Carlos Daniel Paulino (2019-2021), University of Lisbon, Portugal

Arthur Pewsey, University of Extremadura, Spain

Gilbert Saporta, Conservatoire National des Arts et Métiers, France

Alexandra M. Schmidt, McGill University, Canada

Julio Singer, University of Sao Paulo, Brazil

Manuel Scotto, University of Lisbon, Portugal

Lisete Sousa, University of Lisbon, Portugal

Milan Stehlík, University of Valparaíso, Chile and LIT-JK University Linz, Austria

María Dolores Ugarte, Public University of Navarre, Spain

Executive Editor

José A. Pinto Martins, Statistics Portugal

Secretariat

José Cordeiro, Statistics Portugal

Olga Bessa Mendes, Statistics Portugal