



Estimation of a Likelihood Ratio Ordered Family of Distributions

Alexandre Mösching ^{*1,2} and Lutz Dümbgen ^{†1}

¹University of Bern, Department of Mathematics and Statistics, Bern, Switzerland

²F. Hoffmann-La Roche Ltd, Nonclinical Biostatistics, Basel, Switzerland

December 21, 2023

Abstract

Consider bivariate observations $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathbb{R} \times \mathbb{R}$ with unknown conditional distributions Q_x of Y , given that $X = x$. The goal is to estimate these distributions under the sole assumption that Q_x is isotonic in x with respect to likelihood ratio order. If the observations are identically distributed, a related goal is to estimate the joint distribution $\mathcal{L}(X, Y)$ under the sole assumption that it is totally positive of order two. An algorithm is developed which estimates the unknown family of distributions $(Q_x)_x$ via empirical likelihood. The benefit of the stronger regularization imposed by likelihood ratio order over the usual stochastic order is evaluated in terms of estimation and predictive performances on simulated as well as real data.

Keywords: Empirical likelihood, likelihood ratio order, order constraint, quasi-Newton method, stochastic order, total positivity.

AMS 2000 subject classifications: 62G05, 62G08, 62H12.

Acknowledgements: The authors are grateful to Johanna Ziegel, Alexander Jordan and Tilmann Gneiting for stimulating discussions and useful hints. We also thank a reviewer for constructive comments. This work was supported by Swiss National Science Foundation.

1 Introduction

Consider a univariate regression setting with observations $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ in $\mathfrak{X} \times \mathbb{R}$, where \mathfrak{X} is an arbitrary real set. We assume that conditional on $\mathbf{X} := (X_i)_{i=1}^n$, the observations Y_1, Y_2, \dots, Y_n are independent with distributions $\mathcal{L}(Y_i | \mathbf{X}) = Q_{X_i}$, where the distributions $Q_x, x \in \mathfrak{X}$, are unknown. The goal is to estimate the latter under the sole assumption that Q_x is isotonic in x in a certain sense. That means, if (X, Y) denotes a generic observation, the larger (or smaller) the value of X , the larger (or smaller) Y tends to be. An obvious notion of order would be the usual stochastic order, which states that $Q_{x_1} \leq_{\text{st}} Q_{x_2}$ whenever $x_1 \leq x_2$, that is, $Q_{x_1}((-\infty, y]) \geq Q_{x_2}((-\infty, y])$ for all $y \in \mathbb{R}$. This concept has been investigated and generalized by numerous authors, see [Mösching and](#)

*alexandre.moesching@roche.com

†lutz.duembgen@unibe.ch

Dümbgen (2020), Henzi et al. (2021b) and the references cited therein. The latter paper illustrates the application of isotonic distributional regression in weather forecasting, and Henzi et al. (2021a) use it to analyze the length of stay of patients in Swiss hospitals.

The present paper investigates a stronger notion of order, the so-called likelihood ratio order. The usual definition is that for arbitrary points $x_1 < x_2$ in \mathfrak{X} , the distributions Q_{x_1} and Q_{x_2} have densities g_{x_1} and g_{x_2} with respect to some dominating measure such that g_{x_2}/g_{x_1} is isotonic on the set $\{g_{x_1} + g_{x_2} > 0\}$, and this condition will be denoted by $Q_{x_1} \leq_{\text{lr}} Q_{x_2}$. At first glance, this looks like a rather strong assumption coming out of thin air, but it is familiar from mathematical statistics or discriminant analyses and has interesting properties. For instance, $Q_{x_1} \leq_{\text{lr}} Q_{x_2}$ if and only if $Q_{x_1}(\cdot | B) \leq_{\text{st}} Q_{x_2}(\cdot | B)$ for any real interval B such that $Q_{x_1}(B), Q_{x_2}(B) > 0$, where $Q_{x_j}(A | B) := Q_{x_j}(A \cap B)/Q_{x_j}(B)$. Furthermore, likelihood ratio ordering is a frequent assumption or implication of models in mathematical finance, see Beare and Moon (2015), Jewitt (1991). The notion of likelihood ratio order is reviewed thoroughly in Dümbgen and Mösching (2023), showing that it defines a partial order on the set of all probability measures on the real line which is preserved under weak convergence. That material generalizes definitions and results in Shaked and Shanthikumar (2007).

Thus far, estimation of distributions under a likelihood ratio order constraint was mainly limited to settings with two or finitely many samples and populations. First, Dykstra et al. (1995) estimated the parameters of two multinomial distributions that are likelihood ratio ordered via a restricted maximum likelihood approach. After reparametrization, they found that the maximization problem at hand had reduced to a specific bioassay problem treated by Robertson et al. (1988) and which makes use of the theory of isotonic regression. It is then suggested that their approach generalizes well to any two distributions that are absolutely continuous with respect to some dominating measure. Later, Carolan and Tebbs (2005) focused on testing procedures for the equality of two distributions Q_1 and Q_2 versus the alternative hypothesis that $Q_1 \leq_{\text{lr}} Q_2$, in the specific case where the cumulative distribution functions G_i of Q_i , $i = 1, 2$, are continuous. To this end, they made use of the equivalence between likelihood ratio order and the convexity of the ordinal dominance curve $\alpha \mapsto G_2(G_1^{-1}(\alpha))$, $\alpha \in [0, 1]$, which holds in case of G_2 being absolutely continuous with respect to G_1 . The convexity of the ordinal dominance curve was also exploited by Westling et al. (2023) to provide nonparametric maximum likelihood estimators of G_1 and G_2 under likelihood ratio order for discrete, continuous, as well as mixed continuous-discrete distributions using the greatest convex minorant of the empirical ordinal dominance curve. However, this method still necessitates the restrictive assumption that G_2 is absolutely continuous with respect to G_1 . Other attempts at estimating two likelihood ratio ordered distributions include Yu et al. (2017) who treat the estimation problem with a maximum smoothed likelihood approach, requiring the choice of a kernel and bandwidth parameters, and Hu et al. (2023) who suppose absolutely continuous distributions and model the logarithm of the ratio of densities as a linear combination of Bernstein polynomials.

To the best of our knowledge, only Dardanoni and Forcina (1998) considered the problem of estimating an arbitrary fixed number $\ell \geq 2$ of likelihood ratio ordered distributions Q_1, Q_2, \dots, Q_ℓ , all of them sharing the same finite support. They showed that the constrained maximum likelihood problem may be reparametrized to obtain a convex optimization problem with linear inequality constraints, and they propose to solve the latter via a constrained version of the Fisher scoring algorithm. At each step of their procedure, it is necessary to solve a quadratic programming problem.

Within the setting of distributional regression, we follow an empirical likelihood approach (Owen, 1988, 2001) to estimate the family $(Q_x)_{x \in \mathfrak{X}}$ for arbitrary real sets \mathfrak{X} . After a

reparametrization similar to that of [Dardanoni and Forcina \(1998\)](#), we show that the problem of maximizing the (empirical) likelihood under the likelihood ratio order constraint yields again a finite-dimensional convex optimization problem with linear inequality constraints. We did experiments with active set algorithms in the spirit of [Dümbgen et al. \(2021\)](#) which are similar to the algorithms of [Dardanoni and Forcina \(1998\)](#). But, as explained later, the computational burden may become too heavy for large sample sizes n . Alternatively, we devise an algorithm which adapts and extends ideas from [Jongbloed \(1998\)](#) and [Dümbgen et al. \(2006\)](#) for the present setting. It makes use of a quasi-Newton approach, and new search directions are obtained via multiple isotonic weighted least squares regression.

There is an interesting aspect of the present estimation problem. If we assume that the observations (X_i, Y_i) are independent copies of a generic random pair (X, Y) , the new estimation method may also be interpreted as an empirical likelihood estimator of the joint distribution of (X, Y) , hypothesizing that the latter is bivariate totally positive of order two (TP2). That is, for arbitrary intervals A_1, A_2 and B_1, B_2 such that $A_1 < A_2$ and $B_1 < B_2$ element-wise,

$$\mathbb{P}(X \in A_2, Y \in B_1) \mathbb{P}(X \in A_1, Y \in B_2) \leq \mathbb{P}(X \in A_1, Y \in B_1) \mathbb{P}(X \in A_2, Y \in B_2).$$

If the joint distribution of (X, Y) has a density h with respect to Lebesgue measure on $\mathbb{R} \times \mathbb{R}$, or if it is discrete with probability mass function h , then TP2 is equivalent to requiring that

$$h(x_1, y_2)h(y_1, x_2) \leq h(x_1, y_1)h(x_2, y_2) \quad \text{whenever } x_1 < x_2, y_1 < y_2,$$

and this is just a special case of multivariate total positivity of order two ([Karlin, 1968](#)). For further equivalences and results in dimension two, see [Dümbgen and Mösching \(2023\)](#). Interestingly, this TP2 constraint is symmetric in X and Y , and our algorithm exploits this symmetry. A different, more restrictive approach to the estimation of a TP2 distribution is proposed by [Hütter et al. \(2020\)](#). They assume that the distribution of (X, Y) has a smooth density with respect to Lebesgue measure on a given rectangle and devise a sieve maximum likelihood estimator.

The rest of the article is structured as follows. Section 2 explains why empirical likelihood estimation of a family of likelihood ratio ordered distributions is essentially equivalent to the estimation of a discrete bivariate TP2 distribution. In Section 3 we present an algorithm to estimate a bivariate TP2 distribution. In Section 4, a simulation study illustrates the benefits of the new estimation paradigm compared to the usual stochastic order constraint. Proofs and technical details are deferred to the appendix.

2 Two versions of empirical likelihood modelling

With our observations $(X_i, Y_i) \in \mathfrak{X} \times \mathbb{R}$, $1 \leq i \leq n$, let

$$\{X_1, X_2, \dots, X_n\} = \{x_1, \dots, x_\ell\} \quad \text{and} \quad \{Y_1, Y_2, \dots, Y_n\} = \{y_1, \dots, y_m\},$$

with $x_1 < \dots < x_\ell$ and $y_1 < \dots < y_m$. For an index pair (j, k) with $1 \leq j \leq \ell$ and $1 \leq k \leq m$, let

$$w_{jk} := \#\{i : (X_i, Y_i) = (x_j, y_k)\}.$$

That means, the empirical distribution \widehat{R}_{emp} of the observations (X_i, Y_i) can be written as $\widehat{R}_{\text{emp}} = n^{-1} \sum_{j=1}^{\ell} \sum_{k=1}^m w_{jk} \delta_{(x_j, y_k)}$.

2.1 Estimating the conditional distributions Q_x

To estimate $(Q_x)_{x \in \mathfrak{X}}$ under likelihood ratio ordering, we first estimate $(Q_{x_j})_{1 \leq j \leq \ell}$. If that results in $(\widehat{Q}_{x_j})_{1 \leq j \leq \ell}$, we may define

$$\widehat{Q}_x := \begin{cases} \widehat{Q}_{x_1} & \text{if } x < x_1, \\ (1 - \lambda)\widehat{Q}_{x_j} + \lambda\widehat{Q}_{x_{j+1}} & \text{if } x = (1 - \lambda)x_j + \lambda x_{j+1}, \ 1 \leq j < \ell, \ 0 < \lambda < 1, \\ \widehat{Q}_{x_\ell} & \text{if } x > x_\ell. \end{cases}$$

This piecewise linear extension preserves isotonicity with respect to \leq_{lr} , see Lemma A.1.

To estimate $Q_{x_1}, \dots, Q_{x_\ell}$, we restrict our attention to distributions with support $\{y_1, \dots, y_m\}$. That means, we assume temporarily that for $1 \leq j \leq \ell$,

$$Q_{x_j} = \sum_{k=1}^m q_{jk} \delta_{y_k}$$

with weights $q_{j1}, \dots, q_{jm} \geq 0$ summing to one. The empirical log-likelihood for the corresponding matrix $\mathbf{q} = (q_{jk})_{j,k} \in [0, 1]^{\ell \times m}$ equals

$$L_{\text{raw}}(\mathbf{q}) := \sum_{j=1}^{\ell} \sum_{k=1}^m w_{jk} \log q_{jk}. \quad (2.1)$$

Then the goal is to maximize this log-likelihood over all matrices $\mathbf{q} \in [0, 1]^{\ell \times m}$ such that

$$\sum_{k=1}^m q_{jk} = 1 \quad \text{for } 1 \leq j \leq \ell, \quad (2.2)$$

$$q_{j_1 k_2} q_{j_2 k_1} \leq q_{j_1 k_1} q_{j_2 k_2} \quad \text{for } 1 \leq j_1 < j_2 \leq \ell \text{ and } 1 \leq k_1 < k_2 \leq m. \quad (2.3)$$

The latter constraint is equivalent to saying that Q_{x_j} is isotonic in $j \in \{1, \dots, \ell\}$ with respect to \leq_{lr} .

2.2 Estimating the distribution of (X, Y)

Suppose that the observations (X_i, Y_i) are independent copies of a random pair (X, Y) with unknown TP2 distribution R on $\mathbb{R} \times \mathbb{R}$. An empirical likelihood approach to estimating R is to restrict one's attention to distributions

$$R = \sum_{j=1}^{\ell} \sum_{k=1}^m h_{jk} \delta_{(x_j, y_k)}$$

with ℓm weights $h_{jk} \geq 0$ summing to one. The empirical log-likelihood of the corresponding matrix $\mathbf{h} = (h_{jk})_{j,k}$ equals $L_{\text{raw}}(\mathbf{h})$ with the function L_{raw} defined in (2.1). But now the goal is to maximize $L_{\text{raw}}(\mathbf{h})$ over all matrices $\mathbf{h} \in [0, 1]^{\ell \times m}$ satisfying the constraints

$$\sum_{j=1}^{\ell} \sum_{k=1}^m h_{jk} = 1 \quad (2.4)$$

and (2.3). As mentioned in the introduction, requirement (2.3) for \mathbf{h} is equivalent to R being TP2. One can get rid of the constraint (2.4) via a Lagrange trick and maximize

$$L(\mathbf{h}) := L_{\text{raw}}(\mathbf{h}) - nh_{++} + n$$

over all \mathbf{h} satisfying (2.3), where $h_{++} := \sum_j \sum_k h_{jk}$. Indeed, if \mathbf{h} is a matrix in $[0, \infty)^{\ell \times m}$ such that $L_{\text{raw}}(\mathbf{h}) > -\infty$, then $\tilde{\mathbf{h}} := (h_{jk}/h_{++})_{j,k}$ satisfies (2.3) if and only if \mathbf{h} does, and

$$L(\mathbf{h}) = L_{\text{raw}}(\tilde{\mathbf{h}}) + n(\log h_{++} - h_{++} + 1) \leq L_{\text{raw}}(\tilde{\mathbf{h}}) = L(\tilde{\mathbf{h}})$$

with equality if and only if $h_{++} = 1$, that is, $\mathbf{h} = \tilde{\mathbf{h}}$.

2.3 Equivalence of the two estimation problems

For any matrix $\mathbf{a} \in \mathbb{R}^{\ell \times m}$ define the row sums $a_{j+} := \sum_k a_{jk}$ and column sums $a_{+k} := \sum_j a_{jk}$. If \mathbf{h} is an arbitrary matrix in $[0, \infty)^{\ell \times m}$ such that $L_{\text{raw}}(\mathbf{h}) > -\infty$, and if we write

$$h_{jk} = p_j q_{jk} \quad \text{with } p_j := h_{j+} \text{ and } q_{jk} := h_{jk}/h_{j+},$$

then \mathbf{h} satisfies (2.3) if and only if \mathbf{q} does. Furthermore, \mathbf{q} satisfies (2.2), and elementary algebra shows that

$$L(\mathbf{h}) = L_{\text{raw}}(\mathbf{q}) + \sum_{j=1}^{\ell} (w_{j+} \log p_j - np_j + w_{j+}).$$

The unique maximizer $\mathbf{p} = (p_j)_j$ of $\sum_j (w_{j+} \log p_j - np_j + w_{j+})$ is the vector $(w_{j+}/n)_j$, and this implies the following facts:

- If $\hat{\mathbf{h}}$ is a maximizer of $L(\mathbf{h})$ under the constraints (2.3), then $\hat{h}_{j+} = w_{j+}/n$ for all j , and $\hat{q}_{jk} := \hat{h}_{jk}/\hat{h}_{j+}$ defines a maximizer $\hat{\mathbf{q}}$ of $L_{\text{raw}}(\mathbf{q})$ under the constraints (2.2) and (2.3).
- If $\hat{\mathbf{q}}$ is a maximizer of $L_{\text{raw}}(\mathbf{q})$ under the constraints (2.2) and (2.3), then $\hat{h}_{jk} := (w_{j+}/n)\hat{q}_{jk}$ defines a maximizer $\hat{\mathbf{h}}$ of $L(\mathbf{h})$ under the constraints (2.3).

As a final remark, note that the two estimation problems are monotone equivariant in the following sense: If (X, Y) is replaced with $(\tilde{X}, \tilde{Y}) = (\sigma(X), \tau(Y))$ with strictly isotonic functions $\sigma : \mathfrak{X} \rightarrow \mathbb{R}$ and $\tau : \mathbb{R} \rightarrow \mathbb{R}$, then $\mathcal{L}(\tilde{Y}|\tilde{X} = \sigma(x)) = \mathcal{L}(\tau(Y)|X = x)$ for $x \in \mathfrak{X}$. Furthermore, the constraints of likelihood ratio ordered conditional distributions or of a TP2 joint distribution remain valid under such transformations.

2.4 Calibration of rows and columns

The previous considerations motivate to find a maximizer $\hat{\mathbf{h}} \in [0, \infty)^{\ell \times m}$ of $L(\mathbf{h})$ under the constraint (2.3), even if the ultimate goal is to estimate the conditional distributions Q_x , $x \in \mathfrak{X}$. They also indicate two simple ways to improve a current candidate \mathbf{h} for $\hat{\mathbf{h}}$. Let $\tilde{\mathbf{h}}$ be defined via

$$\tilde{h}_{jk} := (w_{j+}/n)h_{jk}/h_{j+},$$

i.e. we rescale the rows of \mathbf{h} such that the new row sums \tilde{h}_{j+} coincide with the empirical weights w_{j+}/n . Then

$$L(\tilde{\mathbf{h}}) - L(\mathbf{h}) = \sum_{j=1}^{\ell} \left(w_{j+} \log \left(\frac{w_{j+}}{nh_{j+}} \right) + nh_{j+} - w_{j+} \right) \geq 0$$

with equality if and only if $\tilde{\mathbf{h}} = \mathbf{h}$. Similarly, one can improve \mathbf{h} by rescaling its columns, i.e. replacing \mathbf{h} with $\tilde{\mathbf{h}}$, where

$$\tilde{h}_{jk} := (w_{+k}/n)h_{jk}/h_{+k}.$$

3 Estimation

3.1 Dimension reduction

The minimization problem mentioned before involves a parameter $\mathbf{h} \in [0, \infty)^{\ell \times m}$ under $\binom{\ell}{2} \binom{m}{2}$ nonlinear inequality constraints. The parameter space and the number of constraints may be reduced as follows.

Lemma 3.1. *Let \mathcal{P} be the set of all index pairs (j, k) such that there exist indices $1 \leq j_1 \leq j \leq j_2 \leq \ell$ and $1 \leq k_1 \leq k \leq k_2 \leq m$ with $w_{j_1 k_2}, w_{j_2 k_1} > 0$.*

- (a) *If $\mathbf{h} \in [0, \infty)^{\ell \times m}$ satisfies (2.3) and $L(\mathbf{h}) > -\infty$, then $h_{jk} > 0$ for all $(j, k) \in \mathcal{P}$.*
- (b) *If such a matrix \mathbf{h} is replaced with $\tilde{\mathbf{h}} := (1_{[(j,k) \in \mathcal{P}]} h_{jk})_{j,k}$, then $\tilde{\mathbf{h}}$ satisfies (2.3), too, and $L(\tilde{\mathbf{h}}) \geq L(\mathbf{h})$ with equality if and only if $\tilde{\mathbf{h}} = \mathbf{h}$.*
- (c) *If $\mathbf{h} \in [0, \infty)^{\ell \times m}$ such that $\{(j, k) : h_{jk} > 0\} = \mathcal{P}$, then constraint (2.3) is equivalent to*

$$h_{j-1,k} h_{j,k-1} \leq h_{j-1,k-1} h_{j,k} \quad \text{for } 1 < j \leq \ell \text{ and } 1 < k \leq m. \quad (3.1)$$

All in all, we may restrict our attention to parameters $\mathbf{h} \in (0, \infty)^{\mathcal{P}}$ satisfying (3.1), where $h_{jk} := 0$ for $(j, k) \notin \mathcal{P}$. Note that (3.1) involves only $(\ell - 1)(m - 1)$ inequalities, and the inequality for one particular index pair (j, k) is nontrivial only if the two pairs $(j - 1, k), (j, k - 1)$ belong to \mathcal{P} .

The set \mathcal{P} consists of all pairs (j, k) such that the support of the empirical distribution \hat{R}_{emp} contains a point (x_{j_1}, y_{k_2}) “northwest” and a point (x_{j_2}, y_{k_1}) “southeast” of (x_j, y_k) . If \mathcal{P} contains two pairs $(j_2, k_1), (j_1, k_2)$ with $j_1 < j_2$ and $k_1 < k_2$, then it contains the whole set $\{j_1, \dots, j_2\} \times \{k_1, \dots, k_2\}$. Figure 1 illustrates the definition of \mathcal{P} . It also illustrates two alternative codings of \mathcal{P} : An index pair (j, k) belongs to \mathcal{P} if and only if $m_j \leq k \leq M_j$, where

$$\begin{aligned} m_j &:= \min\{k : w_{j'k} > 0 \text{ for some } j' \geq j\}, \\ M_j &:= \max\{k : w_{j'k} > 0 \text{ for some } j' \leq j\}. \end{aligned}$$

Note that $m_j \leq M_j$ for all j , $1 = m_1 \leq \dots \leq m_\ell$, and $M_1 \leq \dots \leq M_\ell = m$. Analogously, a pair (j, k) belongs to \mathcal{P} if and only if $\ell_k \leq j \leq L_k$, where

$$\begin{aligned} \ell_k &:= \min\{j : w_{jk'} > 0 \text{ for some } k' \geq k\}, \\ L_k &:= \max\{j : w_{jk'} > 0 \text{ for some } k' \leq k\}. \end{aligned}$$

Here $\ell_k \leq L_k$ for all k , $1 = \ell_1 \leq \dots \leq \ell_m$, and $L_1 \leq \dots \leq L_m = \ell$.

Note that by definition, for any index pair (j, k) ,

$$k \leq M_j \quad \text{if and only if} \quad j \geq \ell_k, \quad (3.2)$$

$$k \geq m_j \quad \text{if and only if} \quad j \leq L_k. \quad (3.3)$$

3.2 Reparametrization and reformulation

If we replace a parameter $\mathbf{h} \in (0, \infty)^{\mathcal{P}}$ with its component-wise logarithm $\boldsymbol{\theta} \in \mathbb{R}^{\mathcal{P}}$, then property (3.1) is equivalent to

$$\theta_{j-1,k-1} + \theta_{j,k} - \theta_{j-1,k} - \theta_{j,k-1} \geq 0 \quad \text{whenever } (j-1, k), (j, k-1) \in \mathcal{P}. \quad (3.4)$$

The set of all $\boldsymbol{\theta} \in \mathbb{R}^{\mathcal{P}}$ satisfying (3.4) is a closed convex cone and is denoted by Θ .

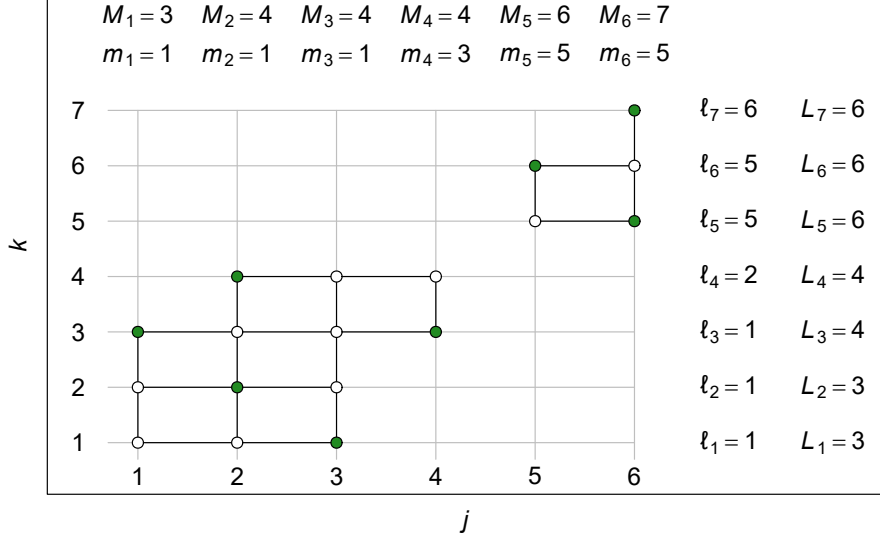


Figure 1: In this specific example, $n \geq 8$ raw observations yielded $\ell = 6$ different values x_j and $m = 7$ different values y_k . The green dots represent those (j, k) with $w_{jk} > 0$. The green dots and black circles represent the set \mathcal{P} .

Now our goal is to minimize

$$f(\boldsymbol{\theta}) := \sum_{(j,k) \in \mathcal{P}} (-w_{jk}\theta_{jk} + n \exp(\theta_{jk})) \quad (3.5)$$

over all $\boldsymbol{\theta} \in \Theta$.

Theorem 3.2. *There exists a unique minimizer $\hat{\boldsymbol{\theta}}$ of $f(\boldsymbol{\theta})$ over all $\boldsymbol{\theta} \in \Theta$.*

Uniqueness follows directly from f being strictly convex, but existence is less obvious, unless $w_{jk} > 0$ for all (j, k) . With $\hat{\boldsymbol{\theta}}$ at hand, the corresponding solution $\hat{\mathbf{h}} \in [0, \infty)^{\ell \times m}$ of the original problem is given by

$$\hat{h}_{jk} = \begin{cases} \exp(\hat{\theta}_{jk}) & \text{if } (j, k) \in \mathcal{P}, \\ 0 & \text{else.} \end{cases}$$

In the proof of Theorem 3.2 and from now on, we view $\mathbb{R}^{\mathcal{P}}$ as a Euclidean space with inner product $\langle \mathbf{x}, \mathbf{y} \rangle := \sum_{(j,k) \in \mathcal{P}} x_{jk} y_{jk}$ and the corresponding norm $\|\mathbf{x}\| := \langle \mathbf{x}, \mathbf{x} \rangle^{1/2}$. For a differentiable function $f : \mathbb{R}^{\mathcal{P}} \rightarrow \mathbb{R}$, its gradient is defined as $\nabla f(\mathbf{x}) := (\partial f(\mathbf{x}) / \partial x_{jk})_{(j,k) \in \mathcal{P}}$.

Let us explain briefly why traditional optimization algorithms may become infeasible for large sample sizes n . Depending on the input data, the set \mathcal{P} may contain more than cn^2 parameters, and the constraint (3.4) may involve at least cn^2 linear inequalities, where $c > 0$ is some generic constant. Even if we restrict our attention to parameters $\boldsymbol{\theta} \in \Theta$ such that a given subset of the inequalities in (3.4) are equalities, they span a linear space of dimension at least $\max(\ell, m)$, because all parameters θ_{jm_j} and $\theta_{\ell_k k}$ are unconstrained, and $\max(\ell, m)$ may be at least cn . Just determining a gradient and Hessian matrix of the target function f within this linear subspace would then require at least cn^4 steps. Consequently, traditional minimization algorithms involving exact Newton steps may be computationally infeasible. Alternatively, we propose an iterative algorithm with quasi Newton steps each of which has running time $O(n^2)$, and the required memory is of this order, too.

3.3 Finding a new proposal

Version 1. To determine whether a given parameter $\boldsymbol{\theta} \in \mathbb{R}^{\mathcal{P}}$ is already optimal and, if not, to obtain a better one, we reparametrize the problem a second time. Let $\tilde{\boldsymbol{\theta}} = T(\boldsymbol{\theta}) \in \mathbb{R}^{\mathcal{P}}$ be given by

$$\tilde{\theta}_{jk} = \begin{cases} \theta_{jm_j} & \text{if } k = m_j, \\ \theta_{jk} - \theta_{j,k-1} & \text{if } m_j < k \leq M_j. \end{cases}$$

Then $\boldsymbol{\theta} = T^{-1}(\tilde{\boldsymbol{\theta}}) = (\sum_{k'=m_j}^k \tilde{\theta}_{jk'})_{j,k}$, and $f(\boldsymbol{\theta})$ is equal to

$$\begin{aligned} \tilde{f}(\tilde{\boldsymbol{\theta}}) &:= \sum_{j=1}^{\ell} \sum_{k=m_j}^{M_j} \left(-w_{jk} \sum_{k'=m_j}^k \tilde{\theta}_{jk'} + n \exp\left(\sum_{k'=m_j}^k \tilde{\theta}_{jk'} \right) \right) \\ &= \sum_{j=1}^{\ell} \sum_{k=m_j}^{M_j} \left(-\underline{w}_{jk} \tilde{\theta}_{jk} + n \exp\left(\sum_{k'=m_j}^k \tilde{\theta}_{jk'} \right) \right) \quad \text{with } \underline{w}_{jk} := \sum_{k'=k}^{M_j} w_{jk'}. \end{aligned}$$

More importantly, we may represent \mathcal{P} as

$$\begin{aligned} \mathcal{P} &= \{(j, m_j) : 1 \leq j \leq \ell\} \cup \{(j, k) : 1 \leq j \leq \ell, m_j < k \leq M_j\} \\ &= \{(j, m_j) : 1 \leq j \leq \ell\} \cup \bigcup_{k=2}^m \{(j, k) : \ell_k \leq j \leq L_{k-1}\}, \end{aligned}$$

where the latter equation follows from (3.2) and (3.3). Now the constraints (3.4) read

$$(\tilde{\theta}_{jk})_{j=\ell_k}^{L_{k-1}} \in \mathbb{R}_{\uparrow}^{L_{k-1}-\ell_k+1} \quad \text{whenever } 2 \leq k \leq m \text{ and } L_{k-1} - \ell_k + 1 \geq 2. \quad (3.6)$$

Here $\mathbb{R}_{\uparrow}^d := \{\boldsymbol{x} \in \mathbb{R}^d : x_1 \leq \dots \leq x_d\}$. The set of $\tilde{\boldsymbol{\theta}} \in \mathbb{R}^{\mathcal{P}}$ satisfying (3.6) is denoted by $\tilde{\Theta}$.

For given $\boldsymbol{\theta}$ and $\tilde{\boldsymbol{\theta}} = T(\boldsymbol{\theta})$, we approximate $\tilde{f}(\tilde{\boldsymbol{x}})$ by the quadratic function

$$\begin{aligned} \tilde{\boldsymbol{x}} &\mapsto \tilde{f}(\tilde{\boldsymbol{\theta}}) + \langle \nabla \tilde{f}(\tilde{\boldsymbol{\theta}}), \tilde{\boldsymbol{x}} - \tilde{\boldsymbol{\theta}} \rangle + 2^{-1} \sum_{(j,k) \in \mathcal{P}} \frac{\partial^2 \tilde{f}}{\partial \tilde{\theta}_{jk}^2}(\tilde{\boldsymbol{\theta}}) (\tilde{x}_{jk} - \tilde{\theta}_{jk})^2 \\ &= \text{const}(\boldsymbol{\theta}) + 2^{-1} \sum_{(j,k) \in \mathcal{P}} \tilde{v}_{jk}(\boldsymbol{\theta}) (\tilde{x}_{jk} - \tilde{\gamma}_{jk}(\boldsymbol{\theta}))^2 \\ &= \text{const}(\boldsymbol{\theta}) + 2^{-1} \sum_{j=1}^{\ell} \tilde{v}_{jm_j}(\boldsymbol{\theta}) (\tilde{x}_{jm_j} - \tilde{\gamma}_{jm_j}(\boldsymbol{\theta}))^2 \\ &\quad + 2^{-1} \sum_{k=2}^m \sum_{\ell_k \leq j \leq L_{k-1}} \tilde{v}_{jk}(\boldsymbol{\theta}) (\tilde{x}_{jk} - \tilde{\gamma}_{jk}(\boldsymbol{\theta}))^2 \end{aligned}$$

with

$$\begin{aligned} \tilde{v}_{jk}(\boldsymbol{\theta}) &:= \frac{\partial^2 \tilde{f}}{\partial \tilde{\theta}_{jk}^2}(\tilde{\boldsymbol{\theta}}) &= n \sum_{k'=k}^{M_j} \exp(\theta_{jk'}), \\ \tilde{\gamma}_{jk}(\boldsymbol{\theta}) &:= \tilde{\theta}_{jk} - \tilde{v}_{jk}(\boldsymbol{\theta})^{-1} \frac{\partial \tilde{f}}{\partial \tilde{\theta}_{jk}}(\tilde{\boldsymbol{\theta}}) &= T_{jk}(\boldsymbol{\theta}) + \tilde{v}_{jk}(\boldsymbol{\theta})^{-1} \underline{w}_{jk} - 1. \end{aligned}$$

This quadratic function of $\tilde{\boldsymbol{x}}$ is easily minimized over $\tilde{\Theta}$ via the pool-adjacent-violators algorithm, applied to the subtuple $(\tilde{x}_{jk})_{j=\ell_k}^{L_{k-1}}$ for each $k = 2, \dots, m$ separately. Then we obtain the proposal

$$\Psi^{\text{row}}(\boldsymbol{\theta}) := T^{-1}(\tilde{\boldsymbol{\theta}}_*(\boldsymbol{\theta})) \quad \text{with} \quad \tilde{\boldsymbol{\theta}}_*(\boldsymbol{\theta}) := \arg \min_{\tilde{\boldsymbol{x}} \in \tilde{\Theta}} \sum_{(j,k) \in \mathcal{P}} \tilde{v}_{jk}(\boldsymbol{\theta}) (\tilde{x}_{jk} - \tilde{\gamma}_{jk}(\boldsymbol{\theta}))^2.$$

Interestingly, if $\boldsymbol{\theta}$ is row-wise calibrated in the sense that $n \sum_{k=m_j}^{M_j} \exp(\theta_{jk}) = w_{j+}$ for $1 \leq j \leq \ell$, then $\tilde{\gamma}_{jm_j}(\boldsymbol{\theta}) = \tilde{\theta}_{jm_j}$ and thus $\Psi_{jm_j}^{\text{row}}(\boldsymbol{\theta}) = \theta_{jm_j}$ for $1 \leq j \leq \ell$.

Version 2. Instead of reparametrizing $\boldsymbol{\theta} \in \Theta$ in terms of its values θ_{jm_j} , $1 \leq j \leq \ell$, and its increments within rows, one could reparametrize it in terms of its values $\theta_{\ell_k k}$, $1 \leq k \leq m$, and its increments within columns, leading to a proposal $\Psi^{\text{col}}(\boldsymbol{\theta})$. Here, $\Psi_{\ell_k k}^{\text{col}}(\boldsymbol{\theta}) = \theta_{\ell_k k}$ for $1 \leq k \leq m$, provided that $\boldsymbol{\theta}$ is column-wise calibrated.

3.4 Calibration

In terms of the log-parametrization with $\boldsymbol{\theta} \in \Theta$, the row-wise calibration mentioned earlier for \mathbf{h} means to replace θ_{jk} with

$$\theta_{jk} - \log\left(\sum_{k'=m_j}^{M_j} \exp(\theta_{jk'})\right) + \log(w_{j+}/n).$$

Analogously, replacing θ_{jk} with

$$\theta_{jk} - \log\left(\sum_{j'=\ell_k}^{L_k} \exp(\theta_{j'k})\right) + \log(w_{+k}/n)$$

leads to a column-wise calibrated parameter $\boldsymbol{\theta}$. Iterating these calibrations alternately, leads to a parameter which is (approximately) calibrated, row-wise as well as column-wise.

3.5 From new proposal to new parameter

Both functions $\Psi = \Psi^{\text{row}}, \Psi^{\text{col}}$ have some useful properties summarized in the next lemma.

Lemma 3.3. *The function Ψ is continuous on Θ with $\Psi(\hat{\boldsymbol{\theta}}) = \hat{\boldsymbol{\theta}}$. For $\boldsymbol{\theta} \in \Theta \setminus \{\hat{\boldsymbol{\theta}}\}$,*

$$\delta(\boldsymbol{\theta}) := \langle \nabla f(\boldsymbol{\theta}), \boldsymbol{\theta} - \Psi(\boldsymbol{\theta}) \rangle > 0,$$

$$f(\boldsymbol{\theta}) - f(\hat{\boldsymbol{\theta}}) \leq \max(2\delta(\boldsymbol{\theta}), \beta_1(\boldsymbol{\theta})\sqrt{\delta(\boldsymbol{\theta})})\|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}\|,$$

and

$$\max_{t \in [0,1]} \left(f(\boldsymbol{\theta}) - f((1-t)\boldsymbol{\theta} + t\Psi(\boldsymbol{\theta})) \right) \geq \min\left(2^{-1}\delta(\boldsymbol{\theta}), \frac{\delta(\boldsymbol{\theta})^2}{\beta_2(\boldsymbol{\theta})\|\boldsymbol{\theta} - \Psi(\boldsymbol{\theta})\|^2}\right)$$

with continuous functions $\beta_1, \beta_2 : \Theta \rightarrow (0, \infty)$.

In view of this lemma, we want to replace $\boldsymbol{\theta} \neq \hat{\boldsymbol{\theta}}$ with $(1-t_*)\boldsymbol{\theta} + t_*\Psi(\boldsymbol{\theta})$ for some suitable $t_* = t_*(\boldsymbol{\theta}) \in [0, 1]$ such that $f(\boldsymbol{\theta})$ really decreases. More specifically, with

$$\rho_{\boldsymbol{\theta}}(t) := f(\boldsymbol{\theta}) - f((1-t)\boldsymbol{\theta} + t\Psi(\boldsymbol{\theta})),$$

our goals are that for some constant $\kappa \in (0, 1]$,

$$\rho_{\boldsymbol{\theta}}(t_*) \geq \kappa \max_{t \in [0,1]} \rho_{\boldsymbol{\theta}}(t),$$

and in case of $\rho_{\boldsymbol{\theta}}$ being (approximately) a quadratic function, t_* should be (approximately) equal to $\arg \max_{t \in [0,1]} \rho_{\boldsymbol{\theta}}(t)$. For that, we proceed similarly as in [Dümbgen et al. \(2006\)](#).

```

 $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta}^{(0)}$ 
 $\delta \leftarrow \infty$ 
 $s \leftarrow 0$ 
while  $\delta \geq \delta_o$  do
   $\boldsymbol{\theta} \leftarrow$  calibration of  $\boldsymbol{\theta}$ 
  if  $s$  is even, do
     $(\boldsymbol{\psi}, \delta) \leftarrow (\Psi^{\text{row}}(\boldsymbol{\theta}), \langle \nabla f(\boldsymbol{\theta}), \boldsymbol{\theta} - \Psi^{\text{row}}(\boldsymbol{\theta}) \rangle)$ 
  else
     $(\boldsymbol{\psi}, \delta) \leftarrow (\Psi^{\text{col}}(\boldsymbol{\theta}), \langle \nabla f(\boldsymbol{\theta}), \boldsymbol{\theta} - \Psi^{\text{col}}(\boldsymbol{\theta}) \rangle)$ 
  end if
   $\rho' \leftarrow \delta$ 
  while  $f(\boldsymbol{\psi}) > f(\boldsymbol{\theta})$  do
     $(\boldsymbol{\psi}, \rho') \leftarrow (2^{-1}(\boldsymbol{\theta} + \boldsymbol{\psi}), 2^{-1}\rho')$ 
  end while
   $t_* \leftarrow \min(1, 2^{-1}\rho' / (\rho' - f(\boldsymbol{\theta}) + f(\boldsymbol{\psi})))$ 
   $\boldsymbol{\theta} \leftarrow (1 - t_*)\boldsymbol{\theta} + t_*\boldsymbol{\psi}$ 
   $s \leftarrow s + 1$ 
end while

```

Table 1: Pseudo code of our algorithm, returning an approximation $\boldsymbol{\theta}$ of $\widehat{\boldsymbol{\theta}}$.

We determine $t_o := 2^{-n_o}$ with n_o the smallest integer such that $\rho_{\boldsymbol{\theta}}(2^{-n_o}) \geq 0$. Then we define a Hermite interpolation of $\rho_{\boldsymbol{\theta}}$:

$$\tilde{\rho}_{\boldsymbol{\theta}}(t) := \rho'_{\boldsymbol{\theta}}(0)t - c_o t^2 \quad \text{with} \quad c_o := t_o^{-1}(\rho'_{\boldsymbol{\theta}}(0) - t_o^{-1}\rho_{\boldsymbol{\theta}}(t_o)) > 0.$$

This new function is such that $\tilde{\rho}_{\boldsymbol{\theta}}(t) = \rho_{\boldsymbol{\theta}}(t)$ for $t = 0, t_o$, and $\tilde{\rho}'_{\boldsymbol{\theta}}(0) = \rho'_{\boldsymbol{\theta}}(0) > 0$. Since $\tilde{\rho}'_{\boldsymbol{\theta}}(t) = \rho'_{\boldsymbol{\theta}}(0) - 2tc_o$, the maximizer of $\tilde{\rho}_{\boldsymbol{\theta}}$ over $[0, t_o]$ is given by

$$t_* := \min(t_o, 2^{-1}\rho'_{\boldsymbol{\theta}}(0)/c_o).$$

As shown in Lemma 1 of [Dümbgen et al. \(2006\)](#), this choice of t_* fulfils the requirements just stated, where $\kappa = 1/4$.

3.6 Complete algorithms

A possible starting point for the algorithm is given by $\boldsymbol{\theta}^{(0)} := (-\log(\#\mathcal{P}))_{(j,k) \in \mathcal{P}}$, but any other parameter $\boldsymbol{\theta}^{(0)} \in \Theta$ would work, too. Suppose we have determined already $\boldsymbol{\theta}^{(0)}, \dots, \boldsymbol{\theta}^{(s)}$ such that $f(\boldsymbol{\theta}^{(0)}) \geq \dots \geq f(\boldsymbol{\theta}^{(s)})$. Let $\Psi(\boldsymbol{\theta}^{(s)})$ be a new proposal with $\Psi = \Psi^{\text{row}}$ or $\Psi = \Psi^{\text{col}}$, and let $\boldsymbol{\theta}^{(s+1)} = (1 - t_*^{(s)})\boldsymbol{\theta}^{(s)} + t_*^{(s)}\Psi(\boldsymbol{\theta}^{(s)})$ with $t_*^{(s)} = t_*(\boldsymbol{\theta}^{(s)}) \in [0, 1]$ as described before. No matter which proposal function Ψ we are using in each step, the resulting sequence $(\boldsymbol{\theta}^{(s)})_{s \geq 0}$ will always converge to $\widehat{\boldsymbol{\theta}}$.

Theorem 3.4. *Let $(\boldsymbol{\theta}^{(s)})_{s \geq 0}$ be the sequence just described. Then $\lim_{s \rightarrow \infty} \boldsymbol{\theta}^{(s)} = \widehat{\boldsymbol{\theta}}$.*

Our numerical experiments showed that a particularly efficient refinement is as follows: Before computing a new proposal $\Psi(\boldsymbol{\theta}^{(s)})$, one should calibrate $\boldsymbol{\theta}^{(s)}$ in the sense that it is row-wise and column-wise calibrated. If s is even, we compute $\Psi^{\text{row}}(\boldsymbol{\theta}^{(s)})$ to determine the next candidate $\boldsymbol{\theta}^{(s+1)}$. If s is odd, we compute $\Psi^{\text{col}}(\boldsymbol{\theta}^{(s)})$ to obtain $\boldsymbol{\theta}^{(s+1)}$. The algorithm stops as soon as $\delta(\boldsymbol{\theta}^{(s)}) = \langle \nabla f(\boldsymbol{\theta}^{(s)}), \boldsymbol{\theta}^{(s)} - \Psi(\boldsymbol{\theta}^{(s)}) \rangle$ is smaller than a prescribed small threshold. Table 1 provides corresponding pseudo code.

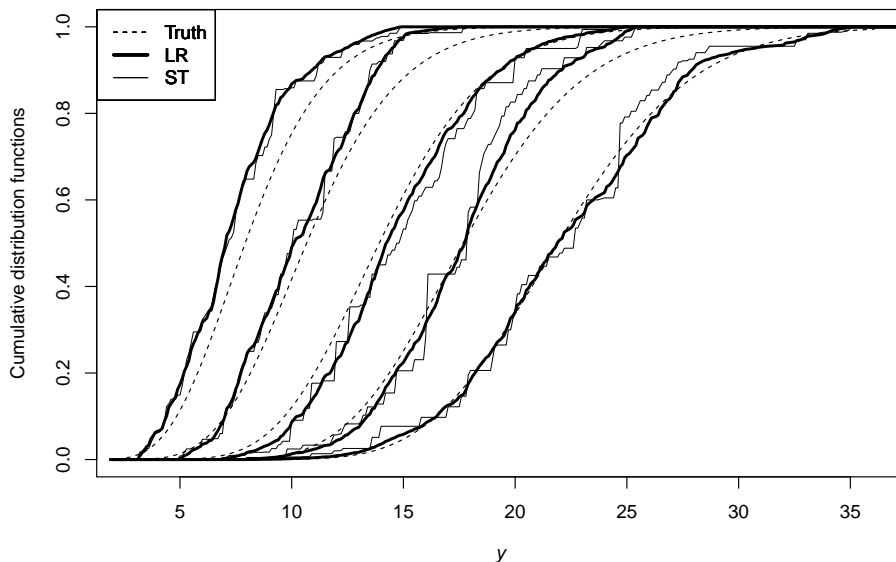


Figure 2: The true conditional Gamma distribution function G_x , the estimate under likelihood ratio (LR) order constraint \hat{G}_x and the estimated under usual stochastic (ST) order constraint \check{G}_x are displayed from left to right for $x \in \{1.5, 2, 2.5, 3, 3.5\}$.

4 Simulation study

In this section, we compare estimation and prediction performances of the likelihood ratio order constrained estimator presented in this article with the estimator under usual stochastic order obtained via isotonic distributional regression. The latter estimator was mentioned briefly in the introduction. It is extensively discussed in [Henzi et al. \(2021b\)](#) and [Mösching and Dümbgen \(2020\)](#).

4.1 A Gamma model

We choose a parametric family of distributions from which we draw observations. We will then use these data to provide distribution estimates which we then compare with the truth. The specific model we have in mind is a family $(Q_x)_{x \in \mathfrak{X}}$ of Gamma distributions with densities

$$g_x(y) := \frac{b(x)^{-a(x)}}{\Gamma(a(x))} y^{a(x)-1} \exp(-y/b(x)),$$

with respect to Lebesgue measure on $(0, \infty)$, with some shape function $a : \mathfrak{X} \rightarrow (0, \infty)$ and scale function $b : \mathfrak{X} \rightarrow (0, \infty)$. Then Q_x is isotonic in $x \in \mathfrak{X}$ with respect to likelihood ratio ordering if and only if both functions a and b are isotonic. Recall that since the family is increasing in likelihood ratio order, it is also increasing with respect to the usual stochastic order.

The specific shape and scale functions used for this study are

$$a(x) := 2 + (x + 1)^2 \quad \text{and} \quad b(x) := 1 - \exp(-10x),$$

defined for $x \in \mathfrak{X} := [1, 4]$. Figure 2 displays corresponding true conditional distribution functions for a selection of x 's.

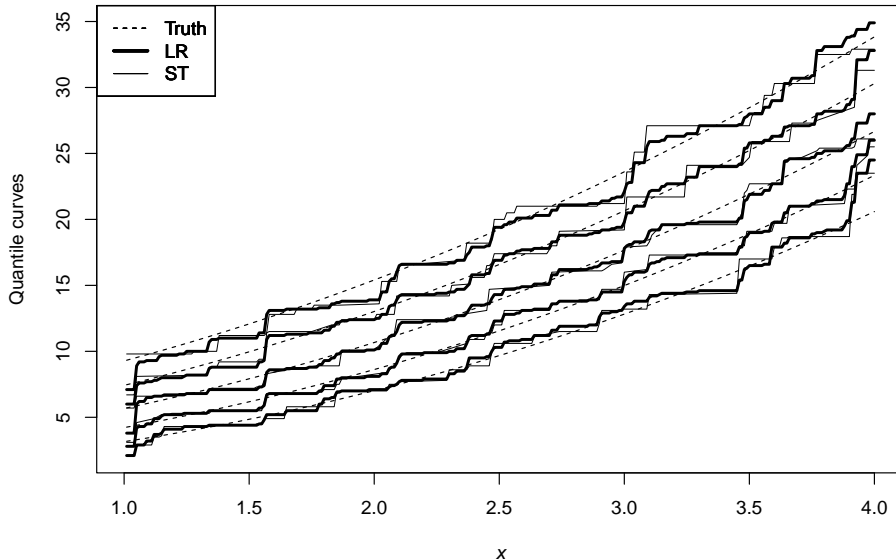


Figure 3: Selection of β -quantile curves. Specifically, a taut-string (Dümbgen and Kovac, 2009) is computed between the lower $\mathfrak{X} \ni x \mapsto \min\{y \in \mathbb{R} : \tilde{G}_x(y) \geq \beta\}$ and upper $\mathfrak{X} \ni x \mapsto \inf\{y \in \mathbb{R} : \tilde{G}_x(y) > \beta\}$ quantile curves for each $\tilde{G} \in \{G, \hat{G}, \check{G}\}$ (corresponding respectively to ‘Truth’, ‘LR’ and ‘ST’) and $\beta \in \{0.1, 0.25, 0.5, 0.75, 0.9\}$.

4.2 Sampling method

Let $\ell_o \in \{50, 1000\}$ be a predefined number and let

$$\mathfrak{X}_o := 1 + \frac{3}{\ell_o} \cdot \{1, 2, \dots, \ell_o\} \subset \mathfrak{X}.$$

For a given sample size $n \in \mathbb{N}$, the sample $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ is obtained as follows: Draw X_1, X_2, \dots, X_n uniformly from \mathfrak{X}_o and sample independently each Y_k from Q_{X_k} . This yields unique covariates $x_1 < \dots < x_\ell$ as well as unique responses $y_1 < \dots < y_m$, for some $1 \leq \ell, m \leq n$.

For each such sample, we compute estimates of $(Q_{x_j})_{j=1}^\ell$ under likelihood ratio order and usual stochastic order constraints. Using linear interpolation, we complete both families of estimates with covariates originally in $\{x_j\}_{j=1}^\ell$ to families of estimates with covariates in the full set \mathfrak{X}_o , see Lemma A.1. We therefore obtain estimates $(\hat{Q}_x)_{x \in \mathfrak{X}_o}$ and $(\check{Q}_x)_{x \in \mathfrak{X}_o}$ under likelihood ratio order and usual stochastic order constraint, respectively. The corresponding families of cumulative distribution functions are written $(\hat{G}_x)_{x \in \mathfrak{X}_o}$ and $(\check{G}_x)_{x \in \mathfrak{X}_o}$, whereas the truth is denoted by $(G_x)_{x \in \mathfrak{X}_o}$. Although the performance of the empirical distribution is worse than those of the two order constrained estimators, it is still useful to study its behaviour, for instance to better understand boundary effects. The family of empirical cumulative distribution functions will be written $(\hat{G}_x)_{x \in \mathfrak{X}_o}$.

4.3 Single sample

Figure 2 provides a visual comparison of a selection of true conditional distribution functions with their corresponding estimates under order constraint for a single sample generated in the setting $\ell_o = 1000$ and $n = 1000$. It shows that the estimates under likelihood ratio order constraint are much smoother than those under usual stochastic order constraint. The former are in general also closer to the truth than the latter. This fact is in reality true on average, as demonstrated in the next paragraph. Smoothness and

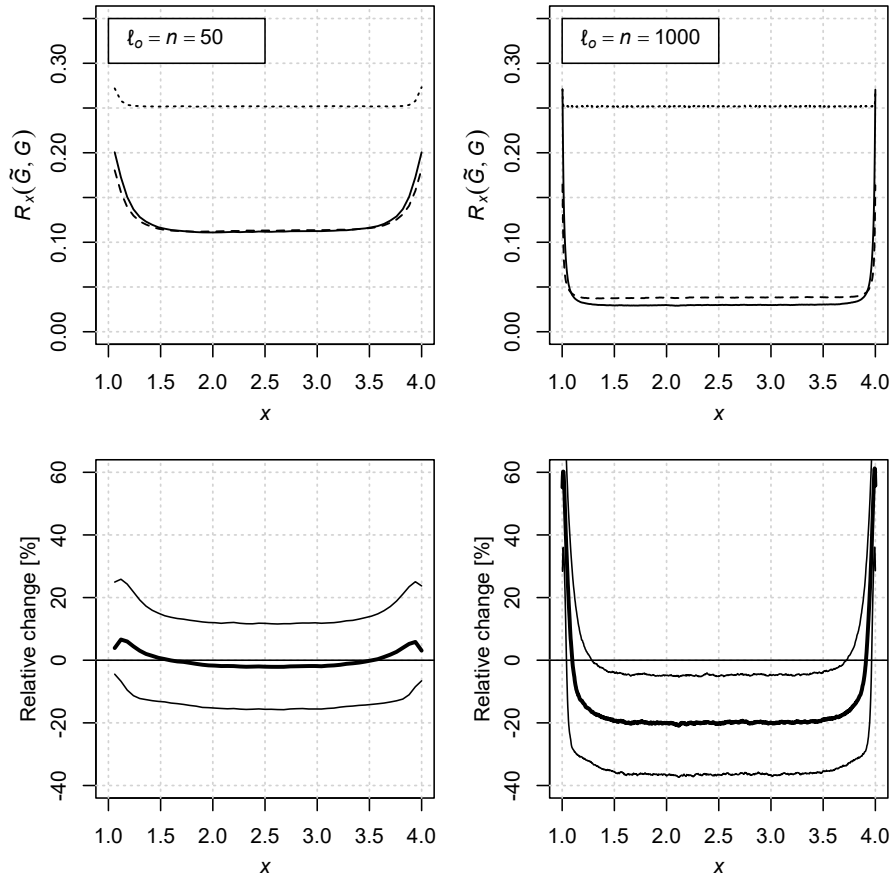


Figure 4: Monte Carlo simulations to evaluate estimation performances with a simple score. First row: Simple scores with \tilde{G} being either \hat{G} (solid line), \check{G} (dashed line) or $\hat{\mathbb{G}}$ (dotted line). Second row: Relative change of score when enforcing a likelihood ratio order constraint over the usual stochastic order constraint. The thicker line is the median variation, whereas the thin lines are the first and third quartiles. Negative values represent an improvement in score.

greater precision in estimation resulting from the likelihood ratio order is also apparent in Figure 3, which displays a selection of quantile curves for each $\tilde{G} \in \{G, \hat{G}, \check{G}\}$.

4.4 A simple score

To assess the ability of each estimator to retrieve the truth, we produce Monte-Carlo estimates of the median of the score

$$R_x(\tilde{G}, G) := \int |\tilde{G}_x(y) - G_x(y)| dQ_x(y),$$

for each estimator $\tilde{G} \in \{\hat{G}, \check{G}, \hat{\mathbb{G}}\}$ and for each $x \in \mathfrak{X}_o$. The above score may be decomposed as a sum of simple expressions involving the evaluation of \tilde{G}_x and G_x on the finite set of unique responses, see Section A.3. We also compute Monte-Carlo quartiles of the relative change in score

$$100 \cdot \frac{R_x(\hat{G}, G) - R_x(\check{G}, G)}{R_x(\check{G}, G)}.$$

The results of the simulations are displayed in Figure 4. A first observation is that the performance of all three estimators decreases towards the boundary points of \mathfrak{X} , and this effect is more pronounced for the two order constrained estimators. This is a known phenomenon from shape constrained inference. However, in the interior of \mathfrak{X} , taking the stochastic ordering into account pays off. The second row of plots in Figure 4 shows the relative change in score when estimating the family of distributions with a likelihood ratio order constraint instead of the usual stochastic order constraint. It is observed that the improvement in score becomes larger and occurs on a wider sub-interval of \mathfrak{X} as ℓ_o and n increase. Only towards the boundary, the usual stochastic order seems to have better performance.

4.5 Theoretical predictive performances

Using the same Gamma model, we evaluate predictive performances of both estimators using the continuous ranked probability score

$$\text{CRPS}(\tilde{G}_x, y) := \int \left(\tilde{G}_x(z) - 1_{[y \leq z]} \right)^2 dz.$$

The CRPS is a strictly proper scoring rule which allows for comparisons of probabilistic forecasts, see [Gneiting and Raftery \(2007\)](#) and [Jordan et al. \(2019\)](#). It can be seen as an extension of the mean absolute error for probabilistic forecasts. The CRPS is therefore interpreted in the same unit of measurement as the true distribution or data.

Because the true underlying distribution is known in the present simulation setting, the expected CRPS score is given by

$$\begin{aligned} S_x(\tilde{G}, G) &:= \int \text{CRPS}(\tilde{G}_x, y) dQ_x(y) \\ &= \sum_{k=0}^m \int_{[y_k, y_{k+1})} (\tilde{G}_x(y_k) - G_x(y))^2 dy + \frac{b(x)}{B(1/2, a(x))}, \end{aligned}$$

where $y_0 := 0$, $y_{m+1} := +\infty$ and $B(\cdot, \cdot)$ is the beta function. As shown in Section A.3, the above sum of integrals may be rewritten as a sum of elementary expressions involving the evaluation of \tilde{G}_x and G_x on the finite set of unique responses, as well as two simple integrals which are computed via numerical integration. Consequently, we compute Monte-Carlo estimates of the median of each score $S_x(\tilde{G}, G)$, $\tilde{G} \in \{\hat{\tilde{G}}, \check{\tilde{G}}, \hat{\mathbb{G}}\}$, as well as estimates of quartiles of the relative change in score when choosing $\hat{\tilde{G}}$ over $\check{\tilde{G}}$.

Figure 5 outlines the results of the simulations. Similar boundary effects as for the simple score are observed. On the interior of \mathfrak{X} , the usual stochastic order improves the naive empirical estimator, and the likelihood ratio order yields the best results. In terms of relative change in score, it appears that imposing a likelihood ratio order constraint to estimate the family of distributions yields an average score reduction of about 0.5% in comparison with the usual stochastic order estimator for a sample of $n = 50$. For $n = 1000$, this improvement occurs on a wider subinterval of \mathfrak{X} and more frequently, as shown by the third quartile curve. Note further that the expected CRPS increases on the interior of \mathfrak{X} . This is due to the fact that the CRPS has the same unit of measurement as the response variable. Since the scale of the response characterized by b increases with x , then so does the corresponding score.

4.6 Empirical predictive performances

We use the weight for age dataset already studied in [Mösching and Dümbgen \(2020\)](#). It comprises the age and weight of $n = 16\,432$ girls whose age in years lies within $\mathfrak{X} :=$

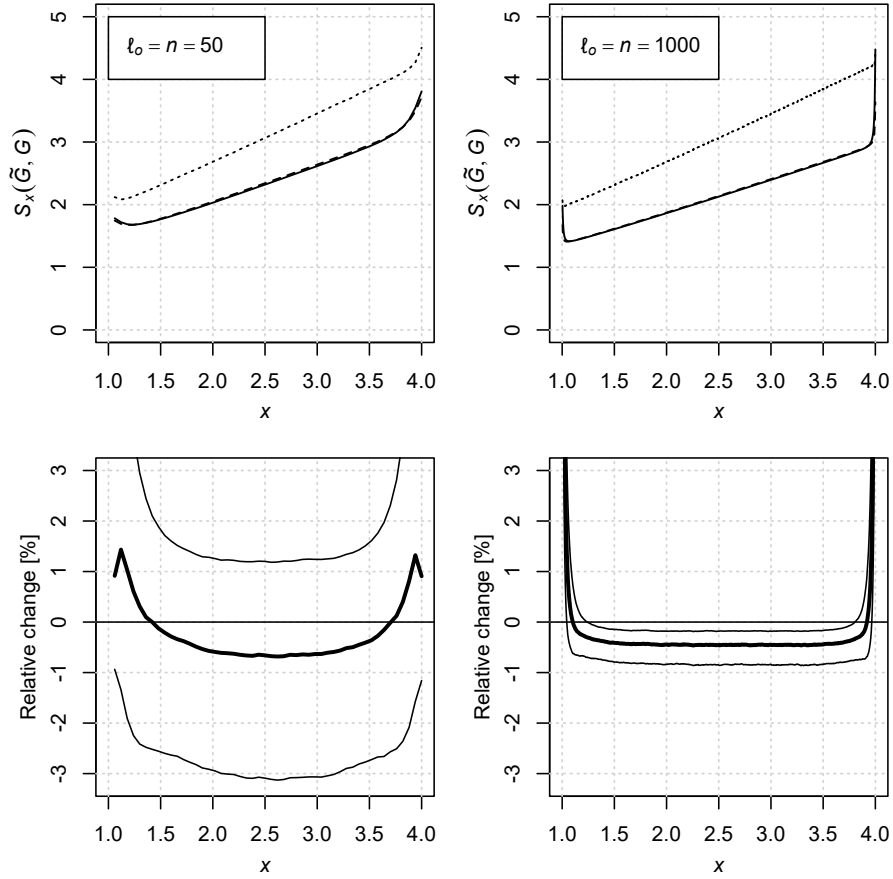


Figure 5: Monte Carlo simulations to evaluate prediction performances using a CRPS-type score. First row: CRPS scores with \tilde{G} being either \hat{G} (solid line), \check{G} (dashed line) or \mathbb{G} (dotted line). Second row: Relative change of score when enforcing a likelihood ratio order constraint over the usual stochastic order constraint.

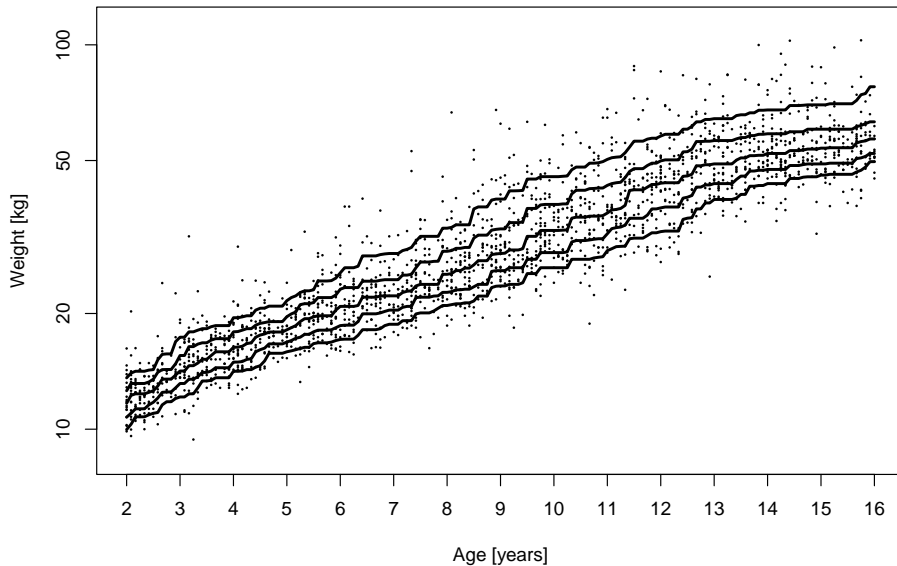


Figure 6: Subsample of the weight for age data and β -quantile curves computed from that sample under likelihood ratio order constraint, $\beta \in \{0.1, 0.25, 0.5, 0.75, 0.9\}$. A logarithmic scale was used for the weight variable.

[2, 16]. A subsample of these data of size 2000 is presented in Figure 6, along with estimated quantile curves under likelihood ratio order using that subsample. The dataset was publicly released as part of the National Health and Nutrition Examination Survey conducted in the US between 1963 and 1991 (data available from www.cdc.gov) and was analyzed by Kuczmarski et al. (2002) with parametric models to produce smooth quantile curves.

Although the likelihood ratio order constraint is harder to justify than the very natural stochastic order constraint, we are interested in the effect of a stronger regularization imposed by the former constraint.

The forecast evaluation is performed using a leave- n_{train} -out cross-validation scheme. More precisely, we choose random subsets $\mathcal{D}_{\text{train}}$ of n_{train} observations which we use to train our estimators. Using the rest of the $n_{\text{test}} := n - n_{\text{train}}$ data pairs in $\mathcal{D}_{\text{test}}$, we evaluate predictive performance by computing the sample median of $\hat{S}_x(\tilde{G}, \mathcal{D}_{\text{test}})$ for each estimator $\tilde{G} \in \{\hat{G}, \check{G}, \hat{G}\}$ and each $x \in \mathfrak{X}_o$, where

$$\hat{S}_x(\tilde{G}, \mathcal{D}_{\text{test}}) := \frac{\sum_{(X,Y) \in \mathcal{D}_{\text{test}}: X=x} \text{CRPS}(\tilde{G}_x, Y)}{\#\{(X, Y) \in \mathcal{D}_{\text{test}} : X = x\}}.$$

Quartile estimates of the relative change in score are also computed.

Figure 7 shows the forecast evaluation results. As expected, the empirical CRPS increases with age, since the spread of the weight increases with age. As to the relative change in score, improvements of about 0.5% can be seen for both training sample sizes. The region of \mathfrak{X} where the estimator under likelihood ratio order constraint shows better predictive performances is the widest for the largest training sample size. These results show the benefit of a stronger regularization.

Code availability

Our procedure is implemented in the R-package LRDistReg and is available from the GitHub of the first author: <https://github.com/AlexandreMoesching/LRDistReg>. Its implementation includes C++ code which is then integrated in R using Rcpp.

References

- BEARE, B. K. and MOON, J.-M. (2015). Nonparametric tests of density ratio ordering. *Econometric Theory* **31** 471–492.
- CAROLAN, C. A. and TEBBS, J. M. (2005). Nonparametric tests for and against likelihood ratio ordering in the two-sample problem. *Biometrika* **92** 159–171.
- DARDANONI, V. and FORCINA, A. (1998). A unified approach to likelihood inference on stochastic orderings in a nonparametric context. *J. Amer. Statist. Assoc.* **93** 1112–1123.
- DÜMBGEN, L., FREITAG-WOLF, S. and JONGBLOED, G. (2006). Estimating a unimodal distribution from interval-censored data. *J. Amer. Statist. Assoc.* **101** 1094–1106.
- DÜMBGEN, L. and KOVAC, A. (2009). Extensions of smoothing via taut strings. *Electron. J. Stat.* **3** 41–75.
- DÜMBGEN, L. and MÖSCHING, A. (2023). On stochastic orders and total positivity. *ESAIM Probab. Stat.* **27** 461–481.

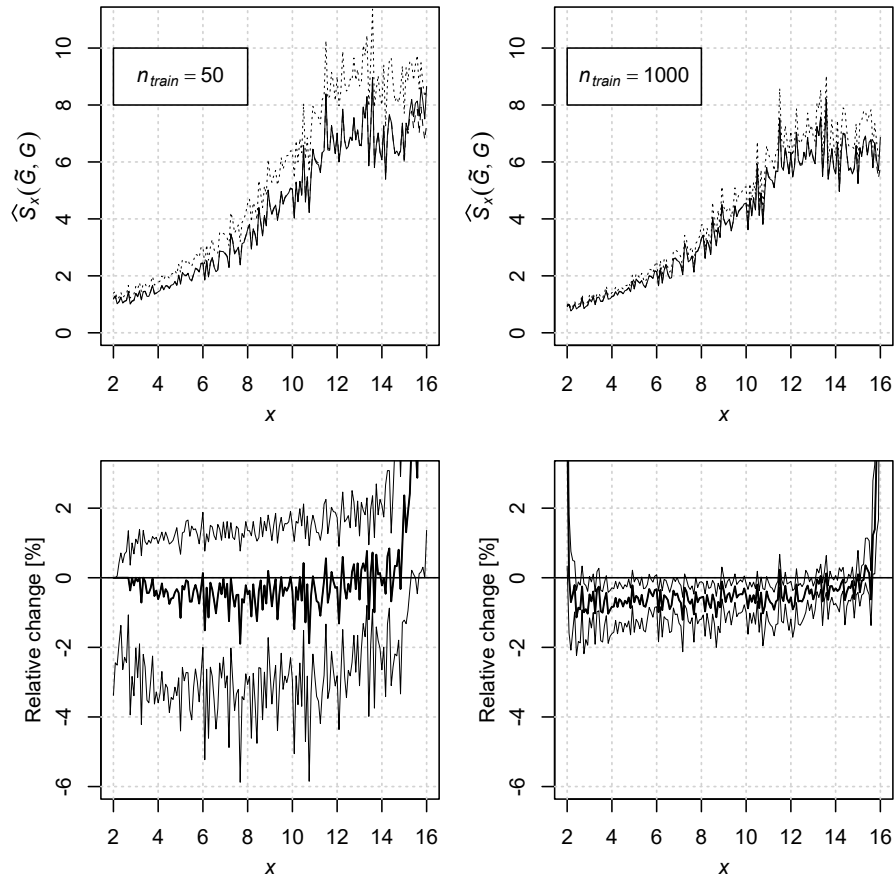


Figure 7: Monte Carlo simulations to evaluate prediction performances using an empirical CRPS score. First row: empirical CRPS scores with \tilde{G} being either \hat{G} (solid line), \tilde{G} (dashed line, hardly distinguishable from solid line) or $\hat{\tilde{G}}$ (dotted line). Second row: Relative change of score when enforcing a likelihood ratio order constraint over the usual stochastic order constraint.

- DÜMBGEN, L., MÖSCHING, A. and STRÄHL, C. (2021). Active set algorithms for estimating shape-constrained density ratios. *Comput. Statist. Data Anal.* **163** Paper No. 107300, 19.
- DYKSTRA, R., KOCHAR, S. and ROBERTSON, T. (1995). Inference for likelihood ratio ordering in the two-sample problem. *J. Amer. Statist. Assoc.* **90** 1034–1040.
- GNEITING, T. and RAFTERY, A. E. (2007). Strictly proper scoring rules, prediction, and estimation. *J. Amer. Statist. Assoc.* **102** 359–378.
- HENZI, A., KLEGER, G.-R., HILTY, M. P., WENDEL GARCIA, P. D. and ZIEGEL, J. F. (2021a). Strictly proper scoring rules, prediction, and estimation. *PLoS ONE* **16** e0247265.
- HENZI, A., ZIEGEL, J. F. and GNEITING, T. (2021b). Isotonic distributional regression. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **83** 963–993.
- HU, D., YUAN, M., YU, T. and LI, P. (2023). Statistical inference for the two-sample problem under likelihood ratio ordering, with application to the ROC curve estimation. *Stat. Med.* **42**(20) 3649–3664.
- HÜTTER, J.-C., MAO, C., RIGOLLET, P. and ROBEVA, E. (2020). Optimal rates for estimation of two-dimensional totally positive distributions. *Electron. J. Stat.* **14**(2) 2600–2652.
- JEWITT, I. (1991). Applications of likelihood ratio orderings in economics. In *Stochastic orders and decision under risk (Hamburg, 1989)*, vol. 19 of *IMS Lecture Notes Monogr. Ser.* Inst. Math. Statist., Hayward, CA, 174–189.
- JONGBLOED, G. (1998). The iterative convex minorant algorithm for nonparametric estimation. *J. Comput. Graph. Statist.* **7** 310–321.
- JORDAN, A., KRÜGER, F. and LERCH, S. (2019). Evaluating probabilistic forecasts with scoringrules. *Journal of Statistical Software* **90** 1–37.
- KARLIN, S. (1968). *Total positivity. Vol. I.* Stanford University Press, Stanford, Calif.
- KUCZMARSKI, R. J., OGDEN, C. L., GUO, S. S., GRUMMER-STRAWN, L. M., FLEGAL, K. M., MEI, Z., WEI, R., CURTIN, L. R., ROCHE, A. F. and JOHNSON, C. L. (2002). CDC Growth Charts for the United States: Methods and development. *Vital Health Stat.* **246**.
- MÖSCHING, A. and DÜMBGEN, L. (2020). Monotone least squares and isotonic quantiles. *Electron. J. Stat.* **14** 24–49.
- OWEN, A. B. (1988). Empirical likelihood ratio confidence intervals for a single functional. *Biometrika* **75** 237–249.
- OWEN, A. B. (2001). *Empirical likelihood*. No. 92 in *Monographs on Statistics and Applied Probability*, Chapman and Hall/CRC.
- ROBERTSON, T., WRIGHT, F. T. and DYKSTRA, R. L. (1988). *Order restricted statistical inference*. Wiley Series in Probability and Mathematical Statistics: Probability and Mathematical Statistics, John Wiley & Sons, Ltd., Chichester.

SHAKED, M. and SHANTHIKUMAR, J. G. (2007). *Stochastic orders*. Springer Series in Statistics, Springer, New York.

WESTLING, T., DOWNES, K. J. and SMALL, D. S. (2023). Nonparametric maximum likelihood estimation under a likelihood ratio order. *Statist. Sinica* **33** in press.

YU, T., LI, P. and QIN, J. (2017). Density estimation in the two-sample problem with likelihood ratio ordering. *Biometrika* **104** 141–152.

A Proofs and technical details

A.1 Proofs for Sections 2 and 3

Lemma A.1. *Let Q_0 and Q_1 be probability distributions on \mathbb{R} such that $Q_0 \leq_{\text{lr}} Q_1$. If we define $Q_t := (1-t)Q_0 + tQ_1$ for $0 < t < 1$, then $Q_s \leq_{\text{lr}} Q_t$ for $0 \leq s < t \leq 1$.*

Proof. By assumption, there exist densities g_0 of Q_0 and g_1 of Q_1 with respect to some dominating measure μ such that g_1/g_0 is isotonic on $\{g_0 + g_1 > 0\}$, and this is equivalent to the property that

$$g_0(y)g_1(x) \leq g_0(x)g_1(y) \quad \text{whenever } x < y.$$

Now, Q_t has density $g_t := (1-t)g_0 + tg_1$ with respect to μ , and elementary algebra reveals that for $0 \leq s < t \leq 1$ and arbitrary $x < y$,

$$g_s(x)g_t(y) - g_s(y)g_t(x) = (t-s)(g_0(x)g_1(y) - g_0(y)g_1(x)) \geq 0,$$

whence $Q_s \leq_{\text{lr}} Q_t$. □

Proof of Lemma 3.1. Let $\mathbf{h} \in [0, \infty)$ satisfy (2.3) and $L(\mathbf{h}) > -\infty$.

As for part (a), it follows from $L(\mathbf{h}) > -\infty$ that $h_{jk} > 0$ whenever $w_{jk} > 0$. We have to show that for arbitrary index pairs $(j_1, k_2), (j_2, k_1)$ with $j_1 \leq j_2$, $k_1 \leq k_2$ and $w_{j_1 k_2}, w_{j_2 k_1} > 0$, also $h_{jk} > 0$ for all $j \in \{j_1, \dots, j_2\}$ and $k \in \{k_1, \dots, k_2\}$.

Since $h_{j_1 k_2}, h_{j_2 k_1} > 0$, it follows from (2.3) that $h_{j_1 k_1}, h_{j_2 k_2} > 0$, too. (If $j_1 = j_2$ or $k_1 = k_2$, this conclusion is trivial.) This type of argument will reappear several times, so we denote it by $A(j_1, j_2, k_1, k_2)$.

Next we show that $h_{j k_1}, h_{j k_2} > 0$ for $j_1 < j < j_2$. Indeed, there exists an index k_* such that $w_{j k_*} > 0$, whence $h_{j k_*} > 0$. If $k_* \leq k_2$, we may conclude from $A(j_1, j, k_*, k_2)$ that $h_{j, k_2} > 0$, and then it follows from $A(j, j_2, k_1, k_2)$ that $h_{j k_1} > 0$. Similarly, if $k_* \geq k_1$, we may conclude from $A(j, j_2, k_1, k_*)$ that $h_{j k_1} > 0$, and then $A(j_1, j, k_1, k_2)$ shows that $h_{j k_2} > 0$.

Analogously, one can show that $h_{j_1 k}, h_{j_2 k} > 0$ for $k_1 < k < k_2$.

Finally, if $j_1 < j < j_2$ and $k_1 < k < k_2$, then we may apply $A(j_1, j, k_1, k)$ or $A(j, j_2, k, k_2)$ to deduce that $h_{jk} > 0$.

As to part (b), since \mathcal{P} contains all pairs (j, k) with $w_{jk} > 0$, we know that $L_{\text{raw}}(\tilde{\mathbf{h}}) = L_{\text{raw}}(\mathbf{h})$, and $n - n\tilde{h}_{++} \geq n - nh_{++}$ with equality if and only if $\tilde{\mathbf{h}} = \mathbf{h}$. This proves the assertions about $L(\tilde{\mathbf{h}})$ and $L(\mathbf{h})$. That $\tilde{\mathbf{h}}$ inherits property (2.3) from \mathbf{h} can be deduced from the fact that for indices $j_1 < j_2$ and $k_1 < k_2$, it follows from $\tilde{h}_{j_1 k_2} \tilde{h}_{j_2 k_1} > 0$, that $(j_1, k_2), (j_2, k_1) \in \mathcal{P}$, so $(j_1, k_1), (j_2, k_2) \in \mathcal{P}$ as well, and $\tilde{h}_{j_1 k_1} \tilde{h}_{j_2 k_2} - \tilde{h}_{j_1 k_2} \tilde{h}_{j_2 k_1}$ is identical to $h_{j_1 k_1} h_{j_2 k_2} - h_{j_1 k_2} h_{j_2 k_1} \geq 0$.

Concerning part (c), we have to show that (3.1) implies (2.3). To this end, let $(j_1, k_2), (j_2, k_1) \in \mathcal{P}$ with $j_1 < j_2$ and $k_1 < k_2$. Since $\{j_1, \dots, j_2\} \times \{k_1, \dots, k_2\} \subset \mathcal{P}$, one can write

$$\frac{h_{j_1-1, k_1} h_{j_2, k_2}}{h_{j_1-1, k_2} h_{j_2, k_1}} = \prod_{k=k_1+1}^{k_2} \frac{h_{j_1-1, k-1} h_{j_2, k}}{h_{j_1-1, k} h_{j_2, k-1}} \geq 1$$

for $j_1 < j \leq j_2$, and

$$\frac{h_{j_1, k_1} h_{j_2, k_2}}{h_{j_1, k_2} h_{j_2, k_1}} = \prod_{j=j_1+1}^{j_2} \frac{h_{j-1, k_1} h_{j, k_2}}{h_{j-1, k_2} h_{j, k_1}} \geq 1,$$

so (2.3) is satisfied as well. \square

Proof of Theorem 3.2. Since f is strictly convex and Θ is convex, f has at most one minimizer in Θ . To prove existence of a minimizer, it suffices to show that

$$f(\boldsymbol{\theta}) \rightarrow \infty \quad \text{as } \boldsymbol{\theta} \in \Theta, \|\boldsymbol{\theta}\| \rightarrow \infty. \quad (\text{A.1})$$

Suppose that (A.1) is false. Then there exists a sequence $(\boldsymbol{\theta}^{(s)})_s$ in Θ such that $\|\boldsymbol{\theta}\| \rightarrow \infty$ but $(f(\boldsymbol{\theta}^{(s)}))_s$ is bounded. With $r_s := \|\boldsymbol{\theta}^{(s)}\|$ and $\mathbf{u}^{(s)} := r_s^{-1} \boldsymbol{\theta}^{(s)}$, we may assume without loss of generality that $\mathbf{u}^{(s)} \rightarrow \mathbf{u}$ as $s \rightarrow \infty$ for some $\mathbf{u} \in \Theta$ with $\|\mathbf{u}\| = 1$. For any fixed $t > 0$ and sufficiently large s , convexity and differentiability of f imply that

$$\begin{aligned} f(\boldsymbol{\theta}^{(s)}) &= f(t\mathbf{u}^{(s)}) + (f(r_s\mathbf{u}^{(s)}) - f(t\mathbf{u}^{(s)})) \\ &\geq f(t\mathbf{u}^{(s)}) + (r_s - t)\partial f(t\mathbf{u}^{(s)})/\partial t. \end{aligned}$$

Since $\lim_{s \rightarrow \infty} f(t\mathbf{u}^{(s)}) = f(t\mathbf{u})$ and $\lim_{s \rightarrow \infty} \partial f(t\mathbf{u}^{(s)})/\partial t = \partial f(t\mathbf{u})/\partial t$, we conclude that

$$\partial f(t\mathbf{u})/\partial t \leq 0 \quad \text{for all } t > 0.$$

But as $t \rightarrow \infty$, the directional derivative $\partial f(t\mathbf{u})/\partial t = \sum_{(j,k) \in \mathcal{P}} (-w_{jk}u_{jk} + u_{jk} \exp(tu_{jk}))$ converges to

$$\begin{cases} \infty & \text{if } u_{jk} > 0 \text{ for some } (j, k) \in \mathcal{P}, \\ -\sum_{(j,k) \in \mathcal{P}} w_{jk}u_{jk} & \text{if } \mathbf{u} \in (-\infty, 0]^{\mathcal{P}}. \end{cases}$$

Consequently, the limiting direction \mathbf{u} lies in $\Theta \cap (-\infty, 0]^{\mathcal{P}}$ and satisfies $u_{jk} = 0$ whenever $w_{jk} > 0$. But as shown below, this implies that $\mathbf{u} = \mathbf{0}$, a contradiction to $\|\mathbf{u}\| = 1$.

The proof of $\mathbf{u} = \mathbf{0}$ is very similar to the proof of Lemma 3.1. If $j_1 \leq j_2$ and $k_1 \leq k_2$ are indices such that $u_{j_1 k_2} = u_{j_2 k_1} = 0$, then it follows from $\mathbf{u} \in (-\infty, 0]^{\mathcal{P}}$ and (3.4) that $u_{j_1 k_1} + u_{j_2 k_2} \geq 0$, whence $u_{j_1 k_1} = u_{j_2 k_2} = 0$. Repeating this argument as in the proof of Lemma 3.1, one can show that for arbitrary $(j_1, k_2), (j_2, k_1) \in \mathcal{P}$ with $j_1 \leq j_2$, $k_1 \leq k_2$, and $w_{j_1 k_2}, w_{j_2 k_1} > 0$, we have $u_{jk} = 0$ for $j_1 \leq j \leq j_2$ and $k_1 \leq k \leq k_2$. By definition of \mathcal{P} , this means that $\mathbf{u} = \mathbf{0}$. \square

Proof of Lemma 3.3. With the linear bijection $T : \mathbb{R}^{\mathcal{P}} \rightarrow \mathbb{R}^{\mathcal{P}}$ and $\tilde{\Theta} = T(\Theta)$, $\tilde{\boldsymbol{\theta}} = T(\boldsymbol{\theta})$, $\tilde{f} = f \circ T^{-1}$, one can show that for arbitrary $\mathbf{x} \in \mathbb{R}^{\mathcal{P}}$ and $\tilde{\mathbf{x}} = T(\mathbf{x})$,

$$\langle \nabla \tilde{f}(\tilde{\boldsymbol{\theta}}), \tilde{\mathbf{x}} - \tilde{\boldsymbol{\theta}} \rangle = \langle \nabla f(\boldsymbol{\theta}), \mathbf{x} - \boldsymbol{\theta} \rangle,$$

so

$$\Psi(\boldsymbol{\theta}) = \arg \min_{\mathbf{x} \in \Theta} (f(\boldsymbol{\theta}) + \langle \nabla f(\boldsymbol{\theta}), \mathbf{x} - \boldsymbol{\theta} \rangle + 2^{-1} \|\mathbf{A}_{\boldsymbol{\theta}}(\mathbf{x}) - \mathbf{A}_{\boldsymbol{\theta}}(\boldsymbol{\theta})\|^2)$$

with

$$\mathbf{A}_\theta(\mathbf{x}) := (\tilde{v}_{jk}(\theta)^{1/2} T_{jk}(\mathbf{x}))_{(j,k) \in \mathcal{P}}$$

and $\tilde{v}_{jk}(\theta) := \partial^2 \tilde{f}(\tilde{\theta}) / \partial \tilde{\theta}_{jk}^2$. It follows from parts (i) and (ii) of Lemma A.2 in Section A.2 that Ψ is continuous on $\mathbb{R}^{\mathcal{P}}$, and that $\delta(\theta) = \langle \nabla f(\theta), \theta - \Psi(\theta) \rangle > 0$ for $\theta \in \Theta \setminus \{\hat{\theta}\}$. Moreover,

$$f(\theta) - f(\hat{\theta}) \leq \max\left(2\delta(\theta), \sqrt{2\delta(\theta)} \|\mathbf{A}_\theta(\theta - \hat{\theta})\|\right).$$

But

$$\|\mathbf{A}_\theta(\mathbf{x})\|^2 \leq \max_{(j,k) \in \mathcal{P}} \tilde{v}_{jk}(\theta) \|T(\mathbf{x})\|^2 \leq 3 \max_{(j,k) \in \mathcal{P}} \tilde{v}_{jk}(\theta) \|\mathbf{x}\|^2,$$

so

$$f(\theta) - f(\hat{\theta}) \leq \max\left(2\delta(\theta), \beta_1(\theta) \sqrt{\delta(\theta)} \|\theta - \hat{\theta}\|\right)$$

with $\beta_1(\theta)$ being the square root of $6 \max_{(j,k) \in \mathcal{P}} \tilde{v}_{jk}(\theta)$. In case of $\Psi = \Psi^{\text{row}}$ and θ being row-wise calibrated, $\beta_1(\theta)^2$ is no larger than $6 \max_{1 \leq j \leq \ell} w_{j+}$, and in case of $\Psi = \Psi^{\text{col}}$ and θ being column-wise calibrated, $\beta_1(\theta)^2 \leq 6 \max_{1 \leq k \leq m} w_{+k}$.

Concerning the lower bound for the maximum of $f(\theta) - f((1-t)\theta' + t\Psi(\theta))$ over all $t \in [0, 1]$, note that for arbitrary $\theta', \theta'' \in \mathbb{R}^{\mathcal{P}}$,

$$\begin{aligned} \frac{d^2}{dt^2} f((1-t)\theta' + t\theta'') &= n \sum_{(j,k) \in \mathcal{P}} \exp((1-t)\theta'_{jk} + t\theta''_{jk})(\theta'_{jk} - \theta''_{jk})^2 \\ &\leq n \max_{(j,k) \in \mathcal{P}} \exp(\max(\theta'_{jk}, \theta''_{jk})) \|\theta' - \theta''\|^2. \end{aligned}$$

Thus part (iii) of Lemma A.2 yields the asserted lower bound with

$$\beta_2(\theta) := 2n \max_{(j,k) \in \mathcal{P}} \exp(\max(\theta_{jk}, \Psi_{jk}(\theta))). \quad \square$$

Proof of Theorem 3.4. It follows from Lemma 3.3 and the construction of the sequence $(\theta^{(s)})_{s \geq 0}$ that

$$f(\theta^{(s)}) - f(\theta^{(s+1)}) \geq \beta(\theta^{(s)})$$

for all $s \geq 0$ with some continuous function $\beta : \Theta \rightarrow [0, \infty)$ such that $\beta > 0$ on $\Theta \setminus \{\hat{\theta}\}$. Note that $f(\theta^{(s)})$ is antitonic in $s \geq 0$, so the sequence $(\theta^{(s)})_{s \geq 0}$ stays in the compact set $R_0 := \{\theta \in \Theta : f(\theta) \leq f(\theta^{(0)})\}$. For each $\theta \in R_0 \setminus \{\hat{\theta}\}$, there exists a $\delta_\theta > 0$ such that the open ball $U(\theta, \delta_\theta)$ with center θ and radius δ_θ satisfies

$$|f - f(\theta)| < \beta(\theta)/3 \quad \text{and} \quad \beta > 2\beta(\theta)/3 \quad \text{on} \quad U(\theta, \delta_\theta).$$

In particular, if $\theta^{(s)} \in U(\theta, \delta_\theta)$ for some $s \geq 0$, then $f(\theta^{(s+1)}) < f(\theta) - \beta(\theta)/3$. Consequently, $\theta^{(s)} \in U(\theta, \delta_\theta)$ for at most one index $s \geq 0$. But for each $\epsilon > 0$, the compact set $\{\theta \in R_0 : \|\theta - \hat{\theta}\| \geq \epsilon\}$ can be covered by finitely many of these balls $U(\theta, \delta_\theta)$. Hence, $\|\theta^{(s)} - \hat{\theta}\| \geq \epsilon$ for at most finitely many indices $s \geq 0$. \square

A.2 Minimizing convex functions via quadratic approximations

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a strictly convex and differentiable function, and let $\Theta \subset \mathbb{R}^d$ be a closed, convex set such that a minimizer

$$\hat{\theta} := \arg \min_{\theta \in \Theta} f(\theta)$$

exists. For $\boldsymbol{\theta}_o \in \Theta$ and some nonsingular matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ consider the quadratic approximation

$$f_o(\mathbf{x}) := f(\boldsymbol{\theta}_o) + \nabla f(\boldsymbol{\theta}_o)^\top (\mathbf{x} - \boldsymbol{\theta}_o) + 2^{-1} \|\mathbf{A}\mathbf{x} - \mathbf{A}\boldsymbol{\theta}_o\|^2$$

of $f(\mathbf{x})$. By construction, $f_o(\boldsymbol{\theta}_o) = f(\boldsymbol{\theta}_o)$ and $\nabla f_o(\boldsymbol{\theta}_o) = \nabla f(\boldsymbol{\theta}_o)$, and there exists a unique minimizer

$$\boldsymbol{\theta}_* := \arg \min_{\boldsymbol{\theta} \in \Theta} f_o(\boldsymbol{\theta}).$$

The next lemma clarifies some connections between $\boldsymbol{\theta}_*$ and $\widehat{\boldsymbol{\theta}}$ in terms of the directional derivative

$$\delta_o := \nabla f(\boldsymbol{\theta}_o)^\top (\boldsymbol{\theta}_o - \boldsymbol{\theta}_*) = -\left. \frac{d}{dt} \right|_{t=0} f(\boldsymbol{\theta}_o + t(\boldsymbol{\theta}_* - \boldsymbol{\theta}_o)).$$

Lemma A.2. (i) The point $\boldsymbol{\theta}_*$ equals $\boldsymbol{\theta}_o$ if and only if $\boldsymbol{\theta}_o = \widehat{\boldsymbol{\theta}}$. Furthermore,

$$2^{-1} \delta_o \leq f_o(\boldsymbol{\theta}_o) - f_o(\boldsymbol{\theta}_*) \leq \delta_o$$

and

$$f(\boldsymbol{\theta}_o) - f(\widehat{\boldsymbol{\theta}}) \leq \nabla f(\boldsymbol{\theta}_o)^\top (\boldsymbol{\theta}_o - \widehat{\boldsymbol{\theta}}) \leq \max\left(2\delta_o, \sqrt{2\delta_o} \|\mathbf{A}\widehat{\boldsymbol{\theta}} - \mathbf{A}\boldsymbol{\theta}_o\|\right).$$

(ii) If f is continuously differentiable, the minimizer $\boldsymbol{\theta}_*$ is a continuous function of $\boldsymbol{\theta}_o \in \Theta$ and \mathbf{A} .

(iii) If f is even twice differentiable such that for some constant $c_o > 0$ and any $t \in [0, 1]$,

$$\frac{d^2}{dt^2} f((1-t)\boldsymbol{\theta}_o + t\boldsymbol{\theta}_*) \leq c_o \|\boldsymbol{\theta}_o - \boldsymbol{\theta}_*\|^2,$$

then in case of $\boldsymbol{\theta}_o \neq \widehat{\boldsymbol{\theta}}$,

$$\max_{t \in [0, 1]} (f(\boldsymbol{\theta}_o) - f((1-t)\boldsymbol{\theta}_o + t\boldsymbol{\theta}_*)) \geq 2^{-1} \min\left(\delta_o, \frac{\delta_o^2}{c_o \|\boldsymbol{\theta}_o - \boldsymbol{\theta}_*\|^2}\right).$$

Proof. By strict convexity of f , $\boldsymbol{\theta}_o = \widehat{\boldsymbol{\theta}}$ if and only if

$$\left. \frac{d}{dt} \right|_{t=0} f(\boldsymbol{\theta}_o + t(\boldsymbol{\theta} - \boldsymbol{\theta}_o)) = \nabla f(\boldsymbol{\theta}_o)^\top (\boldsymbol{\theta} - \boldsymbol{\theta}_o) \geq 0 \quad \text{for all } \boldsymbol{\theta} \in \Theta.$$

But since f_o is strictly convex, too, with $\nabla f_o(\boldsymbol{\theta}_o) = \nabla f(\boldsymbol{\theta}_o)$, the latter displayed condition is also equivalent to $\boldsymbol{\theta}_o = \boldsymbol{\theta}_*$.

Since the asserted inequalities are trivial in case of $\boldsymbol{\theta}_o = \widehat{\boldsymbol{\theta}} = \boldsymbol{\theta}_*$, let us assume in the sequel that $\boldsymbol{\theta}_* \neq \boldsymbol{\theta}_o \neq \widehat{\boldsymbol{\theta}}$. By convexity of f and f_o ,

$$f_o(\boldsymbol{\theta}_o) - f_o(\boldsymbol{\theta}_*) \leq \left. \frac{d}{dt} \right|_{t=1} f_o(\boldsymbol{\theta}_* + t(\boldsymbol{\theta}_o - \boldsymbol{\theta}_*)) = \delta_o$$

and

$$f(\boldsymbol{\theta}_o) - f(\widehat{\boldsymbol{\theta}}) \leq \left. \frac{d}{dt} \right|_{t=1} f(\widehat{\boldsymbol{\theta}} + t(\boldsymbol{\theta}_o - \widehat{\boldsymbol{\theta}})) = \nabla f(\boldsymbol{\theta}_o)^\top (\boldsymbol{\theta}_o - \widehat{\boldsymbol{\theta}}).$$

On the other hand, since $\boldsymbol{\theta}_*$ minimizes f_o over Θ ,

$$0 \leq \left. \frac{d}{dt} \right|_{t=0} f_o(\boldsymbol{\theta}_* + t(\boldsymbol{\theta}_o - \boldsymbol{\theta}_*)) = \nabla f_o(\boldsymbol{\theta}_*)^\top (\boldsymbol{\theta}_o - \boldsymbol{\theta}_*) = \delta_o - \|\mathbf{A}\boldsymbol{\theta}_o - \mathbf{A}\boldsymbol{\theta}_*\|^2,$$

so

$$f_o(\boldsymbol{\theta}_o) - f_o(\boldsymbol{\theta}_*) = \delta_o - 2^{-1} \|\mathbf{A}\boldsymbol{\theta}_o - \mathbf{A}\boldsymbol{\theta}_*\|^2 \geq 2^{-1} \delta_o.$$

Moreover, with $\widehat{\delta} := \nabla f(\boldsymbol{\theta}_o)^\top (\boldsymbol{\theta}_o - \widehat{\boldsymbol{\theta}})$ and $\widehat{\gamma} := \|\mathbf{A}\boldsymbol{\theta}_o - \mathbf{A}\widehat{\boldsymbol{\theta}}\|^2$,

$$\begin{aligned} 2\delta_o &\geq 2(f_o(\boldsymbol{\theta}_o) - f_o(\boldsymbol{\theta}_*)) = 2 \max_{\boldsymbol{\theta} \in \Theta} (f_o(\boldsymbol{\theta}_o) - f_o(\boldsymbol{\theta})) \\ &\geq 2 \max_{t \in [0,1]} (f_o(\boldsymbol{\theta}_o) - f_o((1-t)\boldsymbol{\theta}_o + t\widehat{\boldsymbol{\theta}})) \\ &= \max_{t \in [0,1]} (2t\widehat{\delta} - t^2\widehat{\gamma}) \\ &= 2t_o\widehat{\delta} - t_o^2\widehat{\gamma}, \end{aligned}$$

where $t_o := \min(1, \widehat{\delta}/\widehat{\gamma})$. In case of $\widehat{\delta} \geq \widehat{\gamma}$, we may conclude that $2\delta_o \geq 2\widehat{\delta} - \widehat{\gamma} \geq \widehat{\delta}$, so $\widehat{\delta} \leq 2\delta_o$, and otherwise, $2\delta_o \geq \widehat{\delta}^2/\widehat{\gamma}$, whence $\widehat{\delta} \leq \sqrt{2\delta_o\widehat{\gamma}}$. This proves part (i).

As to part (ii), let $(\boldsymbol{\theta}_o^{(s)})_{s \geq 1}$ be a sequence in Θ with limit $\boldsymbol{\theta}_o$, and let $(\mathbf{A}^{(s)})_{s \geq 1}$ be a sequence of nonsingular matrices in $\mathbb{R}^{d \times d}$ converging to a nonsingular matrix \mathbf{A} . Defining $f_o^{(s)}$ as f_o with $(\boldsymbol{\theta}_o^{(s)}, \mathbf{A}^{(s)})$ in place of $(\boldsymbol{\theta}, \mathbf{A})$, we know that $f_o^{(s)} \rightarrow f_o$ as $s \rightarrow \infty$ uniformly on any bounded subset of \mathbb{R}^d . Consequently, for any fixed $\epsilon > 0$ and $R_\epsilon := \{\boldsymbol{\theta} \in \Theta : \|\boldsymbol{\theta} - \boldsymbol{\theta}_*\| = \epsilon\}$,

$$\gamma_\epsilon^{(s)} := \min_{\boldsymbol{\theta} \in R_\epsilon} f_o^{(s)}(\boldsymbol{\theta}) - f_o^{(s)}(\boldsymbol{\theta}_*) \rightarrow \gamma_\epsilon := \min_{\boldsymbol{\theta} \in R_\epsilon} f_o(\boldsymbol{\theta}) - f_o(\boldsymbol{\theta}_*) > 0$$

as $s \rightarrow \infty$. But as soon as $\gamma_\epsilon^{(s)} > 0$, it follows from convexity of Θ and $f^{(s)}$ that the minimizer $\boldsymbol{\theta}_*^{(s)}$ of $f_o^{(s)}$ satisfies $\|\boldsymbol{\theta}_*^{(s)} - \boldsymbol{\theta}_*\| < \epsilon$.

Part (iii) follows from

$$\begin{aligned} \max_{t \in [0,1]} (f(\boldsymbol{\theta}_o) - f((1-t)\boldsymbol{\theta}_o + t\boldsymbol{\theta}_*)) &= \max_{t \in [0,1]} (f(\boldsymbol{\theta}_o) - f(\boldsymbol{\theta}_o + t(\boldsymbol{\theta}_* - \boldsymbol{\theta}_o))) \\ &\geq \max_{t \in [0,1]} (t\delta_o - 2^{-1}t^2c_o\|\boldsymbol{\theta}_o - \boldsymbol{\theta}_*\|^2) \\ &= t_o\delta_o - 2^{-1}t_o^2c_o\|\boldsymbol{\theta}_o - \boldsymbol{\theta}_*\|^2 \\ &\geq 2^{-1} \min\left(\delta_o, \frac{\delta_o^2}{c_o\|\boldsymbol{\theta}_o - \boldsymbol{\theta}_*\|^2}\right), \end{aligned}$$

where $t_o := \min(1, \delta_o/(c_o\|\boldsymbol{\theta}_o - \boldsymbol{\theta}_*\|^2))$. □

A.3 Technical details for Sections 4

For fixed $\ell_o, n \in \mathbb{N}$, let $(\widehat{G}_x)_{x \in \mathfrak{X}_o}$, $(\check{G}_x)_{x \in \mathfrak{X}_o}$ and $(\widehat{\mathbb{G}}_x)_{x \in \mathfrak{X}_o}$ be estimates of $(G_x)_{x \in \mathfrak{X}_o}$ from a sample $\{(X_i, Y_i)\}_{i=1}^n$ as described in Section 4.2. Then, for all $\tilde{G} \in \{\widehat{G}, \check{G}, \widehat{\mathbb{G}}\}$ and $x \in \mathfrak{X}_o$, the estimate \tilde{G}_x is a step function with jumps in the set $\{y_1, \dots, y_m\}$ of unique observations. For convenience, we further denote $y_0 := 0$, $y_{m+1} := \infty$, and define

$$\tilde{G}_{jk} := \tilde{G}_{x_j}(y_k), \quad 0 \leq k \leq m, \quad \text{and} \quad \tilde{G}_{jm+1} := 1,$$

for all $1 \leq j \leq \ell_o$ and $\tilde{G} \in \{\widehat{G}, \check{G}, \widehat{\mathbb{G}}\}$.

For the remainder of this section, we fix $1 \leq j \leq \ell_o$ and $\tilde{G} \in \{\widehat{G}, \check{G}, \widehat{\mathbb{G}}\}$. Observe that $R_{x_j}(\tilde{G}, G)$ is the sum of the terms

$$R_{x_j}^{(k)}(\tilde{G}, G) = \int_{y_k}^{y_{k+1}} |\tilde{G}_{jk} - G_{x_j}(y)| g_{x_j}(y) dy,$$

defined for $0 \leq k \leq m$, where g_{x_j} is the density of Q_{x_j} with respect to Lebesgue measure. But since

$$\int_\alpha^\beta G_{x_j}(y) g_{x_j}(y) dy = \frac{G_{x_j}(\beta)^2 - G_{x_j}(\alpha)^2}{2},$$

we find that

$$\begin{aligned}
R_{x_j}^{(0)}(\tilde{G}, G) &= G_{j1}^2/2, \\
R_{x_j}^{(k)}(\tilde{G}, G) &= \begin{cases} \rho(\tilde{G}_{jk}, G_{jk+1}) - \rho(\tilde{G}_{jk}, G_{jk}) & \text{if } \tilde{G}_{jk} \geq G_{jk+1}, \\ \rho(\tilde{G}_{jk}, G_{jk}) - \rho(\tilde{G}_{jk}, G_{jk+1}) & \text{if } \tilde{G}_{jk} \leq G_{jk}, \\ \tilde{G}_{jk}^2 - \rho(\tilde{G}_{jk}, G_{jk}) - \rho(\tilde{G}_{jk}, G_{jk+1}) & \text{otherwise,} \end{cases} \\
R_{x_j}^{(m)}(\tilde{G}, G) &= 1/2 - \rho(1, G_{jm}),
\end{aligned}$$

for $1 \leq k < m$, where $\rho(z_1, z_2) := z_1 z_2 - z_2^2/2$.

Similarly, the computation of the CRPS involves the sum of the following integrals

$$S_{x_j}^{(k)}(\tilde{G}, G) := \int_{y_k}^{y_{k+1}} (\tilde{G}_{jk} - G_{x_j}(y))^2 dy,$$

defined for $0 \leq k \leq m$. But integration by parts yields

$$\int_{\alpha}^{\beta} G_{x_j}(y) dy = \beta G_{x_j}(\beta) - \alpha G_{x_j}(\alpha) - c_j(\bar{G}_{x_j}(\beta) - \bar{G}_{x_j}(\alpha))$$

where $c_j := b(x_j)\Gamma(a(x_j) + 1)/\Gamma(a(x_j))$ and \bar{G}_{x_j} denotes the cumulative distribution function of a Gamma distribution with shape $a(x_j) + 1$ and scale $b(x_j)$. In consequence, if we define $\bar{G}_{jk} := \bar{G}_{x_j}(y_k)$ and

$$I_{x_j}^{(k)} := \tilde{G}_{jk}^2(y_{k+1} - y_k) - 2\tilde{G}_{jk}(y_{k+1}G_{jk+1} - y_kG_{jk} - c_j(\bar{G}_{jk+1} - \bar{G}_{jk}))$$

for $1 \leq k < m$, we obtain

$$\sum_{k=0}^m S_{x_j}^{(k)}(\tilde{G}, G) = \int_0^{y_m} G_{x_j}(y)^2 dy + \int_{y_m}^{\infty} (1 - G_{x_j}(y))^2 dy + \sum_{k=1}^{m-1} I_{x_j}^{(k)},$$

where the above two integrals are computed numerically.