

# Localising Strictly Proper Scoring Rules\*

Ramon de Punder<sup>†</sup>    Cees Diks<sup>†</sup>    Roger Laeven<sup>†</sup>    Dick van Dijk<sup>‡</sup>

June 23, 2023

## Abstract

Forecasters are typically not equally interested in all possible realisations of a random variable under scrutiny. Financial risk managers, for instance, usually put relatively more weight on regions of extreme losses. In density forecast comparison, it is common practice to use strictly proper scoring rules to rank a collection of candidate predictive distributions. When focusing on a region of interest, however, weighted scoring rules obtained via conditioning are no longer strictly proper. We develop a general procedure for focusing, i.e., localising, scoring rules in a way that preserves their strict propriety. Our procedure provides a myriad of strictly locally proper scoring rules beyond the censored likelihood score. In particular, the focusing procedure we develop is general enough to handle both univariate and multivariate scoring rules, including the rich class of kernel scores. The one-to-one correspondence between the censored distribution and the original distribution on the region of interest preserves not only strict propriety but also the optimal power properties of the Logarithmic scoring rule. More specifically, our paper generalises the Neyman Pearson lemma, showing that the uniformly most powerful test for a localised version of this lemma’s original hypotheses boils down to a censored likelihood ratio test. Based on a collection of popular scoring rules, including the Logarithmic, Spherical, Quadratic and Continuously Ranked Probability Score (CRPS), Monte Carlo simulations and the results of our empirical illustration align with the intuition that censoring bears, also in general, more desirable power properties than conditioning, especially if the number of expected tail observations is small.

---

\*This working paper is a draft. A previous version titled ‘A General Procedure for Localising Strictly Proper Scoring Rules’ was first published on July 10, 2022.

<sup>†</sup>University of Amsterdam, Tinbergen Institute

<sup>‡</sup>Erasmus University Rotterdam, Tinbergen Institute

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Scoring rules</b>	<b>7</b>
2.1	Regular scoring rules . . . . .	7
2.2	Weighted scoring rules . . . . .	9
<b>3</b>	<b>Localising scoring rules by censoring</b>	<b>11</b>
3.1	Censored scoring rule . . . . .	11
3.2	Generalised censored scoring rule . . . . .	14
3.3	$Z, Q$ -Randomisation . . . . .	15
3.4	Examples . . . . .	16
3.4.1	Semi-local scoring rules . . . . .	16
3.4.2	Distance sensitive scoring rules . . . . .	17
3.5	Localised Neyman–Pearson . . . . .	21
<b>4</b>	<b>Monte Carlo simulation</b>	<b>28</b>
4.1	Size . . . . .	28
4.2	Power . . . . .	30
4.2.1	Laplace tails . . . . .	30
4.2.2	$\mathcal{N}(0, 1)$ versus Student- $t(5)$ : Left-tail . . . . .	32
4.2.3	$\mathcal{N}(0, 1)$ versus Student- $t(5)$ : Centre . . . . .	33
4.2.4	$\mathcal{N}(-0.2, 1)$ versus $\mathcal{N}(-0.2, 1)$ : Centre ( $c = 200$ ) [ $\gamma_F \neq \gamma_P$ ] . . . . .	36
4.2.5	$\mathcal{N}(-1, 1)$ versus $\mathcal{N}(-1, 1)$ : Centre ( $c = 200$ ) [ $\gamma_F \neq \gamma_P$ ] . . . . .	38
<b>5</b>	<b>Empirical Application</b>	<b>40</b>
5.1	Risk management . . . . .	40
5.1.1	Statistical comparison . . . . .	41
5.1.2	Backtesting . . . . .	43
<b>6</b>	<b>Conclusion</b>	<b>47</b>

<b>Appendix</b>	<b>48</b>
A    Proofs . . . . .	48
A.1    Proof Theorem 2 . . . . .	48
A.2    Proof Lemma 1 . . . . .	51
A.3    Proof Corollary 1 . . . . .	51
A.4    Proof Theorem 3 . . . . .	52

# 1 Introduction

Any forecasting application necessitates quantifying the relative performance of different forecasting methods. [Gneiting and Raftery \(2007\)](#) motivated the use of strictly proper scoring rules for this job, which has become the industry standard ([Brehmer and Gneiting, 2020](#); [Patton, 2020](#)). The reason for this is that (strictly) proper scoring rules assign a score to the actual distribution that is (strictly) larger than the score of any other predictive distribution. Although strictly proper scoring rules admit point forecasts (e.g. mean squared error), we concentrate on their use in combination with predictive distributions and densities. Forecasts in the form of predictive distributions have gained interest in many different forecasting fields because they give a complete picture of the stochastic nature of the variable of interest ([Dawid, 1984](#)). At the same time, the specific characteristics of such applications encourage us to zoom in on certain parts of this picture, i.e. to localise the original scoring rule. In this paper, we present a general censoring-based procedure for localising scoring rules that preserves strict propriety. Our framework nests the censored likelihood (csl) scoring rule proposed by [Diks et al. \(2011\)](#) as a special case. We show that the uniformly most powerful test for a localised hypothesis test is based on this strictly locally proper scoring rule.

Motivating examples for local scoring rules can be found in different application areas. In risk management, for example, one is particularly interested in the left tail of the loss distribution, largely driven by regulatory capital requirements, formulated in terms of risk measures such as the Value-at-Risk (VaR) and Expected Shortfall (ES). See e.g. [Diks et al. \(2014\)](#), [Kole et al. \(2017\)](#), [Opschoor et al. \(2017\)](#) and [Diks and Fang \(2020\)](#) for applications. In macroeconomics, policymakers set – whether regulated by law or not – targets for central variables like inflation, nominal GDP and unemployment rates. For such clear targets, it is logical to zoom in on the part of the distribution around the target value. We refer to [Gneiting and Ranjan \(2011\)](#) and [Iacopini et al. \(2022\)](#) (and references therein) for interesting examples in macroeconomics.

The literature on focused scoring rules starts with the weighted likelihood score of

Amisano and Giacomini (2007), which simply multiplies the unweighted logarithmic scoring rule by a weight function. As independently observed by Diks et al. (2011) and Gneiting and Ranjan (2011), this procedure produces improper scoring rules because it favours distributions with more mass assigned to regions with higher weights, independent of the underlying distribution. As proper alternatives, Gneiting and Ranjan (2011) develop the weighted continuously ranked probability scoring (twCRPS) rule, while Diks et al. (2011) propose the conditional (cl) and csl scoring rule. Holzmann and Klar (2017, Theorem 1) observe that the procedure of the cl scoring rule can be generalised to other scoring rules than the logarithmic scoring rule. They propose a general procedure for focusing regular scoring rules that applies the regular scoring rule to a weighted transformation of the original distribution. Their approach differs from ours by the suggested transformation of the original distribution: a conditional vis-à-vis censored distribution. The impact of this difference is that our censoring-based mechanism is the only one guaranteed to deliver strictly locally proper scoring rules. Interestingly, another route leading to the conditioning mechanism of Holzmann and Klar (2017, Theorem 1) is to first generalise the weighted log-likelihood scoring rule proposed by Amisano and Giacomini (2007) and then apply a transformation coined *properisation* by Brehmer and Gneiting (2020, Theorem 1).

Our research also builds on the existing work on strictly proper scoring rules and their associated divergence measures. Although Gneiting and Raftery (2007) are responsible for the formal definition of strict propriety, scoring rules satisfying this property date back to at least the quadratic scoring (QS) rule proposed by Brier (1950). It is useful to know that this research area is dichotomous in the sense that much of the research prior to the rigorous treatment of general probability measures by Gneiting and Raftery (2007) has been conducted relative to discrete distributions on a finite outcome space, while more recent work more often follows the generality of Gneiting and Raftery (2007). For instance, the introduction of the LogS (Good, 1952; Toda, 1963) and spherical scoring (SphS) rule (Roby, 1964; Good, 1971), the initial generalisations of QS and SphS to the  $\text{PowS}_\alpha$  and  $\text{PsSphS}_\alpha$  families, and the axiomatic

characterisations of the LogS,  $\text{PowS}_\alpha$  and  $\text{PsSphS}_\alpha$  rules provided by [Shuford et al. \(1966\)](#), [Savage \(1971\)](#), [Selten \(1998\)](#) and [Jose \(2009\)](#), are all presented in a discrete context. In our analysis, we work with the generalisations of the  $\text{PowS}_\alpha$  and  $\text{PsSphS}_\alpha$  families advocated by [Gneiting and Raftery \(2007\)](#) and [Ovcharov \(2018\)](#).

Moreover, the expected score differences of many scoring rules are recognised as well-known divergence measures, which reduce all together to the class of Bregman divergences ([Bregman, 1967](#)) when solely considering strictly proper scoring rules ([Dawid, 2007](#); [Gneiting and Raftery, 2007](#); [Ovcharov, 2018](#); [Painsky and Wornell, 2019](#)). Consequently, concentrating the score divergences of strictly proper scoring rules excludes all  $f$ -divergences except the Kullback Leibler divergence ([Kullback and Leibler, 1951](#)), which is the unique intersection of the Bregman and  $f$ -divergence families. Due to its favourable properties ([Liese and Vajda, 2006](#)) the Kullback Leibler divergence has become the cornerstone in measuring the discrepancy between densities. For example, it is the divergence that is minimised in the maximum likelihood framework ([Fisher, 1922](#)), which bears optimal properties in the context of testing and estimation. Specifically, the likelihood ratio test is the most powerful test ([Neyman and Pearson, 1933](#)) and maximum likelihood estimators are unbiased estimators reaching the Cramér–Rao lower bound.

Pivotised sample equivalents of the expected score differences are fundamental in hypothesis tests about the relative performance of two candidate predictive distributions. In line with the weighted applications we have in mind, we localise the simple versus simple hypothesis of the Neyman-Pearson lemma into statements about the underlying distribution on the region of interest. By doing so, the hypothesis about the underlying distribution becomes a multiple versus multiple hypothesis, equivalent to the hypothesis studied by [Holzmann and Klar \(2016\)](#). Unlike them, we are still able to derive the uniformly most powerful test for this hypothesis. The test statistic of this test is given by a localised likelihood ratio, where the localisation is performed by censoring, and necessarily not by conditioning.

Power analyses based on localised scoring rules have more frequently been studied

for the [Giacomini and White \(2006\)](#) test ([Diks et al., 2011, 2014](#); [Holzmann and Klar, 2016](#); [Lerch et al., 2017](#)). The null hypothesis of this test entails that the expected score difference between one candidate to the actual distribution is equivalent to the expected score difference between the other candidate and the actual distribution. A great advantage of this test is that all choices underlying the predictor, such as parameter uncertainty, can be seen as an integral part of the candidate, therefore also called prediction methods. For a strictly proper scoring rule, the null implies that both candidates are necessarily misspecified under the null, namely ‘equally misspecified’. Yet, since which distributions are equally off from both candidates is determined by the scoring rule, this means that the null set of the GW test is a function of the selected candidates and the selected scoring rule, complicating comparisons between GW tests based on different scoring rules. To illustrate this interplay, we include a parametric example for which the conditional GW null set coincides with the full parameter space, whereas the censored GW null is a lower-dimensional subspace of the parameter space. We also compare the power properties of the GW test of the censored scoring rules with their conditional counterparts and other commonly used localised scoring rules like the twCRPS of [Gneiting and Raftery \(2007\)](#). In line with [Diks et al. \(2011\)](#), we find that censoring often leads to higher power.

The remainder of this paper is organised as follows. [Section 2](#) describes the fundamental concepts on which the subsequent chapters rely. [Section 3](#) defines the generalised censored scoring rule and includes the assumption under which it is shown to be strictly locally proper. This section also includes a rich collection of examples and a randomisation procedure, called  $Z$ - $Q$ -randomisation, equivalent to the generalised censored scoring rule. It closes with our generalisation of the Neyman-Pearson lemma. [Section 4](#) compares the size and power properties of a test of equal predictive ability of conditional and censored scoring rules. [Section 6](#) concludes.

## 2 Scoring rules

### 2.1 Regular scoring rules

Consider a random variable  $Y : \Omega \rightarrow \mathcal{Y}$  from a complete probability space  $(\Omega, \mathcal{F}, \mathbb{P})$  to the measurable space  $(\mathcal{Y}, \mathcal{G})$ . The goal of a forecaster is to choose a distribution  $F$  from a convex class of candidate distributions  $\mathcal{P}$  on  $(\mathcal{Y}, \mathcal{G})$  that minimises the *score divergence*

$$\mathbb{D}_S(P\|F) := \mathbb{H}_S(P) - \mathbb{E}_P S(F, \cdot)$$

over  $\mathcal{P}$ , where  $\mathbb{H}_S(P) := \mathbb{E}_P S(P, \cdot)$  is the *negative entropy* of  $P$  based on  $S$ . Adhering to [Gneiting and Raftery \(2007\)](#), the selected *scoring rule*  $S$  is restricted to be *strictly proper* to ensure that the forecaster truthfully reports the actual distribution  $P$  as the best candidate from  $\mathcal{P}$ , if  $P \in \mathcal{P}$ . Definitions 1 and 2, adopted from [Holzmann and Klar \(2017\)](#) and [Gneiting and Raftery \(2007\)](#), respectively, formalise both concepts.

**Definition 1** (Scoring rule). *A scoring rule is any extended real-valued ( $\bar{\mathbb{R}} := [-\infty, \infty]$ ) function  $S : \mathcal{P} \times \mathcal{Y} \rightarrow \bar{\mathbb{R}}$  such that  $S(F, \cdot)$  is measurable with respect to  $\mathcal{G}$  and quasi-integrable with respect to all  $P \in \mathcal{P}$ , for all  $F \in \mathcal{P}$ , and for which  $\mathbb{E}_P S(F, \cdot) < \infty$  and  $\mathbb{H}_S(P) \in \mathbb{R}, \forall P, F \in \mathcal{P}$ .*

**Definition 2** ((Strictly) proper scoring rule). *A scoring rule  $S : \mathcal{P} \times \mathcal{Y} \rightarrow \bar{\mathbb{R}}$  is proper relative to  $\mathcal{P}$  if  $\mathbb{D}_S(P\|F) \geq 0, \forall P, F \in \mathcal{P}$ , and strictly proper if, additionally,  $\mathbb{D}_S(P\|F) = 0$  iff  $P = F, \forall P, F \in \mathcal{P}$ .*

If a scoring rule only uses the  $\mu$ -densities  $f \in \mathcal{P}$  of the candidates  $F \in \mathcal{P}$ , it is easier to work with the densities directly, i.e. to define  $S : \mathcal{P} \times \mathcal{Y} \rightarrow \bar{\mathbb{R}}$  and adapt all definitions in this section accordingly. Yet, this is only possible if there exists a  $\sigma$ -finite measure  $\mu$  on  $(\mathcal{Y}, \mathcal{G})$  such that  $F \ll \mu, \forall F \in \mathcal{P}$ . Furthermore, the restrictions on  $S$  in Definition 1 guarantee a meaningful comparison of the expected score of any candidate with the negative entropy of the actual distribution, necessary for identifying (strict) propriety. Since comparisons of candidates  $F \in \mathcal{P}$  are in terms of  $P$ -expectations, the forecaster is,



strictly speaking, only forced to report a member from the P-a.s. equivalence of P when using a strictly proper scoring rule. For clarity, we henceforth suppress technicalities about P-a.s. equivalence. Definition 2 additionally shows that a score divergence is a *divergence measure* (see e.g. Eguchi et al. (1985)) if and only if  $S$  strictly proper. For distributions on  $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$ , where  $\mathcal{B}(\mathcal{Y})$  denotes the Borel  $\sigma$ -algebra on  $\mathcal{Y}$ , the particular form of  $\mathbb{D}_S(P||F)$  makes it a Bregman divergence (Bregman, 1967) under the conditions listed by Ovcharov (2018).

In their review paper, Gneiting and Raftery (2007) provide an abundant list of strictly proper scoring rules, which can be divided into two categories: *local* scoring rules and *distance sensitive* scoring rules (Ehm and Gneiting, 2012). We use the same structure when discussing examples, yet allowing local scoring rules, henceforth called *semi-local*, to also depend on the density via a global norm of the density. Within the class of semi-local rules, we focus on the Logarithmic (LogS) (Good, 1952; Toda, 1963), Quadratic (QS) (Brier, 1950) and (SphS) (Roby, 1964; Good, 1971) scoring rules as well as their generalisations to the Power (PowS $_{\alpha}$ ) and PseudoSpherical (PsSphS $_{\alpha}$ ) families. Our selection of distance-sensitive scoring rules fits into the family of Energy Scores (ES), a subclass of the class of strictly proper scoring rules given by Theorem 5 of Gneiting and Raftery (2007), nesting the real-valued Continuously Ranked Probability Score (CRPS) (Matheson and Winkler, 1976; Hersbach, 2000) as a special case.

The inclusion of the PowS $_{\alpha}$  and PsSphS $_{\alpha}$  families, sharing LogS as a common limiting case for  $\alpha \downarrow 1$ , is partly due to the connection with the expected utility maximisation problems described by Jose et al. (2008). After all, the duality with specific investment problems based on the one-parameter Hyperbolic Absolute Risk Aversion (HARA) utility function family, generated by the absolute risk tolerance function  $\tau_{\alpha}(x) = \beta + \alpha x$ , with  $\beta = 1$  (Merton, 1971, p. 389), gives  $\alpha$  its interpretation as a risk tolerance parameter. Their introduction and axiomatic characterisation are found by Shuford et al. (1966), Savage (1971), Selten (1998) and Jose (2009), though we work with their continuous generalisations provided by Gneiting and Raftery (2007) and Ovcharov (2018).

## 2.2 Weighted scoring rules

In many applications, particular outcomes are of particular importance. To emphasise regions of the outcome space, a forecaster with scoring rule  $S$  is assumed to select a *weight function*  $w \in \mathcal{W}$ , that is, any  $\mathcal{G}$ -measurable map  $w : \mathcal{Y} \rightarrow [0, 1]$ . The forecaster's weight function is zero for outcomes that are of zero interest. Hence, differences in candidates expressed only on  $\{w = 0\} := \{y \in \mathcal{Y} : w(y) = 0\}$  are ideally not accounted for by the scoring rule. Therefore, we restrict the analysis to the class of *localising weighted scoring rules* given by Definition 3, borrowed from [Holzmann and Klar \(2017\)](#).

**Definition 3** (Localising weighted scoring rule). *A weighted scoring rule  $S$ , that is, a map  $S : \mathcal{P} \times \mathcal{Y} \times \mathcal{W} \rightarrow \bar{\mathbb{R}}$  such that  $S_w(\cdot, \cdot)$  is a scoring rule for each  $w \in \mathcal{W}$ , is localising if for any  $P, F \in \mathcal{P}$ ,  $w \in \mathcal{W}$ , it holds that*

$$\forall E \in \mathcal{G} : P(\{w > 0\} \cap E) = F(\{w > 0\} \cap E) \implies S_w(P, y) = S_w(F, y), \forall y \in \mathcal{Y}.$$

Considering the indicator weight function  $w(y) = \mathbb{1}_A(y)$ , taking the value 1 if  $y \in A$ , and 0 otherwise, it is obvious that a localising so-weighted scoring rule cannot be strictly proper. Indeed, any distribution  $\tilde{P}$  equivalent to  $P$  on  $A$  is assigned the same score. Therefore, we instead aim for *strictly locally proper* weighted scoring rules, initially defined by [Holzmann and Klar \(2017\)](#) and included below as Definition 4.

**Definition 4** ((Strictly) locally proper scoring rule). *A weighted scoring rule  $S : \mathcal{P} \times \mathcal{Y} \times \mathcal{W} \rightarrow \bar{\mathbb{R}}$  is locally proper relative to  $(\mathcal{P}, \mathcal{W})$  if it is localising and  $S_w(\cdot, \cdot)$  is proper for each  $w \in \mathcal{W}$ . Furthermore, it is strictly locally proper relative to  $(\mathcal{P}, \mathcal{W})$  if, additionally,*

$$P(\{w > 0\} \cap E) = F(\{w > 0\} \cap E), \forall E \in \mathcal{G} \iff \mathbb{E}_P S_w(P, \cdot) = \mathbb{E}_P S_w(F, \cdot), \forall w \in \mathcal{W}.$$

Before turning to our solution to weighting scoring rules, we first recall two weighting procedures contained in the literature. First, the recipe  $S_w(F, y; w) = w(y)S(F, y)$  proposed by [Amisano and Giacomini \(2007\)](#) for the Logarithmic scoring rule is clearly not strictly locally proper. Indeed, as shown by Example 1 of [Diks et al. \(2011\)](#), it does

not rule out scoring rules that are completely determined by a region where one candidate density dominates the other, yielding a higher expected score for the dominating one, irrespective of the actual distribution.

A second recipe failing to deliver strictly locally proper scoring rules is the *conditional scoring rule*

$$S_w^\sharp(F, y) := w(y)S(F_w^\sharp, y), \quad dF_w^\sharp := \frac{1}{F_w(\mathcal{Y})}dF_w,$$

proposed by [Holzmann and Klar \(2017\)](#). This rule applies the regular scoring rule to the *conditional distribution*  $F_w^\sharp$ , here defined as the *weighted kernel*  $dF_w := w dF$  scaled (sharpened) by a factor  $1/F_w(\mathcal{Y})$ . Again, this scoring rule completely ignores outcomes in  $\{w = 0\}$ . Though now, the distribution is adjusted accordingly, making the scoring rule locally proper. However, since it cannot discriminate between distributions that are proportional to each other on  $\{w > 0\}$ , it is not strictly locally proper ([Holzmann and Klar, 2017](#)). Figure 1a illustrates the potential consequences of this lack of discriminative ability.

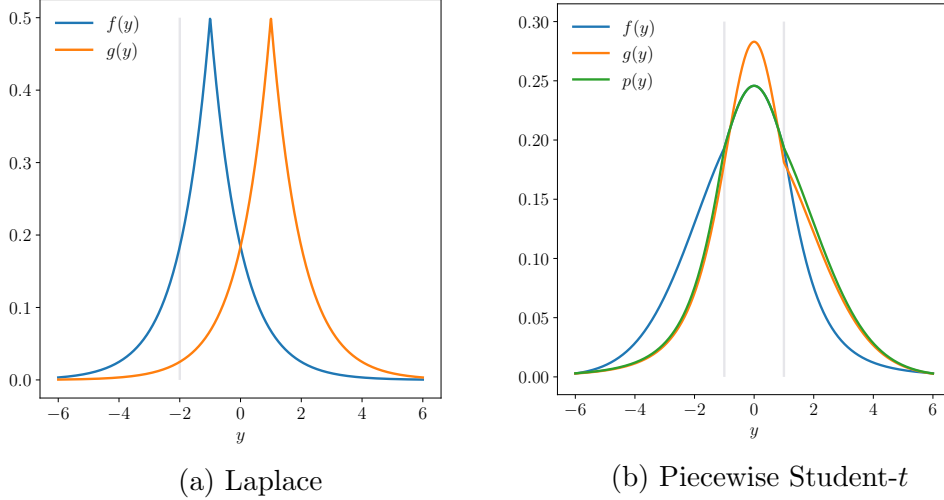
A hint from these examples is that we should not completely forget about  $\{w = 0\}$  when focusing on  $\{w > 0\}$ . Yet, to stay localising, we can only use information about a candidate's distribution on  $\{w = 0\}$  that is implied by the information on  $\{w > 0\}$ . A clear example of a non-localising weighted scoring rule is the twCRPS,

$$\text{twCRPS}(F, y) = \int_{-\infty}^{\infty} w(s) \left( F(s) - \Delta_y(s) \right)^2 ds,$$

for weight functions of the kind  $w(y) = \mathbb{1}_{[r_1, r_2]}(y)$ , where  $r_1 < r_2$  and  $r_1, r_2 \in \mathbb{R}$ . The piecewise Student- $t(\nu_1, \nu_2, \nu_3)$  example displayed in Figure 1b shows that the cumulative character of the CDF destroys the localisation to  $A = [r_1, r_2] = [-1, 1]$ . Here in particular, we selected a distracting candidate G that is similar to P outside A in the sense that  $\nu_{1G} = \nu_{1P} = 3$  and  $\nu_{3G} = \nu_{3P} = 40$ , but different on A, with  $3 = \nu_{2G} < \nu_{2P} = 5$ . In contrast,  $\nu_{2F} = \nu_{2P}$ , while  $\nu_{1F} = \nu_{3P}$  and  $\nu_{1F} = \nu_{3P}$ . The fact that  $\mathbb{D}_{\text{twCRPS}}(P||F) > \mathbb{D}_{\text{twCRPS}}(P||G)$ , while F and P coincide on A, re-

veals that the twCRPS is distracted by the good fit of  $G$  outside  $A$ . Since the bias  $\mathbb{D}_{\text{twCRPS}}(P||F) - \mathbb{D}_{\text{twCRPS}}(P||G) \approx 0.028$  is the sole consequence of the weighted scoring rule being non-localising, it is henceforth referred to as a *localisation bias*.

Figure 1: Non-strictly locally proper scoring rules



Note: (a) The densities  $f$  and  $g$  are both from the Laplace family with common scale parameter  $\theta = 1$ , but different location parameter  $\mu_f = -1$  and  $\mu_g = 1$ . Since Laplace tails are known to be proportional for members of equivalent scale, it follows that  $S_w^\#(f, y) = S_w^\#(g, y)$  on  $A = (-\infty, -2)$ , while  $f \neq g$  on  $A$ .

Therefore,  $S_w^\#$  cannot be strictly locally proper, e.g. consider  $p = f$ . (b)  $f$ ,  $g$  and  $p$  are all piece-wise Student- $t$ , constructed such that  $f = p \neq g$  on  $A = [-1, 1]$ . More specifically, the density  $f(y; \nu_F)$  is the normalisation of the kernel  $\tilde{f}(y; \nu_F) = q(y; \nu_{1F}) \frac{q(-1; \nu_{2F})}{q(-1; \nu_{1F})} \mathbb{1}_{y < -1} + q(y; \nu_{2F}) \mathbb{1}_{-1 \leq y \leq 1} + q(y; \nu_{3F}) \frac{q(1; \nu_{2F})}{q(1; \nu_{3F})} \mathbb{1}_{y > 1}$ . As a result of its non-locality, the twCRPS implies a score divergence indicating  $g$  to be statistically closer to  $p$  on  $A$  than  $f$ . Since  $p = f \neq g$  on  $A$ , the twCRPS is therefore not strictly locally proper.

### 3 Localising scoring rules by censoring

#### 3.1 Censored scoring rule

To overcome issues like the non-strictness and non-locality of the weighted scoring rules discussed above, we propose to use censoring as focusing mechanism. Censoring (Bernoulli, 1760) is a statistical concept that is used in econometrics to model a dependent variable whose value is only partially known (Tobin, 1958). More specifically, for realisations in  $A^c$ , the complement of  $A$ , it is only known that they are not in  $A$ . Events in  $A^c$  are hence indistinguishable after censoring and ‘ $A^c$ ’ could therefore be

viewed as a single outcome of the censored random variable. To avoid confusion, we label observations in  $A^c$  by ‘\*’ rather than ‘ $A^c$ ’ itself, which is nothing but an abstract event for which one can alternatively read ‘NaN’. The censored random variable

$$Y_A^b = \begin{cases} Y, & Y \in A, \\ *, & Y \in A^c, \end{cases}$$

is defined relative to the extended measurable space  $(\mathcal{Y}^*, \mathcal{G}^*)$ , where  $\mathcal{Y}^* = \mathcal{Y} \cup \{*\}$  and  $\mathcal{G}^* = \sigma(\{\mathcal{G}, *\})$ , that is, the smallest  $\sigma$ -algebra containing the collection  $\{\mathcal{G}, *\}$ . Similar to the conditional distribution, we extend the definition of the distribution of  $Y_A^b$  to general weight functions  $w \in \mathcal{W}$ . In particular, we define the *censored distribution* as

$$dF_w^b = dF_w + \bar{F}_w d\delta_*, \quad \bar{F}_w := \int_{\mathcal{Y}} (1 - w) dF, \quad w \in \mathcal{W}, F \in \mathcal{P}, \quad (1)$$

where  $\delta_*$  denotes the Dirac measure at  $*$ , i.e.  $\delta_*(E) = \mathbb{1}_E(*)$ . To make this change of measure well-defined, we consider the original measures  $F \in \mathcal{P}$  relative to the extended measurable space  $(\mathcal{Y}^*, \mathcal{G}^*)$ , by defining  $F(*) = 0$  and taking some value for  $w(*)$ . In case  $F \ll \mu, \forall F \in \mathcal{P}$ , we are invited to work with the  $\mu$ -densities  $f \in \mathcal{P}$  instead, and their associated  $(\mu + \delta_*)$ -densities

$$f_w^b = wf\mathbb{1}_{y \neq *} + \bar{F}_w\mathbb{1}_{y=*}, \quad w \in \mathcal{W}, f \in \mathcal{P}. \quad (2)$$

A detailed proof of this result is deferred to the Online Appendix. Albeit restricted to  $w(y) = \mathbb{1}_A(y)$ , [Borowska et al. \(2020\)](#) also work with an explicit formulation of the censored density, coinciding with  $f_A^b$  in the context of maximum likelihood. Here  $f_A^b$  is preferred notation for  $f_{\mathbb{1}_A}^b$ .

Ideally, the censored scoring rule would be given by  $S_A^b(F, y) = S(F_A^b, y_A^b)$ , as this would fully respect the forecaster’s specific choice of the regular scoring rule  $S$ . The censored scoring rule given by Definition 5 reduces to this definition for the indicator weight function  $w(y) = \mathbb{1}_A(y)$ . The censored scoring rule is also attractive for general

weight functions, but this will be particularly clear from the randomisation perspective taken in Section 3.3, which yields a similar identity for general weight functions; see Equation (4). According to Theorem 1, the censored scoring rule is strictly locally proper. Since Theorem 1 is a corollary of Theorem 2, we have sustainably omitted a proof for this result.

**Definition 5** (Censored scoring rule). *Let  $S : \mathcal{P} \times \mathcal{Y} \rightarrow \bar{\mathbb{R}}$ ,  $\mathcal{P} = \{F_w^b, F \in \mathcal{P}, w \in \mathcal{W}\}$ , denote a scoring rule. Then, the corresponding censored scoring rule is given by the map  $S^b : \mathcal{P} \times \mathcal{Y} \times \mathcal{W} \rightarrow \bar{\mathbb{R}}$ ,*

$$S_w^b(F, y) := w(y)S(F_w^b, y) + (1 - w(y))S(F_w^b, *),$$

where the censored distribution  $F_w^b$  is defined in Equation (1).

**Theorem 1.** *Suppose that the regular scoring rule  $S$  is strictly proper relative to  $\mathcal{P}$ . Then, the censored scoring rule  $S^b$  in Definition 5 is strictly locally proper relative to  $(\mathcal{P}, \mathcal{W})$ .*

The assumption in Theorem 1 ensures that the scoring rule is well-defined with respect to mixed continuous-discrete distributions on  $(\mathcal{Y}^*, \mathcal{G}^*)$ . We will verify that this assumption holds in the examples discussed in Subsection 3.4.

Let us conclude this section by providing some intuition for the result of Theorem 1. Given some weight function  $w \in \mathcal{W}$ , it should be clear that censoring maintains a one-to-one connection with the original distribution on  $\{w > 0\}$ . This relation can be harmed by conditioning due to the additional normalisation of the weighted kernel. This difference is even clearer for indicator weight functions since  $F_A^b = F$ , while  $F_A^\sharp \neq F$ , on  $A$ . Because of this, only the censored scoring rule allows for identifying the original distributions on  $\{w > 0\}$  when comparing two candidates  $F$  and  $G$ . This additionally requires disentanglability of the weighted kernels and discrete probabilities in the censored measures, implied by  $F_w(*) = G_w(*) = 0$ . Consequently, the assumed strict propriety of the original rule localises to  $\{w > 0\}$  for the censored scoring rule.

### 3.2 Generalised censored scoring rule

Given the intuition at the end of the previous section, it is not entirely surprising that one can perform other transformations to the distribution on  $\{w > 0\}$  as long as the transformation is independent of the distribution and traceable when comparing two candidate distributions. The latter requirement is formalised by Assumption 1, under which the *generalised censored scoring rule* in Definition 6 is still strictly locally proper. Appendix A.1 details a proof for this result, summarised by Theorem 2.

**Definition 6** (Generalised censored scoring rule). *Let  $S : \mathcal{P} \times \mathcal{Y} \rightarrow \bar{\mathbb{R}}$  denote a scoring rule. The associated generalised censored scoring rule is given by the map  $S^b : \mathcal{P} \times \mathcal{Y} \times \mathcal{W} \times \mathcal{H} \rightarrow \bar{\mathbb{R}}$ ,*

$$S_{w,H}^b(F, y) = w(y)S(F_{w,H}^b, y) + (1 - w(y))\mathbb{E}_H S(F_{w,H}^b, \cdot), \quad dF_{w,H}^b = dF_w + \bar{F}_w dH,$$

where  $F_{w,H}^b$  is referred to as the *generalised censored distribution* of  $F$ .

**Assumption 1.** *A weight function  $w \in \mathcal{W}$  and nuisance distribution  $H \in \mathcal{H} \subseteq \mathcal{P}$  is such that  $\exists E \in \mathcal{G} : F_w(E) = 0$  and  $H(E) > 0$ ,  $\forall F \in \mathcal{P}, H \in \mathcal{H}$ .*

**Theorem 2.** *Suppose that (i) the regular scoring rule  $S$  in Definition 6 is strictly proper relative to  $\mathcal{P}$ , and (ii)  $\mathcal{W}$  and  $\mathcal{H}$  are such that Assumption 1 is satisfied. Then, the generalised censored scoring rule  $S^b$  in Definition 6 is strictly locally proper relative to  $(\mathcal{P}, \mathcal{W}, \mathcal{H})$ .*

Finally, a corollary of Lemma 3 in the proof of Theorem 2 in Appendix A.1 is that

$$\mathbb{D}_{S_{w,H}^b}(F \| G) = \mathbb{D}_S(F_{w,H}^b \| G_{w,H}^b), \quad (3)$$

i.e. the censored score divergence from  $F$  to  $G$  is the score divergence of the corresponding censored distributions. In particular, this means that we have identified a family of so-called *localised divergence measures*, satisfying the properties of a divergence measure (see Subsection 2.1) on  $\{w > 0\}$ . Indeed, if  $S$  is strictly proper, such that  $\mathbb{D}_S$  is a

divergence measure, it follows that  $\mathbb{D}_{S_{w,H}^b}(F\|G) \geq 0$ , with strict equality if and only if  $F(E \cap \{w > 0\}) = G(E \cap \{w > 0\})$ ,  $\forall E \in \mathcal{G}$ .

### 3.3 $Z, Q$ -Randomisation

The (generalised) censored scoring rule in Definition 5 (6) of the previous section can alternatively be formulated in terms of a randomisation procedure. This procedure relies on an auxiliary random variable  $Z_w$ , indicating, conditional on the realisation  $y$ , whether the observation is censored or not. More specifically, we let

$$y_{Z_w}^b = \varphi(y, Z_w), \quad \varphi(y, Z_w) := \begin{cases} y, & Z_w = 1, \\ *, & Z_w = 0, \end{cases}$$

where  $Z_w|(Y = y) \sim \text{BIN}(1, w(y))$ . By working out the conditional expectation, it is obvious that  $Y_w^b = \mathbb{E}_{Z_w|(Y=Y)}\varphi(Y, Z_w)$  coincides with the specification of the censored random variable in Subsection 3.1. For  $w(y) = \mathbb{1}_A(y)$ , the random variable  $Z_A$  degenerates to being one if  $y \in A$  and zero otherwise, so that  $Y_{Z_A}^b = Y_A^b$  with probability one. Correspondingly, the  $Z$ -randomisation definition of the censored scoring rule reads

$$S_w^b(F, y) = \mathbb{E}_{Z_w|Y=y} S(F_w^b, y_{Z_w}^b), \quad (4)$$

undeniably equivalent to the censored scoring rule defined by Definition 5.

A similar line of reasoning holds for the generalised censored scoring rule. In addition to the auxiliary random variable  $Z_w$ , we introduce an independent random variable  $Q$  with distribution  $H$ . Rather than labelling the observation as censored, we now take a random draw from  $Q$  if  $Z_w = 1$ , i.e. we define

$$y_{H,w}^b := \varphi_{w,H}(y, Z_w, Q), \quad \varphi_{w,H}(y, Z_w, Q) := \begin{cases} Y, & \text{if } Z_w = 1, \\ Q, & \text{if } Z_w = 0. \end{cases}$$

As anticipated, the distribution of  $Y_{H,w}^b = \mathbb{E}_{Z_w|(Y=Y),H}\varphi(Y, Z_w)$  coincides with the



specification of  $F_{w,H}^b$  in Equation (1). Additionally, the generalised censored scoring rule of Definition 6 admits the  $Z, Q$ -randomisation representation

$$S_{H,w}^b(F, y) = \mathbb{E}_{Z_w | (Y=y), H} S(F_{w,H}^b, y_{H,w}^b).$$

The randomisation perspective further clarifies why  $S_{H,w}^b$  generalises  $S_w^b(F, y)$ . Indeed, by choosing a degenerate distribution for  $Q$  at  $*$ , each ‘random draw’ from  $Q$  will be precisely the censoring label  $*$  of the  $Z$ -randomisation procedure. Put differently,  $S_{H,w}^b = S_w^b(F, y)$  for  $H = \delta_*$ .

## 3.4 Examples

### 3.4.1 Semi-local scoring rules

We will now apply our censoring framework to the regular scoring rules introduced in Subsection 2.1. Following the classification into semi-local and distance-sensitive scoring rules, we start with localising the former class. Together with the main characteristics of the LogS, PowS $_\alpha$  and PsSphS $_\alpha$  families, Table 1 presents the localised versions of these families based on conditioning, censoring and generalised censoring. Given the strict propriety classes in Table 1, one can easily verify their strict propriety with respect to  $\ell_\alpha^b$  since  $\|f_w^b\|_\alpha^\alpha \leq 1 + \|f\|_\alpha^\alpha < \infty$ ,  $\forall f \in \mathcal{F}_\alpha$ ,  $\forall w \in \mathcal{W}$ , where  $\alpha = 1$  for LogS. Furthermore, the Bregman generator functions  $\zeta(t)$  refer to the well-known subclass of *separable Bregman divergences*, consisting of the score divergences based on strictly proper scoring rules  $S_\zeta : \mathcal{P}(\mathcal{Y}, \mathcal{G}, \mu) \times \mathcal{Y} \rightarrow \mathbb{R}$  of the form

$$S_\zeta(p, y) = \zeta'(p(y)) - \int_{\mathcal{Y}} \zeta'(p(y))p(y) - \zeta(p(y))\mu(dy).$$

Comparing the censored and conditioned versions of the rules, we notice that the censored variants bear an isolated  $\bar{F}_w$ -dependent term, preserving the coverage probability of  $\{w = 0\}$ . While preserving the likelihood  $\bar{F}_w$  of being censored, Table 1 also shows that the censored scoring rules are independent of  $*$ , the label of a censored

observation. Hence, for this selection of scoring rules, one could alternatively work with an actual number like  $r$  for the location of the residual probability  $\bar{F}_w$ . Strictly speaking, we need to require  $F_w(r) = 0$  in that case, to keep the censored scoring rule strictly locally proper (see Assumption 1), but this is trivially met by restricting to either continuous measures or weight functions satisfying  $w(r) = 0$ , or both. The generalised censored scoring rules in Table 1 show that the invariance with respect to the location of the discrete probability mass holds more generally. In particular, the generalised censored scoring rules turn out to be entirely invariant to the choice of the nuisance density on  $\{w = 0\}$  upon normalisation by the  $\alpha$ -norm of  $h$ , i.e. to the class of densities  $\tilde{h} = h/\|h\|_\alpha$ , where  $\alpha = 1$  for LogS. Since  $\|h\|_1 = 1$ , the latter means that LogS is invariant to the unnormalised choice of  $h$ , as can be seen from Table 1. Finally, Table 1 includes the localised divergence measures  $\mathbb{D}_{S_w^b}$ , which are all localised Bregman divergences since all regular divergences  $\mathbb{D}_S$  in this table are Bregman.

### 3.4.2 Distance sensitive scoring rules

A rich class of distance-sensitive scoring rules is the Energy Score family

$$\text{ES}_\beta(F, y) = \frac{1}{2} \mathbb{E}_F \|\mathbf{Y} - \tilde{\mathbf{Y}}\|_2^\beta - \mathbb{E}_F \|\mathbf{Y} - \mathbf{y}\|_2^\beta, \quad \beta \in (0, 2),$$

known to be strictly proper to the class of Borel probability measures on  $\mathbb{R}^d$  such that  $\mathbb{E}_F \|\mathbf{Y}\|_2^\beta < \infty$  (Gneiting and Raftery, 2007). From this expression, it is immediately clear that the corresponding censored ES family depends, in contrast to the semi-local scoring rules, on  $*$ , or more particularly, the distance  $d(\mathbf{y}) = \|\mathbf{y} - *\|_2$ . Specifically,

$$S_{w,d}^b(F, \mathbf{y}) = \frac{1}{2} \mathbb{E}_{F_w^b} \|\mathbf{Y} - \tilde{\mathbf{Y}}\|_2^\beta - \mathbb{E}_{F_w^b} \left( w(\mathbf{y}) \|\mathbf{Y} - \mathbf{y}\|_2^\beta + (1 - w(\mathbf{y})) d(\mathbf{Y})^\beta \right).$$

Of course, it is not surprising that distance sensitive scoring rules are sensitive to the location of the discrete probability  $\bar{F}_w$ . An easy way to define  $d(\mathbf{y})$  is to simply add the location of  $\bar{F}_w$  by choosing  $* \in \mathbb{R}^d$ . It is important, however, to keep in mind that the censored scoring rule is not invariant with respect to this additional piece of

information. More precisely, the selected value for  $*$ , say  $\mathbf{r}$ , is now not only representing the event of being censored but also the value an observation gets after being censored.

Assuming that the weight function at hand has a finite number of pivotal points  $\mathcal{A} := \{\mathbf{a}_i\}_{i=1}^{n_a}$ , e.g. the edge(s) of an indicator function, the centre of a kernel, etc., we consider the following two choices for the censored distance

$$(i) \quad d_{\text{rand}}(\mathbf{y}) = \left\| \mathbf{y} - \frac{1}{n_a} \sum_{i=1}^{n_a} \mathbf{a}_i \right\|$$

$$(ii) \quad d_{\text{min}}(\mathbf{y}) = \min_{\mathcal{A}} \|\mathbf{y} - \mathbf{a}_i\|$$

The first suggestion is equivalent to taking  $\mathbf{r} = \frac{1}{n_a} \sum_{i=1}^{n_a} \mathbf{a}_i$ , making it straightforward to show that Theorem 1 applies. The second suggestion does not necessarily reduce to a choice for  $\mathbf{r}$ . Though, since this choice is in line with the assumptions of Theorem 1 of Székely and Rizzo (2005), one can still easily verify the assumption of Theorem 1.

We illustrate the role of the censored distance by two concrete examples. For the left-tail indicator function (the same holds for the right-tail indicator), Example 1 shows that both choices for  $d$  lead to a censored scoring rule coinciding with the twCRPS. This is an interesting result since the twCRPS is strictly locally proper for precisely these two types of weight functions (Holzmann and Klar, 2017, Theorem 5). Indeed, for the centre indicator function  $w(y) = \mathbb{1}_{[-r,r]}(y)$ , the twCRPS is knowingly failing to be strictly locally proper since it is non-localising. The corresponding censored scoring rules following from the censored distances in Example 2 are strictly locally proper and hence also still localising. In sharp contrast to the twCRPS, the censoring procedure enforces the weighted scoring rule to be localising by considering anything outside  $[-r, r]$  as the same event. In this way, the censored rules prevent for the localisation bias introduced in Subsection 2.2, illustrated by Figure 1b.

**Example 1.** *Consider the CRPS, i.e. the  $ES_1$  family for  $d = 1$  and take  $w(y) = \mathbb{1}_{(-\infty, r)}(y)$  as weight function. Following the examples of pivotal points of weight functions, we let  $n_a = 1$  and  $a_1 = r$ . The associated censored distances become  $d_{\text{min}}(y) = |y - r|$  and  $d_{\text{rand}}(y) = |y - r|$ . Hence, the choice of both censored distances*

reduce to replacing  $*$  by  $r$ , leading to the scoring rule

$$\begin{aligned}
CRPS_{w,d_{min}}^b(F, y) &= CRPS_{w,d_{rand}}^b(F, y) \\
&= \frac{1}{2} \mathbb{E}_{F_w^b} |Y - \tilde{Y}| - \mathbb{E}_{F_w^b} \left( w(y) |Y - y| + (1 - w(y)) |Y - r| \right) \\
&= twCRPS(F, y).
\end{aligned}$$

**Example 2.** Consider the CRPS and the centre indicator function  $w(y) = \mathbb{1}_{[-r,r]}(y)$ , for which  $a_1 = -r$  and  $a_2 = r$ . The censored distances read  $d_{min}(y) = |y - r| \wedge |y + r|$  and  $d_{rand}(y) = |y|$ . Both censored scoring rules do not coincide with the  $twCRPS$ , at which one arrives if we would put observations below  $-r$  equal to  $-r$  and observations above  $r$  to  $r$ . The latter is clearly not an example of a censored scoring rule, since it uses information outside the region of interest that is not implied by the information within the region of interest. The use of this additional information makes the  $twCRPS$  non-localising and hence prone to the localisation bias illustrated by Figure 1b. The verification of the strict propriety of the CRPS on the extended outcome space  $\mathbb{R}^*$  with censored distance  $d_{min}(y) = |y - r| \wedge |y + r|$  is deferred to the Online Appendix.

Expanding upon the centre indicator example discussed in Example 2, it should be noted that the distances  $d_{min}$  and  $d_{rand}$  sometimes result in distances  $|y - *|$  that are significantly off. For example,  $d_{min}$  can differ greatly from the non-censored distance between an observation  $y \in A$  near  $-r$  and a censored observation located far into the right tail of the distribution before censoring. Although these errors to some extent cancel each other out, this observation can also serve as an inspiration for improvement. In particular, we suggest to alternatively use the generalised censored distribution

$$dF_w^b = dF_w + \bar{F}_w(\gamma d\delta_{a_1} + (1 - \gamma) d\delta_{a_2}), \quad a_1, a_2 \in \mathbb{R}, \gamma \in [0, 1], \quad (5)$$

distributing the residual probability  $\bar{F}_w$  over the pivotal points  $a_1$  and  $a_2$  with proportions  $\gamma$  and  $1 - \gamma$ , respectively.

Unlike the  $twCRPS$ , this measure does not rely on the location of observations

outside the region of interest, except for the fact that they are not in  $A$ . The measure is a generalisation of the censored measure at a single critical point, and can in turn easily be verified to  $n_a$  pivotal points by taking  $dH = \sum_{i=1}^{n_a} \gamma_i d\delta_i$  as reference distribution, with  $\gamma_i$  restricted to the unit simplex  $\Delta(n)$  to maintain  $H$  as a probability measure. The choice of parameter  $\gamma$  depends on the specific application. For the indicator function  $w(y) = \mathbb{1}_{[-r,r]}$  in Example 2, it is appropriate to select  $\gamma = \frac{1}{2}$  when comparing the predictive performance of two candidates that are both symmetric around zero.

The parameter  $\gamma$  depends on the application at hand. For the indicator function  $w(y) = \mathbb{1}_{[-r,r]}$  in Example 2 it makes sense to choose  $\gamma = \frac{1}{2}$  if one aims to compare the predictive ability of two candidates that are both symmetric around zero. In these types of applications, data is typically available to estimate the proportion of residual probabilities of the candidates based on the DGP. It is important to note that using the data instead of the candidates to estimate  $\gamma$ , sets a level playing field for the candidates in terms of their performance on  $A$ . After all, this approach ensures that the relative performance of the candidates on  $A$  is not obscured by the performance outside  $A$  (for the part that is not entirely implied by the distribution on  $A$ ).

Mathematically, we can illustrate the difference between the generalised censored scoring rule based on the censored measure in Equation (5) and the twCRPS as follows. For the centre indicator function  $w(y) = \mathbb{1}_A(y)$ , where  $A = [a_1, a_2]$ , we have the following equality

$$\text{twCRPS}(F, y) = \text{CRPS}(F_w^\dagger, y_w^\dagger), \quad dF_w^\dagger = dF_w + \bar{F}_w(\gamma_F d\delta_{a_1} + (1 - \gamma_F) d\delta_{a_2})$$

where  $\gamma_F = F_{wL}/\bar{F}_w$ ,  $F_{wL} = F(A_L)$ ,  $A_L^c = (-\infty, a_1)$ . Furthermore,  $y_w^\dagger = y\mathbb{1}_A(y) + a_1\mathbb{1}_{A_L^c}(y) + a_2\mathbb{1}_{A_R^c}(y)$ , with  $A_R^c = (a_2, \infty)$ , allowing the twCRPS to assign different scores to observations in  $A_L^c$  and  $A_R^c$ . One crucial distinction between the generalised censored measure and  $F_w^\dagger$  is that the latter candidate's reference distribution depends on the candidate itself, namely through the dependence of the proportion parameter on  $F$ . In expectation, the difference between the twCRPS and the generalised censored

scoring rule reduces to precisely this difference between  $\gamma = P_{wL}/\bar{P}_w$ , where  $P$  denotes the underlying distribution of  $Y$ , and  $\gamma_F$ . Specifically,

$$\mathbb{E}_{P_{\text{twCRPS}}}(\mathbf{F}, Y) = \mathbb{E}_P \text{CRPS}_w^\dagger(\mathbf{F}, Y),$$

where the only difference between  $\text{CRPS}_w^\dagger$  and  $\text{CRPS}_w^\flat$  is the dependence on  $\mathbf{F}_w^\dagger$  rather than  $\mathbf{F}_w^\flat$ , i.e.

$$\text{CRPS}_w^\dagger(\mathbf{F}, y) = \begin{cases} \text{CRPS}(\mathbf{F}_w^\dagger, y), & \text{if } y \in A \\ \gamma \text{CRPS}(\mathbf{F}_w^\dagger, a_1) + (1 - \gamma)(\mathbf{F}_w^\dagger, a_2), & \text{if } y \in A^c, \end{cases}$$

which, contrary to the twCRPS, does not depend on whether an observation is in  $A_L$  or  $A_R$ .

For centre indicator case, for which the twCRPS is not strictly locally proper and hence not a generalised censored scoring rule, we have now derived the alternative (close to censoring) procedure, which is helpful in two ways. (i) By revealing the recipe for obtaining the twCRPS, we uncovered the multivariate twCRPS for practioners that are despite the localisation-bias still willing to use the twCRPS in a multivariate setting. (ii) We have uncovered precisely the difference between the twCRPS and the generalised censored scoring rule, i.e.  $\gamma$  versus  $\gamma_F$  in the definition of the focused measure.

### 3.5 Localised Neyman–Pearson

In anticipation of our favourite applications, we now switch to an explicit time series context. In particular, consider a stochastic process  $\{Y_t : \Omega \rightarrow \mathcal{Y}\}_{t=1}^T$  from a complete probability space  $(\Omega, \mathcal{F}, \mathbb{P})$  to a measurable space  $(\mathcal{Y}^T, \mathcal{G}^T)$ , where  $\mathcal{Y}^T$  and  $\mathcal{G}^T$  denote the product outcome space and  $\sigma$ -algebra of the individual outcome spaces  $\mathcal{Y}$  and  $\sigma$ -algebras  $\mathcal{G}$ , respectively. The process generates the filtration  $\{\mathcal{F}_t\}_{t=1}^T$ , in which  $\mathcal{F}_t = \sigma(Y_1, \dots, Y_t)$  is the information set at time  $t$ , satisfying  $\mathcal{F}_t \subseteq \mathcal{F}_{t+1} \subseteq \mathcal{F}$ ,  $\forall t$ . We denote

predictive distributions of  $Y_{t+1}$  based on  $\mathcal{F}_t$  by  $F_t$ , predictive distribution functions by  $F_t$  and predictive  $\mu_t$ -densities by  $f_t$ . The existence of the sequence of densities  $f_t$  is implied by the existence of a sequence of measures  $\{\mu_t\}$  such that  $F_t \ll \mu_t, \forall t$ . Furthermore, the regions of interest  $A_t \subseteq \mathcal{Y}$  are always assumed to be  $\mathcal{F}_t$ -measurable.

The aim of this section is to derive a uniformly most powerful (UMP) test for the following null and alternative hypothesis

$$\mathbb{H}_0 : p_{0t} \mathbb{1}_{A_t} = f_{0t} \mathbb{1}_{A_t}, \quad \forall t \quad \text{vs} \quad \mathbb{H}_1 : p_{1t} \mathbb{1}_{A_t} = f_{1t} \mathbb{1}_{A_t}, \quad \forall t. \quad (6)$$

Although the predictive densities  $f_{jt} = \frac{F_{jt}}{d\mu_t}$ ,  $j \in \{0, 1\}$ , are assumed to be known, the testing problem remains a multiple versus multiple hypothesis test due to the lacking specification of the density outside the regions of interest  $A_t$ . Yet, since the densities  $p_{jt}$  must integrate to one on  $A_t \cup A_t^c$ , the null hypothesis does imply that these densities integrate to  $F_{jt}(A_t^c)$  on  $A_t^c$ . Therefore, the implied specification on  $A_t^c$  can be summarised as

$$\frac{F_{jt}(A^c)}{H_{jt}(A^c)} h_{jt} \mathbb{1}_{A_t^c} = F_{jt}(A^c) [h_{jt}]_{A_t^c}^\# \mathbb{1}_{A_t^c}, \quad j \in \{0, 1\},$$

where the unknown densities  $h_{jt} = \frac{H_{jt}}{d\mu_t}$  can be seen as infinite dimensional nuisance parameters.

Explicitising the implied assumption on  $A_t^c$  in the hypotheses and phrasing them in terms of a statement about the whole sample distribution leads to the following equivalent hypotheses

$$\mathbb{H}_j : p_j(\mathbf{y}) = \prod_{t=0}^{T-1} \left( f_{jt}(y_{t+1}) \mathbb{1}_{A_t}(y_{t+1}) + F_{jt}(A^c) [h_{jt}]_{A_t^c}^\#(y_{t+1}) \mathbb{1}_{A_t^c}(y_{t+1}) \right), \quad j \in \{0, 1\}.$$

Since the densities  $f_{jt}$  are fixed, and the densities  $h_{jt}$  are unrestricted under both hypothesis, the class of densities satisfying hypothesis  $\mathbb{H}_j$  can alternatively be written

as

$$\mathcal{P}_j = \left\{ \prod_{t=0}^{T-1} \left( f_j(y_{t+1}) \mathbb{1}_{A_t}(y_{t+1}) + F_{jt}(A^c) [h_{jt}]_{A_t^c}^\#(y_{t+1}) \mathbb{1}_{A_t^c}(y_{t+1}) \right), h_j \in \mathcal{H} \right\}, \quad j \in \{0, 1\},$$

in which  $\mathcal{H}$  denotes the space of all densities on  $A^c = \prod_{t=0}^{T-1} A_t^c$ .

Let  $\phi : \mathcal{Y}^T \rightarrow [0, 1]$  denote a test function determining which values should be included in the critical region. In terms of the index set of all observations  $\mathcal{I} = \{1, \dots, T\}$ , this space can also be denoted as  $\mathcal{Y}(\mathcal{I}) = \prod_{t \in \mathcal{I}} \mathcal{Y}_t$ . The aim of this section is to find a uniformly most powerful (UMP) test  $\phi^*$  of size  $\alpha$  for testing problem (6), i.e. a solution to the maximisation problem

$$\max_{\phi \in \Phi(\alpha)} \mathbb{E}_{p_1} \phi, \quad \Phi(\alpha) = \{\phi : \sup_{p_0 \in \mathcal{P}_0} \mathbb{E}_{p_0} \phi \leq \alpha\}. \quad (7)$$

As a first step toward the solution, given by Theorem 3, let us fix an  $h_1 \in \mathcal{H}$  so that the distribution under the alternative is completely known. Given the fact that the hypotheses are, in the end, silent about the shape of the density on  $A^c$ , we conjecture that a UMP test neglects the information about the shape of the density on  $A^c$ . If  $T = 2$ , for example, and we consider the optimal test on  $A_1 \times A_2^c$ , our intuition is that an optimal test does not care about the shape of  $[h_2]_{A_2^c}^\#$ , that is, the specific values  $[h_2]_{A_2^c}^\#(y_2)$  for all  $y_2 \in A_2^c$ , but just about the total probability of an outcome falling into  $A_2^c$ . In other words, we expect that a solution to problem (7) has integrated out the dependence on the nuisance densities.

Although it is obvious that marginalising out the still assumed to be fixed density  $h_1 \in \mathcal{H}$  is harmless in terms of power, it is non-trivial that this is an affordable strategy in terms of size for all  $h_0 \in \mathcal{H}$ . Lemma 1 and its proof show that the subclass of tests disregarding information about the shape of  $h_1$  is guaranteed to be size correct. In our search for the UMP test, Corollary 1 then formalises the idea that we can restrict our attention to tests of the conjectured kind.

**Lemma 1.** *Consider testing problem (6) and suppose that the outcomes  $(y_t)_{t \in \mathcal{I}_A}$  are*



in  $A_t$ , and the remaining  $n - k$ , with  $k = |\mathcal{I}_A|$ , observations  $(y_t)_{t \in \mathcal{I}_{A^c}}$  in  $A_t^c$ . For an arbitrary but fixed density  $h_1 \in \mathcal{H}$ , the test

$$\psi_{h_1} : \mathcal{Y}^T \rightarrow [0, 1], \quad \psi_{h_1} = \int_{\mathcal{Y}(\mathcal{I}_{A^c})} \phi_{h_1}^* \prod_{t \in \mathcal{I}_{A^c}} [h_{1t}]_{A_t^c}^\# \mathbb{1}_{A_t^c} d\mu_t$$

where  $\phi_{h_1}^*$  denotes a solution to problem (7), is such that  $\psi_{h_1} \in \Phi(\alpha)$ .

**Corollary 1.** Consider testing problem (6) and suppose that the outcomes  $(y_t)_{t \in \mathcal{I}_A}$  are in  $A_t$ , and the remaining  $T - k$ , with  $k = |\mathcal{I}_A|$ , observations  $(y_t)_{t \in \mathcal{I}_{A^c}}$  in  $A_t^c$ . Let  $\Psi(\alpha) \subseteq \Phi(\alpha)$  denote the class of size  $\alpha$  tests on  $\mathcal{Y}^T$  that are constant in arguments varying in  $\mathcal{Y}(\mathcal{I}_{A^c})$ . Then,

$$\max_{\phi \in \Phi(\alpha)} \mathbb{E}_{p_1} \phi = \max_{\psi \in \Psi(\alpha)} \mathbb{E}_{p_1} \psi, \quad \forall h_1 \in \mathcal{H}.$$

For any fixed  $h_1 \in \mathcal{H}$ , the reduced optimisation problem resulting from Corollary 1, simplifies to a simple versus simple hypothesis in terms of the censored measures  $d[F_{jt}]_{A_t}^b = \mathbb{1}_{A_t} dF_{jt} + F_{jt}(A_t^c) d\delta_*$ , enabling us to formalise a localised version of the Fundamental Lemma of [Neyman and Pearson \(1933\)](#), included below as Theorem 3.

**Theorem 3** (Localised Neyman-Pearson). *The UMP test for testing problem (6) is given by*

$$\phi_A^b(\mathbf{y}) = \begin{cases} 1, & \text{if } \lambda(\mathbf{y}) > c \\ \gamma & \text{if } \lambda(\mathbf{y}) = c \\ 0, & \text{if } \lambda(\mathbf{y}) < c \end{cases} \quad \lambda(\mathbf{y}) = \frac{[f_1]_A^b(\mathbf{y})}{[f_0]_A^b(\mathbf{y})}, \quad [f_j]_A^b(\mathbf{y}) = \prod_{t=0}^{T-1} [f_{jt}]_{A_t}^b(y_{t+1}),$$

where  $j \in \{0, 1\}$  and  $c$  is the largest constant such that  $[F_0]_A^b(\lambda(\mathbf{y}) \geq c) \geq \alpha$  and  $[F_0]_A^b(\lambda(\mathbf{y}) \leq c) \geq 1 - \alpha$ , and  $\gamma \in [0, 1]$  is such that  $\alpha = [F_0]_A^b(\lambda(\mathbf{y}) > c) + \gamma[F_0]_A^b(\lambda(\mathbf{y}) = c)$ .

It is worth emphasising that the obtained equivalence between testing problem (6)

and  $\mathbb{H}_j : p_j = [f_j]_{A_t}^b, j \in \{0, 1\}$ , is a priori unobvious, since

$$p_j = \left( f_j \mathbb{1}_A + F_j(A^c)[h_j]_{A^c}^\# \mathbb{1}_{A^c} \right) \implies p_j = [f_j]_A^b,$$

but not the other way around. Formulated differently,

$$\mathbb{H}_j : [p_j]_A^b = [f_j]_A^b$$

is a multiple versus multiple hypothesis about  $p_j$  (for example satisfied if  $p_j = [f_j]_A^b$ ), but a simple versus simple hypothesis about  $[p_j]_A^b$ . Furthermore, we have included an example for the special case that  $T = 1$  in the Online Appendix, showing that we arrive at the same solution as [Holzmann and Klar \(2016\)](#) for this special case.

We close this section with two corollaries of Theorem 3, the proofs of which are deferred to the Online Appendix. Corollary 2 reveals that, unsurprisingly, the localised NP test given by Theorem 3 can alternatively be formulated by the censored likelihood score of [Diks et al. \(2011\)](#). Corollary 3 ensures that the conditional operator does not bear a UMP test too, making the censored operator strictly preferred over the conditional one in the current setting.

**Corollary 2.** *Another formulation of the UMP test for testing problem (6) is given by the test defined in Theorem 3, with  $\lambda(\mathbf{y})$  replaced by  $\tilde{\lambda}(\mathbf{y}) = \sum_{t=0}^{T-1} (S_{A_t}^{csl}(f_{1t}, y_{t+1}) - S_{A_t}^{csl}(f_{0t}, y_{t+1}))$ , where  $S_{A_t}^{csl}$  denotes the censored likelihood score (csl) proposed by [Diks et al. \(2011\)](#).*

**Corollary 3.** *For testing problem (6), the test*

$$\phi_A^\#(\mathbf{y}) = \begin{cases} 1, & \text{if } \lambda^\#(\mathbf{y}) > c \\ \gamma & \text{if } \lambda^\#(\mathbf{y}) = c \\ 0, & \text{if } \lambda^\#(\mathbf{y}) < c \end{cases} \quad \lambda_A^\#(\mathbf{y}) = \frac{[f_1]_A^\#(\mathbf{y})}{[f_0]_A^\#(\mathbf{y})} \mathbb{1}_A(\mathbf{y}), \quad [f_j]_A^\#(\mathbf{y}) = \prod_{t=1}^T [f_{jt}]_{A_t}^\#(y_t),$$

where  $j \in \{0, 1\}$  and  $c$  is the largest constant such that  $[F_0]_A^b(\lambda(\mathbf{y}) \geq c) \geq \alpha$  and

$[F_0]_A^b(\lambda(\mathbf{y}) \leq c) \geq 1 - \alpha$ , and  $\gamma \in [0, 1]$  is such that  $\alpha = [F_0]_A^b(\lambda(\mathbf{y}) > c) + \gamma[F_0]_A^b(\lambda(\mathbf{y}) = c)$ , is not UMP.

Table 1: Examples semi-local scoring rules

Name	Logarithmic	Power family	PseudoSpherical family
		Regular	
$S(f, y)$	$\text{LogS}(f, y) = \log f(y)$	$\text{PowS}_\alpha(f, y) = \alpha f(y)^{\alpha-1} - (\alpha-1)\ f\ _\alpha^\alpha$	$\text{PsSphS}_\alpha(f, y) = \frac{f(y)^{\alpha-1}}{\ f\ _\alpha^{\alpha-1}}$
Special cases		$\text{QS}(f, y) = \text{PowS}_2(f, y)$	$\text{SphS}(f, y) = \text{PsSphS}_2(f, y)$
		$\text{LogS}(f, y) = \lim_{\alpha \downarrow 1} \text{PowS}_\alpha(f, y)$	$\text{LogS}(f, y) = \lim_{\alpha \downarrow 1} \text{PsSphS}_\alpha(f, y)$
$\mathbb{H}_S(f)$	$\mathbb{E}_f \log f$	$\ f\ _\alpha^\alpha$	$\ f\ _\alpha$
$\mathbb{D}_S(f\ g)$	$\text{KL}(f\ g) = \mathbb{E}_f \log \left( \frac{f}{g} \right)$	$\ f\ _\alpha^\alpha - \alpha \int g^{\alpha-1}(f-g) d\mu - \ g\ _\alpha^\alpha$	$\ f\ _\alpha - \frac{\int f g^{\alpha-1} d\mu}{\ g\ _\alpha^{\alpha-1}}$
$\alpha = 2$	–	$\ f - g\ _2^2$	$\ f\ _2(1 - C(f, g))$
SP class	$\hat{f}_1$	$\hat{f}_\alpha$	$\hat{f}_\alpha$
$\zeta(t)$	$t \log t$	$t^\alpha$	–
$S(\tilde{f}, \tilde{y})$	$\log f(y) - \log  b $	$\left(\frac{1}{ b }\right)^{\alpha-1} \text{PowS}_\alpha(f, y)$	$\left(\frac{1}{ b }\right)^{\frac{\alpha-1}{\alpha}} \text{PsSphS}_\alpha(f, y)$
		Focused	
$S_w^\#(f, y)$	$w(y) \log \left( \frac{f(y)}{1 - \bar{F}_w} \right)$	$w(y) \left( \alpha \left( \frac{f_w(y)}{1 - \bar{F}_w} \right)^{\alpha-1} - (\alpha-1) \left\  \frac{f_w(y)}{1 - \bar{F}_w} \right\ _\alpha^\alpha \right)$	$w(y) \frac{f_w(y)^{\alpha-1}}{\ f_w\ _\alpha^{\alpha-1}}$
$S_w^b(f, y)$	$w(y) \log f(y) + (1 - w(y)) \log \bar{F}_w$	$w(y) \alpha f_w(y)^{\alpha-1} + (1 - w(y)) \alpha \bar{F}_w^{\alpha-1} - (\alpha-1) (\ f_w\ _\alpha^\alpha + \bar{F}_w^\alpha)$	$\frac{w(y) f_w(y)^{\alpha-1} + (1 - w(y)) \bar{F}_w^{\alpha-1}}{(\ f_w\ _\alpha^\alpha + \bar{F}_w^\alpha)^{\frac{\alpha-1}{\alpha}}}$
$S_{w,h}^b(f, y)$	$w(y) \log f(y) + (1 - w(y)) \log \bar{F}_w$	$w(y) \alpha f_w(y)^{\alpha-1} + (1 - w(y)) \alpha \bar{F}_w^{\alpha-1} \ h\ _\alpha^\alpha - (\alpha-1) (\ f_w\ _\alpha^\alpha + \bar{F}_w^\alpha \ h\ _\alpha^\alpha)$	$\frac{w(y) f_w(y)^{\alpha-1} + (1 - w(y)) \bar{F}_w^{\alpha-1} \ h\ _\alpha^\alpha}{(\ f_w\ _\alpha^\alpha + \bar{F}_w^\alpha \ h\ _\alpha^\alpha)^{\frac{\alpha-1}{\alpha}}}$
$\mathbb{H}_{S_w}^b(f)$	$\int \log(f) f_w d\mu + \log(\bar{F}_w) \bar{F}_w + \int \log(w) f_w d\mu$	$\ f_w\ _\alpha^\alpha + \bar{F}_w^\alpha$	$(\ f_w\ _\alpha^\alpha + \bar{F}_w^\alpha)^{\frac{1}{\alpha}}$
$\mathbb{D}_{S_w^b}(f\ g)$	$\int \log \left( \frac{f}{g} \right) f_w d\mu + \log \left( \frac{\bar{F}_w}{\bar{G}_w} \right) \bar{F}_w$	$\ f_w\ _\alpha^\alpha + \bar{F}_w^\alpha - \int f_w g_w^{\alpha-1} d\mu - \bar{F}_w \bar{G}_w^{\alpha-1} - (\alpha-1) (\ g_w\ _\alpha^\alpha + \bar{G}_w^\alpha)$	$(\ f_w\ _\alpha^\alpha + \bar{F}_w^\alpha)^{\frac{1}{\alpha}} - \frac{\int f_w g_w^{\alpha-1} d\mu + \bar{F}_w \bar{G}_w^{\alpha-1}}{(\ g_w\ _\alpha^\alpha + \bar{G}_w^\alpha)^{\frac{\alpha-1}{\alpha}}}$

Note:  $C(f, g) = \frac{\int f g d\mu}{\sqrt{\int f^2 d\mu \int g^2 d\mu}}$  and  $\mathbb{D}_S(f\|g)$  denote the cosine similarity and the score divergence between  $f$  and  $g$ , respectively.  $\mathbb{H}_S(f)$  denotes the negative entropy and the  $\text{PowS}_\alpha$  and  $\text{PsSphS}_\alpha$  families are restricted to  $\alpha > 1$ . Furthermore,  $\hat{f}_\alpha$  denotes the space for which the  $L^\alpha$ -norm is finite, where  $\mu$  is measure relative to which the densities  $p$  and  $f$  are defined, i.e.  $\frac{dF}{d\mu}$ , with  $F \ll \mu$ . The common limiting case of the  $\text{PowS}_\alpha$  and  $\text{PsSphS}_\alpha$  remains to hold for conditioning and censoring, i.e.  $\lim_{\alpha \downarrow 1} \frac{1}{\alpha-1} \text{PowS}_{\alpha,w}^x(f, y) = \lim_{\alpha \downarrow 1} \frac{1}{\alpha-1} \text{PowS}_{\alpha,w}^x(f, y) = \text{LogS}^x(f, y)$ ,  $x \in \{\#, b\}$ . The generalised censored distribution  $S_{w,h}^b$  departs from a density  $h$  which support is a subset of  $\{w = 0\} \subseteq \mathcal{Y}$ . The weight function is restricted accordingly.

## 4 Monte Carlo simulation

In this section, we compare the size and power properties of the conditional and censored scoring rules of a selection of regular scoring rules. The selected simulation design is similar to the ones conducted by [Diks et al. \(2011\)](#), [Lerch et al. \(2017\)](#) and [Holzmann and Klar \(2016\)](#). In particular, we use the score difference series of two candidates F and G, that is, realisations of  $D = S(F, Y) - S(G, Y)$  to employ the [Giacomini and White \(2006\)](#) test, for the null hypothesis

$$\mathbb{H}_0 : \mathbb{E}_P S(F, Y) = \mathbb{E}_P S(G, Y),$$

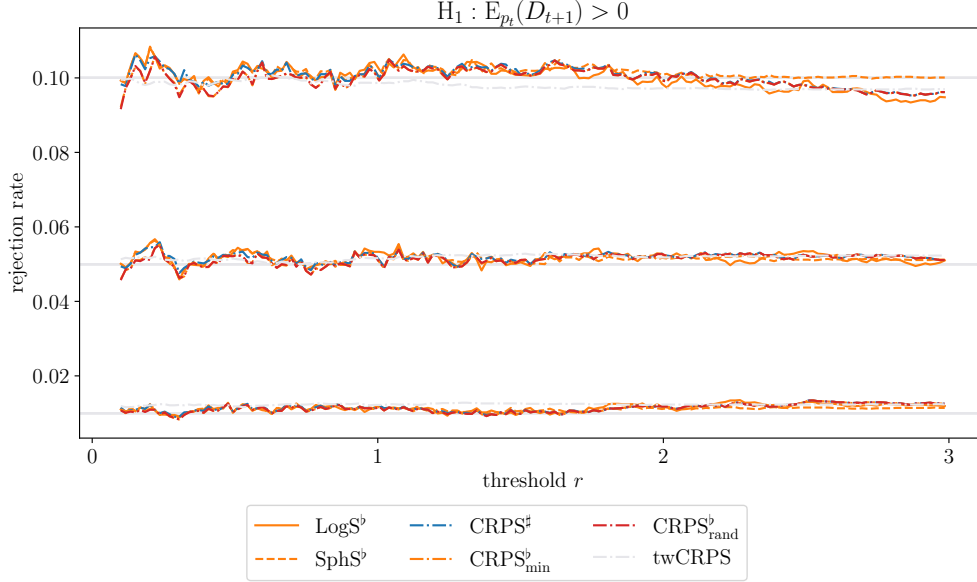
by means of the Diebold Mariano-type statistic  $t_T = \frac{\frac{1}{T} \sum_{t=1}^T d_t}{\sqrt{\hat{\sigma}_t^2/T}}$ , where  $\hat{\sigma}_t$  should be a heteroskedasticity and autocorrelation-consistent (HAC) variance estimator in non-i.i.d. settings. This null hypothesis, which is equivalent to  $\mathbb{H}_0 : \mathbb{D}_S(P \| F) = \mathbb{D}_S(F \| G)$ , is rejected if it is unlikely enough that quoting F instead of P leads to the same information loss as quoting G instead of P.

A natural conjecture is that strictly locally proper scoring rules generally lead to higher power since they are sensible with respect to all measurable aspects of the distribution. When comparing conditional and censored scoring rules, we thus expect the censored scoring rules to have more power. Though, it is worth mentioning that the simulation study of [Diks et al. \(2011\)](#) already presents some examples in which the conditional scoring rule indicates higher power. This can be explained by understanding that some candidates do not have aspects for which the conditional scoring rule has a blind spot. In such cases, the conditioning transformation can either enhance or alleviate the discriminative ability of the scoring rule when compared to censoring.

### 4.1 Size

As carefully explained by [Diks et al. \(2011\)](#), the null hypothesis of the GW test forces a particularly symmetric design. We adopt the design of [Diks et al. \(2011\)](#), using a centre-

Figure 2: Size properties GW test



indicator weight function  $w(y) = \mathbb{1}_{[-r,r]}(y)$  combined with and i.i.d. standard normal DGP and normal candidates with unit variance and means  $\mu_f = -0.2$  and  $\mu_g = 0.2$ . Due to the symmetry, the norms and  $\bar{F}_w$ -probabilities of the candidates are equivalent, leading to coinciding DM statistics based on QS and SphS scoring rules. Additionally, the equal norms and discrete probabilities also imply the censoring and conditioning rules to be equivalent within a semi-local scoring rule family since observations outside the region of interest obtain the same scores under both candidates in this case.

Figure 2 displays the rejection rates for rejection the null of equal predictive ability against the one-sided alternative that candidate  $f$  is statistically closer to  $p$  than  $g$ . The rejection rates are given at nominal significance levels of 0.01, 0.05 and 0.10, for focused versions of the LogS, SphS and CRPS scoring rules, based on 10,000 simulations. Given the discussion above, this gives a complete picture of the selection  $\{\text{LogS, SphS, QS, CRPS}\} \times \{\#, b\}$ . The twCRPS is added since it will also be included as benchmark in the power studies based on weight functions for which the censored CRPS variants do not reduce to the twCRPS (see Example 2). None of the displayed rejection rates give reason to doubt the size correctness of the tests.

## 4.2 Power

### 4.2.1 Laplace tails

Our first experiment studies the consequences of the lack of the conditional rule to disentangle two proportional tails when using the left tail indicator function  $w(y) = \mathbb{1}_{(-\infty, r)}(y)$ . In particular, we follow up on the Laplace example given by Figure 1a in Subsection 2.2, analysing two Laplace candidates with different location  $\mu_f = -1$  and  $\mu_g = 1$  but equivalent scale  $\theta_f = \theta_g = 1$ . Interestingly, even if  $\mu_p \rightarrow \mu_f$ , the conditional scoring rule does not have any power against the null of the candidates being statically equally far away from  $p$ , that is, for thresholds  $r < \mu_f$ , for which the conditional distributions on  $(-\infty, r)$  coincide. Since movements of  $p$  in terms of  $\mu_p$  are invisible through the lens of a conditional score divergence, this is essentially not a lack of power against  $\mathbb{H}_0$ , which is based on the conditional scoring rule. Yet, it is a lack of power against the distributions being statistically equally far away from the actual density on  $\{w > 0\}$  through the lens of the regular score divergence and, therefore, still a lack of local discriminative ability. More fundamentally, the GW test degenerates in this case, as the score differences are exactly zero.

Leaving this extreme case, we analyse what happens if the scale parameters are not exactly the same, but close. Specifically, we let  $\theta_f = 1$  and  $\theta_g = 1.1$ . Figure 3 shows the rejection rates of the GW test if the DGP is  $f$  (left-hand side) or  $g$  (right-hand side) in favour of  $f$  (top) or  $g$  (bottom). Since both candidates are now also different through the lens of the conditional rule, the subfigures on the diagonal display actual power while the off-diagonal ones show spurious power. Concerning the selection of scoring rules, it is good to remember that both censored CRPS rules (based on  $d_{\min}$  and  $d_{\text{rand}}$ ) coincide with each other and with the twCRPS, see Example 1 for details.

Three observations strike us. First of all, the increase in power from the conditional operator to the censoring operator is immense for all four scoring rules and thresholds  $r < \mu_f$ . The difference decreases over the interval  $r \in (\mu_f, \mu_g)$ , after which both conditioning and censoring have close to unit power. This observation is in line with the

Figure 3: Laplace experiment ( $c = 20$ )

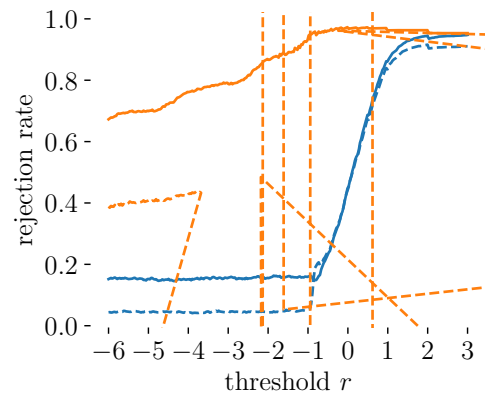
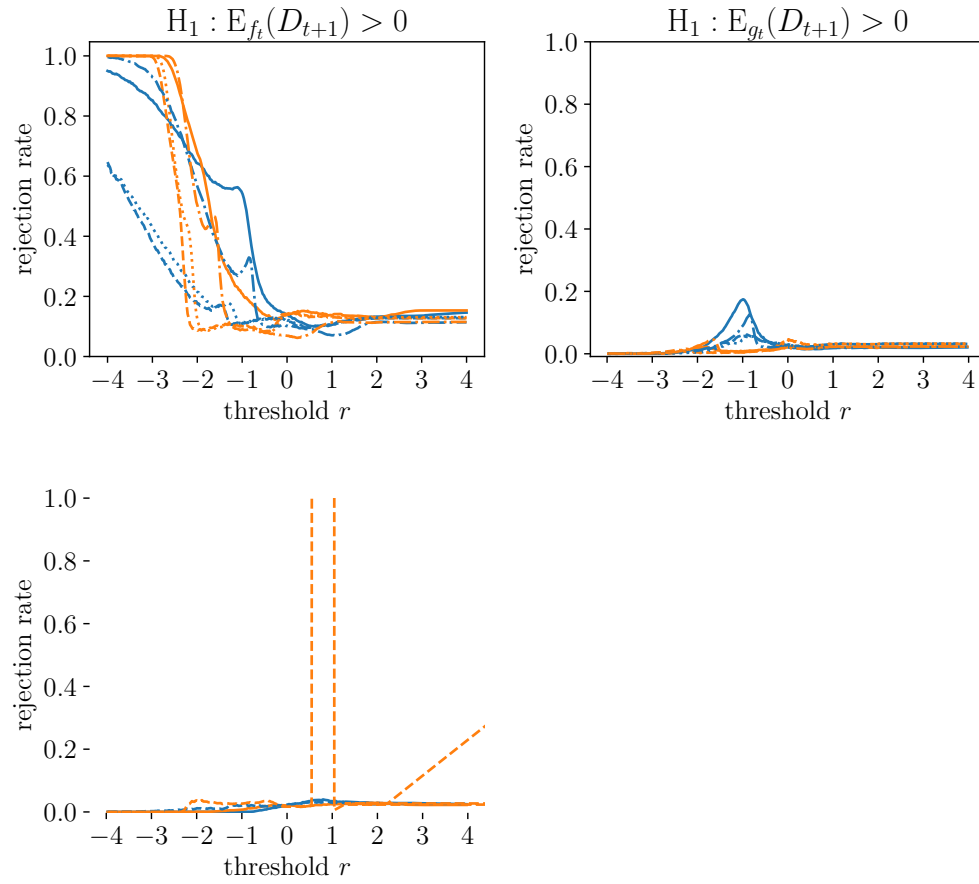




Figure 4:  $\mathcal{N}(0, 1)$  versus Student- $t(5)$ : Left-tail ( $c = 20$ )



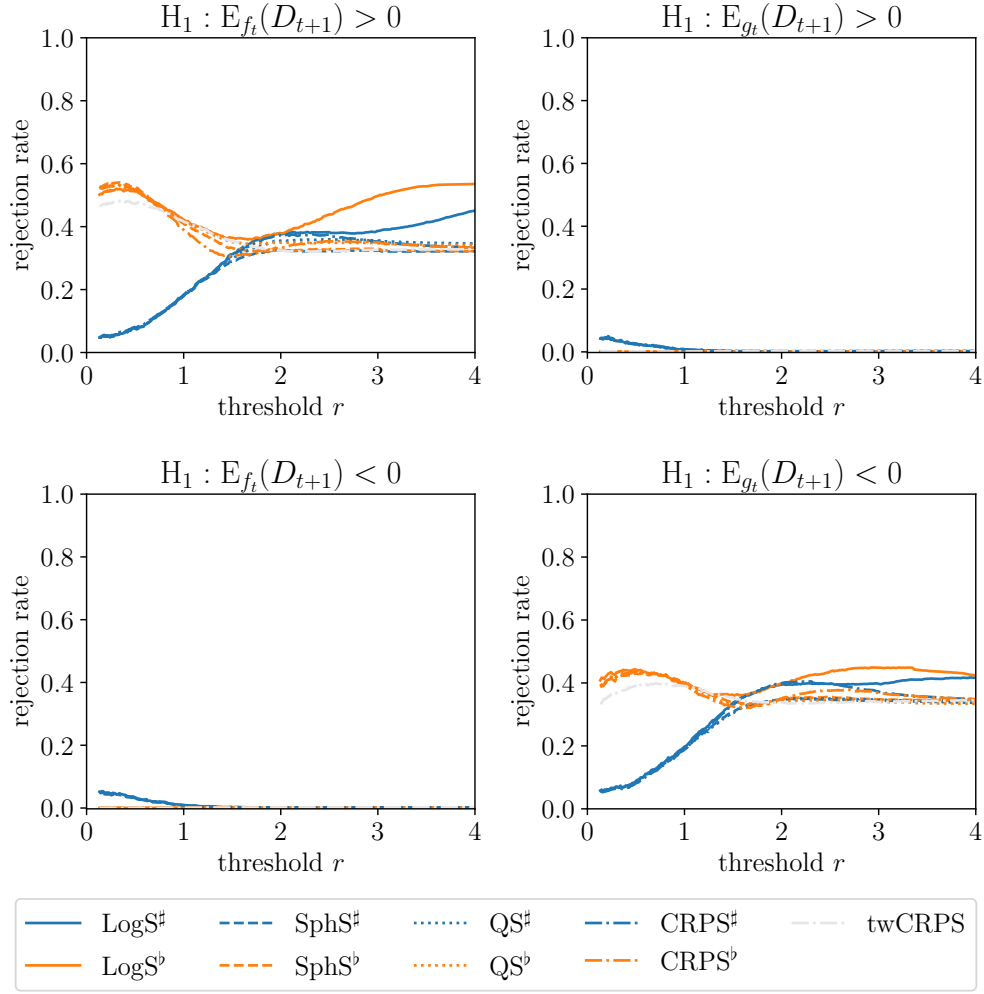
less monotonic, intersecting the graph of the competing focusing operator rejection rates. The latter occurs by construction since the densities of the candidates intersect as well, see [Diks et al. \(2011\)](#) for a discussion. Starting with the clearest differences, we note the spurious power humps of the conditional rules if the Student- $t(5)$  distribution is the DGP. In contrast, the censored scoring rules have almost no spurious power. Staying in the right column of [Figure 4](#), the rejection rates in the bottom row reveal a clear preference for the censoring operator. Indeed, the exceptions of higher conditional power are rather weak, while the difference between the rejection rates (far) into the left-tail is particularly large for the Logarithmic and Spherical scoring rule. On the other hand, if the standard normal distribution is the DGP, then there is hardly (a difference in) spurious power. The differences between the rejection rates representing power are more extreme when the data is generated from the standard Normal distribution, yet so is their drop between  $r = -2$  and  $r = -1$ , clouding a clear preference for one of the focusing operators for these intermediate tail values of  $r$ .

#### 4.2.3 $\mathcal{N}(0, 1)$ versus Student- $t(5)$ : Centre

In our third Monte Carlo experiment, we focus on the centre of the candidate distributions by implementing the weight function  $w(y) = \mathbb{1}_{[-r, r]}(y)$ . [Figure 5](#) displays the rejection rates for the same selection of regular scoring rules as in the previous experiments. Based on [Figure 5](#), the added value of censoring relative to conditioning is overwhelming: Censoring leads to higher power and lower spurious power, in particular for values smaller than  $r = 1$ , which are of particular interest in applications.

As discussed in [Example 2](#), the variants of the censored CRPS no longer coincide with each other, nor with the twCRPS. The  $\text{CRPS}_w^b$  displayed in the [Figure 5](#) is the generalised censored scoring rule based on the generalised censored measure in [Equation \(5\)](#). Due to the symmetry of the set up, there is visually no difference between using the suggested value  $\gamma = \frac{1}{2}$  or the estimated proportion  $\hat{\gamma}$ . We have also calculated the  $\text{CRPS}^\dagger(F, y)$  introduced in [Subsection 3.4.2](#), which visually coincides with the twCRPS in this case. As can be seen from [Figure 6](#), the alternative-distance based

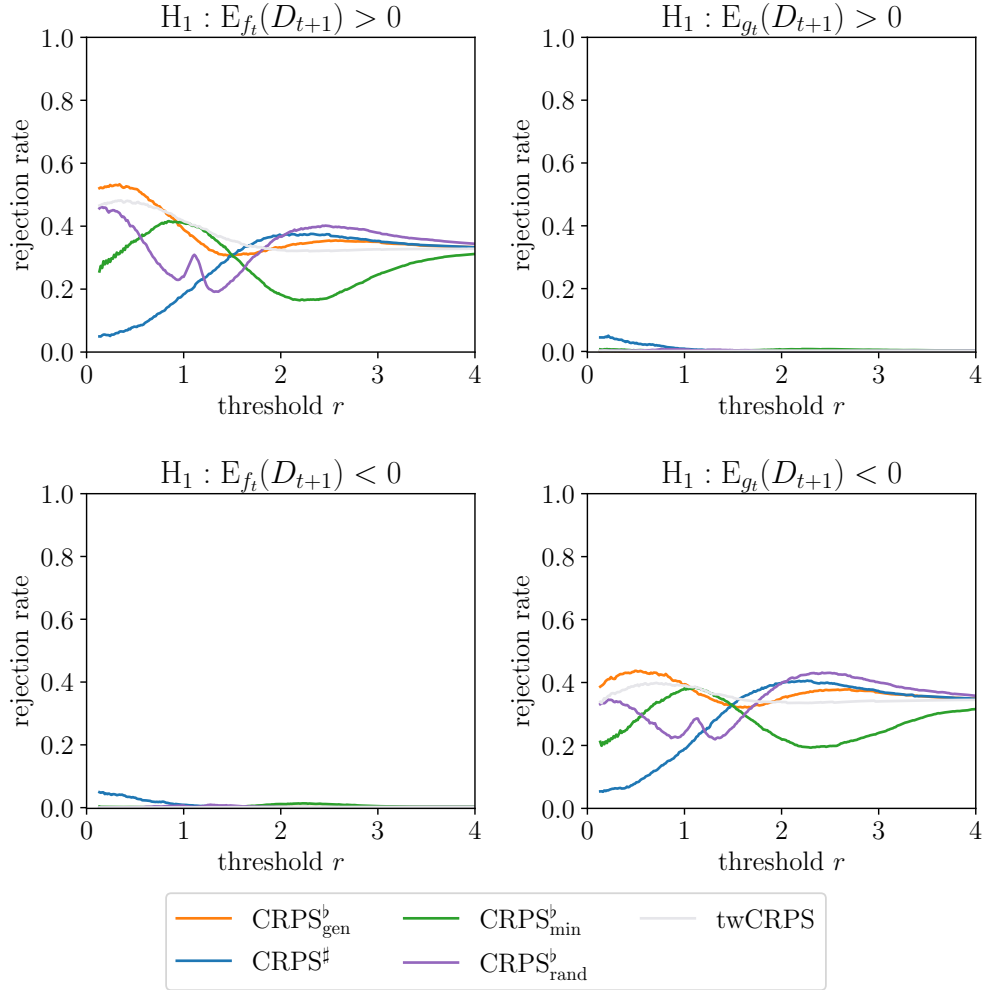
Figure 5:  $\mathcal{N}(0, 1)$  versus Student- $t(5)$ : Centre ( $c = 200$ )



One-sided rejection rates of the GW-test of equal predictive ability of the candidates  $f_t$  (standard normal) and  $g_t$  (Student- $t(5)$ ) at a nominal significance level of 0.05 based on 10,000 simulations. The DGP is either  $f_t$  (left-hand side) or  $g_t$  (right-hand side). Moreover, rejections in the top panels are in favour of  $f_t$ , while rejections in the bottom panels are in favour of  $g_t$ . The incorporated weight function is  $w(y) = 1_{(-r, r)}(y)$  and the number of expected observations in the region of interest is kept constant at  $c = 200$ .

CRPS variants do not coincide with the generalised censored CRPS, but lead to lower power than the generalised censored CRPS. Reconsidering the discussion of Subsection 3.4.2, this is unsurprising, since the generalised censored scoring rule incorporates more of the available information.

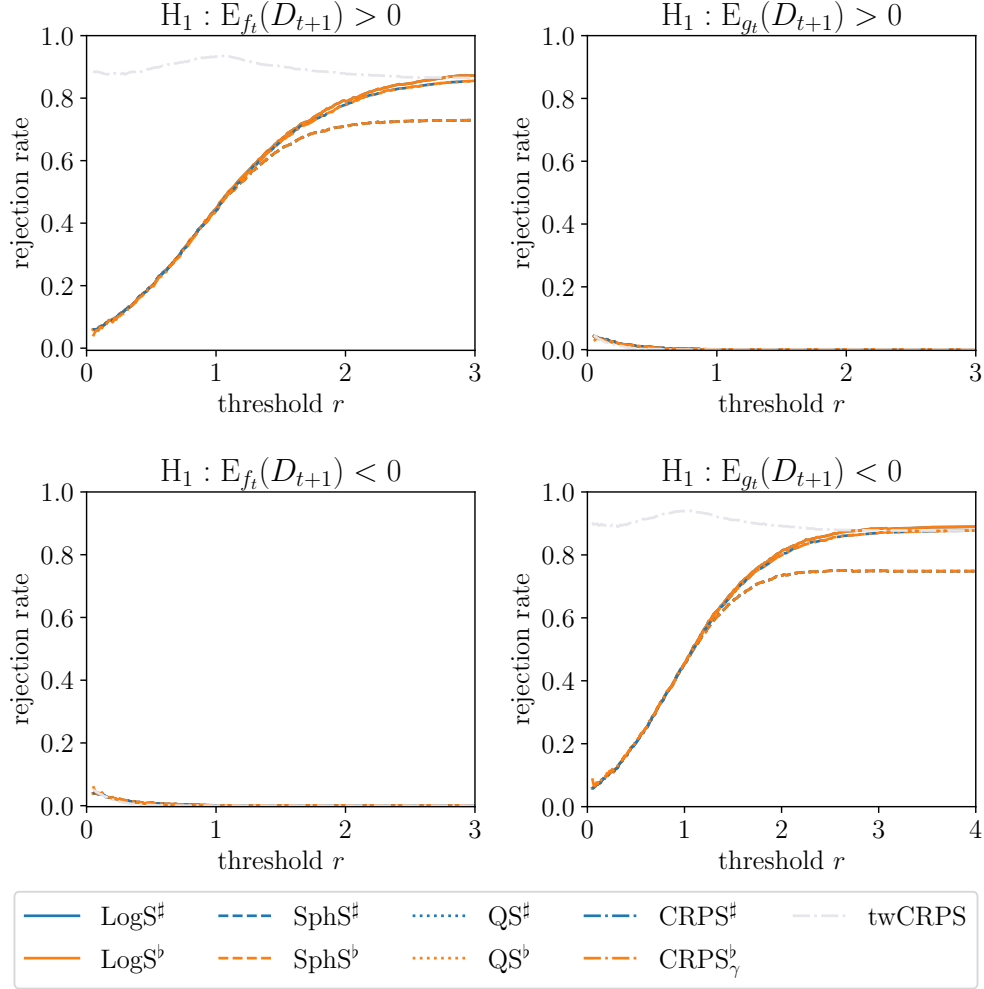
Figure 6:  $\mathcal{N}(0, 1)$  versus Student- $t(5)$ : Centre ( $c = 200$ )



One-sided rejection rates of the GW test of equal predictive ability of the candidates  $f_t$  (standard normal) and  $g_t$  (Student- $t(5)$ ) at a nominal significance level of 0.05 based on 10,000 simulations. The DGP is either  $f_t$  (left-hand side) or  $g_t$  (right-hand side). Moreover, rejections in the top panels are in favour of  $f_t$ , while rejections in the bottom panels are in favour of  $g_t$ . The incorporated weight function is  $w(y) = 1_{[-r, r]}(y)$  and the number of expected observations in the region of interest is kept constant at  $c = 200$ .

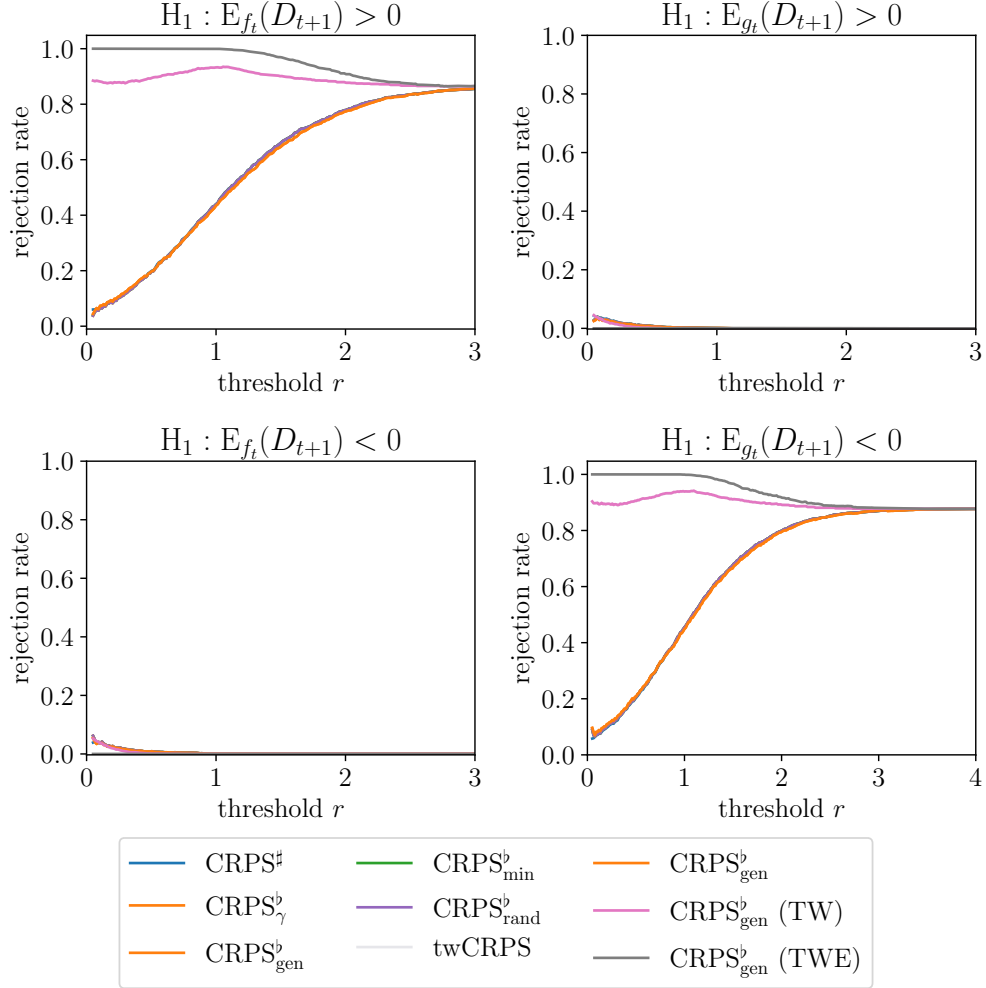
#### 4.2.4 $\mathcal{N}(-0.2, 1)$ versus $\mathcal{N}(-0.2, 1)$ : Centre ( $c = 200$ ) [ $\gamma_F \neq \gamma_P$ ]

Figure 7:  $\mathcal{N}(-0.2, 1)$  versus  $\mathcal{N}(-0.2, 1)$ : Centre ( $c = 200$ )



One-sided rejection rates of the GW test of equal predictive ability of the candidates  $f_t$  ( $\mathcal{N}(-0.2, 1)$ ) and  $g_t$  ( $\mathcal{N}(0.2, 1)$ ) at a nominal significance level of 0.05 based on 10,000 simulations. The DGP is either  $f_t$  (left-hand side) or  $g_t$  (right-hand side). Moreover, rejections in the top panels are in favour of  $f_t$ , while rejections in the bottom panels are in favour of  $g_t$ . The incorporated weight function is  $w(y) = 1_{[-r, r]}(y)$  and the number of expected observations in the region of interest is kept constant at  $c = 200$ .

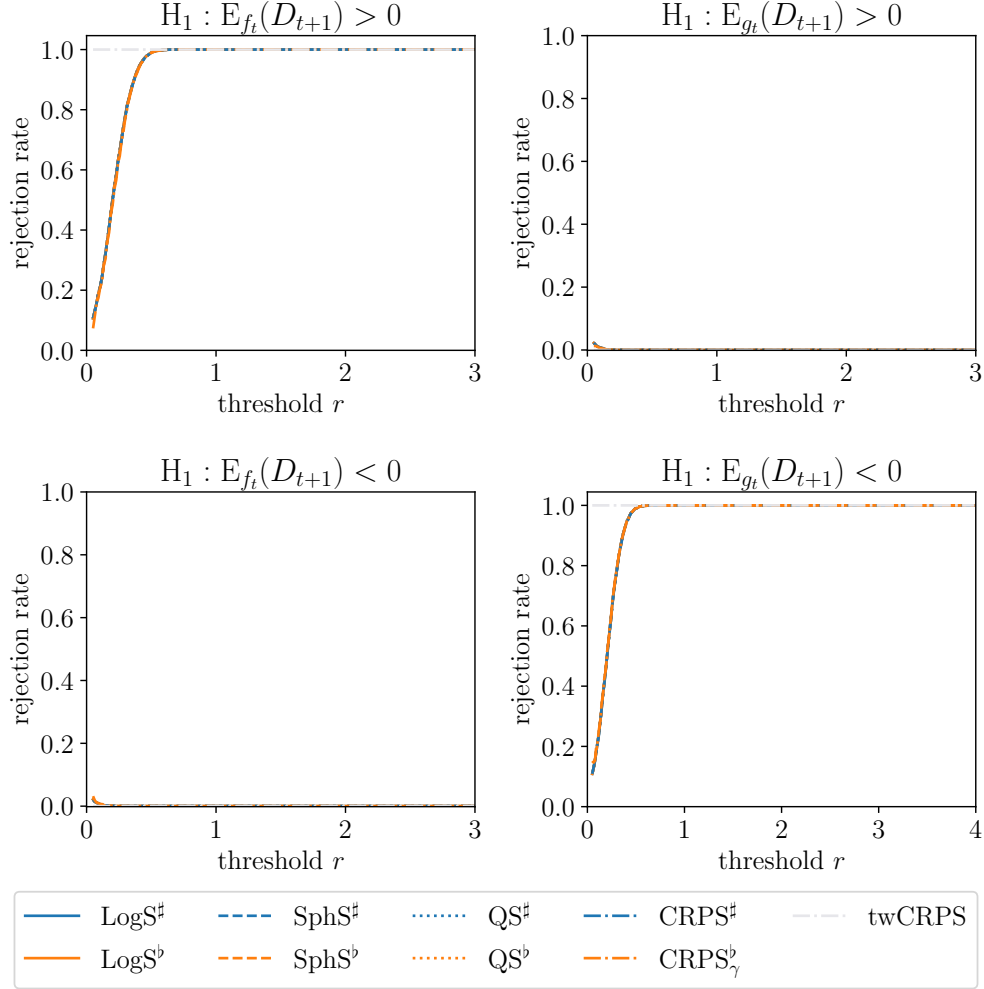
Figure 8:  $\mathcal{N}(-0.2, 1)$  versus  $\mathcal{N}(-0.2, 1)$ : Centre ( $c = 200$ )



One-sided rejection rates of the GW test of equal predictive ability of the candidates  $f_t$  ( $\mathcal{N}(-0.2, 1)$ ) and  $g_t$  ( $\mathcal{N}(0.2, 1)$ ) at a nominal significance level of 0.05 based on 10,000 simulations. The DGP is either  $f_t$  (left-hand side) or  $g_t$  (right-hand side). Moreover, rejections in the top panels are in favour of  $f_t$ , while rejections in the bottom panels are in favour of  $g_t$ . The incorporated weight function is  $w(y) = 1_{[-r, r]}(y)$  and the number of expected observations in the region of interest is kept constant at  $c = 200$ .

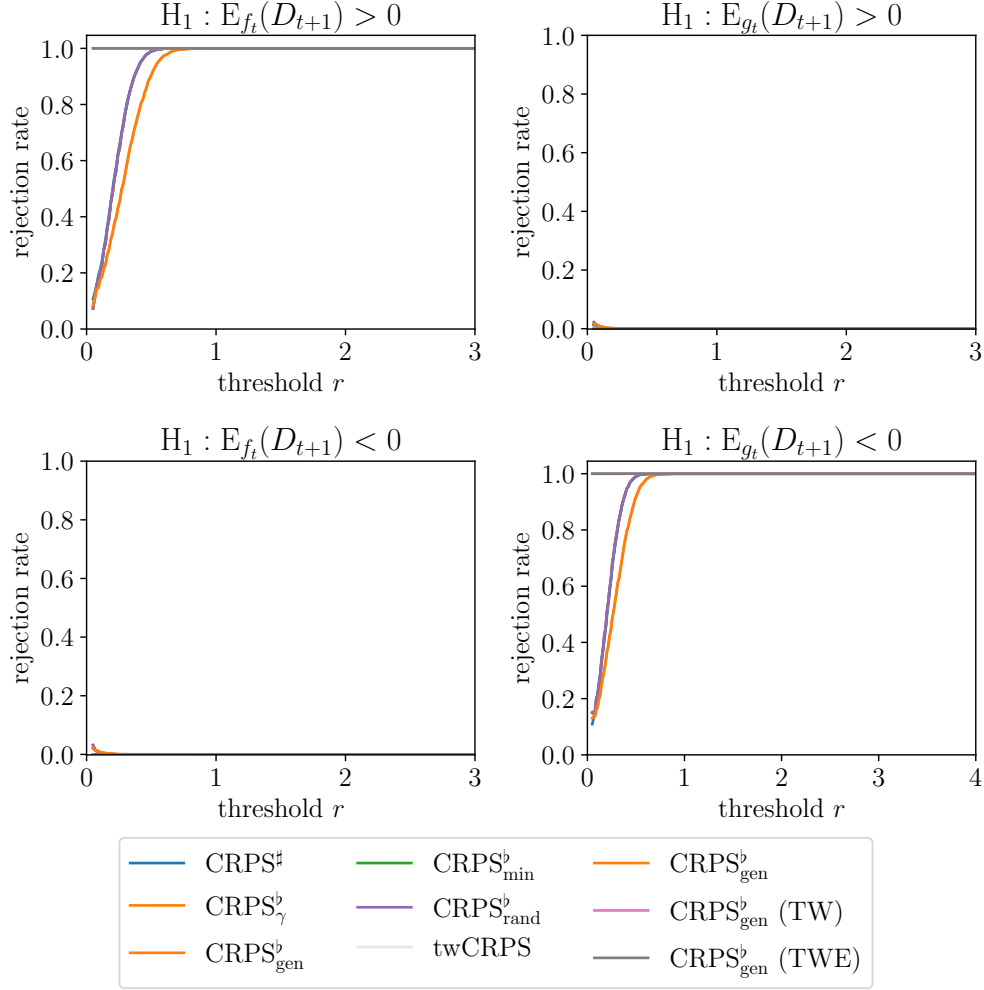
#### 4.2.5 $\mathcal{N}(-1, 1)$ versus $\mathcal{N}(-1, 1)$ : Centre ( $c = 200$ ) [ $\gamma_F \neq \gamma_P$ ]

Figure 9:  $\mathcal{N}(-0.2, 1)$  versus  $\mathcal{N}(-0.2, 1)$ : Centre ( $c = 200$ )



One-sided rejection rates of the GW test of equal predictive ability of the candidates  $f_t$  ( $\mathcal{N}(-1, 1)$ ) and  $g_t$  ( $\mathcal{N}(1, 1)$ ) at a nominal significance level of 0.05 based on 10,000 simulations. The DGP is either  $f_t$  (left-hand side) or  $g_t$  (right-hand side). Moreover, rejections in the top panels are in favour of  $f_t$ , while rejections in the bottom panels are in favour of  $g_t$ . The incorporated weight function is  $w(y) = 1_{[-r, r]}(y)$  and the number of expected observations in the region of interest is kept constant at  $c = 200$ .

Figure 10:  $\mathcal{N}(-1, 1)$  versus  $\mathcal{N}(-1, 1)$ : Centre ( $c = 200$ )



One-sided rejection rates of the GW test of equal predictive ability of the candidates  $f_t$  ( $\mathcal{N}(-1, 1)$ ) and  $g_t$  ( $\mathcal{N}(1, 1)$ ) at a nominal significance level of 0.05 based on 10,000 simulations. The DGP is either  $f_t$  (left-hand side) or  $g_t$  (right-hand side). Moreover, rejections in the top panels are in favour of  $f_t$ , while rejections in the bottom panels are in favour of  $g_t$ . The incorporated weight function is  $w(y) = 1_{[-r, r]}(y)$  and the number of expected observations in the region of interest is kept constant at  $c = 200$ .



## 5 Empirical Application

### 5.1 Risk management

Evaluating the downside risk of asset returns is a crucial task in risk management, particularly for meeting regulatory requirements related to risk measures like the Value-at-Risk ( $\text{VaR}_{\hat{f}_t}^q$ ), which represents the  $q$ -th quantile of the model-based estimated density forecast  $\hat{f}_t$  and the more recently required Expected Shortfall  $\text{ES}_{\hat{f}_t}^q$ , being the expected loss conditional on the the loss being below its  $\text{VaR}_{\hat{f}_t}^q$ . To fulfil this objective, we opt for a weight function of  $w_t(y_t) = \mathbb{1}_{(-\infty, r_t^q]}(y_t)$ , and choose for the variable of interest  $y_t$  the log-returns of the S&P500, that is,  $y_t = \log(P_t/P_{t-1})$ , where  $P_t$  is the adjusted closing price on day  $t$ . The dataset used for this study consists of 6,777 observations in total, spanning from January 2, 1996, to December 30, 2022, and obtained from Yahoo Finance.

Each of our selected forecast methods can be represented as

$$Y_t | \mathcal{F}_{t-1} \sim \mathcal{D}(\mu, \sigma_t^2, \boldsymbol{\vartheta}),$$

where  $\mathcal{D}(\mu, \sigma_t^2, \boldsymbol{\vartheta})$  denotes a parametric family of distributions with mean  $\mu$ , variance  $\sigma_t^2$  and other parameters  $\boldsymbol{\vartheta}$ . We have also considered AR(1) and AR(5) models for the specification of the conditional mean, but did not find substantial improvements relative to the constant mean specification. In our analysis, we consider two conditional variance models: the GARCH(1,1) model proposed by [Bollerslev \(1987\)](#) and defined by the equation

$$\sigma_t^2 = \omega + \alpha(y_t - \mu)^2 + \beta\sigma_{t-1}^2,$$

and the RGARCH(1,1) model proposed by Hansen et al. (2012), which is given by

$$\begin{aligned}\sigma_t^2 &= \omega + \alpha x_{t-1} + \beta \sigma_{t-1}^2, \\ x_t &= \xi + \phi \sigma_t^2 + \tau z_t + \kappa(z_t^2 - 1) + u_t,\end{aligned}$$

where  $x_t$  represents the realised measure,  $z_t = (y_t - \mu)/\sigma_t$  and  $u_t$  denotes a white noise process with variance  $\sigma_u^2$ . The realised measure is downloaded from the Risklab page of Dacheng Xiu's website: <https://dachxiu.chicagobooth.edu/#risklab>. Notably, the equation for the realised measure is necessary for predictions beyond one step ahead, and the term  $\tau z_t + \kappa(z_t^2 - 1)$  captures the leverage effect. Furthermore, we use either a normal distribution or a Student- $t$  distribution with  $\nu$  degrees of freedom. Although the heavy tails of the Student- $t$  distribution have proven useful for daily stock return data, the normal distribution remains a common choice. We estimate the parameters using full-fledged maximum likelihood estimation, based on a rolling estimation window of length  $T_{\text{est}}$ .

### 5.1.1 Statistical comparison

As an empirical extension of our power analysis from Chapter 4, we compare the relative performance of the forecasting methods via the Model Confidence Set (MCS), as described by Hansen et al. (2011). The MCS procedure expands the GW hypothesis to larger sets of  $H_0$ -equivalent methods, employing an iterative elimination procedure based on the equivalence tests  $\text{Tmax}_k$  or  $\text{TR}_k$ , with  $k$  indicating the block length of the block bootstrap required for estimating the non-standard asymptotic distribution of these test statistics. Optimal power properties of censoring in the GW environment intuitively accelerate elimination in the MCS procedure, resulting in smaller MCS  $p$ -values and consequently, lower MCS cardinality.

Table 2 delineates the MCS  $p$ -values and the deduced  $\text{MCS}_{0.90}$  and  $\text{MCS}_{0.75}$  for the six previously mentioned forecast methods, based on their 1-step and 5-step ahead density forecasts. We present the outcomes for  $q \in \{0.01, 0.1, 0.2\}$ , using the  $\text{TR}_{20}$

statistic with a block bootstrap size of  $B = 10,000$  and an estimation window of  $T_{\text{est}} = 1,000$ , and later validate the stability of these findings. When examining the censored ( $b$ ) and conditional ( $\sharp$ ) columns, the reported results overwhelmingly support enhanced power via censoring. The cardinality of the censored MCS never surpasses its conditional counterpart, the censored MCS often being a (strict) subset of the conditional MCS. Notable differences, especially for  $\text{MCS}_{0.75}$  of QS or CRPS at  $q = 0.10$ , underline the overall relative increase in cardinality of  $\text{MCS}_{0.90}$  ( $\text{MCS}_{0.75}$ ) by 78% (75%) when opting for conditioning over censoring. The significance of these reductions is further emphasised by the fact that the resulting MCS encompass more complex model specifications, which would be the optimal choices in the absence of parameter and forecasting uncertainty.

For  $h = 5$ , reductions through censoring also occur, albeit less frequently. Excluding the CRPS rule, the table suggests that censoring leads to a smaller  $\text{MCS}_{0.90}$  ( $\text{MCS}_{0.75}$ ) 3.0 (1.7) times more often. With the CRPS inclusion, this frequency equilibrates due to the stronger eliminative capacity of the conditional CRPS. Again, the differences can become quite sizeable, here on average leading to a relative increase of cardinality of 31% (41%), or even 50% (62%) excluding the CRPS, when using conditioning instead.

To ascertain the robustness of our findings, we expand our study across different estimation window lengths  $T_{\text{est}}$ , block lengths  $k$ , and the  $T_{\text{max}}$  statistics for a more comprehensive set of quantiles, also including  $q = 0.05, 0.15$  and  $0.25$ . Table 3 highlights the consistency of these results, revealing only one occurrence of increased  $\text{MCS}^b$ -cardinality at both confidence levels and even no exception for the average relative increase of the size of from  $\text{MCS}^b$  to  $\text{MCS}^\sharp$ . These findings align with reported differences between forecasting horizons in Table 2, more generally influenced by the relatively strong performance of the  $\text{CRPS}^\sharp$  rule on  $h = 5$  outcomes. Furthermore, the first row of Table 3 extends the results from Table 2 to all quantiles, showing that censoring typically reduces MCS cardinality 71% of the time and only infrequently, 4% of the time, increases it. The TR statistic typically intensifies elimination, bearing smaller MCS  $p$ -values, translating into movements outside the  $\text{MCS}_{0.90}$  or  $\text{MCS}_{0.75}$  for

the censored scoring rules, where MCS  $p$ -values were typically closer to the boundaries. Differences in outcomes due to varying block lengths ( $k = 20$  and  $k = 100$ ) are minor, and the influence of the estimation window lengths on the MCS  $p$ -values is not substantial.

### 5.1.2 Backtesting

Beyond the statistical assessment of forecast methods, we compute their 1- and 5-step ahead Value at Risk ( $\text{VaR}_{\hat{f}_t}^q$ ) and Expected Shortfall ( $\text{ES}_{\hat{f}_t}^q$ ). These measures provide only partial insight into the forecasts, since the tail component of the density forecast carries more comprehensive information than a single quantile ( $\text{VaR}_{\hat{f}_t}^q$ ) or conditional moment  $\text{ES}_{\hat{f}_t}^q = \mathbb{E}_{\hat{f}_t} \left( Y_{t+h} | Y_{t+h} \leq \text{VaR}_{\hat{f}_t}^q \right)$ . Notably, the conditioning in  $\text{ES}_{\hat{f}_t}^q$  is a quantile of the density forecast itself rather than  $\hat{r}_t^q$ , a.s. implying a discrepancy between the operational region of  $\text{ES}_{\hat{f}_t}^q$  and the focused scoring rules introduced above. Additionally, if the  $\text{VaR}_{\hat{f}_t}^q$  is quite off, then the ‘risk’ indicated by  $\text{ES}_{\hat{f}_t}^q$  can become quite detached from the true risk  $\text{ES}_p^q$ , where  $p$  denotes the density of the DGP. Hence, the  $\text{ES}_{\hat{f}_t}^q$  is (particularly) useful when the  $\text{VaR}_{\hat{f}_t}^q$  is accurate, i.e. we preferably have a good fit for the pair  $(\text{VaR}_{\hat{f}_t}^q, \text{ES}_{\hat{f}_t}^q)$ , rather than just  $\text{ES}_{\hat{f}_t}^q$  itself.

We highlight a corollary before discussing results. Given a fixed level  $q$ , let  $r$  be such that  $\text{VaR}_{\hat{f}_t}^q \vee \text{VaR}_p^q \leq r$ . A property of the censored scoring rule is its ability to render the true  $(\text{VaR}_p^q, \text{ES}_p^q)$  pair, since

$$\mathbb{D}_{S_w^b}(p||f) = 0 \implies (\text{VaR}_p^q, \text{ES}_p^q) = (\text{VaR}_{\hat{f}_t}^q, \text{ES}_{\hat{f}_t}^q), \quad (8)$$

where  $w(y) = \mathbb{1}_{(-\infty, r)}(y)$ . This is a direct consequence of (3), i.e. another corollary of Lemma 2, and holds also more generally for any functional on distributions on  $\{w > 0\}$ . In (sharp) contrast,  $\mathbb{D}_{S_w^\#}(p||f) = 0$  implies that  $p \propto f$  on  $(-\infty, r)$  and hence  $(\text{VaR}_p^q, \text{ES}_p^q) \neq (\text{VaR}_{\hat{f}_t}^q, \text{ES}_{\hat{f}_t}^q)$ , unless  $\bar{F}_w = \bar{P}_w$ . Therefore, model selection based on censored scoring rules aligns more effectively with backtesting of functionals of the distribution compared to model selection based on conditional scoring rules.

Table 2 reports the backtesting results, demonstrating a consistent preference for Student- $t_\nu$  models by Expected Shortfall (ES) measures, with VaR's top models fluctuating between Student- $t_\nu$  and Normal. Interestingly, at  $q = 0.10$  for  $h = 1$ , there's a significant discrepancy in methods preferred by the risk measures. The MCS accommodates this by retaining at least one model from both camps, the censored MCS favouring accurate  $\text{VaR}_{\hat{f}_t}^q$  coverage slightly more than conditional ones. Table 3 supports this, revealing that censoring typically leads to greater alignment with the scoring rule in the VaR column, while conditioning excels in the ES column. Notably, these outcomes are often reached with smaller MCS under censoring. Hence, the small percentages of mismatches particularly underscore the usefulness of censoring, as it frequently identifies a compact set of models yielding reliable risk measures. For  $h = 5$ , the percentages are generally smaller. However, this is not because forecasting risk measures further into the future is easier. On the contrary, the MCS are less reduced for  $h = 5$ , likely due to the greater difficulty in distinguishing between models as a result of increased forecasting noise.

As discussed earlier, it is sensible to examine the pair  $(\text{VaR}_{\hat{f}_t}^q, \text{ES}_{\hat{f}_t}^q)$  rather than  $\text{ES}_{\hat{f}_t}^q$  in isolation. Moreover, Equation (8) suggests that censoring may generate Model Confidence Sets (MCS) containing forecast models that produce  $(\text{VaR}_{\hat{f}_t}^q, \text{ES}_{\hat{f}_t}^q)$  pairs closer to the true pair. Support for this conjecture is found in Table 3. Despite often being smaller, the censored MCS contains well-fitted  $(\text{VaR}_{\hat{f}_t}^q, \text{ES}_{\hat{f}_t}^q)$  pairs, defined as 0% mismatches for both VaR and ES, more than twice as often (9 versus 4). If we accept up to 4% mismatches, the comparison remains favourable: 14 versus 6, endorsing censored MCS as a superior selection mechanism prior to VaR and ES calculations.

Table 2: Evaluating forecast methods in the left-tail

$q$	Method	LogS		QS		SphS		CRPS		Backtesting	
		$b$	$\sharp$	$b$	$\sharp$	$b$	$\sharp$	$b$	$\sharp$	hitrate	ES
$h = 1$											
0.01	RGARCH- $t_\nu$	<b>1.00</b>	<b>0.60</b>	<b>0.45</b>	<b>0.95</b>	<b>0.65</b>	<b>0.88</b>	<b>0.73</b>	<b>0.97</b>	<b>0.016</b>	<b>0.002*</b>
	TGARCH- $t_\nu$	<b>0.99</b>	<b>1.00</b>	<b>0.63</b>	<b>1.00</b>	<b>0.88</b>	<b>1.00</b>	<b>0.81</b>	<b>1.00</b>	<b>0.014</b>	<b>0.008*</b>
	GARCH- $t_\nu$	<b>0.53</b>	<b>0.69</b>	<b>0.34</b>	<b>0.84</b>	<b>0.65</b>	<b>0.88</b>	<b>0.81</b>	<b>0.91</b>	<b>0.015</b>	<b>0.044*</b>
	RGARCH- $\mathcal{N}$	0.09	0.19	<b>1.00</b>	<b>0.95</b>	<b>1.00</b>	<b>0.66</b>	<b>0.81</b>	<b>0.97</b>	0.020	0.221
	TGARCH- $\mathcal{N}$	0.03	0.09	<b>0.63</b>	<b>0.95</b>	<b>0.88</b>	<b>0.38</b>	<b>1.00</b>	<b>0.97</b>	0.021	0.257
	GARCH- $\mathcal{N}$	0.01	0.09	<b>0.45</b>	<b>0.84</b>	<b>0.65</b>	<b>0.45</b>	<b>0.81</b>	<b>0.64</b>	0.023	0.294
0.10	RGARCH- $t_\nu$	<b>1.00</b>	<b>0.73</b>	<b>1.00</b>	<b>0.70</b>	<b>0.35</b>	<b>0.88</b>	0.12	<b>0.74</b>	0.115	<b>0.074*</b>
	TGARCH- $t_\nu$	0.05	<b>1.00</b>	0.16	<b>1.00</b>	0.02	<b>1.00</b>	0.05	<b>1.00</b>	0.113	<b>0.045*</b>
	GARCH- $t_\nu$	0.00	<b>0.40</b>	0.01	<b>0.43</b>	0.00	<b>0.37</b>	0.01	0.20	0.118	<b>0.048*</b>
	RGARCH- $\mathcal{N}$	0.05	0.03	<b>0.49</b>	<b>0.70</b>	<b>1.00</b>	0.21	<b>1.00</b>	<b>0.74</b>	<b>0.103*</b>	0.160
	TGARCH- $\mathcal{N}$	0.00	0.01	0.06	<b>0.39</b>	0.00	0.03	0.12	0.04	<b>0.100*</b>	0.168
	GARCH- $\mathcal{N}$	0.00	0.00	0.00	0.06	0.00	0.01	0.03	0.01	<b>0.103*</b>	0.206
0.20	RGARCH- $t_\nu$	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>0.40</b>	<b>1.00</b>	<b>0.26</b>	0.10	<b>0.42</b>	<b>0.203*</b>	<b>0.111</b>
	TGARCH- $t_\nu$	0.02	0.23	0.10	0.06	0.02	0.02	0.01	0.14	<b>0.204*</b>	<b>0.076</b>
	GARCH- $t_\nu$	0.00	0.06	0.06	0.00	0.00	0.00	0.00	0.01	<b>0.210*</b>	<b>0.084</b>
	RGARCH- $\mathcal{N}$	0.01	0.04	0.00	<b>1.00</b>	0.02	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	0.183	0.144
	TGARCH- $\mathcal{N}$	0.00	0.00	0.00	0.01	0.00	0.00	0.01	0.09	0.180	0.134
	GARCH- $\mathcal{N}$	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.01	0.186	0.152
$h = 5$											
0.01	RGARCH- $t_\nu$	<b>0.37</b>	<b>0.87</b>	0.12	<b>1.00</b>	0.24	<b>1.00</b>	<b>0.57</b>	<b>1.00</b>	<b>0.015</b>	<b>0.227</b>
	TGARCH- $t_\nu$	<b>0.83</b>	<b>1.00</b>	<b>0.86</b>	<b>0.45</b>	<b>1.00</b>	<b>0.40</b>	<b>0.65</b>	<b>0.63</b>	<b>0.017</b>	<b>0.128*</b>
	GARCH- $t_\nu$	<b>1.00</b>	<b>0.96</b>	0.17	<b>0.38</b>	<b>0.47</b>	<b>0.40</b>	<b>0.65</b>	<b>0.44</b>	<b>0.017</b>	<b>0.054*</b>
	RGARCH- $\mathcal{N}$	0.01	0.05	0.12	<b>0.81</b>	0.18	<b>0.40</b>	<b>0.57</b>	0.18	0.019	0.521
	TGARCH- $\mathcal{N}$	0.01	0.05	<b>1.00</b>	<b>0.75</b>	<b>1.00</b>	0.23	<b>1.00</b>	0.16	0.022	0.461
	GARCH- $\mathcal{N}$	0.01	0.04	0.17	<b>0.75</b>	<b>0.41</b>	<b>0.27</b>	<b>0.65</b>	0.09	0.022	0.480
0.10	RGARCH- $t_\nu$	<b>0.46</b>	0.15	<b>0.35</b>	0.16	<b>0.55</b>	0.11	<b>0.47</b>	<b>0.43</b>	<b>0.100*</b>	<b>0.156</b>
	TGARCH- $t_\nu$	<b>1.00</b>	<b>0.36</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>0.41</b>	<b>1.00</b>	<b>1.00</b>	0.108	<b>0.096</b>
	GARCH- $t_\nu$	<b>0.56</b>	<b>1.00</b>	<b>0.35</b>	<b>0.70</b>	<b>0.55</b>	<b>1.00</b>	<b>0.60</b>	<b>0.43</b>	0.113	<b>0.075*</b>
	RGARCH- $\mathcal{N}$	0.00	0.00	0.01	0.16	<b>0.26</b>	0.03	<b>0.53</b>	0.01	0.090	0.242
	TGARCH- $\mathcal{N}$	0.00	0.00	0.01	0.16	0.21	0.00	<b>0.60</b>	0.00	<b>0.095*</b>	0.242
	GARCH- $\mathcal{N}$	0.00	0.00	0.01	0.16	0.20	0.00	<b>0.46</b>	0.00	<b>0.099*</b>	0.247
0.20	RGARCH- $t_\nu$	0.15	<b>0.36</b>	0.00	<b>0.89</b>	0.02	<b>0.97</b>	<b>0.37</b>	<b>0.67</b>	0.175	<b>0.155</b>
	TGARCH- $t_\nu$	<b>1.00</b>	<b>0.36</b>	<b>0.77</b>	<b>1.00</b>	<b>1.00</b>	<b>0.97</b>	<b>1.00</b>	<b>1.00</b>	<b>0.192*</b>	<b>0.114</b>
	GARCH- $t_\nu$	<b>0.78</b>	<b>1.00</b>	<b>1.00</b>	<b>0.89</b>	<b>0.50</b>	<b>0.97</b>	<b>0.37</b>	<b>0.67</b>	<b>0.201*</b>	<b>0.086*</b>
	RGARCH- $\mathcal{N}$	0.00	0.00	0.00	<b>0.89</b>	0.00	<b>1.00</b>	<b>0.32</b>	<b>0.67</b>	0.160	0.193
	TGARCH- $\mathcal{N}$	0.00	0.00	0.00	<b>0.79</b>	0.00	<b>0.90</b>	<b>0.37</b>	0.02	0.172	0.167
	GARCH- $\mathcal{N}$	0.00	0.00	0.00	<b>0.79</b>	0.00	<b>0.76</b>	0.23	0.01	<b>0.176</b>	0.185

Note: Columns titled by scoring rules present MCS  $p$ -values for assessing predictive equality of the forecast methods listed,  $h$ -steps ahead, based on LogS, QS, SphS and CRPS scoring rule variants, both censored ( $b$ ) and conditional ( $\sharp$ ). The emphasis is on the left-tail, incorporated by weight function  $w_t(y_t) = \mathbb{1}_{(-\infty, r_q]}(y_t)$ , for various  $q$  values. Bold (underlined)  $p$ -values signify the forecast method's inclusion in MCS<sub>0.75</sub> (MCS<sub>0.90</sub>). All MCS  $p$ -values utilise the TR20 statistic, implementing  $B = 10,000$  block bootstrap simulations with blocklength  $k = 20$ . The estimation window is  $T_{\text{est}} = 1,000$ . The backtesting columns present the hit rate, instances beneath a method's VaR at level  $q$ , and the absolute difference between actual and density forecast-implied expected shortfall at level  $q$ . The top three models are bolded; an asterisk denotes non-rejection of null for correct coverage (hit rate) or difference (ES) at a 0.05 significance level.

Table 3: Overview and robustness of left-tail application

$T_{\text{est}}$	Stat.	MCS <sub>0.90</sub>							MCS <sub>0.75</sub>						
		MCS			VaR		ES		MCS			VaR		ES	
		<	>	%	b	#	b	#	<	>	%	b	#	b	#
$h = 1$															
1000	TR <sub>20</sub>	<b>71</b>	4	128	17	17	17	<b>8</b>	<b>63</b>	8	104	<b>17</b>	21	25	<b>13</b>
	TR <sub>100</sub>	<b>71</b>	4	130	17	17	17	<b>8</b>	<b>63</b>	8	104	<b>17</b>	17	17	<b>0</b>
	Tmax <sub>20</sub>	<b>46</b>	13	58	<b>0</b>	13	4	4	<b>71</b>	8	107	<b>4</b>	13	17	<b>4</b>
	Tmax <sub>100</sub>	<b>46</b>	13	48	<b>0</b>	13	4	4	<b>67</b>	13	102	<b>4</b>	13	17	<b>4</b>
750	TR <sub>20</sub>	<b>54</b>	13	81	<b>13</b>	25	13	13	<b>50</b>	13	80	<b>13</b>	29	29	<b>17</b>
	Tmax <sub>20</sub>	<b>33</b>	29	19	<b>0</b>	17	4	4	<b>63</b>	17	100	<b>4</b>	17	21	<b>8</b>
1250	TR <sub>100</sub>	<b>63</b>	8	122	17	17	17	<b>0</b>	<b>50</b>	4	99	17	17	21	<b>17</b>
	Tmax <sub>20</sub>	<b>58</b>	17	74	8	8	13	<b>0</b>	<b>50</b>	17	83	<b>8</b>	13	21	<b>17</b>
$h = 5$															
1000	TR <sub>20</sub>	<b>38</b>	25	69	17	<b>0</b>	4	<b>0</b>	<b>50</b>	42	72	<b>17</b>	21	4	<b>0</b>
	TR <sub>100</sub>	<b>38</b>	25	69	17	<b>0</b>	4	<b>0</b>	<b>50</b>	42	66	<b>17</b>	21	4	<b>0</b>
	Tmax <sub>20</sub>	<b>42</b>	25	43	0	0	0	0	<b>50</b>	42	59	<b>0</b>	17	0	0
	Tmax <sub>100</sub>	<b>38</b>	17	44	0	0	0	0	<b>46</b>	38	59	<b>0</b>	17	0	0
750	TR <sub>100</sub>	33	<b>46</b>	44	<b>0</b>	8	0	0	33	<b>50</b>	54	<b>0</b>	25	0	0
	Tmax <sub>20</sub>	<b>29</b>	25	53	<b>0</b>	4	0	0	<b>42</b>	38	53	<b>0</b>	17	0	0
1250	TR <sub>20</sub>	<b>38</b>	33	61	4	4	0	0	<b>50</b>	33	61	<b>8</b>	25	0	0
	Tmax <sub>20</sub>	<b>42</b>	29	46	<b>0</b>	4	0	0	<b>46</b>	42	46	<b>4</b>	17	0	0

Note: The table summarises MCS and backtesting results using varying values for estimation window  $T_{\text{est}}$ , equivalence test statistics TR <sub>$k$</sub>  and Tmax <sub>$k$</sub> , blocklength  $k$  across forecast horizons  $h = 1$  and  $h = 5$ , based on  $B = 10,000$  bootstrap replications. All values are percentages. Columns labelled with |MCS| refer to MCS cardinality. Across 24 combinations of  $q \in \{0.01, 0.05, 0.10, 0.15, 0.20, 0.25\}$  and  $S \in \{\text{LogS}, \text{QS}, \text{SphS}, \text{CRPS}\}$ , the < (>)-column displays the frequency of  $|\text{MCS}_{1-\alpha}^b| < (>) |\text{MCS}_{1-\alpha}^\#|$  and the %-column indicates average relative cardinality increase from  $\text{MCS}_{1-\alpha}^\#$  to  $\text{MCS}_{1-\alpha}^b$ . The VaR (ES) column shows the frequency of the  $\text{MCS}_{1-\alpha}$  containing one of the top three models based on VaR (ES) backtesting results. Bolded numbers indicate strictly more smaller (<) or larger (>)  $\text{MCS}_{1-\alpha}^b$  as well as strictly less times the  $\text{MCS}_{1-\alpha}$  contains a top 3 VaR or ES method.

## 6 Conclusion

In many applications, forecasters are not equally interested in all possible outcomes of the random variable of interest. For such cases, we have motivated the use of censoring as focusing device. In particular, we have shown that focusing scoring rules by applying them to censored distributions leads to strictly locally proper scoring rules. To the best of our knowledge, we are first in deriving a transformation of the original scoring rule that preserves strict propriety.

Our approach offers considerable flexibility in terms of the original scoring rule, the weight function, and the outcome space. For specific choices, the (generalised) censored scoring rule delivers intuitively sound scoring rules that can easily be implemented by practitioners. When applied to the logarithmic scoring rule, our focusing technique produces the well-established censored likelihood score. Our framework also facilitates the derivation of weighted versions of the Energy Score family, which can be considered the multivariate equivalent of the twCRPS for left- or right-tail indicator functions. For other weight functions, the censored CRPS rule is strictly locally proper, while the twCRPS is not. By uncovering a procedure close to censoring, we established to clarify the localisation bias of the twCRPS and how one can derive a multivariate version of the twCRPS for the centre indicator function.

The censored likelihood score also appears in a second important result of this paper. In particular, we have shown that the UMP test for a localised version of the standard simple versus simple Neyman Pearson testing problem is based on the censored likelihood ratio. Furthermore, the results of our Monte Carlo study suggest that our theoretical findings spill over to the finite sample properties of other forecast evaluation tests. In our experiments, striking differences in power almost always favour censoring. In the empirical application, we identified an enhanced elimination of poorly performing models during the MCS process when using censoring instead of conditioning. Moreover, models that resulted in superior VaR and ES backtesting outcomes were more frequently incorporated in the MCS process when based on censoring.



# Appendix

## A Proofs

### A.1 Proof Theorem 2

For clarity of exposition, we first prove the main ingredients of the proof via two isolated lemmas and a corollary.

**Lemma 2.** *Consider the censored scoring rule defined in Definition 6.  $\forall w \in \mathcal{W}$  and  $H \in \mathcal{H}$ , the following identity holds  $\int_{\mathcal{Y}} S_{w,H}^b(F, y)P(dy) = \int_{\mathcal{Y}} S(F_{w,H}^b, y)P_{w,H}^b(dy)$ .*

*Proof.*

$$\begin{aligned}
\int_{\mathcal{Y}} S_{w,H}^b(F, y)P(dy) &= \int_{\mathcal{Y}} \left( w(y)S(F_{w,H}^b, y) + (1 - w(y)) \int_{\mathcal{Y}} S(F_{w,H}^b, q)H(dq) \right) P(dy), \\
&= \int_{\mathcal{Y}} w(y)S(F_{w,H}^b, y)P(dy) + \int_{\mathcal{Y}} S(F_{w,H}^b, q) \int_{\mathcal{Y}} (1 - w(y))P(dy)H(dq), \\
&= \int_{\mathcal{Y}} S(F_{w,H}^b, y)P_w(dy) + \int_{\mathcal{Y}} S(F_{w,H}^b, y)\bar{P}_w H(dy), \\
&= \int_{\mathcal{Y}} S(F_{w,H}^b, y)(P_w(dy) + \bar{P}_w H(dy)), \\
&= \int_{\mathcal{Y}} S(F_{w,H}^b, y)P_{w,H}^b(dy).
\end{aligned}$$

□

**Lemma 3.** *Consider two distributions  $P$  and  $F$  on the same measurable space  $(\mathcal{Y}, \mathcal{G})$ . On the same space, let their censored counterparts  $P_{w,H}^b$  and  $F_{w,H}^b$  be given by Definition 6. Then,*

$$F_{w,H}^b(E) = G_{w,H}^b(E), \quad \forall E \in \mathcal{G} \iff F(E \cap \{w > 0\}) = G(E \cap \{w > 0\}), \quad \forall E \in \mathcal{G}.$$

*Proof.* “ $\implies$ ” We start with the most challenging direction, for which Assumption 1

is of critical importance. First, note that

$$\begin{aligned}
F_{w,H}^b(E) &= G_{w,H}^b(E), \quad \forall E \in \mathcal{G} \\
\implies F_{w,H}^b(E \cap \{w = c\}) &= G_{w,H}^b(E \cap \{w = c\}), \quad \forall E \in \mathcal{G} \\
\implies \int_{\mathcal{Y}} (1-w) dF_H(E \cap \{w = c\}) &= \int_{\mathcal{Y}} (1-w) dG_H(E \cap \{w = c\}), \quad \forall E \in \mathcal{G} \\
\implies \int_{\mathcal{Y}} (1-w) dF_H(\{w = c\}) &= \int_{\mathcal{Y}} (1-w) dG_H(\{w = c\}), \\
\implies \int_{\mathcal{Y}} (1-w) dF &= \int_{\mathcal{Y}} (1-w) dG,
\end{aligned}$$

where  $c$  denotes a constant such that Assumption 1 is satisfied. Then, exploit this equality to conclude

$$\begin{aligned}
F_{w,H}^b(E) &= G_{w,H}^b(E), \quad \forall E \in \mathcal{G} \\
\implies \int_{\mathcal{Y}} w(y) \mathbb{1}_{y \in E} dF &= \int_{\mathcal{Y}} w(y) \mathbb{1}_{y \in E} dG, \quad \forall E \in \mathcal{G} \\
\implies F(E \cap \{w > 0\}) &= G(E \cap \{w > 0\}), \quad \forall E \in \mathcal{G}.
\end{aligned}$$

“ $\Leftarrow$ ” The other direction is somewhat trivial. Indeed,

$$\begin{aligned}
F(E \cap \{w > 0\}) &= G(E \cap \{w > 0\}), \quad \forall E \in \mathcal{G} \\
\implies \int_{\mathcal{Y}} w(y) \mathbb{1}_{y \in E} dF &= \int_{\mathcal{Y}} w(y) \mathbb{1}_{y \in E} dG, \quad \forall E \in \mathcal{G} \\
\implies \int_{\mathcal{Y}} (1-w) dF &= \int_{\mathcal{Y}} (1-w) dG,
\end{aligned}$$

and the two implied results jointly imply  $F_{w,H}^b(E) = G_{w,H}^b(E)$ ,  $\forall E \in \mathcal{G}, \forall H \in \mathcal{H}$ .  $\square$

**Corollary 4.** *The censored scoring rule defined in Definition 6 is localising  $\forall H \in \mathcal{H}$ .*

*Proof.* Suppose that  $F(E \cap \{w > 0\}) = G(E \cap \{w > 0\})$ ,  $\forall E \in \mathcal{G}$ . Then, by Lemma 3,  $F_{w,H}^b(E) = G_{w,H}^b(E)$ ,  $\forall E \in \mathcal{G}$ , whence it follows that  $S_{w,H}^b(P, y) = S_{w,H}^b(F, y)$ ,  $\forall y \in \mathcal{Y}$ .  $\square$

We now turn to the main body of the proof. The definition of a strictly locally

proper scoring rule (Definition 4) and the definitions on which this definition is built, that is, the definition of a locally proper scoring rule (Definition 4) and a localising weighted scoring rule (Definition 3), reveal that we need to prove a list of three things  $\forall H \in \mathcal{H}$ : (i)  $S_{w,H}^b(P, y)$  must be localising relative to  $\mathcal{W}$ , (ii)  $S_{w,H}^b(P, y)$  must be proper relative to  $\mathcal{P}$ ,  $\forall w \in \mathcal{W}$  and (iii) the if and only if statement in Definition 4. We prove them one by one.

(i)  $S_{w,H}^b(P, y)$  is localising relative to  $\mathcal{W}$ ,  $\forall H \in \mathcal{H}$ , by Corollary 4.

(ii) Fix an arbitrary  $w \in \mathcal{W}$  and  $H \in \mathcal{H}$ . Since  $\mathcal{P}_{w,H}^b \subseteq \mathcal{P}$ ,  $S$  is strictly proper relative to  $\mathcal{P}_{w,H}^b$ , i.e.

$$\int_{\mathcal{Y}} S(P_{w,H}^b, y) P_{w,H}^b(dy) \geq \int_{\mathcal{Y}} S(F_{w,H}^b, y) P_{w,H}^b(dy), \quad \forall P_{w,H}^b, F_{w,H}^b \in \mathcal{P}_{w,H}^b, \quad (\text{A.1})$$

which is by definition of the class  $\mathcal{P}_{w,H}^b \equiv \{[P]_{w,H}^b, P \in \mathcal{P}\}$  equivalent to

$$\int_{\mathcal{Y}} S([P]_{w,H}^b, y) [P]_{w,H}^b(dy) \geq \int_{\mathcal{Y}} S([F]_{w,H}^b, y) [F]_{w,H}^b(dy), \quad \forall P, F \in \mathcal{P}, \quad (\text{A.2})$$

and hence, by Lemma 2, also

$$\int_{\mathcal{Y}} S_{w,H}^b(P, y) P(dy) \geq \int_{\mathcal{Y}} S_{w,H}^b(F, y) P(dy), \quad \forall P, F \in \mathcal{P}. \quad (\text{A.3})$$

Therefore,  $S_{w,H}^b(P, y)$  is proper relative to  $\mathcal{P}$  by Definition 2.

(iii) Since  $S$  is strictly proper relative to  $\mathcal{P}^b$  and hence  $\mathcal{P}_{w,H}^b$ , it also follows that,  $\forall w \in \mathcal{W}$  and  $H \in \mathcal{H}$ ,

$$\int_{\mathcal{Y}} S(P_{w,H}^b, y) P_{w,H}^b(dy) = \int_{\mathcal{Y}} S(F_{w,H}^b, y) P_{w,H}^b(dy) \iff P_{w,H}^b = F_{w,H}^b,$$

and thus, by Lemma 3,

$$\int_{\mathcal{Y}} S(P_{w,H}^b, y) P_w^b(dy) = \int_{\mathcal{Y}} S(F_{w,H}^b, y) P_w^b(dy) \iff P(E \cap \{w > 0\}) = F(E \cap \{w > 0\}),$$

$\forall E \in \mathcal{G}$ , and hence, by Lemma 2, also

$$\int_{\mathcal{Y}} S_{w,H}^\flat(P, y) P(dy) = \int_{\mathcal{Y}} S_{w,H}^\flat(F, y) P(dy) \iff P(E \cap \{w > 0\}) = F(E \cap \{w > 0\}),$$

which is the desired if and only if statement of Definition 4.

But then, as we have verified each of the listed conditions (i) to (iii), we have shown that  $S_{w,H}^\flat(P, y)$  is strictly locally proper relative to  $(\mathcal{P}, \mathcal{W})$ ,  $\forall H \in \mathcal{H}$ .  $\square$

## A.2 Proof Lemma 1

Due to the integral over  $\mathcal{Y}(\mathcal{I}_{A^c})$ , any test  $\psi_{h_1}$  is constant in arguments varying in  $\mathcal{Y}(\mathcal{I}_{A^c})$ . We can use this observation to simplify the size of a test  $\psi_{h_1}$ . In particular,  $\forall h_1 \in \mathcal{H}$ , we have that

$$\begin{aligned} \sup_{p_0 \in \mathcal{P}_0} \mathbb{E}_{p_0} \psi_{h_1} &= \left( \prod_{t \in \mathcal{I}_{A^c}} F_0(A_t^c) \right) \sup_{h_0 \in \mathcal{H}} \int_{\mathcal{Y}^T} \psi_{h_1} \prod_{t \in \mathcal{I}_A} f_{0t} \mathbb{1}_{A_t} \prod_{t \in \mathcal{I}_{A^c}} [h_{0t}]_{A_t^c}^\# \mathbb{1}_{A_t^c} d\mu_t \\ &= \left( \prod_{t \in \mathcal{I}_{A^c}} F_0(A_t^c) \right) \int_{\mathcal{Y}(\mathcal{I}_A)} \psi_{h_1} \prod_{t \in \mathcal{I}_A} f_{0t} \mathbb{1}_{A_t} d\mu_t \\ &= \left( \prod_{t \in \mathcal{I}_{A^c}} F_0(A_t^c) \right) \int_{\mathcal{Y}^T} \phi_{h_1}^* \prod_{t \in \mathcal{I}_{A^c}} [h_{1t}]_{A_t^c}^\# \mathbb{1}_{A_t^c} d\mu_t \prod_{t \in \mathcal{I}_A} f_{0t} \mathbb{1}_{A_t} d\mu_t \\ &\leq \left( \prod_{t \in \mathcal{I}_{A^c}} F_0(A_t^c) \right) \sup_{h_0 \in \mathcal{H}} \int_{\mathcal{Y}^T} \phi_{h_1}^* \prod_{t \in \mathcal{I}_{A^c}} [h_{0t}]_{A_t^c}^\# \mathbb{1}_{A_t^c} d\mu_t \prod_{t \in \mathcal{I}_A} f_{0t} \mathbb{1}_{A_t} d\mu_t \\ &= \sup_{p_0 \in \mathcal{P}_0} \mathbb{E}_{p_0} \phi_{h_1}^* \\ &\leq \alpha, \end{aligned}$$

since  $\phi_{h_1}^* \in \Phi(\alpha)$ . Hence,  $\psi_{h_1} \in \Phi(\alpha)$ .  $\square$

## A.3 Proof Corollary 1

Fix an arbitrary  $h_1 \in \mathcal{H}$ . Since  $\Psi(\alpha) \subseteq \Phi(\alpha)$ , we trivially have that  $\max_{\phi \in \Phi(\alpha)} \mathbb{E}_{p_1} \phi \geq \max_{\psi \in \Psi(\alpha)} \mathbb{E}_{p_1} \psi$ . Now suppose that  $\max_{\phi \in \Phi(\alpha)} \mathbb{E}_{p_1} \phi < \max_{\psi \in \Psi(\alpha)} \mathbb{E}_{p_1} \psi$ . Then, we can always define the test  $\tilde{\psi} = \int_{\mathcal{Y}(\mathcal{I}_{A^c})} \phi^* \prod_{t \in \mathcal{I}_{A^c}} [h_{1t}]_{A_t^c}^\# \mathbb{1}_{A_t^c} d\mu_t$ , with  $\phi^* \in \arg \max_{\phi \in \Phi(\alpha)} \mathbb{E}_{p_1} \phi$ ,

satisfying  $\mathbb{E}_{p_1} \phi^* = \mathbb{E}_{p_1} \tilde{\psi}$ . But, by Lemma 1,  $\tilde{\psi} \in \Psi(\alpha)$ , in which case  $\max_{\phi \in \Phi(\alpha)} \mathbb{E}_{p_1} \phi = \max_{\psi \in \Psi(\alpha)} \mathbb{E}_{p_1} \tilde{\psi}$ , contradicting the assumed strict inequality.  $\square$

#### A.4 Proof Theorem 3

For any fixed  $h_1 \in \mathcal{H}$ , the most powerful test of size  $\alpha$  is a solution to the following restricted maximisation problem

$$\begin{aligned}
\max_{\phi \in \Phi(\alpha)} \mathbb{E}_{p_1} \phi &= \max_{\alpha \in \Delta_{\bar{T}}(\alpha)} \sum_{k=0}^T \sum_{s=1}^{\binom{T}{k}} \max_{\phi_{k,s} \in \Phi(\alpha_{k,s})} \mathbb{E}_{p_1} (\phi_{k,s} | y_t \in A_t, \forall i \in \mathcal{I}_A(k, s) \wedge y_t \in A_t^c, \forall i \in \mathcal{I}_{A^c}(k, s)) \\
&= \max_{\alpha \in \Delta_{\bar{T}}(\alpha)} \sum_{k=0}^T \sum_{s=1}^{\binom{T}{k}} \max_{\phi_{k,s} \in \Psi(\alpha_{k,s})} \mathbb{E}_{p_1} (\phi_{k,s} | y_t \in A_t, \forall i \in \mathcal{I}_A(k, s) \wedge y_t \in A_t^c, \forall i \in \mathcal{I}_{A^c}(k, s)) \\
&= \max_{\alpha \in \Delta_{\bar{T}}(\alpha)} \sum_{k=0}^T \sum_{s=1}^{\binom{T}{k}} \max_{\phi_{k,s} \in \Psi(\alpha_{k,s})} \left( \prod_{t \in \mathcal{I}_{A^c}} F_1(A_t^c) \right) \int_{\mathcal{Y}(\mathcal{I}_A)} \phi_{k,s} \prod_{t \in \mathcal{I}_A} f_{1t} \mathbb{1}_{A_t} d\mu_t \\
&= \max_{\alpha \in \Delta_{\bar{T}}(\alpha)} \sum_{k=0}^T \sum_{s=1}^{\binom{T}{k}} \max_{\phi_{k,s} \in \Psi(\alpha_{k,s})} \int_{\mathcal{Y}^T} \phi_{k,s} \prod_{t=0}^{T-1} d[F_t]_{A_t}^b \\
&= \max_{\alpha \in \Delta_{\bar{T}}(\alpha)} \sum_{k=0}^T \sum_{s=1}^{\binom{T}{k}} \max_{\phi_{k,s} \in \Phi(\alpha_{k,s})} \int_{\mathcal{Y}^T} \phi_{k,s} \prod_{t=0}^{T-1} d[F_t]_{A_t}^b,
\end{aligned}$$

where  $\bar{T} = \sum_{k=0}^T \binom{T}{k}$  and  $\Delta_{\bar{T}}(\alpha_0) = \{\alpha_0 \in [0, \alpha_0]^{\bar{T}} : \iota'_{\bar{T}} \alpha_0 = \alpha_0\}$ , with  $\iota_{\bar{T}}$  denoting column vector of ones of length  $\bar{T}$ . The first equality exploits that the test function can be decomposed into test functions operating on a single part of the partitioning of the outcome space  $\mathcal{Y}^T$ , in which case the maximisation problem can be split into finding an optimal test on each of the partitioned parts conditional on the amount of size spent on each part and the optimal distribution of size over the partition of the outcome space.

The second equality holds by Corollary 1, the third equality uses that the optimal test is constant in arguments varying in  $A^c$ , the fourth equality holds by definition of the censored measure and the fifth equality uses that all tests that are non-constant in arguments varying in  $A^c$  map under the censored measure onto tests that are constant in arguments varying in  $A^c$ .

Finally, the result follows by observing that the final maximisation problem is equivalent to finding the optimal test  $\phi_A^b$  for the testing problem  $\mathbb{H}_j : p_j = \prod_{t=0}^{T-1} [f_j]_{A_t}^b$ ,  $j \in \{0, 1\}$ , for which  $\phi_A^b$  is the UMP test by the Fundamental Lemma of [Neyman and Pearson \(1933\)](#). By the equivalence,  $\phi_A^b$  is, for any  $h_1 \in \mathcal{H}$ , also the most powerful test for testing problem (6). But, since the test  $\phi^b$  is independent of  $h_1$ , it is the UMP test for testing problem (6). □

## References

- Amisano, G. and R. Giacomini (2007). Comparing density forecasts via weighted likelihood ratio tests. *Journal of Business & Economic Statistics* 25(2), 177–190.
- Bernoulli, D. (1760). Essai d’une nouvelle analyse de la mortalité causée par la petite vérole, et des avantages de l’inoculation pour la prévenir. *Histoire de l’Acad., Roy. Sci.(Paris) avec Mem*, 1–45.
- Bollerslev, T. (1987). A conditionally heteroskedastic time series model for speculative prices and rates of return. *The review of economics and statistics*, 542–547.
- Borowska, A., L. Hoogerheide, S. Koopman, and H. van Dijk (2020). Partially censored posterior for robust and efficient risk evaluation. *Journal of Econometrics* 217(2), 335–355.
- Bregman, L. M. (1967). The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR Computational Mathematics and Mathematical Physics* 7(3), 200–217.
- Brehmer, J. and T. Gneiting (2020). Properization: constructing proper scoring rules via bayes acts. *Annals of the Institute of Statistical Mathematics* 72(3), 659–673.
- Brier, G. W. (1950). Verification of forecasts expressed in terms of probability. *Monthly weather review* 78(1), 1–3.
- Dawid, A. P. (1984). Statistical theory: The prequential approach. *Journal of the Royal Statistical Society: Series A (General)* 147(2), 278–290.
- Dawid, A. P. (2007). The geometry of proper scoring rules. *Annals of the Institute of Statistical Mathematics* 59(1), 77–93.
- Diks, C. and H. Fang (2020). Comparing density forecasts in a risk management context. *International Journal of Forecasting* 36(2), 531–551.
- Diks, C., V. Panchenko, O. Sokolinskiy, and D. van Dijk (2014). Comparing the accuracy of multivariate density forecasts in selected regions of the copula support. *Journal of Economic Dynamics and Control* 48, 79–94.
- Diks, C., V. Panchenko, and D. van Dijk (2011). Likelihood-based scoring rules for comparing density forecasts in tails. *Journal of Econometrics* 163(2), 215–230.
- Eguchi, S. et al. (1985). A differential geometric approach to statistical inference on the basis of contrast functionals. *Hiroshima mathematical journal* 15(2), 341–391.
- Ehm, W. and T. Gneiting (2012, FEB). Local proper scoring rules of order two. *Annals of Statistics* 40(1), 609–637.
- Fisher, R. A. (1922). On the mathematical foundations of theoretical statistics. *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character* 222(594-604), 309–368.

- Giacomini, R. and H. White (2006). Tests of conditional predictive ability. *Econometrica* 74(6), 1545–1578.
- Gneiting, T. and A. Raftery (2007). Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association* 102(477), 359–378.
- Gneiting, T. and R. Ranjan (2011). Comparing density forecasts using threshold- and quantile-weighted scoring rules. *Journal of Business & Economic Statistics* 29(3), 411–422.
- Good, I. (1952). Rational decisions. *Journal of the Royal Statistical Society* 14, 107–114.
- Good, I. (1971). Comment on “measuring information and uncertainty.”. *Foundation of Statistical Inference*, 265–273.
- Hansen, P. R., Z. Huang, and H. H. Shek (2012, SEP-OCT). Realized garch: a joint model for returns and realized measures of volatility. *Journal of Applied Econometrics* 27(6, SI), 877–906.
- Hansen, P. R., A. Lunde, and J. M. Nason (2011). The model confidence set. *Econometrica* 79(2), 453–497.
- Hersbach, H. (2000, OCT). Decomposition of the continuous ranked probability score for ensemble prediction systems. *Weather and Forecasting* 15(5), 559–570.
- Holzmann, H. and B. Klar (2016). Weighted scoring rules and hypothesis testing. Available at arXiv:1611.07345v2.
- Holzmann, H. and B. Klar (2017). Focusing on regions of interest in forecast evaluation. *Annals of applied statistics* 11(4), 2404–2431.
- Iacopini, M., F. Ravazzolo, and L. Rossini (2022). Proper scoring rules for evaluating density forecasts with asymmetric loss functions. *Journal of Business & Economic Statistics*, 1–15.
- Jose, V. R. (2009). A characterization for the spherical scoring rule. *Theory and Decision* 66(3), 263–281.
- Jose, V. R. R., R. F. Nau, and R. L. Winkler (2008). Scoring rules, generalized entropy, and utility maximization. *Operations research* 56(5), 1146–1157.
- Kole, E., T. Markwat, A. Opschoor, and D. Van Dijk (2017). Forecasting value-at-risk under temporal and portfolio aggregation. *Journal of Financial Econometrics* 15(4), 649–677.
- Kullback, S. and R. A. Leibler (1951). On information and sufficiency. *The annals of mathematical statistics* 22(1), 79–86.
- Lerch, S., T. L. Thorarinsdottir, F. Ravazzolo, and T. Gneiting (2017). Forecaster’s dilemma: Extreme events and forecast evaluation. *Statistical Science*, 106–127.



- Liese, F. and I. Vajda (2006). On divergences and informations in statistics and information theory. *IEEE Transactions on Information Theory* 52(10), 4394–4412.
- Matheson, J. and R. Winkler (1976). Scoring rules for continuous probability distributions. *Management Science* 22(10), 1087–1096.
- Merton, R. C. (1971). Optimum consumption and portfolio rules in a continuous-time model. pp. 621–661. Elsevier.
- Neyman, J. and E. S. Pearson (1933). Ix. on the problem of the most efficient tests of statistical hypotheses. *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character* 231(694-706), 289–337.
- Opschoor, A., D. van Dijk, and M. van der Wel (2017). Combining density forecasts using focused scoring rules. *Journal of Applied Econometrics* 32(7), 1298–1313.
- Ovcharov, E. (2018). Proper scoring rules and bregman divergence. *Bernoulli* 24(1), 53–79.
- Painsky, A. and G. W. Wornell (2019). Bregman divergence bounds and universality properties of the logarithmic loss. *IEEE Transactions on Information Theory* 66(3), 1658–1673.
- Patton, A. J. (2020). Comparing possibly misspecified forecasts. *Journal of Business & Economic Statistics* 38(4), 796–809.
- Roby, T. B. (1964). Belief states: A preliminary empirical study. Technical report, Tufts University Medford MA.
- Savage, L. J. (1971). Elicitation of personal probabilities and expectations. *Journal of the American Statistical Association* 66(336), 783–801.
- Selten, R. (1998). Axiomatic characterization of the quadratic scoring rule. *Experimental Economics* 1(1), 43–61.
- Shuford, E. H., A. Albert, and H. E. Massengill (1966). Admissible probability measurement procedures. *Psychometrika* 31(2), 125–145.
- Székely, G. J. and M. L. Rizzo (2005). A new test for multivariate normality. *Journal of Multivariate Analysis* 93(1), 58–80.
- Tobin, J. (1958). Estimation of relationships for limited dependent variables. *Econometrica*, 24–36.
- Toda, M. (1963). Measurement of subjective probability distribution. rep. no. 3, 1963, state college, pennsylvania. *Institute for Research*.

# Supplementary Material: A General Procedure for Localising Strictly Proper Scoring Rules

Ramon de Punder\*      Cees Diks\*      Roger Laeven\*      Dick van Dijk<sup>†</sup>

June 23, 2023

## **Abstract**

This document contains additional proofs, derivations and examples of the paper ‘Localising Strictly Proper Scoring Rules’.

---

\*University of Amsterdam, Tinbergen Institute

<sup>†</sup>Erasmus University Rotterdam, Tinbergen Institute

# Contents

A	Additional Proofs . . . . .	2
A.1	Proof censored density . . . . .	2
A.2	Proof Corollary 2 . . . . .	3
A.3	Proof Corollary 3 . . . . .	3
B	Derivations Table 1 . . . . .	4
B.1	LogS . . . . .	5
B.2	PsSphS $_{\alpha}$ . . . . .	5
B.3	PowS $_{\alpha}$ . . . . .	8
C	Examples . . . . .	10
C.1	Propriety CRPS* for Example 2 . . . . .	10
C.2	Localised NP for $T = 1$ . . . . .	11
C.3	CRPS . . . . .	12

## A Additional Proofs

### A.1 Proof censored density

We defined the measures  $\mu$  and  $F$  to the extended measurable space  $(\mathcal{Y}^*, \mathcal{G}^*)$  by putting  $\mu^*(E) = \mu(E \setminus \{*\})$ , if  $* \in E$  and  $\mu^*(E) = \mu(E)$ , otherwise,  $\forall E \in \mathcal{G}^*$ . To simplify the notation, we drop the subscript  $*$  in the notation of the extended measures, while still considering all measures with respect to the extended measurable space  $(\mathcal{Y}^*, \mathcal{G}^*)$

Since  $(\mu + \delta_*)(E) = 0$  implies that both  $\mu(E) = 0$  and  $\delta_*(E) = 0$ ,  $\forall E \in \mathcal{G}^*$ , we have that both  $\mu \ll \mu + \delta_*$  and  $\delta_* \ll \mu + \delta_*$ . As a consequence,

$$f_{w,h}^b := \frac{dF_w^b}{d(\mu + \delta_*)} = w \frac{dF}{d(\mu + \delta_*)} + \bar{F}_w \frac{d\delta_*}{d(\mu + \delta_*)}$$

is the censored  $(\mu + \delta_*)$ -density of  $F_w^b$ .

We can simplify this density as follows. Understand

Put together, we arrive at

$$f_{w,h}^b(y) = w(y) \frac{dF}{d\mu}(y) \mathbb{1}_{\mathcal{Y} \setminus \{*\}}(y) + \bar{F}_w \mathbb{1}_*(y) = w(y) f(y) \mathbb{1}_{y \neq *} + \bar{F}_w \mathbb{1}_{y=*}, \quad y \in \mathcal{Y},$$

where  $f$  denotes the  $\mu$ -density of  $F$ . □

## A.2 Proof Corollary 2

The test based on  $\tilde{\lambda}(\mathbf{y})$  is equivalent to the UMP test in Theorem 3, since

$$\begin{aligned} \tilde{\lambda}(\mathbf{y}) &= \sum_{t=0}^{T-1} (S_{A_t}^{\text{csl}}(f_{1t}, y_{t+1}) - S_{A_t}^{\text{csl}}(f_{0t}, y_{t+1})) \\ &= \sum_{t=0}^{T-1} \left( \log \left( [f_{1t}]_{A_t}^b(y_{t+1}) \right) - \log \left( [f_{0t}]_{A_t}^b(y_{t+1}) \right) \right) \\ &= \log \lambda(\mathbf{y}) \end{aligned}$$

and hence  $\lambda(\mathbf{y}) \underset{<}{\geq} c \iff \tilde{\lambda}(\mathbf{y}) \underset{<}{\geq} \tilde{c}$ , with  $\tilde{c} = \log c$ .

## A.3 Proof Corollary 3

We show that  $\phi_A^\sharp$  is not UMP by a specific counterexample in which the power of  $\phi_A^\sharp$  is strictly smaller than the power of  $\phi_A^b$ . In particular, suppose that  $T = 1$  and consider two densities  $f_0$  and  $f_1$  that are different on  $A = [r, \infty)$ , for some constant  $r > 0$ . Furthermore, assume that

$$\int_{\{y: \lambda(y) > r\}} F_0(dy) > \alpha, \quad \lambda(y) = \frac{f_1(y)}{f_0(y)}. \quad (\text{A.1})$$

For  $T = 1$ , the likelihood ratios of the conditional and censored test simplify to

$$\begin{aligned} \lambda_A^\sharp(y) &= \frac{[f_1]_A^\sharp(y)}{[f_0]_A^\sharp(y)} = \frac{\frac{f_1(y)}{F_1(A)}}{\frac{f_0(y)}{F_0(A)}} \mathbb{1}_A(y) = \frac{F_0(A^c)}{F_1(A^c)} \frac{f_1(y)}{f_0(y)} \mathbb{1}_A(y) \\ \lambda_A^b(y) &= \frac{[f_1]_A^b(y)}{[f_0]_A^b(y)} = \frac{f_1(y)}{f_0(y)} \mathbb{1}_A(y) + \frac{F_1(A^c)}{F_0(A^c)} \mathbb{1}_{A^c}(y). \end{aligned}$$

Due to restriction (A.1), the corresponding critical regions  $C^\sharp = [c^\sharp, \infty)$  and  $C^\flat = [c^\flat, \infty)$  are both contained in  $A$ . Hence, an example in which  $\sharp$  has higher power than  $\flat$ , would not only be a counterexample to Theorem 3 but also to the fundamental lemma of Neyman and Pearson (1933).

There exist many examples for which the power of the censored test is strictly larger than the power of the conditional test. For instance, suppose that  $y \sim \text{Exp}(\theta_j)$ ,  $j \in \{0, 1\}$ , with  $\theta_0 > \theta_1$ . Then, the critical regions follow from the equation

$$\alpha = \int_{\{y: \lambda(y) > c^*\}} \theta_0 e^{-\theta_0 y} dy = \int_{\{y: a^* \left(\frac{\theta_1}{\theta_0}\right) e^{-(\theta_1 - \theta_0)y} > c^*\}} \theta_0 e^{-\theta_0 y} dy = 1 - F_0 \left( \frac{1}{\theta_0 - \theta_1} \log \left( \frac{\theta_0}{\theta_1} \right) \frac{c^*}{a^*} \right),$$

where  $a^\sharp = \frac{1 - F_0(r)}{1 - F_1(r)} = e^{-(\theta_0 - \theta_1)r}$  and  $a^\flat = 1$ . Isolating  $c^*$ , gives

$$c^* = ba^*, \quad b = \frac{\theta_0}{\theta_1} e^{(\theta_0 - \theta_1)F_0^{-1}(1-\alpha)} > 0.$$

Now, the power of the conditional test is only weakly larger than the power of the censored test, if

$$\int_{\{y: \lambda(y) > c^\sharp\}} \theta_1 e^{-\theta_1 y} dy \geq \int_{\{y: \lambda(y) > c^\flat\}} \theta_1 e^{-\theta_1 y} dy \iff c^\sharp \geq c^\flat \iff (\theta_0 - \theta_1)r \leq 0.$$

But then, as  $\theta_0 > \theta_1$  and  $r > 0$ , it follows that the power of the conditional test is always strictly smaller than the power of the censored test. Consequently, the conditional test  $\phi_A^\sharp$  is not UMP.  $\square$

## B Derivations Table 1

The results in Table 1 hold under the assumption that all  $F \in \mathcal{P}$  are Borel measures on  $\mathbb{R}^d$  satisfying  $F(r) = 0$ , with  $r \in \mathbb{R}^d$ . Furthermore, the assumption on  $h$  and  $w$  in the generalised censored scoring rule examples can predominantly be simplified due to the observation that  $w(y)h(y) = 0$ ,  $\forall y \in \mathcal{Y}$  and  $f_w(y) = 0$ ,  $\forall y \in \{w = 0\}$ .

## B.1 LogS

$$\text{LogS}(\tilde{f}, \tilde{y}) = \log \tilde{f}(\tilde{y}) = \log f(y) - \log |b| \stackrel{\text{eqv.}}{=} \log f(y).$$

$$\begin{aligned} \text{LogS}_w^\sharp(f, y) &= w(y) \log \left( \frac{w(y)f(y)}{1 - \bar{F}_w} \right) \\ &\stackrel{\text{eqv.}}{=} w(y) \log \left( \frac{f(y)}{1 - \bar{F}_w} \right) \\ &= S_w^{\text{cl}}(f, y), \end{aligned}$$

$$\begin{aligned} \text{LogS}_w^\flat(f, y) &= w(y) \log (w(y)f(y)\mathbb{1}_{y \neq r} + \bar{F}_w\mathbb{1}_{y=r}) + (1 - w(y)) \log \bar{F}_w \\ &\stackrel{\mu\text{-a.e.}}{=} w(y) \log (w(y)f(y)) + (1 - w(y)) \log \bar{F}_w \\ &\stackrel{\text{eqv.}}{=} w(y) \log (f(y)) + (1 - w(y)) \log \bar{F}_w \\ &= S_w^{\text{csl}}(f, y). \end{aligned}$$

$$\begin{aligned} \text{LogS}_{w,h}^\flat(f, y) &= w(y) \log f_{w,h}^\flat(y) + (1 - w(y)) \int_{\mathcal{Y}} \log f_{w,h}^\flat(q) h(q) \mathrm{d}q, \\ &= w(y) \left( \log (f_w(y)) \mathbb{1}_{w>0} + \log (\bar{F}_w h(y)) \mathbb{1}_{w=0} \right) \\ &\quad + (1 - w(y)) \int_{\{w=0\}} \left( \log (f_w(q)) \mathbb{1}_{w>0} + \log (\bar{F}_w h(q)) \mathbb{1}_{w=0} \right) h(q) \mathrm{d}q, \\ &= w(y) \log f_w(y) + (1 - w(y)) \int_{\{w=0\}} \log (\bar{F}_w h(q)) h(q) \mathrm{d}q, \\ &\stackrel{\text{eqv.}}{=} w(y) \log f(y) + (1 - w(y)) \log \bar{F}_w, \\ &= S^{\text{csl}}(f, y). \end{aligned}$$

## B.2 PsSphS $_\alpha$

$$\text{PsSphS}_\alpha(\tilde{f}, \tilde{y}) = \frac{\tilde{f}(\tilde{y})}{\|\tilde{f}\|_\alpha^{\alpha-1}} = \frac{\left(\frac{1}{|b|}\right)^{\alpha-1} f(y)^{\alpha-1}}{\left(\frac{1}{|b|}\right)^{\frac{(\alpha-1)^2}{\alpha}} \|f\|_\alpha^{\alpha-1}} = \left(\frac{1}{|b|}\right)^{\frac{\alpha-1}{\alpha}} \text{PsSphS}_\alpha(f, y).$$

Next, we show the limit. Rescaling the  $\text{PsSphS}_\alpha$  family by a factor  $\frac{1}{\alpha-1}$ , we obtain

$$\begin{aligned}
\lim_{\alpha \downarrow 1} \frac{1}{\alpha-1} \left( \frac{f(y)}{\|f\|_\alpha} \right)^{\alpha-1} &= \lim_{\alpha \downarrow 1} \frac{(\alpha-1) \left( \frac{f(y)}{\|f\|_\alpha} \right)^{\alpha-1}}{(\alpha-1)^2} \\
&= \lim_{\alpha \downarrow 1} \frac{\left( \frac{f(y)}{\|f\|_\alpha} \right)^{\alpha-1} + (\alpha-1) \left( \log \left( \frac{f(y)}{\|f\|_\alpha} \right) + (\alpha-1) \left( \frac{f(y)}{\|f\|_\alpha} \right)^{-1} \frac{\partial}{\partial \alpha} \frac{f(y)}{\|f\|_\alpha} \right) \left( \frac{f(y)}{\|f\|_\alpha} \right)^{\alpha-1}}{2(\alpha-1)} \\
&= \frac{1}{2} \lim_{\alpha \downarrow 1} \frac{1}{\alpha-1} \left( \frac{f(y)}{\|f\|_\alpha} \right)^{\alpha-1} + \frac{1}{2} \lim_{\alpha \downarrow 1} \log \left( \frac{f(y)}{\|f\|_\alpha} \right) \left( \frac{f(y)}{\|f\|_\alpha} \right)^{\alpha-1} \\
&\quad + \frac{1}{2} \lim_{\alpha \downarrow 1} (\alpha-1) \left( \frac{f(y)}{\|f\|_\alpha} \right)^{\alpha-2} \frac{\partial}{\partial \alpha} \frac{f(y)}{\|f\|_\alpha},
\end{aligned}$$

and hence

$$\lim_{\alpha \downarrow 1} \frac{1}{\alpha-1} \left( \frac{f(y)}{\|f\|_\alpha} \right)^{\alpha-1} = \log f(y), \tag{A.2}$$

since  $\|f\|_1 = 1$ . It might be helpful to note that the second equality in the first display follows from L'Hôpital's rule combined with the following derivative

$$\frac{\partial}{\partial \alpha} \left( \frac{f(y)}{\|f\|_\alpha} \right)^{\alpha-1} = \log \left( \left( \frac{f(y)}{\|f\|_\alpha} \right) + (\alpha-1) \left( \frac{f(y)}{\|f\|_\alpha} \right)^{-1} \frac{\partial}{\partial \alpha} \frac{f(y)}{\|f\|_\alpha} \right) \left( \frac{f(y)}{\|f\|_\alpha} \right)^{\alpha-1}.$$

For the conditional  $\text{PsSphS}_\alpha$  family, we find

$$\begin{aligned}
\text{PsSphS}_{\alpha,w}^\sharp(f, y) &= w(y) \frac{\left( \frac{f_w(y)}{1-F_w} \right)^{\alpha-1}}{\left( \int_{\mathcal{Y}} \left( \frac{f_w}{1-F_w} \right)^\alpha d\mu \right)^{\frac{\alpha-1}{\alpha}}} \\
&= w(y) \frac{f_w(y)^{\alpha-1}}{\|f_w\|_\alpha^{\alpha-1}} \\
&= w(y) \left( \frac{f_w(y)}{\|f_w\|_\alpha} \right)^{\frac{\alpha-1}{\alpha}}
\end{aligned}$$

By the close similarity with Equation (A.2), it is uncomplicated to obtain the following



limit

$$\begin{aligned}
\lim_{\alpha \downarrow 1} \frac{1}{\alpha - 1} \text{PsSphS}_{\alpha, w}^{\sharp}(f, y) &= w(y) \lim_{\alpha \downarrow 1} \frac{1}{\alpha - 1} \left( \frac{f_w(y)}{\|f_w\|_{\alpha}} \right)^{\alpha-1} \\
&= w(y) \log \left( \frac{f_w(y)}{\|f_w\|_1} \right) \\
&= w(y) \log f_w^{\sharp}(y) \\
&= \text{LogS}_w^{\sharp}(f, y),
\end{aligned}$$

since  $\|f_w\|_1 = \int_{\mathcal{Y}} w f d\mu = 1 - \bar{F}_w$ . Clearly, this result also follows directly from the linearity of limits, as

$$\lim_{\alpha \downarrow 1} \frac{1}{\alpha - 1} \text{PsSphS}_{\alpha}^{\sharp}(f, y) = w(y) \lim_{\alpha \downarrow 1} \frac{1}{\alpha - 1} \text{PsSphS}_{\alpha}(f_w^{\sharp}, y) = w(y) \log f_w^{\sharp}(y) = \text{LogS}_w^{\sharp}(f, y). \tag{A.3}$$

Moreover, for the censored  $\text{PsSphS}_{\alpha}$  family, it follows that

$$\begin{aligned}
\text{PsSphS}_w^{\flat}(f, y) &= \frac{w(y) (f_w(y) \mathbb{1}_{y \neq r} + \bar{F}_w \mathbb{1}_{y=r})^{\alpha-1} + (1 - w(y)) \bar{F}_w^{\alpha-1}}{\left( \int_{\mathcal{Y}} (f_w(y) \mathbb{1}_{y \neq r} + \bar{F}_w \mathbb{1}_{y=r})^{\alpha} (\mu + \delta_r)(dy) \right)^{\frac{\alpha-1}{\alpha}}} \\
&= \frac{w(y) (f_w(y)^{\alpha-1} \mathbb{1}_{y \neq r} + \bar{F}_w^{\alpha-1} \mathbb{1}_{y=r}) + (1 - w(y)) \bar{F}_w^{\alpha-1}}{\left( \int_{\mathcal{Y}} (f_w(y))^{\alpha} dy + \bar{F}_w^{\alpha} \right)^{\frac{\alpha-1}{\alpha}}} \\
&\stackrel{\mu\text{-a.e.}}{=} \frac{w(y) f_w(y)^{\alpha-1} + (1 - w(y)) \bar{F}_w^{\alpha-1}}{\left( \|f_w(y)\|_{\alpha}^{\alpha} + \bar{F}_w^{\alpha} \right)^{\frac{\alpha-1}{\alpha}}}.
\end{aligned}$$

For the limit of  $\alpha \downarrow 1$ , we cannot directly apply Equation (A.2) as we did for the

conditional case. Nevertheless, we obtain a similarly satisfying result, namely

$$\begin{aligned}
\lim_{\alpha \downarrow 1} \frac{1}{\alpha - 1} \text{PsSphS}_w^b(f, y) &= w(y) \lim_{\alpha \downarrow 1} \frac{1}{\alpha - 1} \left( \frac{f_w(y)}{(\|f_w\|_\alpha^\alpha + \bar{F}_w^\alpha)^{\frac{1}{\alpha}}} \right)^{\alpha-1} \\
&\quad + (1 - w(y)) \lim_{\alpha \downarrow 1} \frac{1}{\alpha - 1} \left( \frac{\bar{F}_w}{(\|f_w\|_\alpha^\alpha + \bar{F}_w^\alpha)^{\frac{1}{\alpha}}} \right)^{\alpha-1} \\
&= w(y) \left( \lim_{\alpha \downarrow 1} \log \left( \frac{f_w(y)}{(\|f_w\|_\alpha^\alpha + \bar{F}_w^\alpha)^{\frac{1}{\alpha}}} \right) \left( \frac{f_w(y)}{(\|f_w\|_\alpha^\alpha + \bar{F}_w^\alpha)^{\frac{1}{\alpha}}} \right)^{\alpha-1} \right. \\
&\quad \left. + \lim_{\alpha \downarrow 1} (\alpha - 1) \left( \frac{f_w(y)}{(\|f_w\|_\alpha^\alpha + \bar{F}_w^\alpha)^{\frac{1}{\alpha}}} \right)^{\alpha-2} \frac{\partial}{\partial \alpha} \frac{f_w(y)}{(\|f_w\|_\alpha^\alpha + \bar{F}_w^\alpha)^{\frac{1}{\alpha}}} \right) \\
&\quad + (1 - w(y)) \left( \lim_{\alpha \downarrow 1} \log \left( \frac{\bar{F}_w}{(\|f_w\|_\alpha^\alpha + \bar{F}_w^\alpha)^{\frac{1}{\alpha}}} \right) \left( \frac{\bar{F}_w}{(\|f_w\|_\alpha^\alpha + \bar{F}_w^\alpha)^{\frac{1}{\alpha}}} \right)^{\alpha-1} \right. \\
&\quad \left. + \lim_{\alpha \downarrow 1} (\alpha - 1) \left( \frac{\bar{F}_w}{(\|f_w\|_\alpha^\alpha + \bar{F}_w^\alpha)^{\frac{1}{\alpha}}} \right)^{\alpha-2} \frac{\partial}{\partial \alpha} \frac{\bar{F}_w}{(\|f_w\|_\alpha^\alpha + \bar{F}_w^\alpha)^{\frac{1}{\alpha}}} \right) \\
&= w(y) \log f_w(y) + (1 - w(y)) \log \bar{F}_w \\
&= \text{LogS}_w^b(f, y),
\end{aligned}$$

where we have used that  $\|f_w\|_1 + \bar{F}_w = 1 - \bar{F}_w + \bar{F}_w = 1$ .

### B.3 PowS $_\alpha$

$$\begin{aligned}
\text{PowS}_\alpha(\tilde{f}, \tilde{y}) &= \alpha (\tilde{f}(\tilde{y}))^{\alpha-1} - (\alpha - 1) \|\tilde{f}\|_\alpha^\alpha \\
&= \alpha \left( \frac{1}{|b|} \right)^{\alpha-1} f(y) - (\alpha - 1) \left( \frac{1}{|b|} \right)^{\alpha-1} \|f\|_\alpha^\alpha \\
&= \left( \frac{1}{|b|} \right)^{\alpha-1} \text{PowS}_\alpha(f, y),
\end{aligned}$$

as

$$\begin{aligned}
\|\tilde{f}\|_\alpha^\alpha &= \int_{\tilde{\mathcal{Y}}} \tilde{f}(\tilde{y})^\alpha \mu(d\tilde{y}) \\
&= \left(\frac{1}{|b|}\right)^{\alpha-1} \int_{\tilde{\mathcal{Y}}} \left(f\left(\frac{\tilde{y}-a}{b}\right)\right)^\alpha \frac{1}{|b|} \mu(d\tilde{y}) \\
&= \left(\frac{1}{|b|}\right)^{\alpha-1} \int_{\mathcal{Y}} (f(y))^\alpha \mu(dy) \\
&= \left(\frac{1}{|b|}\right)^{\alpha-1} \|f\|_\alpha^\alpha.
\end{aligned}$$

Next, we verify the limit for the non-focused family. Specifically,

$$\begin{aligned}
\lim_{\alpha \downarrow 1} \frac{1}{\alpha-1} \text{PowS}_\alpha &= \lim_{\alpha \downarrow 1} \frac{1}{\alpha-1} (\alpha f(y)^{\alpha-1} - (\alpha-1) \|f\|_\alpha^\alpha) \\
&= \lim_{\alpha \downarrow 1} \frac{(\alpha-1) \alpha f(y)^{\alpha-1}}{(\alpha-1)^2} - 1 \\
&= \lim_{\alpha \downarrow 1} \frac{\alpha f(y)^{\alpha-1} + (\alpha-1) f(y)^{\alpha-1} (1 + \alpha \log f(y))}{2(\alpha-1)} - 1 \\
&= \frac{1}{2} \left( \lim_{\alpha \downarrow 1} \frac{1}{\alpha-1} \alpha f(y)^{\alpha-1} - 1 \right) + \frac{1}{2} \left( \lim_{\alpha \downarrow 1} f(y)^{\alpha-1} (1 + \alpha \log f(y)) - 1 \right)
\end{aligned}$$

and hence

$$\lim_{\alpha \downarrow 1} \frac{1}{\alpha-1} \text{PowS}_\alpha(f, y) = \log f(y).$$

Furthermore, the conditional version of the  $\text{PowS}_\alpha$  family displayed in Table 1 is nothing but a direct application of the conditioning procedure. For the limit of the  $\text{PowS}_{\alpha,w}^\sharp$ , we recall Equation (A.3) and immediately conclude that  $\lim_{\alpha \downarrow 1} \frac{1}{\alpha-1} \text{PowS}_{\alpha,w}^\sharp(f, y) = \text{LogS}_w^\sharp(f, y)$ .

Turning to the censored focusing method, we recall from the analysis in Appendix B.2 that  $\|f_w^b\|_\alpha^\alpha = \|f_w(y)\|_\alpha^\alpha + \bar{F}_w^\alpha$ . Using this result, we obtain

$$\begin{aligned}
\text{PowS}_{\alpha,w}^b(f, y) &= w(y) \alpha (f_w(y) \mathbb{1}_{y \neq r} + \bar{F}_w \mathbb{1}_{y=c})^{\alpha-1} + (1-w(y)) \alpha \bar{F}_w^{\alpha-1} - (\alpha-1) \|f_w^b\|_\alpha^\alpha \\
&\stackrel{\mu\text{-a.e.}}{=} w(y) \alpha f_w(y)^{\alpha-1} + (1-w(y)) \alpha \bar{F}_w^{\alpha-1} - (\alpha-1) (\|f_w\|_\alpha^\alpha + \bar{F}_w^\alpha),
\end{aligned}$$

which bears the following limit

$$\begin{aligned}
\lim_{\alpha \downarrow 1} \frac{1}{\alpha - 1} \text{PowS}_{\alpha, w}^b(f, y) &= w(y) \lim_{\alpha \downarrow 1} \frac{(\alpha - 1) \alpha f_w(y)^{\alpha-1}}{(\alpha - 1)^2} + (1 - w(y)) \lim_{\alpha \downarrow 1} \frac{(\alpha - 1) \alpha \bar{F}_w^{\alpha-1}}{(\alpha - 1)^2} - 1 \\
&= \frac{1}{2} \lim_{\alpha \downarrow 1} \left( w(y) \left( \frac{1}{\alpha - 1} \alpha f_w(y)^{\alpha-1} - 1 \right) + (1 - w(y)) \left( \frac{1}{\alpha - 1} \alpha \bar{F}_w^{\alpha-1} - 1 \right) \right) \\
&\quad + \frac{1}{2} \lim_{\alpha \downarrow 1} \left( w(y) \left( f_w(y)^{\alpha-1} (1 + \alpha \log f_w(y)) - 1 \right) \right. \\
&\quad \left. + (1 - w(y)) \left( \bar{F}_w^{\alpha-1} (1 + \alpha \log \bar{F}_w) - 1 \right) \right).
\end{aligned}$$

Therefore,

$$\lim_{\alpha \downarrow 1} \frac{1}{\alpha - 1} \text{PowS}_{\alpha, w}^b(f, y) = w(y) \log f_w(y) + (1 - w(y)) \log \bar{F}_w = \text{LogS}_w^b(f, y).$$

## C Examples

### C.1 Propriety CRPS\* for Example 2

The extended metric  $d^* : \mathbb{R} \cup \{*\} \times \mathbb{R} \cup \{*\} \rightarrow \mathbb{R}$  is given by

$$d^*(x, y) = \begin{cases} |x - y|, & \text{if } x \in \mathbb{R}, y \in \mathbb{R} \\ d(x), & \text{if } x \in \mathbb{R}, y = * \\ d(y), & \text{if } x = *, y \in \mathbb{R} \\ 0, & \text{if } x = y = *. \end{cases}$$

It is worth mentioning that Theorem 1 of [Székely and Rizzo \(2005\)](#) does not require  $d(y)$  to be a metric. In particular, the triangle inequality does not need to hold. Therefore, the selected distance  $d^*(y) = |y - r_1| \vee |y - r_2|$ , as this implies a continuous  $d^*(x_w^b, y_w^b)$  on  $\{w > 0\} \cup \{*\} \times \{w > 0\} \cup \{*\}$ . Indeed, for any  $x < r$ , it follows that  $\lim_{y \downarrow r_1} d(x, y) = |x - r_1| = \lim_{y \rightarrow *} d(x, y)$ , while for any  $x > r$ , it follows that  $\lim_{y \uparrow r_2} d(x, y) = |x - r_2| = \lim_{y \rightarrow *} d(x, y)$ . Furthermore, let  $x_1, \dots, x_n$  be  $n$  distinct but arbitrary outcomes in  $\{w > 0\}$ . Consider the case, where one outcome is  $*$ , wlog

the last one. Then

$$\sum_{i=1}^n \sum_{j=1}^n a_i a_j d^*(x_i, x_j) = \sum_{i=1}^n \sum_{j=1}^n a_i a_j d(\tilde{x}_i, \tilde{x}_j) \leq 0, \quad \tilde{x}_i = x_i \mathbb{1}_{\{w>0\}}(x_i) + r_k \mathbb{1}_{\{*\}},$$

with strict equality if and only if  $a_i = 0$ ,  $\forall i$  since  $d$  is a strictly negative definite kernel on  $\mathbb{R} \times \mathbb{R}$  and  $x_1, \dots, x_{n-1}, r_k | x_i \in \{w > 0\} \forall i$  is a particular set of outcomes in  $\mathbb{R}$ ,  $\forall k \in \{1, 2\}$ . If none of the outcomes is  $*$ , we can borrow the strictly negative definiteness of  $d$  on  $\mathbb{R} \times \mathbb{R}$  in a similar way, by understanding that any set of arbitrary outcomes in  $\{w > 0\}$  is a particular set of arbitrary outcomes in  $\mathbb{R}$ . But then, Theorem 1 of [Székely and Rizzo \(2005\)](#) applies, whence the CRPS is also proper with respect to  $(\mathbb{R}^*, d^*)$ .

## C.2 Localised NP for $T = 1$

Consider the special case  $T = 1$ . For one observation, it is straightforward to derive a most powerful test on  $A^c$ . For any  $h_1 \in \mathcal{H}$ , the relevant maximisation problem simplifies to

$$\begin{aligned} \max_{\phi \in \Phi(\alpha)} \mathbb{E}_{p_1} \phi(y) &= \max_{\phi \in \Phi(\alpha)} \left\{ \mathbb{E}_{f_1} \phi_A(y) + F_1(A^c) \mathbb{E}_{[h_1]_{A^c}^\#} \phi_{A^c}(y) \right\} \\ &= \max_{\alpha_A \leq \alpha} \left\{ \max_{\phi_A \in \Phi_A(\alpha_A)} \{ \mathbb{E}_{f_1} \phi_A(y) \} + F_1(A^c) \max_{\phi_{A^c} \in \Phi_{A^c}(\alpha - \alpha_A)} \{ \mathbb{E}_{[h_1]_{A^c}^\#} \phi_{A^c}(y) \} \right\} \\ &= \max_{\alpha_A \leq \alpha} \left\{ \max_{\phi_A \in \Phi_A(\alpha_A)} \{ \mathbb{E}_{f_1} \phi_A(y) \} + F_1(A^c) \frac{\alpha - \alpha_A}{F_0(A^c)} \mathbb{1}_{A^c} \right\}. \end{aligned}$$

After all, rejecting with probability  $\frac{\alpha - \alpha_A}{F_0(A^c)}$  if  $y \in A^c$  is optimal since this is size correct and any more complicated test function  $\phi_{A^c}$  has lower power. This can be verified as follows. For all level  $\alpha - \alpha_A$  tests  $\phi_{A^c}$ , i.e.  $\phi_{A^c} \in \Phi_{A^c}(\alpha - \alpha_A)$ , we have that

$$\begin{aligned} F_1(A^c) \mathbb{E}_{[h_1]_{A^c}^\#} \phi_{A^c}(y) &\leq F_1(A^c) \sup_{h_1 \in \mathcal{H}} \{ \mathbb{E}_{[h_1]_{A^c}^\#} \phi_{A^c}(y) \} \\ &= F_1(A^c) \sup_{h_0 \in \mathcal{H}} \{ \mathbb{E}_{[h_0]_{A^c}^\#} \phi_{A^c}(y) \} \\ &\leq F_1(A^c) \frac{\alpha - \alpha_A}{F_0(A^c)}. \end{aligned}$$

Consequently, the test

$$\phi_{A^c}^*(y) = \frac{\alpha - \alpha_A}{F_0(A^c)}, \quad y \in A^c,$$

is most powerful against any other test  $\phi_{A^c}^*(y)$  of size  $\alpha - \alpha_A$ .

This solution, also documented by [Holzmann and Klar \(2016\)](#), coincides with the UMP test given by Theorem 3. Indeed, suppose that the size  $\alpha$  is such that  $\frac{F_1(A^c)}{F_0(A^c)} = c$ , i.e. not all of the size is spent on  $A$ , then the randomisation probability  $\gamma$  in Theorem 3 is such that

$$\alpha = \alpha_A + \gamma F_0 \left( \lambda(y) = \frac{F_1(A^c)}{F_0(A^c)} \right) = \alpha_A + \gamma F_0(A^c) \implies \gamma = \frac{\alpha - \alpha_A}{F_0(A^c)}.$$

### C.3 CRPS

$$\begin{aligned} \text{CRPS}_{w_1}^b(F, y) &= w_1(y) \text{CRPS}(F_{w_1}^b, y) + (1 - w_1(y)) \text{CRPS}(F_{w_1}^b, r) \\ &= \mathbb{1}_{(-\infty, r)}(y) \left( \int_{-\infty}^r (F(s) - \Delta_y(s))^2 ds + \int_r^\infty (1 - 1)^2 ds \right) \\ &\quad + (1 - \mathbb{1}_{(-\infty, r)}(y)) \left( \int_{-\infty}^r (F(s) - \Delta_r(s))^2 ds + \int_r^\infty (1 - 1)^2 ds \right) \\ &= \mathbb{1}_{(-\infty, r)}(y) \left( \int_{-\infty}^r (F(s) - \Delta_y(s))^2 ds + \int_r^\infty (1 - 1)^2 ds \right) \\ &\quad + (1 - \mathbb{1}_{(-\infty, r)}(y)) \left( \int_{-\infty}^r (F(s) - \Delta_y(s))^2 ds + \int_r^\infty (1 - 1)^2 ds \right) \\ &= \int_{-\infty}^\infty w_1(s) (F(s) - \Delta_y(s))^2 ds \end{aligned}$$

$$\begin{aligned}
\text{CRPS}_{w_2}^b(F, y) &= w_2(y) \text{CRPS}(F_{w_2}^b, y) + (1 - w_2(y)) \text{CRPS}(F_{w_2}^b, r) \\
&= \mathbb{1}_{(r, \infty)}(y) \left( \int_{-\infty}^r (0 - 0)^2 ds + \int_r^\infty (F(s) - \Delta_y(s))^2 ds \right) \\
&\quad + (1 - \mathbb{1}_{(r, \infty)}(y)) \left( \int_{-\infty}^r (0 - 0)^2 ds + \int_r^\infty (F(s) - \Delta_r(s))^2 ds \right) \\
&= \mathbb{1}_{(r, \infty)}(y) \left( \int_{-\infty}^r (0 - 0)^2 ds + \int_r^\infty (F(s) - \Delta_y(s))^2 ds \right) \\
&\quad + (1 - \mathbb{1}_{(r, \infty)}(y)) \left( \int_{-\infty}^r (0 - 0)^2 ds + \int_r^\infty (F(s) - \Delta_y(s))^2 ds \right) \\
&= \int_{-\infty}^\infty w_2(s) (F(s) - \Delta_y(s))^2 ds
\end{aligned}$$

## References

- Holzmann, H. and B. Klar (2016). Weighted scoring rules and hypothesis testing. Available at [arXiv:1611.07345v2](https://arxiv.org/abs/1611.07345v2).
- Neyman, J. and E. S. Pearson (1933). Ix. on the problem of the most efficient tests of statistical hypotheses. *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character* 231 (694-706), 289–337.
- Székely, G. J. and M. L. Rizzo (2005). A new test for multivariate normality. *Journal of Multivariate Analysis* 93(1), 58–80.