

No. 2011-061

**ROBUST SOLUTIONS OF OPTIMIZATION PROBLEMS  
AFFECTED BY UNCERTAIN PROBABILITIES**

By Aharon Ben-Tal, Dick den Hertog, Anja De Waegenare,  
Bertrand Melenberg, Gijs Rennen

May 25, 2011

ISSN 0924-7815

# Robust solutions of optimization problems affected by uncertain probabilities

Aharon Ben-Tal \*

*Department of Industrial Engineering and Management,  
Technion – Israel Institute of Technology, Haifa 32000, Israel  
CentER Extramural Fellow, CentER, Tilburg University, The Netherlands*

Dick den Hertog, Anja De Waegenare, Bertrand Melenberg, Gijs Rennen

*Department of Econometrics and Operations Research, Tilburg University,  
P.O. Box 90153, 5000 LE Tilburg, The Netherlands*

May 25, 2011

## Abstract

In this paper we focus on robust linear optimization problems with uncertainty regions defined by  $\phi$ -divergences (for example, chi-squared, Hellinger, Kullback-Leibler). We show how uncertainty regions based on  $\phi$ -divergences arise in a natural way as confidence sets if the uncertain parameters contain elements of a probability vector. Such problems frequently occur in, for example, optimization problems in inventory control or finance that involve terms containing moments of random variables, expected utility, etc. We show that the robust counterpart of a linear optimization problem with  $\phi$ -divergence uncertainty is tractable for most of the choices of  $\phi$  typically considered in the literature. We extend the results to problems that are nonlinear in the optimization variables. Several applications, including an asset pricing example and a numerical multi-item newsvendor example, illustrate the relevance of the proposed approach.

**Keywords:** robust optimization,  $\phi$ -divergence, goodness-of-fit statistics.

**JEL codes:** C61.

## 1 Introduction

Several papers in the late 1990s ([29], [3], [4], [19], [20]) started a revival of Robust Optimization (RO), both in terms of theoretical aspects, as well as practical applications. For a survey we refer to [5] or [8]. Consider, for example, a linear constraint with uncertain parameters. The idea of robust optimization is to define a so-called *uncertainty region* for the uncertain parameters, and then to require that the constraint should hold for all parameter values in this uncertainty region. The optimization problem modeling this requirement is called the *Robust Counterpart Problem* (RCP). Although the RCP typically has an infinite number of constraints, it is still tractable (polynomially solvable) for several optimization problems and several choices of the uncertainty region. In particular, the robust counterpart for a linear programming problem

---

\*Part of this work was done during a visit of the first author at CWI, Amsterdam, The Netherlands.

with polyhedral or ellipsoidal uncertainty regions reduces to a linear programming and a conic quadratic programming problem, respectively.

When applying the RO methodology to a practical problem, a major modeling decision concerns the choice of the uncertainty region  $U$ . Such a choice should fulfil three basic requirements. First,  $U$  should be consistent with whatever data (and information) is available on the uncertain parameters. Second,  $U$  should be statistically meaningful. Third,  $U$  should be such that the corresponding RCP is tractable. The latter requirement is essential when confronting an optimization problem having a large-scale design dimension and/or large scale parameter space.

In this paper we are concerned mainly with optimization problems where the uncertain parameters are probabilities. This is the case when the objective function and/or the constraint functions involve terms with expectations (such as moments of random variables, or expected utility, etc.). For such problems we advocate the use of uncertainty regions that are constructed as confidence sets using  $\phi$ -divergence functionals. Such functionals include the Hellinger distance, the Kullback Leibler, the Burg, and the chi-squared divergence, and many others. We choose  $\phi$ -divergences because these play a fundamental role in statistics (see [30] and [35]). The main contribution of this paper is showing that the choice of  $U$  as such uncertainty sets indeed fulfils the above three requirements:

- $U$  is based on empirical probability estimates obtained from historical data.
- $U$  is shown to relate to a statistical confidence region based on asymptotic theory.
- $U$  is such that the corresponding RCP is shown to be tractable: for basically all significant  $\phi$ -divergence functionals, the resulting robust counterpart problem is polynomially solvable. In fact, in many cases it reduces to a linear, or a conic quadratic problem.

Using (smooth)  $\phi$ -divergences, uncertainty regions can easily be constructed as (approximate) confidence sets, when the probabilities can be estimated from historical data. This follows from applying asymptotic theory. Moreover,  $\phi$ -divergences also allow the construction of confidence sets when the probabilities are calculated using additional information, represented by some underlying statistical model. In this way, smaller confidence sets can be obtained without reducing the confidence level. This is a consequence of the so-called information processing theorem, valid for  $\phi$ -divergences, see [30]. The size of the uncertainty region can be controlled by the confidence level of the confidence set. For example, the choice of a 95% confidence level will result in an uncertainty set which is (statistically) significant. Combined with the tractability of the RCP with these uncertainty sets,  $\phi$ -divergences therefore present an appealing approach in robust optimization.

We illustrate the relevance of the proposed approach by applying it first to an investment problem. We show that our approach yields a natural link with standard asset pricing theory. We also present a numerical illustration in terms of a multi-item newsvendor problem. Using our robust optimization approach leads to solutions that are quite robust, while at the same time exhibiting good average optimal performance.

We now discuss related papers. In Chapter 2 of [5] probabilistic arguments are used to construct an uncertainty region by using partial *a priori* knowledge on the underlying distribution of the uncertain parameters. Klabjan et al. [28] use the well-known chi-squared statistic, which is a special case of a  $\phi$ -divergence statistic, to define uncertainty regions for the unknown demand distribution in an inventory control problem. In their approach a robust dynamic programming problem has to be solved. Calafiore [11] studies portfolio selection problems in which the true

distribution of asset returns is unknown. He assumes that the true distribution is only known to lie within a certain distance from an estimated one and uses the Kullback-Leibler divergence to measure the distance. Our analysis includes the Kullback-Leibler divergence as a special case. Moreover, whereas Calafiore’s approach requires solving a “nested” optimization problem, our approach allows for a tractable reformulation of the robust counterpart. Wang et al. [39] studied robust optimization for data-driven newsvendor problems, in which the uncertainty set for the unknown distribution is defined as a “likelihood region”. Bertsimas and Brown [7] interpret robust optimization in terms of coherent risk measures. Ben-Tal et al. [6] consider the soft robust optimization approach and establish for such optimization problems a link with convex risk measures. Related research on robust optimal portfolio choice with uncertainty sets based on confidence sets include [16], [23], and [24] (for an overview, see [22]). These papers typically use mean or covariance matrix-based confidence sets, while we use confidence sets based on  $\phi$ -divergences.

The remainder of this paper is organized as follows. We start with an introduction to robust linear optimization in Section 2, and to  $\phi$ -divergences in Section 3. In Section 3 we also discuss the construction of uncertainty sets as confidence sets using  $\phi$ -divergences. In Section 4 we study the robust counterparts for problems with  $\phi$ -divergence uncertainty regions. In Section 5 we show that for different choices of the  $\phi$ -divergence, the robust counterpart can be reformulated as a tractable problem. In Section 6 we present some applications, including a numerical multi-item newsvendor example, and Section 7 concludes the paper with topics for further research.

## 2 Introduction to robust linear optimization

In this paper the main focus is on robust linear optimization. Without loss of generality, we focus on robust counterpart problems of the form

$$\min \{c^T x \mid Ax \leq b, \forall A \in \widehat{U}\},$$

where  $x \in \mathbb{R}^n$  is the optimization vector,  $c \in \mathbb{R}^n$  and  $b \in \mathbb{R}^l$  are given (known) parameters,  $A \in \mathbb{R}^{l \times n}$  is a matrix with uncertain parameters, and  $\widehat{U}$  is a given uncertainty region for  $A$ . Indeed, as shown in [5], for robust linear optimization we can, without loss of generality, assume that the objective and the right-hand-side of the constraints are certain.

Moreover, as also shown in [5], for robust linear optimization, we can without loss of generality assume constraint-wise uncertainty. Hence, we focus on a single constraint, which we assume to be of the form

$$(a + Bp)^T x \leq \beta, \forall p \in U, \tag{1}$$

where  $x \in \mathbb{R}^n$  is the design vector,  $a \in \mathbb{R}^n$ ,  $B \in \mathbb{R}^{n \times m}$ , and  $\beta \in \mathbb{R}$  are given (known) parameters,  $p \in \mathbb{R}^m$  is the uncertain parameter, and  $U$  the uncertainty region for  $p$ .

In Table 1 the tractability results for several standard choices of  $U$  are given. For a detailed treatment, see [5]. The last line in the table is a new result, and will be proved in this paper. Here, we briefly discuss the derivation of the robust counterparts of the standard choices of  $U$ , illustrating the general principles. We shall apply these general principles in our approach as well. To start, the results for the box and ball uncertainty region can easily be obtained by finding the worst-case solution with respect to  $p$ , i.e., by solving

$$\max \{p^T B^T x \mid p \in U\}.$$

For the polyhedral and cone uncertainty region we can use duality. Specifically, under the assumption that the uncertainty region is a cone  $K$  which contains a strictly feasible solution (i.e., there exists a  $\bar{p}$  such that  $C\bar{p} + d \in \text{int}K$ ) it holds that:

$$\max\{p^T B^T x \mid Cp + d \in K\} = \min\{d^T y \mid C^T y = -B^T x, y \in K^*\},$$

where  $K^*$  denotes the dual cone of  $K$ . This means that  $x$  satisfies (1) if and only if  $x$  satisfies

$$a^T x + \min\{d^T y \mid C^T y = -B^T x, y \in K^*\} \leq \beta.$$

Hence, we have that  $x$  satisfies (1) if and only if  $(x, y)$  satisfies

$$\begin{cases} a^T x + d^T y \leq \beta \\ C^T y = -B^T x \\ y \in K^*. \end{cases}$$

Moreover, in [5] it is shown that if  $U$  is the intersection of different “tractable cones” then the robust counterpart can also be reformulated as a tractable problem.

In this paper we shall show that if the uncertainty region  $U$  is based on a  $\phi$ -divergence, the robust counterpart can also be reformulated as a tractable optimization problem.

### 3 Introduction to $\phi$ -divergence

In this section, we first define the concept of  $\phi$ -divergence, and discuss some properties that will be useful in obtaining tractable reformulations of the robust counterpart of problem (1), when the uncertainty region  $U$  is defined in terms of a  $\phi$ -divergence. Next, we discuss how to construct uncertainty regions as (approximate) confidence sets based on a  $\phi$ -divergence.

#### 3.1 Definition and some characteristics

The  $\phi$ -divergence (“distance”) between two vectors<sup>1</sup>  $p = (p_1, \dots, p_m)^T \geq 0, q = (q_1, \dots, q_m)^T \geq 0$  in  $\mathbb{R}^m$  is defined as

$$I_\phi(p, q) = \sum_{i=1}^m q_i \phi\left(\frac{p_i}{q_i}\right), \quad (2)$$

where  $\phi(t)$  is convex for  $t \geq 0$ ,  $\phi(1) = 0$ ,  $0\phi(a/0) := a \lim_{t \rightarrow \infty} \phi(t)/t$ , for  $a > 0$ , and  $0\phi(0/0) := 0$ . We refer to the function  $\phi$  as the  $\phi$ -divergence function. We shall mainly focus on probability vectors  $p$  and  $q$  that satisfy the additional constraint  $p^T e = 1$  and  $q^T e = 1$ , where  $e$  denotes a column vector of ones of the same dimension as  $p$  and  $q$ . However, some of our results are also valid more generally for  $p \geq 0$  and  $q \geq 0$ .

Different choices for  $\phi$  have been proposed in the literature. For a good overview, see [32]. Table 2 contains the most known and used choices for  $\phi$ . The power divergence class presented in the bottom row of Table 2 was proposed by Cressie and Read [15] to be used in case of multinomial data, and is since then extensively studied, see for example [27]. The expression for  $\theta = 0$  is obtained by taking the limit  $\theta \rightarrow 0$ , and the expression for  $\theta = 1$  by taking the limit  $\theta \rightarrow 1$ .

---

<sup>1</sup>In case of a vector, we interpret the inequality componentwise, i.e.,  $p = (p_1, \dots, p_m)^T \geq 0$  means  $p_i \geq 0$  for  $i = 1, \dots, m$ . Similarly,  $p > 0$  means  $p_i > 0$  for  $i = 1, \dots, m$ .

Uncertainty region	$U$	Robust Counterpart	Tractability
Box	$\ p\ _\infty \leq 1$	$a^T x + \ B^T x\ _1 \leq \beta$	LP
Ball	$\ p\ _2 \leq 1$	$a^T x + \ B^T x\ _2 \leq \beta$	CQP
Polyhedral	$Cp + d \geq 0$	$\begin{cases} a^T x + d^T y \leq \beta \\ C^T y = -B^T x \\ y \geq 0 \end{cases}$	LP
Cone (closed, convex, pointed)	$Cp + d \in K$	$\begin{cases} a^T x + d^T y \leq \beta \\ C^T y = -B^T x \\ y \in K^* \end{cases}$	Conic Opt.
Separable functions	$\sum_i f_{\ell i}(p_i) \leq 0, \forall \ell \in \{1, \dots, L\}$	$\begin{cases} a^T x + \sum_\ell \sum_i \lambda_\ell f_{\ell i}^* \left( \frac{s_{\ell i}}{\lambda_\ell} \right) \leq \beta \\ \sum_\ell s_{\ell i} = b_i^T x, \quad i \in \{1, \dots, m\} \\ \lambda \geq 0. \end{cases}$	Convex Opt.

Table 1: Robust linear optimization for different choices for the uncertainty region in terms of  $p = (p_1, \dots, p_m)^T$ . The functions  $f_{\ell i}$  are assumed to be convex,  $f_{\ell i}^*$  is the conjugate function of  $f_{\ell i}$ , and  $K^*$  denotes the dual cone of  $K$ .

Table 3 shows that the Cressie and Read class contains several well-known  $\phi$ -divergence functions proposed in the literature (up to normalization). Notice that when the  $\phi$ -divergence function corresponding to a  $\phi$ -divergence is differentiable at  $t = 1$ , the function  $\varphi(t) = \phi(t) - \phi'(t)(t - 1)$  also yields a  $\phi$ -divergence, satisfying (for probability vectors)  $I_\varphi(p, q) = I_\phi(p, q)$ , with  $\varphi(1) = \varphi'(1) = 0$  and  $\varphi(t) \geq 0$ , see [32].

Given some  $\phi$  with corresponding  $\phi$ -divergence  $I_\phi(p, q)$ , the so-called adjoint of  $\phi$  is defined for  $t \geq 0$  as (see [2]):

$$\tilde{\phi}(t) := t\phi\left(\frac{1}{t}\right). \quad (3)$$

It holds that  $\tilde{\phi}$  satisfies the conditions for  $\phi$ -divergence functions, and  $I_{\tilde{\phi}}(p, q) = I_\phi(q, p)$ . Later in this paper we will also use other properties of  $\tilde{\phi}$ . For example, it is easy to see that the adjoint of the adjoint function is the function itself, i.e.,  $\tilde{\tilde{\phi}} = \phi$ . Moreover, the function  $\phi$  is called self-adjoint if  $\tilde{\phi} = \phi$ . As can be seen from Table 2, the  $J$ -divergence and the variation distance are self-adjoint. For other interesting properties of  $\tilde{\phi}$  we refer to [2].

We will show in Section 4 that the robust counterpart of a linear constraint with  $\phi$ -divergence uncertainty can be reformulated in terms of the so-called conjugate of  $\phi$ . The conjugate is a function  $\phi^* : \mathbb{R} \rightarrow \mathbb{R} \cup \{\infty\}$  which is defined as follows:

$$\phi^*(s) = \sup_{t \geq 0} \{st - \phi(t)\}. \quad (4)$$

In Table 2 we only present the expressions of  $\phi^*$  on its effective domain  $\text{dom}(\phi^*)$ , i.e., the part of the domain where  $\phi^*(s) < \infty$ .<sup>2</sup>

In some cases  $\phi^*$  does not exist in a (known) closed form. This is, for example, the case for the  $J$ -divergence distance measure (see Table 2). In the sequel we will use the following two propositions to determine tractable reformulations of the robust counterpart in cases where  $\phi^*$  does not exist in closed form. The first proposition applies when  $\phi$  can be written as the sum of two  $\phi$ -divergence functions  $\phi_1$  and  $\phi_2$ . The conditions required in the proposition are fulfilled in case the functions  $f_1$  and  $f_2$  are  $\phi$ -divergence functions.

**Proposition 3.1** [36] *Assume that  $f_1$  and  $f_2$  are convex, and the intersection of the relative interiors of the effective domains of  $f_1$  and  $f_2$  is nonempty, i.e.,  $\text{ri}(\text{dom } f_1) \cap \text{ri}(\text{dom } f_2) \neq \emptyset$ . Then*

$$(f_1 + f_2)^*(s) = \inf_{s_1 + s_2 = s} (f_1^*(s_1) + f_2^*(s_2)),$$

and the inf is attained for some  $s_1, s_2$ . □

The following proposition relates the conjugate of the adjoint function to the conjugate of the original function.

**Proposition 3.2** [26] *For the conjugate of a  $\phi$ -divergence function and the conjugate of its adjoint, we have*

$$\phi^*(s) = \inf\{y \in \mathbb{R} : (\tilde{\phi})^*(-y) \leq -s\}. \quad \square$$

---

<sup>2</sup>These  $\phi^*$  correspond to  $\phi$  with effective domain  $\text{dom}(\phi) = (0, \infty)$ . Thus, we set  $\phi(t) = \infty$  for  $t \leq 0$ .

Divergence	$\phi(t)$	$\phi(t), t \geq 0^a$	$I_\phi(p, q)$	$\phi^*(s)$	$\tilde{\phi}(t)$	RCP
Kullback-Leibler	$\phi_{kl}(t)$	$t \log t - t + 1$	$\sum p_i \log \left( \frac{p_i}{q_i} \right)$	$e^s - 1$	$\phi_b(t)$	S.C.
Burg entropy	$\phi_b(t)$	$-\log t + t - 1$	$\sum q_i \log \left( \frac{q_i}{p_i} \right)$	$-\log(1 - s), s < 1$	$\phi_{kl}(t)$	S.C.
J-divergence	$\phi_j(t)$	$(t - 1) \log t$	$\sum (p_i - q_i) \log \left( \frac{p_i}{q_i} \right)$	no closed form	$\phi_j(t)$	S.C.
$\chi^2$ -distance	$\phi_c(t)$	$\frac{1}{t}(t - 1)^2$	$\sum \frac{(p_i - q_i)^2}{p_i}$	$2 - 2\sqrt{1 - s}, s < 1$	$\phi_{mc}(t)$	CQP
Modified $\chi^2$ -distance	$\phi_{mc}(t)$	$(t - 1)^2$	$\sum \frac{(p_i - q_i)^2}{q_i}$	$\begin{cases} -1, & s < -2 \\ s + s^2/4, & s \geq -2 \end{cases}$	$\phi_c(t)$	CQP
Hellinger distance	$\phi_h(t)$	$(\sqrt{t} - 1)^2$	$\sum (\sqrt{p_i} - \sqrt{q_i})^2$	$\frac{s}{1-s}, s < 1$	$\phi_h(t)$	CQP
$\chi$ divergence of order $\theta > 1$	$\phi_{ca}^\theta(t)$	$ t - 1 ^\theta$	$\sum q_i \left  1 - \frac{p_i}{q_i} \right ^\theta$	$s + (\theta - 1) \left( \frac{ s }{\theta} \right)^{\frac{\theta}{\theta-1}}$	$t^{1-\theta} \phi_{ca}^\theta(t)$	CQP
Variation distance	$\phi_v(t)$	$ t - 1 $	$\sum  p_i - q_i $	$\begin{cases} -1, & s \leq -1 \\ s, & -1 \leq s \leq 1 \end{cases}$	$\phi_v(t)$	LP
Cressie and Read	$\phi_{cr}^\theta(t)$	$\frac{1-\theta+\theta t-t^\theta}{\theta(1-\theta)}, \quad \theta \neq 0, 1^b$	$\frac{1}{\theta(1-\theta)}(1 - \sum p_i^\theta q_i^{1-\theta})$	$\frac{1}{\theta}(1 - s(1 - \theta))^{\frac{\theta}{\theta-1}} - \frac{1}{\theta}$ $s < \frac{1}{1-\theta}$	$\phi_{cr}^{1-\theta}(t)$	CQP

Table 2: Some  $\phi$ -divergence examples, with their conjugates and adjoints. The last column indicates the tractability of (1); S.C. means “admits self-concordant barrier”.

<sup>a</sup> $\phi(t) = \infty$ , for  $t < 0$

<sup>b</sup>Note that  $\phi_{cr}^1(t) = \phi_b(t)$  and  $\phi_{cr}^0(t) = \phi_{kl}(t)$ .



$\theta$	$\phi_{cr}^\theta(t)$	Equivalent with
2	$\frac{1}{2}(t^2 - 2t + 1) = \frac{1}{2}(t - 1)^2$	modified $\chi^2$ -distance
1	$t(\log t - 1) + 1$	Kullback-Leibler
$\frac{1}{2}$	$4\left(\frac{1}{2} + \frac{1}{2}t - \sqrt{t}\right) = 2(1 - \sqrt{t})^2$	Hellinger distance
-1	$\frac{1}{2}\left(-2 + t + \frac{1}{t}\right) = \frac{1}{2}\left(\sqrt{t} - \frac{1}{\sqrt{t}}\right)^2$	$\chi^2$ -distance

Table 3: Some specific choices for  $\theta$  for the Cressie and Read  $\phi$ -divergence class.

To choose between different  $\phi$ -divergences one might use some representation theorem for  $\phi$ -divergences, as given in, for example, Reid and Williamson [35], see also [30]. For instance, a useful representation theorem states that  $\phi$ -divergences can be represented by a weighted average of basic  $\phi$ -divergences, where the weights are exclusively determined by the second order derivative of  $\phi$  (possibly considered as a generalized function).

### 3.2 Construction of uncertainty regions

In this subsection we describe how to construct uncertainty regions for probability vectors  $p$  as (approximate) confidence sets using  $\phi$ -divergences. We consider settings in which there is a fixed number  $m$  of given scenarios for a random variable  $Z$ , where the components of the probability vector  $p = (p_1, \dots, p_m)^T$  are given by  $p_i \equiv \mathbb{P}(Z \in C_i)$ ,  $i = 1, \dots, m$ . Here,  $p_i$  represents the probability that scenario  $i$  will occur, where  $C_i$ ,  $i = 1, \dots, m$ , form a partition (of measurable sets) of the outcome space of  $Z$ . As *basic case* we take the case where we only observe  $Z \in C_i$ ,  $i = 1, \dots, m$ . In this situation we can assume without loss of generality that  $Z \in \{1, \dots, m\}$ , where  $Z = i$  in case of scenario  $i$ . But we shall also consider cases where  $Z$  contains more information than just which of the  $m$  scenarios occurs. To capture both the basis case and more general cases, we assume the existence of a (measurable) transformation  $G$ , such that  $G(Z) = i$  if  $Z \in C_i$ ,  $i = 1, \dots, m$ . The basis case then corresponds to the situation where  $G$  is a one-to-one transformation.

Denote by  $\mathbb{P}_Z$  the probability distribution of  $Z$ . We shall assume that  $\mathbb{P}_Z$  belongs to a parameterized set of probability distributions  $\{\mathbb{P}_\theta \mid \theta \in \Theta \subset \mathbb{R}^d\}$ , i.e., there exists some  $\theta_0 \in \Theta$ , such that  $\mathbb{P}_Z = \mathbb{P}_{\theta_0}$ . We write  $p_\theta = (p_{1,\theta}, \dots, p_{m,\theta})^T$ , with  $p_{i,\theta} = \mathbb{P}_\theta(G(Z) = i)$ , and we write  $p_0 = p_{\theta_0}$ . We consider the case where the probability distributions  $\mathbb{P}_\theta$  are dominated by a common  $\sigma$ -finite measure  $\mu$ . The density of  $\mathbb{P}_\theta$  with respect to  $\mu$  is denoted by  $f_\theta$ , where we shall write  $f_0 = f_{\theta_0}$ . In the basic case, when we only observe the scenarios, we have  $Z \in \{1, \dots, m\}$ , and we can take, for example,  $\Theta = \mathbb{R}^{m-1}$ ,  $\mu$  the counting measure, and

$$\mathbb{P}_\theta(Z = i) = f_\theta(i) \times 1 = \exp(\theta_i) / \sum_{j=1}^m \exp(\theta_j), \quad i = 1, \dots, m,$$

for  $\theta = (\theta_1, \dots, \theta_m)^T$ , with normalization  $\theta_m \equiv 0$ . We then have

$$\mathbb{P}_\theta(Z = i) = p_{i,\theta} = f_\theta(i), \quad i = 1, \dots, m,$$

so that there is a one-to-one correspondence between the sets  $\mathcal{P} := \{p \in \mathbb{R}^m \mid p \geq 0, p^T e = 1\}$  and  $\mathcal{P}_\Theta := \{p_\theta \mid \theta \in \Theta\}$ .

Given this setting, we shall discuss the construction of uncertainty sets as confidence sets, under the assumption that a sample  $Z_1, \dots, Z_N$ , randomly drawn from  $\mathbb{P}_Z$ , is given. The  $\phi$ -divergence

between the densities  $f_\theta$  and  $f_0$  is given by

$$I_\phi(f_\theta, f_0) = \int \phi\left(\frac{f_\theta}{f_0}\right) f_0 d\mu.$$

We shall first construct a confidence set in terms of  $f_\theta$  which we then use to construct a confidence set in terms of  $p_\theta$ . Let  $\hat{\theta}$  denote the Maximum Likelihood estimator of  $\theta$ , and denote  $\hat{f}_0 = f_{\hat{\theta}}$ . In the basic case we get  $\hat{f}_0 = q_N$ , where  $q_N = (q_{1,N}, \dots, q_{m,N})^T$  is the  $m$ -dimensional vector containing as components the sample frequencies of the  $m$  scenarios based on the random sample  $Z_1, \dots, Z_N$ . We shall use  $I_\phi(f_\theta, \hat{f}_0)$  as estimator for  $I_\phi(f_\theta, f_0)$ .<sup>3</sup> Pardo [32] presents the characteristics of this estimator under the assumption that  $\phi$  is twice continuously differentiable in a neighborhood of 1, with  $\phi''(1) > 0$ . Most  $\phi$ -divergences reported in Table 2 satisfy this condition. Under the probability distribution  $\mathbb{P}_Z = \mathbb{P}_{\theta_0}$  and under appropriate additional regularity conditions, he shows that the normalized estimated  $\phi$ -divergence

$$\frac{2N}{\phi''(1)} I_\phi(f_\theta, \hat{f}_0) \tag{5}$$

asymptotically (i.e., for  $N \rightarrow \infty$ ) follows a  $\chi_d^2$ -distribution, with degrees of freedom determined by the dimension of the parameter set  $\Theta$ . In terms of the densities  $f_\theta$  we therefore have the following (approximate)  $(1 - \alpha)$ -confidence set around  $f_0$ :

$$\left\{ f_\theta \mid I_\phi(f_\theta, \hat{f}_0) \leq \rho \right\}, \tag{6}$$

where<sup>4</sup>

$$\rho = \rho_\phi(N, d, \alpha) := \frac{\phi''(1)}{2N} \chi_{d,1-\alpha}^2. \tag{7}$$

Based on the *information (or data) processing theorem*, see [30] for a new proof, we have

$$I_\phi(p_\theta, \hat{p}_0) \leq I_\phi(f_\theta, \hat{f}_0),$$

where  $\hat{p}_0 = p_{\hat{\theta}}$  is the estimator of  $p_0$ , using  $\hat{\theta}$  as estimator for  $\theta$ . Thus, we have

$$\{\theta \in \Theta \mid I_\phi(p_\theta, \hat{p}_0) \leq \rho\} \supset \left\{ \theta \in \Theta \mid I_\phi(f_\theta, \hat{f}_0) \leq \rho \right\}.$$

We also have

$$\{p \in \mathcal{P} \mid I_\phi(p, \hat{p}_0) \leq \rho\} \supset \{p \in \mathcal{P}_\Theta \mid I_\phi(p_\theta, \hat{p}_0) \leq \rho\}. \tag{8}$$

This implies the left hand side of (8) as (approximate) confidence set of confidence level at least  $(1 - \alpha)$  for  $p \in \mathcal{P}$  around  $\hat{p}_0$ . In the basic case, i.e., when we only observe the scenarios, we have that the dimension of  $\Theta$  equals  $m - 1$ , so that  $d = m - 1$  in (7). But with additional information we might be able to parameterize  $f_\theta$  by means of  $\Theta \subset \mathbb{R}^d$  with  $d < m - 1$ . Then, using (7) with  $d < m - 1$ , we get a smaller confidence set, but of the same confidence level.

The confidence set (8) is based on asymptotics ( $N \rightarrow \infty$ ), and therefore only approximately valid. In order to improve the approximation, several possibilities exist, see [32]. One possibility in the basic case is to consider the statistic

$$\frac{1}{\sqrt{\delta_\phi}} \left( \frac{2N}{\phi''(1)} I_\phi(p, q_N) - \gamma_\phi \right), \tag{9}$$

<sup>3</sup>It is also possible to avoid the use of the Maximum Likelihood estimator, see, for example, [10] or [30].

<sup>4</sup>In this expression  $\chi_{m-1,1-\alpha}^2$  is the  $1 - \alpha$  percentile of the  $\chi_{m-1}^2$ -distribution, i.e.,  $P(X \geq \chi_{m-1,1-\alpha}^2) = \alpha$ , with  $X$  following a  $\chi_{m-1}^2$ -distribution.

Divergence	$h(t)$	$\phi$
Renyi	$\frac{1}{\theta(\theta-1)} \log(\theta(\theta-1)t+1); \quad \theta \neq 0, 1$	$\phi_{cr}^\theta; \quad \theta \neq 0, 1$
Sharma-Mittal	$\frac{1}{v-1} \left( (1 + \theta(\theta-1)x)^{\frac{v-1}{\theta-1}} - 1 \right); \quad \theta, v \neq 1$	$\phi_{cr}^\theta; \quad \theta \neq 0, 1$
Bhattacharyya	$-\log(1 - \frac{1}{4}t)$	$\phi_{cr}^{\frac{1}{2}}$

Table 4: Examples of  $(h, \phi)$ -divergence statistics.

instead of (5). The ‘‘correction parameters’’  $\delta_\phi$  and  $\gamma_\phi$ , satisfying  $\delta_\phi \rightarrow 1$  and  $\gamma_\phi \rightarrow 0$  for  $N \rightarrow \infty$ , are defined at p. 190 of [32]. These corrections ensure that the test statistic has the same mean and variance as the limiting  $\chi^2$ -distribution, up to order  $1/N$ . We can use (9) to construct an approximate confidence interval, similar to (6), but due to the correction terms the approximation might be better for smaller sample sizes.

In the literature also several so-called  $(h, \phi)$ -divergence statistics have been proposed. Such a  $(h, \phi)$ -divergence between two probability vectors  $p \geq 0$  and  $q \geq 0$  in  $\mathbb{R}^m$  is defined as  $h(I_\phi(p, q))$ , for some appropriately chosen  $h$ . Some examples, taken from [32], are given in Table 4. Let  $h$  be increasing and continuously differentiable in a neighborhood of 0. Then, under  $\mathbb{P}_Z$ , the statistic

$$\frac{2N}{h'(0)\phi''(1)} h(I_\phi(f_\theta, \hat{f}_0)),$$

follows the same distribution as the statistic in (5). Therefore, the uncertainty regions in (6) and (8), with

$$\rho = \rho_{(h,\phi)}(N, d, \alpha) := h^{-1} \left( \frac{h'(0)\phi''(1)}{2N} \chi_{d,1-\alpha}^2 \right), \quad (10)$$

are approximate  $(1 - \alpha)$ -confidence intervals. Thus, (8) with this choice of  $\rho$  yields a  $(h, \phi)$ -divergence based uncertainty region.

## 4 Robust counterpart with $\phi$ -divergence uncertainty

In this section we derive the robust counterpart (RCP) for (1) with a  $\phi$ -divergence based uncertainty region. We consider the following robust linear constraint:

$$(a + Bp)^T x \leq \beta, \quad \forall p \in U, \quad (11)$$

where  $x \in \mathbb{R}^n$  is the optimization vector,  $a \in \mathbb{R}^n$ ,  $B \in \mathbb{R}^{n \times m}$ , and  $\beta \in \mathbb{R}$  are given parameters,  $p \in \mathbb{R}^m$  is the uncertain parameter, and

$$U = \{p \in \mathbb{R}^m \mid p \geq 0, Cp \leq d, I_\phi(p, q) \leq \rho\}, \quad (12)$$

where  $q \in \mathbb{R}^m$  (with  $q \geq 0$ ),  $\rho > 0$ ,  $d \in \mathbb{R}^k$ , and  $C \in \mathbb{R}^{k \times m}$  are given. As discussed in the previous section, when the uncertainty region is constructed as confidence set, we will have  $q = \hat{p}_0$ , the empirical or data based estimate. Formulation (12) is somewhat more general than we considered in the previous section. To deal with  $p$  as a probability vector we include the constraints  $e^T p \leq 1$  and  $e^T p \geq 1$ . But if some additional information concerning  $p$  is available that can be expressed in terms of linear (in)equalities, these can also be included in the uncertainty region as given by (12). We shall assume that these additional constraints are such that  $q \in U$ .

We prove the following theorem.

**Theorem 4.1** A vector  $x \in \mathbb{R}^n$  satisfies (11) with uncertainty region  $U$  given by (12) such that  $q \in U$  if and only if there exist  $\eta \in \mathbb{R}^k$  and  $\lambda \in \mathbb{R}$  such that  $(x, \lambda, \eta)$  satisfies

$$\begin{cases} a^T x + d^T \eta + \rho \lambda + \lambda \sum_i q_i \phi^* \left( \frac{b_i^T x - c_i^T \eta}{\lambda} \right) \leq \beta \\ \eta \geq 0, \lambda \geq 0, \end{cases} \quad (13)$$

where  $b_i$  and  $c_i$  are the  $i$ -th columns of  $B$  and  $C$ , respectively, and  $\phi^*$  is the conjugate function given by (4), with  $0\phi^* \left( \frac{s}{0} \right) := 0$  if  $s \leq 0$  and  $0\phi^* \left( \frac{s}{0} \right) := +\infty$  if  $s > 0$ .

**Proof:** We have that (11) holds if and only if:

$$\beta \geq \max_p \{(a + Bp)^T x \mid p \in U\} = \max_{p \geq 0} \left\{ (a + Bp)^T x \mid Cp \leq d, \sum_{i=1}^m q_i \phi \left( \frac{p_i}{q_i} \right) \leq \rho \right\}. \quad (14)$$

The Lagrange function for the optimization problem on the right-hand-side of (14) is given by:

$$L(p, \lambda, \eta) = (a + Bp)^T x + \rho \lambda - \lambda \sum_{i=1}^m q_i \phi(p_i/q_i) + \eta^T (d - Cp),$$

and the dual objective function is:

$$g(\lambda, \eta) = \max_{p \geq 0} L(p, \lambda, \eta).$$

Since  $q \in U$ , it follows that  $U$  is regular in the sense that  $Cq \leq d$  and  $I_\phi(q, q) = 0 < \rho$ . Due to this regularity of  $U$  strong duality holds. Hence, it follows that  $x$  satisfies (11) if and only if  $\min_{\lambda, \eta \geq 0} g(\lambda, \eta) \leq \beta$ , where the min is attained for some  $\lambda \geq 0$ ,  $\eta \geq 0$ . Equivalently,  $x$  satisfies (11) if and only if  $g(\lambda, \eta) \leq \beta$  for some  $\lambda \geq 0$  and  $\eta \geq 0$ . The dual objective function satisfies:

$$\begin{aligned} g(\lambda, \eta) &= a^T x + d^T \eta + \rho \lambda + \max_{p \geq 0} \sum_{i=1}^m (p_i (b_i^T x) - p_i (c_i^T \eta) - \lambda q_i \phi(p_i/q_i)) \\ &= a^T x + d^T \eta + \rho \lambda + \sum_{i=1}^m \max_{p_i \geq 0} (p_i (b_i^T x - c_i^T \eta) - \lambda q_i \phi(p_i/q_i)) \\ &= a^T x + d^T \eta + \rho \lambda + \sum_{i=1}^m q_i \max_{t \geq 0} \{t (b_i^T x - c_i^T \eta) - \lambda \phi(t)\} \\ &= a^T x + d^T \eta + \rho \lambda + \sum_{i=1}^m q_i (\lambda \phi)^* (b_i^T x - c_i^T \eta). \end{aligned} \quad (15)$$

Finally, we have  $(\lambda \phi)^* (s) = \lambda \phi^* \left( \frac{s}{\lambda} \right)$  for  $\lambda \geq 0$ , where we define  $0\phi^* \left( \frac{s}{0} \right) := (0\phi)^* (s)$ , which equals 0 if  $s \leq 0$  and  $+\infty$  if  $s > 0$ .  $\square$

In the RCP (13) we need  $\phi^*$ , the conjugate function of  $\phi$ . These conjugates are given in Table 2. However, for the  $J$ -divergence, the conjugate function is not available in a closed form expression. Nevertheless, in the next section, where we discuss the tractability aspects of (13), we also derive a tractable representation of the RCP for this case.

We present four corollaries. The first corollary specializes the theorem to a probability vector without additional constraints.

**Corollary 4.2** A vector  $x \in \mathbb{R}^n$  satisfies (11) with uncertainty region  $U$  given by

$$U = \{p \in \mathbb{R}^m \mid p \geq 0, e^T p = 1, I_\phi(p, q) \leq \rho\},$$

such that  $q \in U$  if and only if there exist  $\eta \in \mathbb{R}$  and  $\lambda \in \mathbb{R}$  such that  $(x, \lambda, \eta)$  satisfies

$$\begin{cases} a^T x + \eta + \rho\lambda + \lambda \sum_i q_i \phi^* \left( \frac{b_i^T x - \eta}{\lambda} \right) \leq \beta \\ \lambda \geq 0. \end{cases} \quad (16)$$

Consider next the following nonlinear constraint in  $x \in \mathbb{R}^n$ :

$$(a + Bp)^T f(x) \leq \beta, \quad \forall p \in U, \quad (17)$$

where  $a \in \mathbb{R}^k$ ,  $B \in \mathbb{R}^{k \times m}$ ,  $x \in \mathbb{R}^n$ , and  $f : \mathbb{R}^n \rightarrow \mathbb{R}^k$ . In the sequel we shall assume that  $b_i^T f(\cdot) : \mathbb{R}^n \rightarrow \mathbb{R}$  is convex for all  $i$  (with  $b_i$  the  $i$ -th column of  $B$ ). Constraints such as (17) may occur if  $p$  is a probability vector and if the objective and/or constraints of a nonlinear programming problem depend on moments of a random variable. One example is the class of expected utility maximization (see Section 6.1). We have the following corollary.

**Corollary 4.3** A vector  $x \in \mathbb{R}^n$  satisfies (17) with uncertainty region  $U$  given by (12) such that  $q \in U$  if and only if there exist  $\eta \in \mathbb{R}^k$  and  $\lambda \in \mathbb{R}$  such that  $(x, \lambda, \eta)$  satisfies

$$\begin{cases} a^T f(x) + d^T \eta + \rho\lambda + \lambda \sum_i q_i \phi^* \left( \frac{b_i^T f(x) - c_i^T \eta}{\lambda} \right) \leq \beta \\ \eta \geq 0, \lambda \geq 0. \end{cases} \quad (18)$$

**Proof:** The dual objective is given by (15) with  $x$  replaced by  $f(x)$ . Therefore, it follows from (13) that (17) is equivalent to (18).  $\square$

In case we are not sure which  $\phi$ -divergence to use, we might combine several  $\phi$ -divergences. For example, we can take the uncertainty region as an intersection of (a finite number of)  $\phi$ -divergences, given by

$$U = \{p \in \mathbb{R}^m \mid p \geq 0, Cp \leq d, I_{\phi_\ell}(p, q) \leq \rho_\ell, \ell \in \{1, \dots, L\}\}, \quad (19)$$

where  $\phi_\ell$  are the corresponding  $\phi$ -divergence functions and  $\rho_\ell > 0$  are given. Again, we assume  $q \in U$ . We have the following corollary.

**Corollary 4.4** A vector  $x \in \mathbb{R}^n$  satisfies (11) with uncertainty region  $U$  given by (19) such that  $q \in U$  if and only if there exist  $\eta \in \mathbb{R}^k$  and  $\lambda = (\lambda_1, \dots, \lambda_L)^T \in \mathbb{R}^L$  such that  $(x, \lambda, \eta)$  satisfies

$$\begin{cases} a^T x + d^T \eta + \sum_\ell \lambda_\ell \rho_\ell + \sum_\ell \lambda_\ell \sum_i q_i \phi_\ell^* \left( \frac{s_{\ell i}}{\lambda_\ell} \right) \leq \beta \\ \sum_\ell s_{\ell i} = b_i^T x - c_i^T \eta, \quad i \in \{1, \dots, m\} \\ \eta \geq 0, \lambda \geq 0. \end{cases}$$

**Proof:** In case of (15) we get  $\sum_{\ell} \lambda_{\ell} \rho_{\ell}$  instead of  $\lambda \rho$  and  $(\sum_{\ell} \lambda_{\ell} \phi_{\ell})^* (b_i^T x - c_i^T \eta)$  instead of  $(\lambda \phi)^* (b_i^T x - c_i^T \eta)$ . Using Proposition 3.1 we get

$$\left( \sum_{\ell} \lambda_{\ell} \phi_{\ell} \right)^* (b_i^T x - c_i^T \eta) = \min_{\sum_{\ell} s_{\ell i} = b_i^T x - c_i^T \eta} \sum_{\ell} (\lambda_{\ell} \phi_{\ell})^* (s_{\ell i}).$$

Since this expression appears in the “ $\leq$ ”-inequality in (13), we may ignore the “min”. Finally, using  $(\lambda_{\ell} \phi_{\ell})^* (s) = \lambda_{\ell} \phi_{\ell}^* \left( \frac{s}{\lambda_{\ell}} \right)$  we arrive at the result of the corollary.  $\square$

In the derivation of the RCP we did not exploit the special structure of the  $\phi$ -divergence functions. Therefore, suppose that the uncertainty region in (11) is defined by separable constraint functions:

$$U = \{p \in \mathbb{R}^m \mid \sum_i f_{\ell i}(p_i) \leq 0, \forall \ell \in \{1, \dots, L\}\}, \quad (20)$$

where  $f_{\ell i}$  are convex functions such that for each  $i$  we have  $\cap_{\ell=1}^L \text{ri}(\text{dom } f_{\ell i}) \neq \emptyset$ . Then the following corollary gives a tractable reformulation of the RCP. This result extends the classes of uncertainty regions for which tractable RCPs are derived in the literature. See also Table 1.

**Corollary 4.5** *A vector  $x \in \mathbb{R}^n$  satisfies (11) with uncertainty region  $U$  given by (20) such that for some  $\bar{p} = (\bar{p}_1 \dots, \bar{p}_m)^T \in U$*

$$\sum_i f_{\ell i}(\bar{p}_i) < 0, \forall \ell \in \{1, \dots, L\} \quad (21)$$

*if and only if there exist  $\lambda = (\lambda_1, \dots, \lambda_L)^T \in \mathbb{R}^L$  such that  $(x, \lambda)$  satisfies*

$$\begin{cases} a^T x + \sum_{\ell} \sum_i \lambda_{\ell} f_{\ell i}^* \left( \frac{s_{\ell i}}{\lambda_{\ell}} \right) \leq \beta \\ \sum_{\ell} s_{\ell i} = b_i^T x, \quad i \in \{1, \dots, m\} \\ \lambda \geq 0, \end{cases}$$

*where  $f_{\ell i}^*$  denotes the conjugate of  $f_{\ell i}$ ,  $\ell \in \{1, \dots, L\}$ ,  $i \in \{1, \dots, m\}$ .*

**Proof:** The proof follows from the proof of Theorem 4.1 combined with that of Corollary 4.5, where (21) is Slater’s condition guaranteeing that strong duality holds in this case.  $\square$

Finally, we briefly consider the case where we lack sufficient data or information to determine the nominal value  $q$ . If that is the case we might add a second robustness layer by replacing the constraint in (13) by:

$$a^T x + d^T \eta + \rho \lambda + \lambda \sum_i q_i \phi^* \left( \frac{b_i^T x - c_i^T \eta}{\lambda} \right) \leq \beta, \quad \forall q \in \bar{U}, \quad (22)$$

where  $\bar{U}$  is the uncertainty region for  $q$ . Note that the left-hand-side of (22) is an affine function in  $q$ , so this constraint is a special case of (17). Therefore, if  $\bar{U}$  is again a  $\phi$ -divergence based uncertainty region we can use the results obtained in this paper to determine tractable reformulations, and if  $\bar{U}$  is polyhedral or ellipsoidal we can use the results in [5]. Suppose, for example, that

$$\bar{U} = \{q \geq 0 \mid \sum_i |q_i - \bar{q}_i| \leq \varrho, \sum_i q_i = 1\},$$

for some  $\varrho > 0$ . Then it can easily be proven that  $(x, \eta, \lambda)$  satisfies (22) if and only if  $(x, \eta, \lambda, \mu)$  satisfies

$$\begin{cases} a^T x + d^T \eta + \rho \lambda + \mu + \varrho \|\gamma'\|_\infty + \bar{q}^T \gamma' \leq \beta \\ \gamma'_i = \lambda \phi^* \left( \frac{b_i^T x - c_i^T \eta}{\lambda} \right) \leq \mu, & i \in \{1, \dots, m\} \\ \eta \geq 0, \lambda \geq 0. \end{cases} \quad (23)$$

## 5 Tractability of the robust counterpart

In this section we investigate a number of questions related to the RCP (13) and by answering these questions we illustrate tractable reformulations for a selection of  $\phi$ -divergence functions, including the Burg-, the Kuhlback-Leibler-, and the  $J$ -divergence. We present the tractability results of the other  $\phi$ -divergences, which can be treated in a similar way, in the Appendix. The last column of Table 2 summarizes the tractability results.

Questions that need to be addressed so as to derive tractable RCPs (13), are:

1. What to do if  $\phi^*$  is not differentiable?
2. What to do if  $\phi^*$  does not exist in a closed form?
3. What is the convexity status of the first constraint function in (13)?
4. Does the constraint set (13) admit a self-concordant barrier?

The first question is relevant since some  $\phi^*$  functions presented in Table 2 are not differentiable. However, for all these cases we can reformulate the problem as a differentiable problem by adding extra variables and constraints.

Question 2 will be addressed below, when we discuss uncertainty regions based on the  $J$ -divergence.

To answer question 3, concerning the convexity issue, observe that for a  $\phi$ -divergence function  $\phi$  its conjugate  $\phi^*$  is also convex. Moreover, since

$$\lambda \phi^* \left( \frac{b_i^T x - c_i^T \eta}{\lambda} \right) = \sup_{t \geq 0} \{ (b_i^T x - c_i^T \eta) t - \lambda \phi(t) \}, \quad (24)$$

and the supremum over linear functions is convex, we obtain that the left hand side of (24) is jointly convex in  $\lambda$ ,  $x$ , and  $\eta$ , which means that the constraint function in (13) is convex. In a similar way we find that the constraint function in (18) is also convex, since we assume that  $b_i^T f(\cdot)$  is convex for all  $i$ .

An affirmative answer to question 4 is very desirable since it implies the possibility to use polynomial-time interior point algorithms (see [31]). We shall address this question for the Burg- and Kuhlback-Leibler-divergences. As we shall see later, after answering question 2 in case of the  $J$ -divergence, we will also be able to answer question 4 for the  $J$ -divergence.<sup>5</sup>

---

<sup>5</sup>In case of the other  $\phi$ -divergences presented in Table 2 we find that the RCP can even be reformulated as a CQP or LP problem. See the Appendix.

To investigate question 4, we reformulate the constraint set (13) as follows:

$$\begin{cases} a^T x + d^T \eta + \rho \lambda + q^T z \leq \beta, \\ \lambda \phi^* \left( \frac{b_i^T x - c_i^T \eta}{\lambda} \right) \leq z_i, \quad \forall i \\ \eta \geq 0, \lambda \geq 0, \end{cases} \quad (25)$$

with  $z = (z_1, \dots, z_m)^T$ . In case of the Burg-divergence (like many others) we have  $\text{dom}(\phi^*) = (-\infty, u)$ , for  $u < \infty$ . As a consequence, the middle inequalities of (25) can be reformulated as

$$\lambda f \left( \frac{s_i}{\lambda} \right) \leq z_i, \quad s_i = \lambda u - (b_i^T x - c_i^T \eta) \geq 0, \quad \forall i, \quad (26)$$

with  $f(s) := \phi^*(u - s)$ .

Reformulation (26) cannot be used in case of the Kuhlback-Leibler-divergence, since the effective domain of the conjugate of its  $\phi$ -divergence function equals the real line. In this case we apply Proposition 3.2 and obtain that the RCP (13) is equivalent to

$$\begin{cases} a^T x + d^T \eta + \rho \lambda + \lambda \sum_i q_i \inf \left[ y \in \mathbb{R} : (\tilde{\phi})^*(-y) \leq \frac{-b_i^T x + c_i^T \eta}{\lambda} \right] \leq \beta \\ \eta \geq 0, \lambda \geq 0, \end{cases} \quad (27)$$

which, in turn, is equivalent to

$$\begin{cases} a^T x + d^T \eta + \rho \lambda + q^T z \leq \beta \\ \lambda (\tilde{\phi})^* \left( \frac{-z_i}{\lambda} \right) \leq -b_i^T x + c_i^T \eta, \quad \forall i \\ \eta \geq 0, \lambda \geq 0, \end{cases} \quad (28)$$

with again  $z = (z_1, \dots, z_m)^T$ . In case of the Kuhlback-Leibler-divergence (like many others) we have  $\text{dom}((\tilde{\phi})^*) = (-\infty, \tilde{u})$ , for  $\tilde{u} < \infty$ . As a consequence, the middle inequalities of (28) can be reformulated as

$$\lambda f \left( \frac{s_i}{\lambda} \right) \leq -b_i^T x + c_i^T \eta, \quad s_i = \lambda \tilde{u} + z_i \geq 0, \quad \forall i, \quad (29)$$

with now  $f(s) := (\tilde{\phi})^*(\tilde{u} - s)$ .

Our aim is to establish self-concordance for the logarithmic barrier function for the constraint set (25) combined with (26) (in case of the Burg-divergence) and for the constraint set (28) combined with (29) (in case of the Kuhlback-Leibler-divergence). We first recall the definition of a self-concordant function.

**Definition:** Let  $F \subset \mathbb{R}^n$  be an open and convex set. A function  $\varphi : F \rightarrow \mathbb{R}$  is called  $\kappa$ -**self-concordant on**  $F$ , with  $\kappa \geq 0$ , if  $\varphi$  is  $C^3(F)$ , and  $\forall y \in F$  and  $\forall h \in \mathbb{R}^n$  the following inequality holds:

$$|\nabla^3 \varphi(y)[h, h, h]| \leq 2\kappa (h^T \nabla^2 \varphi(y) h)^{\frac{3}{2}},$$

where  $\nabla^3 \varphi(y)[h, h, h]$  denotes the third order differential of  $\varphi$  at  $y$  and  $h$ .

We shall use the next theorem.

**Theorem 5.1** *If for a convex function  $f : \mathbb{R}^+ \rightarrow \mathbb{R}$  it holds that:*

$$\left| f'''(s) \right| \leq \kappa f''(s)/s, \quad \text{for some } \kappa > 0, \quad (30)$$

*then the logarithmic barrier function for*

$$\{y f(s/y) \leq z, s \geq 0, y \geq 0\} \quad (31)$$

*is  $(2 + \frac{\sqrt{2}}{3}\kappa)$ -self-concordant.*



**Proof:** Define  $g(s, y) = yf(s/y)$ . According to Lemma A.2 in [17] it holds that if there exists a  $\beta$  such that

$$|\nabla^3 g(s, y)[h, h, h]| \leq \beta h^T \nabla^2 g(s, y) h \sqrt{\frac{h_1^2}{s^2} + \frac{h_2^2}{y^2}}, \quad \forall h \in \mathbb{R}^2, \quad (32)$$

in which  $\nabla^3 g(s, y)[h, h, h]$  is the third order differential, then the logarithmic barrier function for (31), given by

$$-\ln(z - g(s, y)) - \ln s - \ln y, \quad (33)$$

is  $(1 + \frac{1}{3}\beta)$ -self-concordant. We now prove that (32) holds for  $\beta = 3 + \kappa\sqrt{2}$ . It can easily be verified that for the second order differential we have

$$\nabla^2 g(s, y)[h, h] = h^T \nabla^2 g(s, y) h = f''(s/y) \left( \frac{h_1^2}{y} - \frac{2sh_1 h_2}{y^2} + \frac{s^2 h_2^2}{y^3} \right). \quad (34)$$

Moreover, for the third order differential we have

$$\begin{aligned} \nabla^3 g(s, y)[h, h, h] &= f''(s/y) \left( -\frac{3h_1^2 h_2}{y^2} + \frac{6sh_1 h_2^2}{y^3} - \frac{3s^2 h_2^3}{y^4} \right) + \\ &f'''(s/y) \left( \frac{h_1^3}{y^2} - \frac{3sh_1^2 h_2}{y^3} + \frac{3s^2 h_1 h_2^2}{y^4} - \frac{s^3 h_2^3}{y^5} \right). \end{aligned}$$

Using  $|f'''(s)| \leq \kappa f''(s)/s$ , for some  $\kappa > 0$ , we have

$$\begin{aligned} |\nabla^3 g(s, y)[h, h, h]| &\leq f''(s/y) \left| -\frac{3h_1^2 h_2}{y^2} + \frac{6sh_1 h_2^2}{y^3} - \frac{3s^2 h_2^3}{y^4} \right| + \\ &\kappa \frac{y}{s} f''(s/y) \left| \frac{h_1^3}{y^2} - \frac{3sh_1^2 h_2}{y^3} + \frac{3s^2 h_1 h_2^2}{y^4} - \frac{s^3 h_2^3}{y^5} \right| \\ &\leq 3f''(s/y) \left( \frac{h_1^2}{y} - \frac{2sh_1 h_2}{y^2} + \frac{s^2 h_2^2}{y^3} \right) \frac{|h_2|}{y} + \\ &\kappa f''(s/y) \left( \frac{h_1^2}{y} - \frac{2sh_1 h_2}{y^2} + \frac{s^2 h_2^2}{y^3} \right) \left( \frac{|h_1|}{s} + \frac{|h_2|}{y} \right) \\ &\leq (3 + \kappa\sqrt{2}) h^T \nabla^2 g(s, y) h \sqrt{\frac{h_1^2}{s^2} + \frac{h_2^2}{y^2}}. \end{aligned}$$

This proves that (32) holds for  $\beta = 3 + \kappa\sqrt{2}$  and hence that the corresponding logarithmic barrier function is  $(2 + \frac{\sqrt{2}}{3}\kappa)$ -self-concordant.  $\square$

To apply this theorem, notice that we reformulated the relevant parts of both the constraint set (25) combined with (26) (in case of the Burg-divergence) and the constraint (28) combined with (29) (in case of the Kuhlback-Leibler-divergence) as (31). Thus, if  $f$  in (26) or in (29) satisfies condition (30), then the theorem implies that the logarithmic barrier function for the corresponding constraint set is self-concordant. In case of the **Burg-divergence** we have  $\text{dom}(\phi^*) = (-\infty, 1)$ , resulting in  $f(s) = -\log(s)$ . This function  $f$  satisfies condition (30) with  $\kappa = 2$ . Therefore, it follows that in case of the Burg-divergence (25) combined with (26) is tractable. As an immediate consequence we also have that (28) combined with (29) is tractable in case of the **Kuhlback-Leibler-divergence**.

Finally, we return to question 2: what to do if  $\phi^*$  does not exist in closed form? For these cases Propositions 3.1 and 3.2 may be of help.

First, consider the case where  $\phi^*$  is not available in closed form expression, but there exist  $\phi$ -divergences  $\phi_1$  and  $\phi_2$  such that  $\phi = \phi_1 + \phi_2$ , and  $\phi_1^*$  and  $\phi_2^*$  are available in closed form. Then, applying Proposition 3.1, we obtain that the RCP (13) is equivalent to

$$\begin{cases} a^T x + d^T \eta + \rho \lambda + \lambda \sum_i q_i \min_{s_{1i} + s_{2i} = b_i^T x - c_i^T \eta} [\phi_1^*(s_{1i}/\lambda) + \phi_2^*(s_{2i}/\lambda)] \leq \beta \\ \eta \geq 0, \lambda \geq 0. \end{cases}$$

Since the first inequality is a “ $\leq$ ” one, we may delete the “min” and get the following system of inequalities in  $(x, \eta, \lambda, s_1, s_2)$  to represent the RCP (13):

$$\begin{cases} a^T x + d^T \eta + \lambda \rho + \lambda \sum_i q_i [\phi_1^*(s_{1i}/\lambda) + \phi_2^*(s_{2i}/\lambda)] \leq \beta \\ s_{1i} + s_{2i} = b_i^T x - c_i^T \eta, \quad \forall i \\ \eta \geq 0, \lambda \geq 0. \end{cases} \quad (35)$$

This is a tractable problem if, loosely speaking, the corresponding problems for  $\phi_1$  and  $\phi_2$  are tractable.

We can apply this approach to the  **$J$ -divergence** for which there is no closed form expression available for  $\phi^*$ . The crucial observation is that in case of the  $J$ -divergence we have  $\phi_j(t) = (t - 1) \log t = t \log t - \log t = \phi_{kl}(t) + \phi_b(t)$ , where  $\phi_{kl}(t)$  is the Kullback-Leibler  $\phi$ -divergence function and  $\phi_b(t)$  the Burg  $\phi$ -divergence function.

To complete our analysis, we give an example of a  $\phi$ -divergence function, for which a closed form expression for its conjugate is not available, but for which still a tractable RCP can be derived by using Proposition 3.2. Suppose that

$$\phi(t) = |t - 1|^\theta t^{1-\theta}.$$

It can be verified that this is a  $\phi$ -divergence function corresponding to a well-defined  $\phi$ -divergence. However,  $\phi^*$  is not available in a closed form expression, and hence (13) cannot be used directly. To overcome this problem we observe that  $\tilde{\phi} = \phi_{ca}^\theta$ , i.e., the adjoint of  $\phi$  is the  $\chi$ -divergence function of order  $\theta$ , for which a closed form expression for its conjugate is available (see Table 2). Therefore, one can obtain a tractable RCP for this choice of  $\phi$  by using (28), based on an application of Proposition 3.2.

## 6 Applications

In this section we first present an expected utility maximization framework in general terms, which we then specialize to an investment problem and to the newsvendor example. Next, to illustrate the performance of  $\phi$ -divergence based robust optimization, we present as numerical example a multi-item newsvendor optimization problem.

### 6.1 Expected utility maximization

We consider a decision maker who faces a problem in which the outcome of the decision is uncertain, and depends on which scenario will be realized. Let  $x \in \mathbb{R}^n$  denote the decision variable, let  $\bar{r}(x, i)$  denote the payoff from decision  $x$  if scenario  $i = 1, \dots, m$  occurs, and let  $u(r)$  denote the utility that the decision maker attaches to payoff  $r \in \mathbb{R}$ . Then, the optimization problem is given by:

$$\max_{x \in X} \sum_i p_i \times u(\bar{r}(x, i)), \quad (36)$$

where  $X \subset \mathbb{R}^n$  denotes the feasible region for the decision variable  $x$ , and where  $p_i$  is the probability of scenario  $i$  to occur. In case the probability vector  $p = (p_1, \dots, p_m)^T$  is not known, the robust counterpart problem is:<sup>6</sup>

$$\max_{x \in X} \min_{p \in U} \sum_i p_i \times u(\bar{r}(x, i)), \quad (37)$$

or, equivalently,

$$\begin{aligned} & \max z \\ & \text{s.t. } \sum_i p_i \times u(\bar{r}(x, i)) \geq z, \quad p \in U \\ & \quad x \in X. \end{aligned}$$

**Corollary 6.1** *The RCP (37) with uncertainty region as given in Corollary 4.2 is equivalent to:*

$$\max_{x \in X, \lambda \geq 0, \eta} \left\{ -\eta - \rho\lambda - \lambda \sum_i q_i \phi^* \left( \frac{-u(\bar{r}(x, i)) - \eta}{\lambda} \right) \right\}. \quad (38)$$

**Proof:** The proof follows from combining Corollaries 4.2 and 4.3, with  $a = 0$ ,  $B = I^{m \times m}$ , and  $f_i(x) = u(\bar{r}(x, i))$ .  $\square$

Optimization problem (38) is a concave optimization problem if  $u(\bar{r}(x, i))$  is concave in  $x$  for all  $i$  and the feasible set  $X$  is convex.

## 6.2 Investment example

As a special case we consider an investment problem. Let  $R_i \in \mathbb{R}^n$  be an  $n$ -dimensional vector of gross returns in case of scenario  $i$ . Investors can choose portfolios represented by a vector of weights  $x$  belonging to the set  $X \equiv \{x \in \mathbb{R}^n \mid x^T e = 1\}$ . If scenario  $i$  occurs, the portfolio with weights  $x$  yields as gross return  $\bar{r}(x, i) = x^T R_i$ . Let  $R_i = (R_{1i}, \tilde{R}_i^T)^T$ , with  $\tilde{R}_i$  the  $(n-1)$ -dimensional subvector of  $R_i$ , containing the gross returns of assets 2 to  $n$ . Similarly, let  $x = (x_1, \tilde{x}^T)^T$  and  $e = (1, \tilde{e}^T)^T$ . Then optimization problem (36) becomes

$$\max_{\tilde{x}} \sum_i p_i \times u \left( R_{1i} + \tilde{x}^T \tilde{R}_i^e \right), \quad (39)$$

with  $\tilde{R}_i^e = \tilde{R}_i - R_{1i} \tilde{e}$ . Similarly, the RCP becomes (37), with  $\bar{r}(x, i) = R_{1i} + \tilde{x}^T \tilde{R}_i^e$  and  $U$  as given in Corollary 4.2. We shall assume that  $u$  is differentiable and that its derivative satisfies  $u'(\cdot) > 0$ .

The first order optimality conditions for problem (39) are given by

$$\sum_i p_i \times u' \left( R_{1i} + \tilde{x}^T \tilde{R}_i^e \right) \times (R_{ji} - R_{1i}) = 0, \quad j = 2, \dots, n. \quad (40)$$

---

<sup>6</sup>See [25] for an axiomatization of this utility.

This equation is a special case of the “basic equation of asset pricing,” see [12], p. 1517.<sup>7</sup> It is an equilibrium condition, stating that the weighted average of the excess return of any asset  $j$  in excess of a reference asset (in our case asset 1) equals zero. The positive random variable realizing these weights  $u' \left( R_{1i} + \tilde{x}^T \tilde{R}_i^e \right)$  is a so-called Stochastic Discount Factor (SDF). Without risk, the SDF would be constant, and the equilibrium condition (40) becomes the condition that all (nonrandom) returns are equal. As its name suggests, the “basic equation of asset pricing” is heavily used in finance, particularly in equilibrium pricing. But also when estimating and testing a particular asset pricing model, one typically makes use of the implied SDF.

A natural question is whether the first order conditions of the RCP (37) or its equivalence (38) are also a special case of the “basic equation of asset pricing.” To obtain the first order conditions of the RCP, we shall assume that  $\phi^*$  is differentiable, with  $(\phi^*)'(\cdot) > 0$ ,<sup>8</sup> and we shall assume that  $\lambda > 0$ . Then we find as first order conditions

$$\begin{aligned} (\text{w.r.t. } \eta) \quad & -1 + \sum_i \tilde{q}_i = 0 \\ (\text{w.r.t. } \lambda) \quad & -\rho - \sum_i q_i \times \phi^* \left( \frac{-u \left( R_{1i} + \tilde{x}^T \tilde{R}_i^e \right) - \eta}{\lambda} \right) \\ & - \frac{1}{\lambda} \sum_i \tilde{q}_i \times \left( R_{1i} + \tilde{x}^T \tilde{R}_i^e + \eta \right) = 0 \\ (\text{w.r.t. } \tilde{x}) \quad & \sum_i \tilde{q}_i \times u' \left( R_{1i} + \tilde{x}^T \tilde{R}_i^e \right) \times (R_{ji} - R_{1i}) = 0, \quad j = 2, \dots, n, \end{aligned}$$

where

$$\tilde{q}_i \equiv q_i \times (\phi^*)' \left( \frac{-u \left( R_{1i} + \tilde{x}^T \tilde{R}_i^e \right) - \eta}{\lambda} \right), \quad i = 1, \dots, m.$$

If we combine the first order conditions with respect to  $\eta$  and  $\tilde{x}$ , we see that we have to solve the same system of equations as in case of (40). The difference is that the probabilities  $p_i$  are replaced by  $\tilde{q}_i$ ,  $i = 1, \dots, m$ .<sup>9</sup> To become a special case of the “basic equation of asset pricing,” we consider the equations as expectations with respect to  $q_{i,N}$ , the empirical counterparts of  $p_i$ . We then find as SDF<sup>10</sup>

$$(\phi^*)' \left( \frac{-u \left( R_{1i} + \tilde{x}^T \tilde{R}_i^e \right) - \eta}{\lambda} \right) \times u' \left( R_{1i} + \tilde{x}^T \tilde{R}_i^e \right).$$

<sup>7</sup>We assume rational expectations, i.e., the probabilities  $p_i$  represent the “true” probabilities. The derived SDF is up to normalization, since equation (40) is in terms of excess returns.

<sup>8</sup>It follows from the assumptions that  $\text{dom} \phi = \mathbb{R}^+$ , and hence  $(\phi^*)'(\cdot) \geq 0$ . However, from Table 2 we see that for some choices of  $\phi$ , like the modified  $\chi^2$ -distance or the variation distance,  $(\phi^*)'(\cdot)$  may be zero. Such choices of  $\phi$  are excluded in the sequel, as they do not result in a strictly positive SDF, required in the “basic equation of asset pricing.”

<sup>9</sup>Moreover, the expectation of the optimal RCP portfolio return with respect to these probabilities  $\tilde{q}_i$  has to equal the maximum value of the objective function (38). This latter requirement follows from a reformulation of the first order conditions with respect to  $\lambda$ :

$$\sum_i \tilde{q}_i \left( R_{1i} + \tilde{x}^T \tilde{R}_i^e \right) = -\eta - \rho\lambda - \lambda \sum_i q_i \phi^* \left( \frac{-u \left( R_{1i} + \tilde{x}^T \tilde{R}_i^e \right) - \eta}{\lambda} \right).$$

<sup>10</sup>Again up to normalization, see footnote 7. Moreover, this SDF is the relevant one from an empirical point of view in case  $q = q_N$ , which is consistent for the true probability vector.

Thus, our reformulation (38), specialized to the investment problem, allows a straightforward way to retrieve the SDF in case of the robust optimization problem. This makes reformulation (38) also relevant from the point of view of equilibrium pricing and empirical finance.

A special case is obtained when  $u(r) = r$ . Then the RCP can be reformulated as

$$\min_{\tilde{x}} \max_{p \in U} \sum_i p_i \times \left( - \left( R_{1i} + \tilde{x}^T \tilde{R}^e \right) \right).$$

The inner maximization represents a coherent risk measure (see [1]). Thus, in this special case of the RCP the portfolio weights are determined by minimizing a coherent risk measure. The paper [33] provides an application.

### 6.3 Newsvendor example

In this subsection, we consider as application of utility maximization the single-item **newsvendor** problem. The newsvendor's problem is how many units of a product (item) to order, taking into account that the demand for the product is stochastic. Due to uncertainty, the newsvendor can face both unsold items or unmet demand. The unsold items will return a loss because their salvage value is lower than the purchase price. In the case of unmet demand the newsvendor incurs a cost of lost sales, which may include a penalty for the lost customer goodwill.

Let  $u(r)$  denote the newsvendor's utility from net profit  $r \in \mathbb{R}$ . His objective is to choose the order quantity  $x = Q$  in order to maximize the expected utility (36) of his net profit

$$\bar{r}(Q, i) = v \min(d_i, Q) + s(Q - d_i)^+ - l(d_i - Q)^+ - cQ,$$

where  $d_i \geq 0$  is the uncertain demand in scenario  $i$ ,  $v$  is the unit selling price,  $s$  is the salvage value per unsold item,  $l$  is the shortage cost per unit of unsatisfied demand, and  $c$  is the purchasing price per unit. A standard assumption for this problem is  $v + l \geq r$ .

In case the probability distribution of the demand is unknown, the RCP is given by:

$$\max_Q \min_{p \in U} \left\{ \sum_i p_i \times u \left( v \min(d_i, Q) + s(Q - d_i)^+ - l(d_i - Q)^+ - cQ \right) \right\}.$$

It follows immediately from Corollary 6.1 that with a  $\phi$ -divergence uncertainty region  $U$  as given by Corollary 4.2 this problem is equivalent to:

$$\min_{Q, \lambda \geq 0, \eta} \left\{ \eta + \rho\lambda + \lambda \sum_i q_{i,N} \phi^* \left( \frac{-u(v \min(d_i, Q) + s(Q - d_i)^+ - l(d_i - Q)^+ - cQ) - \eta}{\lambda} \right) \right\}.$$

With a concave utility function  $u(\cdot)$ , the assumption  $v + l \geq s$  ensures that

$$-u(v \min(d_i, Q) + s(Q - d_i)^+ - l(d_i - Q)^+ - cQ)$$

is convex in  $Q$  for all  $i$ .

Several papers study risk aversion in the newsvendor model by using as objective function expected utility ([18]), mean-variance ([13]), or conditional value-at risk ([14]). Still, all these papers assume that the entire demand distribution is known. Our approach can be used to add risk aversion with respect to the unknown demand distribution in cases where only some historical data is given.

Item ( $j$ )	1	2	3	4	5	6	7	8	9	10	11	12
$c$	4	5	6	4	5	6	4	5	6	4	5	6
$v$	6	8	9	5	9	8	6	8	9	6.5	7	8
$s$	2	2.5	1.5	1.5	2.5	2	2.5	1.5	2	2	1.5	1
$l$	4	3	5	4	3.5	4.5	3.5	3	5	3.5	3	5
$q_{1,N}^{(j)}$	0.375	0.250	0.375	0.127	0.958	0.158	0.485	0.142	0.679	0.392	0.171	0.046
$q_{2,N}^{(j)}$	0.375	0.250	0.250	0.786	0.007	0.813	0.472	0.658	0.079	0.351	0.484	0.231
$q_{3,N}^{(j)}$	0.250	0.500	0.375	0.087	0.035	0.029	0.043	0.200	0.242	0.257	0.345	0.723

Table 5: Parameter values for the multi-item newsvendor example.

**Remark.** Our approach can also be applied to regret approaches for the newsvendor model. Perakis and Roels [34] study regret in newsvendor models in which only partial information is given, for example, mean, variance, symmetry, or unimodality. Our result here can be used to minimize robustly the regret when only some historical demand data is available.

#### 6.4 Numerical illustration: multi-item newsvendor example

As numerical illustration, we consider a multi-item newsvendor problem (see, for example, [21], [37]). This problem deals with optimizing the inventory of several items which can only be sold in one period. Due to the uncertain demand, this newsvendor can face both unsold items or unmet demand. As in the single-item case, the unsold items will return a loss, and unmet demand generates a cost of lost sales. For each item  $j$ , we define the purchase cost  $c_j$ , the selling price  $v_j$ , the salvage value of unsold items  $s_j$ , and the cost of lost sales  $l_j$ . Furthermore, we denote  $\gamma$  for the budget that is available for the purchase of the items.

We assume that demand for item  $j$  is a random variable that can take on  $m$  values, denoted as  $d_i$ ,  $i = 1, \dots, m$  (i.e., for simplicity, the same possible outcomes for all items  $j$ ). We denote  $p_i^{(j)}$  for the unknown probability that the demand for item  $j$  equals  $d_i$ , and we let the uncertainty region for  $p^{(j)} = (p_1^{(j)}, \dots, p_m^{(j)})^T$  be given by:

$$U^{(j)} := \left\{ p^{(j)} \in \mathbb{R}^m \mid p^{(j)} \geq 0, (p^{(j)})^T e = 1, I_\phi \left( p^{(j)}, q_N^{(j)} \right) \leq \rho \right\}, \quad (41)$$

where  $q_N^{(j)}$  represents the sample-based estimated probability distribution for item  $j$ .

Denote by  $Q_j$  the order quantity for item  $j$ . We consider two types of multi-item newsvendor problems. The first is to maximize the sum of the profits:

$$\max_Q \sum_j \sum_i p_i^{(j)} \bar{r}_j(Q_j, i),$$

and the second is to maximize the minimal profit:

$$\max_Q \min_j \sum_i p_i^{(j)} \bar{r}_j(Q_j, i).$$

The robust versions of these problems can be stated as:

$$\begin{aligned} & \max \|z\| \\ & \text{s.t.} \quad -c_j Q_j + \sum_i p_i^{(j)} f_{i,j}(Q_j) \geq z_j, \quad \forall j, \forall p^{(j)} \in U^{(j)} \\ & \quad \sum_j c_j Q_j \leq \gamma, \end{aligned}$$

for the case where the norm in the objective is either the 1-norm, or the  $\infty$ -norm, respectively, and with

$$f_{i,j}(Q_j) = v_j \min\{d_i, Q_j\} + s_j \max\{0, Q_j - d_i\} - l_j \max\{0, d_i - Q_j\}.$$

It follows from (18) that this problem can be reformulated as:

$$\begin{aligned} & \max \|z\| \\ & \text{s.t. } -c_j Q_j - \eta_j - \lambda_j \rho - \lambda_j \sum_i q_{i,N}^{(j)} \phi^* \left( \frac{-f_{i,j}(Q_j) - \eta_j}{\lambda_j} \right) \geq z_j, \quad \forall j \\ & \sum_j c_j Q_j \leq \gamma \\ & \lambda \geq 0. \end{aligned}$$

Our numerical results apply to the case with  $n = 12$  different items, and  $m = 3$  scenarios for the demand for each item: low demand (4), medium demand (8), and high demand (10), denoted as  $d_1 = 4$ ,  $d_2 = 8$ , and  $d_3 = 10$ , respectively. The parameter values of the revenue functions, as well as the values of  $q_{i,N}^{(j)}$ , are as given in Table 6.4. Furthermore, the budget is set at  $\gamma = 1000$ .

We solve the RCP for the Burg-divergence (or the Kullback-Leibler-divergence in terms of  $\tilde{\phi}$ ) and for the Cressie and Read  $\phi$ -divergence function with  $\theta = 0.5$ . For both  $\phi$ -divergence functions, we consider the case where  $\rho = \rho^a$  is the test statistic (5) and the case where  $\rho = \rho^c$  is the corrected test statistic (9). In each case, the confidence level is set at  $\alpha = 0.05$ , and we determine the robust optimal solutions for different sample sizes  $N = 10, 20, \dots, 1000$ .

Using the solutions of the RCP problems and the solution of the non-robust problem (i.e., assuming that  $q_N$  is the true probability vector), we make several comparisons. First, we compare the performance of the robust versus the non-robust solutions for the different values of the sample size  $N$  (which in turn yields different values for  $\rho_\phi^a$  and  $\rho_\phi^c$ ). Second, we compare the results of the two  $\phi$ -divergence measures. Third, for each  $\phi$ -divergence measure, we look at the effect of using the corrected test-statistic instead of the approximate test-statistic, i.e., using  $\rho = \rho_\phi^c$  instead of  $\rho = \rho_\phi^a$ .

To make comparisons, we proceed as follows. First, we sample 10,000 hypothetically true  $p$ -vectors. Next, for each sampled probability vector  $p$ , we calculate the value of the objective function for the non-robust as well as for the robust optimal solutions. We then compare the performance of the different solutions by determining the mean and the range (i.e., the minimum and the maximum value) of the objective values corresponding to the sampled  $p$ -vectors.

The  $p$ -vectors are sampled such that approximately 95 percent of the sample satisfies  $I_{\phi_{mc}}(p, q_N) \leq \bar{\rho} := \rho_{\phi_{mc}}^a$ , where  $\phi_{mc}$  denotes the modified  $\chi^2$ -divergence. Specifically, we sample  $p_i$ , for  $i = 1, \dots, m-1$ , from a normal distribution  $N(q_{i,N}, \sigma_i)$ , and set  $p_m = 1 - \sum_{i=1}^{m-1} p_i$ . If this sampling returns a probability vector (i.e.,  $p_i \geq 0$  for  $i = 1, \dots, m$ ) we accept the vector. Otherwise, we repeat the sampling until a valid  $p$ -vector is found. To satisfy  $I_{\phi_{mc}}(p, q_N) \leq \bar{\rho}$  for approximately 95 percent of the sampled  $p$ -vectors, we determine the value of  $\sigma_i$  of the normal distribution as follows. We know that the condition  $I_{\phi_{mc}}(p, q_N) \leq \bar{\rho}$  is satisfied if (but not only if) the following holds:

$$\frac{(p_i - q_{i,N})^2}{q_{i,N}} \leq \frac{\bar{\rho}}{m} \Leftrightarrow q_{i,N} - \sqrt{\frac{\bar{\rho}}{m} q_{i,N}} \leq p_i \leq q_{i,N} + \sqrt{\frac{\bar{\rho}}{m} q_{i,N}}.$$

For the normal distribution about 95 percent of the values are within two standard deviations from the mean. Therefore, we take  $\sigma_i = \frac{1}{2}\sqrt{\frac{\bar{p}}{m}q_{i,N}}$ . Because  $\bar{p}$  can be relatively large for small values of  $N$  and to avoid too many invalid samples, we put an upper bound of  $\frac{1}{2}q_{i,N}$  on  $\sigma_i$ . Figures 1 and 2 display the range and the mean of the objective values corresponding to the

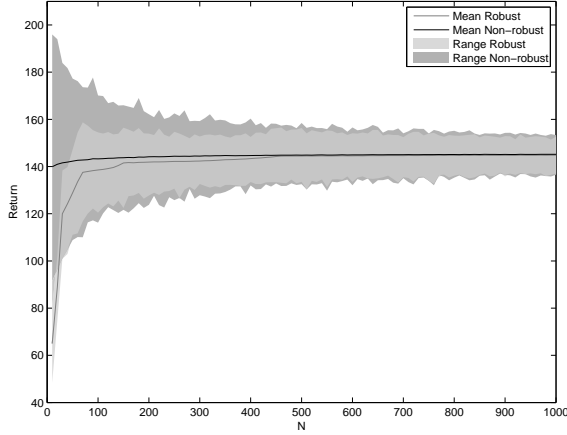


Figure 1: Cressie-Read for  $\theta = 0.5$ , and  $\rho_\phi^c$ , and the 1-norm.

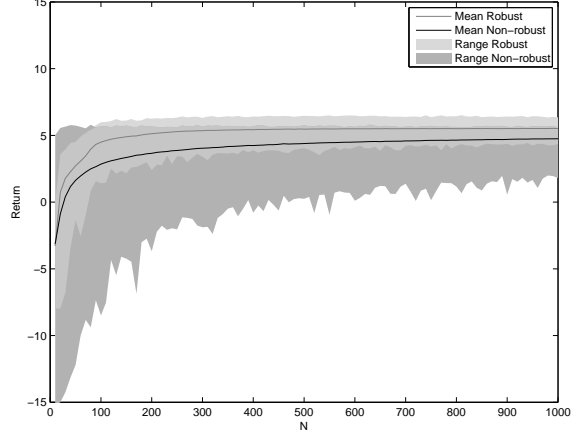


Figure 2: Cressie-Read for  $\theta = 0.5$ , and  $\rho_\phi^c$ , and the  $\infty$ -norm.

sampled  $p$ -vectors for the Cressie and Read-divergence function with  $\rho = \rho_\phi^c$ , for the 1-norm and the  $\infty$ -norm, respectively. The results for the Burg-divergence function are essentially similar.

Concerning the value of the objectives of the robust optimizations, there is a significant difference between using the 1-norm and the  $\infty$ -norm. For the 1-norm (Figure 1), it holds that for small values of  $N$  the mean of the objective values for the robust solution is lower than the mean of the objective values for the non-robust solution, but as  $N$  grows the two methods have practically the same mean profit. In contrast, for the  $\infty$ -norm (Figure 2), the mean of the objective values for the robust solution is higher than the mean for the non-robust solution. Moreover, the dispersion of objective values for the robust solution is significantly smaller than the range of objective values for the non-robust solution for the  $\infty$ -norm. In particular, the robust solution avoids substantial losses.

Concerning the effect of  $N$ , the effect of using  $\rho_\phi^c$  versus  $\rho_\phi^a$ , and the differences between the two  $\phi$ -divergence measures, we observe the following:

**Effect of  $N$ .** Because 95 percent of the sampled  $p$ -vectors needs to satisfy  $I_{\phi_{mc}}(p, q_N) \leq \bar{p}$ , and because  $\bar{p}$  is decreasing in  $N$ , the range of the expected returns becomes smaller as  $N$  increases. However, because 5 percent of the sampled  $p$ -vectors does not need to satisfy  $I_{\phi_{mc}}(p, q_N) \leq \bar{p}$ , the range does not converge to a single value.

**Effect of  $\rho_\phi^c$  versus  $\rho_\phi^a$ .** With regard to the differences between the robust solutions in case  $\rho_\phi^a$  is used (i.e., the uncertainty region is based on the approximate test statistic) and when  $\rho_\phi^c$  is used (i.e., the uncertainty region is based on the corrected test statistic), we observe that there are significant differences only for relatively small values for  $N$ . This occurs of course since the effect of the correction becomes smaller as  $N$  increases.

**Comparison of different  $\phi$ -divergence measures.** The different  $\phi$ -divergence measures lead to different optimal quantities, but the structure of the solutions is similar. The mean expected



utility as well as the range of the expected utilities over the sampled  $p$ -vectors is similar for the two  $\phi$ -divergence measures.

## 7 Concluding remarks

In this paper we have shown that the robust counterpart of linear and nonlinear optimization problems with uncertainty regions defined by  $\phi$ -divergence distance measures can be reformulated as tractable optimization problems. Thus, these uncertainty regions are useful alternatives to uncertainty regions considered in the existing literature, particularly so when the uncertainty is associated with probabilities. In this latter case, we have shown that uncertainty regions based on  $\phi$ -divergence test statistics have a natural interpretation in terms of statistical confidence sets. This allows for an approach that is fully data-driven.

Our approach also has other applications. For example,  $\phi$ -divergence distances can be directly used as the distance in the so-called Globalized Robust Counterpart methodology (see Chapter 3 in [5]).

Let us now mention some directions for further research. First, different choices of  $\phi$  have been proposed in the literature [32], each of them with different statistical properties. It could be interesting to study the differences in performance of optimal solutions of robust counterpart problems such as (13) for different choices of  $\phi$ .

Next, in the classical statistical literature many goodness-of-fit statistics are considered that do not belong to the  $\phi$ -divergence class. It is an interesting topic for further research to analyze whether the corresponding robust counterparts are tractable.

In terms of practical applications, it may be useful to extensively study the applicability of the proposed approach to, for example, asset liability management problems and other inventory control problems. In particular, with respect to inventory control problems it may be interesting to extend the work of Wagner [38]. In that paper the Wagner-Whitin model with backlogged demand and period-dependent costs is analyzed in settings in which the demand distribution is not known. Our analysis can most likely be used to extend the analysis to the more practical case where only some historical demand data is given.

Finally, several commonly used risk measures in finance (for example, mean-variance, expected shortfall) are nonlinear in probabilities. It is a challenging question whether the proposed approach can be extended to problems in which the unknown probability vector  $p$  appears nonlinearly. In [5] techniques are described to deal with certain types of nonlinear uncertainty, and maybe similar techniques can be used in this case.

## References

- [1] P. Artzner, F. Delbaen, J.-M. Eber, and D. Heath. Coherent Measures of Risk. *Mathematical Finance*, 9(3):203–228, 1999.
- [2] A. Ben-Tal, A. Ben-Israel, and M. Teboulle. Certainty Equivalents and information Measures: Duality and Extremal Principles. *Journal of Mathematical Analysis and Applications*, 157:211–236, 1991.
- [3] A. Ben-Tal and A. Nemirovsky. Stable Truss Topology Design via Semidefinite Programming. *SIAM Journal on Optimization*, 7(4):991–1016, 1997.

- [4] A. Ben-Tal and A. Nemirovsky. Robust Convex Optimization. *Mathematics of Operations Research*, 23(4):769–805, 1998.
- [5] A. Ben-Tal, L. El Ghaoui, and A. Nemirovsky. *Robust Optimization*. Princeton Press, Princeton, 2009.
- [6] A. Ben-Tal, D. Bertsimas, and D. Brown. A Soft Robust Model for Optimizing Under ambiguity. *Operations Research*, 58(4-part-2):1220–1234, 2010.
- [7] D. Bertsimas and D. Brown. Constructing Uncertainty Sets for Robust Linear Optimization. *Operations Research*, 57(6):1483–1495, 2009.
- [8] D. Bertsimas, D. Brown, and C. Caramanis. Theory and Applications of Robust Optimization. Working paper, 2008.
- [9] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, Cambridge, 2004.
- [10] M. Broniatowski and A. Keziou. Parametric Estimation and Tests through Divergences and the Duality Technique. *Journal of Multivariate Analysis*, 100:16–36, 2009.
- [11] G.C. Calafiore. Ambiguous Risk Measures and Optimal Robust Portfolios. *SIAM Journal on Optimization*, 18:853–877, 2007.
- [12] J.Y. Campbell. Asset Pricing at the Millennium, *Journal of Finance*, 55(4):1515–1567, 2000.
- [13] Y. Chen and A. Federgruen. Mean-Variance Analysis of Basic Inventory Models. Working paper, Columbia University, New York, 2000.
- [14] Y. Chen, M. Xu, and Z.G. Zhang. A Risk-Averse Newsvendor Model under the CVar Criterion. *Operations Research*, 57(4):1040–1044, 2009.
- [15] N. Cressie, and T.R. Read. Multinomial Goodness-of-Fit Tests. *Journal of the Royal Statistical Society, Series B*, 46:440–464, 2009.
- [16] E. Delage and Y. Ye. Distributionally Robust Optimization under Moment Uncertainty with Application to Data-Driven Problems. *Operations Research*, forthcoming, 2010.
- [17] D. den Hertog. *Interior Point Approach to Linear, Quadratic and Convex Programming*. Kluwer Academic Publishers, Dordrecht, 1994.
- [18] L. Eeckhoudt, C. Gollier, and H. Schlesinger. Risk Averse (and Prudent) Newsboy. *Management Science*, 41(5):786–794, 2009.
- [19] L. El Ghaoui and H. Lebr et. Robust Solution to Least-Squares Problems with Uncertain Data. *SIAM Journal on Matrix Analysis and Applications*, 18:1035–1064, 1997.
- [20] L. El Ghaoui, F. Oustry, and H. Lebret. Robust Solutions to Uncertain Semidefinite Programs. *SIAM Journal on Optimization*, 9:33–52, 1998.
- [21] S.J. Erlebacher. Optimal and Heuristic Solutions for the Multi-Item Newsvendor Problem with a Single Capacity Constraint. *Production and Operations Management*, 9(3):303–318, 2000.

- [22] F. Fabozzi, D. Huang, and G. Zhou. Robust Portfolios: Contributions from Operations Research and Finance. *Annals of Operations Research*, 176:191–220, 2010.
- [23] L. Garlappi, R. Uppal, and T. Wang. Portfolio Selection with Parameter and Model Uncertainty: A Multi-Prior Approach. *Review of Financial Studies*, 20:41–81, 2007.
- [24] D. Goldfarb and G. Iyengar. Robust Portfolio Selection Problems. *Mathematics of Operations Research*, 28, 1–38, 2003.
- [25] I. Gilboa and D. Schmeidler. Maxmin Expected Utility with Non-Unique Prior. *Journal of Mathematical Economics*, 18:141–153, 1989.
- [26] A. A. Gushchin. On an Extension of the Notion of  $f$ -divergence. *Theory of Probability and its Applications*, 52(3):439–455, 2008.
- [27] L. Jager and J.A. Wellner. Goodness-of-Fit Tests via phi-divergences. *Annals of Statistics*, 35:2018–2053, 2007.
- [28] D. Klabjan, D. Simchi-Levi, and M. Song. Robust Stochastic Lot-Sizing by means of Histograms. Working paper.
- [29] P. Kouvelis and G. Yu. *Robust Discrete Optimization and its Application*. Kluwer Academic Publishers, London, 1997.
- [30] F. Liese and I. Vajda. On Divergences and Information in Statistics and Information Theory. *IEEE Transactions on Information Theory*, 52(10):4394–4412, 2006.
- [31] Y. Nesterov and A. Nemirovsky. *Interior Point Polynomial Algorithms in Convex Programming*. SIAM Studies in Applied Mathematics, Philadelphia, 1993.
- [32] L. Pardo. *Statistical Inference Based on Divergence Measures*. Chapman & Hall/CRC, Boca Raton, 2006.
- [33] A. Pelsser. Robustness, Model Uncertainty and Pricing. Working paper, Maastricht University, 2010.
- [34] G. Perakis and G. Roels. Regret in the Newsvendor Model with Partial Information. *Operations Research*, 56(1):188–203, 2008.
- [35] M. Reid and R. Williamson. Information, Divergence and Risk for Binary Experiments. Working Paper, Australian National University, 2009.
- [36] R. T. Rockafeller. *Convex Analysis*. Princeton Press, Princeton, 1970.
- [37] I. Moon and E. A. Silver. The Multi-Item Newsvendor Problem with a Budget Constraint and Fixed Ordering Costs. *Journal of the Operational Research Society*, 51(5):602–608, 2000.
- [38] M.R. Wagner. Fully Distribution-Free Profit Maximization: The Inventory Management Case. *Mathematics of Operations Research*, 35(4):728–741, 2010.
- [39] Z. Wang, P.W. Glynn, and Y. Ye. Likelihood Robust Optimization for Data-Driven Newsvendor Problems. Working paper, 2009.

## Appendix: Tractable reformulations

In this appendix we give the final tractable reformulations for (13) for different choices of  $\phi$ . The tractable reformulations for Kullback-Leibler, Burg entropy, and  $J$ -divergence are already derived in Section 5.

$\chi^2$ -distance (CQP)

$$\begin{cases} a^T x + d^T \eta + \lambda \rho + 2\lambda e^T q - 2q^T y \leq \beta \\ \sqrt{y_i^2 + \frac{1}{4}(b_i^T x - c_i^T \eta)^2} \leq \frac{1}{2}(2\lambda - b_i^T x + c_i^T \eta), & \forall i \\ b_i^T x - c_i^T \eta \leq \lambda, & \forall i \\ \eta \geq 0, \lambda \geq 0. \end{cases}$$

Modified  $\chi^2$  distance (CQP)

$$\begin{cases} a^T x + d^T \eta + \lambda(\rho - e^T q) + \frac{1}{4}q^T y \leq \beta \\ \sqrt{z_i^2 + \frac{1}{4}(\lambda - \mu_i)^2} \leq \frac{1}{2}(\lambda + \mu_i), & \forall i \\ z_i \geq 0, & \forall i \\ z_i \geq b_i^T x - c_i^T \eta + 2\lambda, & \forall i \\ \eta \geq 0, \lambda \geq 0. \end{cases}$$

Hellinger distance (CQP)

$$\begin{cases} a^T x + d^T \eta + \lambda \rho - \lambda e^T q + q^T y \leq \beta \\ \sqrt{\lambda^2 + \frac{1}{4}(y_i - \lambda + b_i^T x - c_i^T \eta)^2} \leq \frac{1}{2}(y_i + \lambda - b_i^T x + c_i^T \eta), & \forall i \\ b_i^T x - c_i^T \eta \leq \lambda, & \forall i \\ \eta \geq 0, \lambda \geq 0. \end{cases}$$

$\chi$  divergence of order  $\theta$  (CQP)

$$\begin{cases} a^T x + d^T \eta + \lambda \rho + \sum_i q_i (b_i^T x - c_i^T \eta) + \lambda(\theta - 1) \sum_i q_i \left(\frac{z_i}{\theta \lambda}\right)^{\frac{\theta}{\theta-1}} \leq \beta \\ z_i \geq -b_i^T x + c_i^T \eta, & \forall i \\ z_i \geq b_i^T x - c_i^T \eta, & \forall i. \\ \eta \geq 0, \lambda \geq 0. \end{cases}$$

Variation distance (LP)

$$\begin{cases} a^T x + d^T \eta + \lambda \rho + q^T y \leq \beta \\ y_i \geq -\lambda, & \forall i \\ y_i \geq b_i^T x - c_i^T \eta, & \forall i \\ b_i^T x - c_i^T \eta \geq \lambda, & \forall i \\ \eta \geq 0, \lambda \geq 0. \end{cases}$$

**Cressie and Read (CQP)**

$$\begin{cases} a^T x + d^T \eta + \lambda \rho + \frac{\lambda}{\theta} \sum q_i \left( \left( \frac{y_i}{\lambda} \right)^{\frac{\theta}{\theta-1}} - 1 \right) \leq \beta \\ y_i = \lambda - (1 - \theta)(b_i^T x - c_i^T \eta), \quad \forall i \\ \eta \geq 0, \lambda \geq 0. \end{cases}$$