# Robust Hypothesis Testing with $\alpha$ -Divergence

Gökhan Gül, Student Member, IEEE, Abdelhak M. Zoubir, Fellow, IEEE,

Abstract—A robust minimax test for two composite hypotheses, which are determined by the neighborhoods of two nominal distributions with respect to a set of distances - called  $\alpha$ -divergence distances, is proposed. Sion's minimax theorem is adopted to characterize the saddle value condition. Least favorable distributions, the robust decision rule and the robust likelihood ratio test are derived. If the nominal probability distributions satisfy a symmetry condition, the design procedure is shown to be simplified considerably. The parameters controlling the degree of robustness are bounded from above and the bounds are shown to be resulting from a solution of a set of equations. The simulations performed evaluate and exemplify the theoretical derivations.

*Index Terms*—Detection, hypothesis testing, robustness, least favorable distributions, minimax optimization, likelihood ratio test.

#### I. INTRODUCTION

Decision theory has been an active field of research benefiting from contributions from several disciplines, such as economics, engineering, mathematics, or statistics. A decision maker (or a detector) chooses a course of action from several possibilities. A detector is said to be optimal or to be giving the best decision for a particular problem if the decision rule of interest minimizes (or maximizes) a well defined cost function, e.g., the error probability (or the probability of detection) [1]. In addition to the fact that decision theory is truly an interdisciplinary subject of research, there are many areas of engineering, where decision theory finds applications, e.g., radar, sonar, seismology, communications and biomedicine. For some applications, such as image and speech classification or pattern recognition, interest is in a statistical test that performs well on average. However, for safety oriented applications such as seismology or forest fire detection, as well as for biomedical applications such as early cancer detection from magnetic resonance images or X-ray images, interest is in maximizing the worst case performance because the consequences of an incorrect decision can be severe [1].

In general, any practical application of decision theory can be formulated as a hypothesis testing problem. For binary hypothesis testing, it is assumed that under each hypothesis  $\mathcal{H}_i$ , the received data  $y = (y_1, \ldots, y_n) \in \Omega$  follows a particular distribution  $F_i$  corresponding to a density function  $f_i$ ,  $i \in \{0, 1\}$ . A decision rule  $\delta$  partitions the whole observation space  $\Omega$  into non-overlapping regions corresponding to each hypothesis. The optimality of the decision rule  $\delta$  depends on the correctness of the assumption that the data y follows  $F_i$ . However, in many practical applications either  $F_0$  and/or  $F_1$ are partially known or are affected by some secondary physical effects that go unmodeled [2].

Imprecise knowledge of  $F_0$  or  $F_1$  leads, in general, to performance degradation and a useful approach is to extend the known model by accepting a set of distributions  $\mathcal{F}_i$ , under each hypothesis  $\mathcal{H}_i$ , that are populated by probability distributions  $G_i$ , which are at the neighborhood of the nominal distribution  $F_i$  based on some distance D [1]. Under some mild conditions on D, it can be shown that the best (error minimizing) decision rule  $\hat{\delta}$  for the worst case (error maximizing) pair of probability distributions  $(\hat{G}_0, \hat{G}_1) \in \mathcal{F}_0 \times \mathcal{F}_1$  accepts a saddle value. Therefore, such a test design guarantees a certain level of detection at all times. This type of optimization is known as minimax optimization and the corresponding worst case distributions  $(\hat{G}_0, \hat{G}_1)$  are called least favorable distributions (LFD)s [3].

The literature in this field is unfortunately not rich. One of the earliest and probably the most crucial work goes back to Huber, who proposed a robust version of the probability ratio test for the  $\epsilon$ -contamination and total variation classes of distributions [4]. He proved the existence of least favorable distributions and showed that the corresponding robust test was a censored version of the nominal likelihood ratio for both uncertainty classes. In a later work, Huber and Strassen extended the  $\epsilon$ -contamination neighborhood to a larger class, which includes five different distances as special cases [5]. It was also shown that the robust test resulting from this new neighborhood was still a censored likelihood ratio test. Although it was found to be less engineering oriented by Levy [1], the largest classes for which similar conclusions have been made was for the 2-alternating capacities proposed by Huber and Strassen [6].

Another approach for robust hypothesis testing was proposed by Dabak and Johnson based on the fact that the choice of measures defining the contamination neighborhoods was arbitrary [7]. They chose the relative entropy (KL-divergence) because it is a natural distance between probability measures and therefore a natural way to define the contamination neighborhoods. Somewhat surprisingly, the robust test which minimizes the KL-divergence between the LFDs obtained from the closed balls with respect to the relative entropy distance was not a clipped likelihood ratio test, but a nominal likelihood ratio test with a modified threshold. It was noted that their approach was not robust for all sample sizes but when Kullback's theorem is valid, that is for a large number of observations [7]. The difference in the robust tests for  $\epsilon$ -contamination and relative entropy neighborhoods lies in the fact that all the densities in the class of distributions based on relative entropy are absolutely continuous with respect to the nominal distributions, but not for the case of the  $\epsilon$ -contamination class.

A question left open by Dabak and Johnson was the design of

G. Gül and A. M. Zoubir are with the Signal Processing Group, Institute of Telecommunications, Technische Universitt Darmstadt, 64283, Darmstadt, Germany (e-mail: ggul@spg.tu-darmstadt.de; zoubir@spg.tu-darmstadt.de)

Manuscript received April 19, 2005; revised January 11, 2007.

a robust test for a finite number of samples. Levy answered this question under two assumptions; monotone increasing nominal likelihood ratio and symmetric nominal density functions  $(f_0(y) = f_1(-y))$ , where  $y \in \mathbb{R}$ . He implied that a robust test based on the relative entropy would be more suitable for modeling errors rather than outliers, due to the smoothness property (absolute continuity). He also showed that the resulting robust test was neither equivalent to Huber's

nor to Dabak's robust test; it was a completely different test

[2]. Although KL-divergence is a smooth and a natural distance between probability measures, it is not clear why KL divergence should be considered to build uncertainty sets, especially since there are many other divergences, which are also smooth and have nice theoretical properties, e.g. the symmetry property, which KL-divergence does not have. Besides, theoretically nice properties do not always lead to preferable engineering applications, see for example [8, p.7]. In this respect, KLdivergence can be replaced by the  $\alpha$ -divergence because  $\alpha$ -divergence includes uncountably many distances as special cases, e.g.  $\chi^2$  distance for  $\alpha = 2$  [9], it reduces to the KLdivergence as  $\alpha \rightarrow 1$  and shares similar theoretical properties with the KL-divergence such as smoothness, convexity or satisfiability of (generalized) Pythagorean inequality [10]. Moreover, the flexibility provided by the choice of  $\alpha$  results in performance improvements in various signal processing applications and implies the sub-optimality of the KL-divergence. For example, in the design of distributed detection networks with power constraints,  $\alpha$ -divergence is considered as the distance between the probability measures, and error exponents of both kinds are maximized over all  $\alpha \in (0, 1)$  [11]. In nonnegative matrix factorization [12], and indexing and retrieval [13], the optimal value of  $\alpha$  (with respect to some objective function) is found to be 1/2 corresponding to the squared Hellinger distance. In medical applications; e.g. in medical image segmentation [14], restoration [15] and registration [16], the  $\alpha$ -divergence is considered and the optimal value of  $\alpha$  is found to be a non-standard value, i.e. a value which does not correspond to any known distance. There are also theoretical works which take advantage of the  $\alpha$ -divergence in the design of statistical tests. It is reported for parametric models [17], [18] as well as for non parametric models [19] that the use of  $\alpha$ -divergence as the distance between probability measures, again with some non-standard values of  $\alpha$ , e.g.  $\alpha = 1.6$  in [18] and  $\alpha = 1.3$  or  $\alpha = 1.5$  in [19], leads to promising results. However, non of these aforementioned works have the property of minimax robustness. Furthermore, in non of the aforementioned works, it is possible to adjust the tradeoff between robustness and detection performance. Additionally, the parametric models have a possibly invalid assumption that the actual probability distributions can be represented by a parametric model. This motivates the work in this paper: a minimax robust design of hypothesis testing with the  $\alpha$ -divergence distance, where the robustness is adjustable with respect to the detection performance by the choice of two robustness parameters,  $\epsilon_0$  and  $\epsilon_1$ .

The related literature can be summarized as follows: In [3], the symmetry constraint that was imposed in [2] was removed,

considering the squared Hellinger distance. In [20], the number of non-linear equations that needs to be solved to be able to design the robust test was reduced and a formula from where the maximum robustness parameters could be obtained was derived. In [21], robust approaches were extended to distributed detection problems where communication from the sensors to the fusion center is constrained. In a recent work [22], based on the KL-divergence, the monotone increasing likelihood ratio constraint was removed.

In this paper, A minimax robust test for two composite hypotheses, which are formed by the neighborhoods of two nominal distributions with respect to the  $\alpha$ -divergence, is designed. It is shown that for any  $\alpha$ , the corresponding robust test is the same and unique. There is no constraint on the choice of nominal distributions. Therefore, our design generalizes [2]. Since the  $\alpha$ -divergence includes the KL-divergence or the squared Hellinger distance as a special case, cf. [9], our work also generalizes the works in [3], [20] and [22]. The advantage of considering the  $\alpha$ -divergence for modeling errors is that it allows the designer to choose a single parameter that accounts for the distance without carrying out tedious steps of derivations for the design of a robust test. Additionally, the a priori probabilities in our work are not required to be equal, which was assumed in all previous works on model mismatch. An example is cognitive radio where the primary user may be idle for most of the time, i.e.  $P(\mathcal{H}_0) \gg P(\mathcal{H}_1)$ [23]. Last but not least, the work in this paper allows vector valued observations.

The organization of this paper is as follows. In the following section, some background to the minimax optimization problem is given and characterization the saddle value condition is detailed, before the problem definition is stated. Section III is divided into three parts. In the first part, the minimax optimization problem is solved and the least favorable distributions, the robust decision rule as well as the robust likelihood ratio, which are later shown to be determined via solving two nonlinear equations, are obtained. The second part shows how the problem is simplified if the nominal probability density functions satisfy the symmetry condition. In the third part, the maximum of the robustness parameters, above which a minimax robust test cannot be designed, are derived. In Section IV simulation results that illustrate the validity of the theoretical derivations are detailed. Finally, the paper is concluded in Section V.

#### **II. PROBLEM FORMULATION**

## A. Background

Let  $(\Omega, \mathscr{A})$  be a measurable space with the probability measures  $F_0$ ,  $F_1$ ,  $G_0$ ,  $G_1$ , and G on it, having the density functions  $f_0$ ,  $f_1$ ,  $g_0$ ,  $g_1$  and g respectively, with respect to some dominating measure  $\mu$ , i.e.,  $F_i, G_i, G \ll \mu$ ,  $i \in \{0, 1\}$ . It is assumed that the nominal measures are distinct, i.e. the condition  $F_0 = F_1 \mu$ -almost everywhere is not true. Consider the binary composite hypothesis testing problem

$$\mathcal{H}_0^c : G = G_0$$
  
$$\mathcal{H}_1^c : G = G_1 \tag{1}$$

where the measures  $G_i$  are defined whenever their corresponding density functions  $g_i$  belong to the closed ball

$$\mathcal{G}_{i} = \{g_{i} : D(g_{i}, f_{i}) \le \epsilon_{i}\} \quad i \in \{0, 1\},$$
(2)

where D is a distance between the density functions. In other words, every density function  $g_i$  which is at least  $\epsilon_i$  close to the nominal density  $f_i$  is a member of the uncertainty class  $\mathcal{G}_i$  and defines  $G_i$ ,  $i \in \{0, 1\}$ . We choose D to be the  $\alpha$ -divergence i.e.,

$$D(g, f; \alpha) := \frac{1}{\alpha(1 - \alpha)} \left( 1 - \int_{\Omega} g^{\alpha} f^{1 - \alpha} \mathrm{d}\mu \right), \alpha \in \mathbb{R} \setminus \{0, 1\}$$
(3)

since it is a convex distance for every  $\alpha$  and it includes various distances as special cases [9, p.1536].<sup>1</sup>. Given that  $y \in \Omega$  has been observed, a randomized decision rule  $\delta : \Omega \mapsto [0, 1]$ maps each y to a real number in the unit interval. Let  $\Delta$  be the set of all decision rules (functions). Then, for any possible choice of  $\delta \in \Delta$ , the following error types are well defined: first, the false alarm probability

$$P_F(\delta, f_0) = \int_{\Omega} \delta f_0 \mathrm{d}\mu, \qquad (4)$$

second, the miss detection probability

$$P_M(\delta, f_1) = \int_{\Omega} (1 - \delta) f_1 \mathrm{d}\mu, \qquad (5)$$

and third, the overall error probability

$$P_E(\delta, f_0, f_1) = P(\mathcal{H}_0)P_F(\delta, f_0) + P(\mathcal{H}_1)P_M(\delta, f_1).$$
 (6)

It is well known that  $P_E$  is minimized if the decision rule is chosen to be the likelihood ratio test

$$\delta(y) = \begin{cases} 0, & l(y) < \rho \\ \kappa(y), & l(y) = \rho \\ 1, & l(y) > \rho \end{cases}$$
(7)

where  $\rho = P(\mathcal{H}_0)/P(\mathcal{H}_1)$  is some threshold,  $l(y) := f_1/f_0(y)$  is the likelihood ratio at observation y and  $\kappa : \Omega \to [0, 1]$ .

#### B. Saddle value specification

In this section, the existence of a saddle value condition due to the functional topology of the minimax optimization problem is shown. Minimax theorem, which is attributed to John von Neumann, gives the necessary conditions such that the existence of a saddle value is guaranteed [24]. However, it is applicable if and only if both sets over which the maximization and minimization is performed are compact. Note that the closed balls ( $\mathcal{G}_0$  and  $\mathcal{G}_1$ ) with respect to the  $\alpha$ -divergence distance are not compact, therefore Von Neumann's minimax theorem is not applicable in our case. Here, we adopt Sion's minimax theorem [25],

$$\sup_{\substack{(g_0,g_1)\in\mathcal{G}_0\times\mathcal{G}_1\\ =\min\\ \delta\in\Delta}} \min_{\substack{g_0,g_1)\in\mathcal{G}_0\times\mathcal{G}_1}} P_E(\delta,g_0,g_1),$$
(8)

<sup>1</sup>Notice that  $\alpha$ -divergence is preferred against the Rényi's  $\alpha$ -divergence because Rényi's  $\alpha$ -divergence is convex only for  $\alpha \in [0, 1][9, p.1540]$ 

which removes the compactness constraint on the set over which maximization is performed. In order for (8) to be valid the following conditions must hold:

- The objective function P<sub>E</sub>(δ, ·) is real valued, upper semicontinuous and quasi-concave on G<sub>0</sub> × G<sub>1</sub> for all δ ∈ Δ
- The objective function  $P_E(\cdot, (g_0, g_1))$  is lower semicontinuous and quasi-convex on  $\Delta$  for all  $(g_0, g_1) \in \mathcal{G}_0 \times \mathcal{G}_1$
- $\Delta$  is a compact convex subset of a linear topological space
- $\mathcal{G}_0 \times \mathcal{G}_1$  is a convex subset of a linear topological space

The first two conditions hold true because  $P_E$  is a real valued continuous function, and linear on all three terms  $\delta, g_0, g_1$ , therefore both convex and concave. The last condition is also true because, all convex combinations of  $g_i^0 \in \mathcal{G}_i$  and  $g_i^1 \in \mathcal{G}_i$ are in  $\mathcal{G}_i$  since D is a convex distance and the Cartesian product of convex sets is again a convex set. Similarly,  $\Delta$  is a convex set because for any  $t \in [0,1]$  and for all  $\delta_0, \delta_1 \in \Delta$ ,  $t\delta_0 + (1-t)\delta_1 \in \Delta$ . Note that any continuous function is also upper or lower semi-continuous and any convex function is also quasi-convex. Lastly,  $\Delta$ , which is equivalent to  $[0,1]^{\Omega}$ in infinite dimensional vector space, is the product of uncountably many compact sets [0, 1]. According to Tychonoff's theorem,  $\Delta$  is compact with respect to the product topology [26], [27]. Note that any finitely supported discretization of  $g_0$  and  $g_1$  makes both  $\mathcal{G}_0 \times \mathcal{G}_1$  and  $\Delta$  compact with respect to the standard topology. This is a straightforward result of Heine-Borel theorem [28, Theorem 2.41].

Accordingly, based on Sions's minimax theorem, there exists a saddle value for the objective function  $P_E$ , i.e.,

$$P_{E}(\delta, \hat{g}_{0}, \hat{g}_{1}) \ge P_{E}(\hat{\delta}, \hat{g}_{0}, \hat{g}_{1}) \ge P_{E}(\hat{\delta}, g_{0}, g_{1}).$$
(9)

Since  $P_E$  is distinct in  $g_0$  and  $g_1$ , we also have

$$P_F(\hat{\delta}, g_0) \le P_F(\hat{\delta}, \hat{g}_0)$$
$$P_M(\hat{\delta}, g_1) \le P_M(\hat{\delta}, \hat{g}_1).$$
(10)

#### C. Problem definition

Based on (10), the minimax optimization problem (8) can be solved considering the Karush-Kuhn-Tucker (KKT) multipliers. Hence, the problem formulation can be restated as

Maximization:  

$$\hat{g}_{0} = \arg \sup_{g_{0} \in \mathcal{G}_{0}} P_{F}(\delta, g_{0})$$
s.t.  $g_{0} > 0, \Upsilon(g_{0}) = \int_{\mathbb{R}} g_{0} d\mu = 1$ 

$$\hat{g}_{1} = \arg \sup_{g_{1} \in \mathcal{G}_{1}} P_{M}(\delta, g_{1})$$
s.t.  $g_{1} > 0, \Upsilon(g_{1}) = \int_{\mathbb{R}} g_{1} d\mu = 1$ 

Minimization:  $\hat{\delta} = \arg\min_{\delta \in \Delta} P_E(\delta, \hat{g}_0, \hat{g}_1).$  (11)

In (11), there are two separate maximization problems, which are coupled with the minimization problem through the deci-

4

sion rule  $\delta^2$ 

#### III. ROBUST DETECTION WITH $\alpha$ -DIVERGENCE

The following theorem provides a solution for (11), which is composed of the least favorable densities  $\hat{g}_0$  and  $\hat{g}_1$ , the robust decision rule  $\hat{\delta}$ , the robust likelihood ratio function  $\hat{l} = \hat{g}_1/\hat{g}_0$ in parametric forms, as well as two non-linear equations from which the parameters can be obtained. Before the statement of the theorem, let  $l_l$  and  $l_u$  be two real numbers with  $0 < l_l \le 1 \le l_u < \infty$ . Furthermore, let

$$k(l_l, l_u) = \frac{\int_{\mathcal{I}_1} (l - l_l) f_0 \mathrm{d}\mu}{\int_{\mathcal{I}_3} (l_u - l) f_0 \mathrm{d}\mu},$$
(12)

implying the robust likelihood ratio function

$$\hat{l} = \frac{\hat{g}_1}{\hat{g}_0} = \begin{cases} l_l^{-1}l, & l < \rho l_l \\ \rho, & \rho l_l \le l \le \rho l_u \\ l_u^{-1}l, & l > \rho l_u \end{cases}$$
(18)

provide a unique solution to (11). Furthermore, the parameters  $l_l$  and  $l_u$  can be determined by solving

$$\frac{1}{z(l_l, l_u; \alpha, \rho)^{\alpha}} \left( l_l^{\alpha} \int_{\mathcal{I}_1} f_0 d\mu + \int_{\mathcal{I}_2} \Phi'_0(l_l, l_u; \alpha, \rho)^{\alpha} f_0 d\mu + (k(l_l, l_u)l_u)^{\alpha} \int_{\mathcal{I}_3} f_0 d\mu \right) = x(\alpha, \varepsilon_0)$$
(19)

and

$$\frac{1}{z(l_l, l_u; \alpha, \rho)^{\alpha}} \left( \int_{\mathcal{I}_1} f_1 d\mu + \int_{\mathcal{I}_2} \Phi'_1(l_l, l_u; \alpha, \rho)^{\alpha} f_1 d\mu \right)$$
$$f_1 d\mu, + k(l_l, l_u)^{\alpha} \int_{\mathcal{I}_3} f_1 d\mu \right) = x(\alpha, \varepsilon_1)$$
(20)

where 
$$\Phi'_{j}(l_{l}, l_{u}; \alpha, \rho) = z(l_{l}, l_{u}; \alpha, \rho)\Phi_{j}$$
, and  $x(\alpha, \varepsilon) = 1 - \alpha(1 - \alpha)\varepsilon$ .

A proof of Theorem III.1 is given in three stages. In the maximization stage, the Karush-Kuhn-Tucker (KKT) multipliers are used to determine the parametric forms of the LFDs,  $\hat{g}_0$  and  $\hat{g}_1$ , and the robust likelihood ratio function  $\hat{l}$ . In the minimization stage, the LFDs and the robust decision rule  $\hat{\delta}$  are made explicit. Finally, in the optimization stage, four parameters that are needed to design the test are reduced to two parameters without loss of generality.

Proof:

## A. Derivation of LFDs and the robust decision rule

1) Maximization step: Consider the Lagrangian function

$$L(g_0, \lambda_0, \mu_0) = P_F(\delta, g_0) + \lambda_0(\epsilon_0 - D(g_0, f_0; \alpha)) + \mu_0(1 - \Upsilon(g_0)))$$
(21)

where  $\mu_0$  and  $\lambda_0 \ge 0$  are the KKT multipliers. It can be seen that L is a strictly concave functional of  $g_0$ , as  $\partial^2 L/\partial g_0^2 < 0$ for every  $\lambda_0 > 0$ . Therefore, there exists a unique solution to (21), in case all KKT conditions are met [29, Chapter 5]. More explicitly the Lagrangian can be stated as

$$L(g_0, \lambda_0, \mu_0) = \int_{\mathbb{R}} \delta g_0 - \mu_0 g_0 + \frac{\lambda_0}{\alpha(1-\alpha)} \left( (1-\alpha) f_0 + \alpha g_0 - \left(\frac{g_0}{f_0}\right)^{\alpha} f_0 \right) + \lambda_0 \epsilon_0 + \mu_0 d\mu.$$
(22)

Note that similar to [2], the positivity constraint  $g_0 \ge 0$  (or  $g_1 \ge 0$ ) is not imposed, because for some  $\alpha$ , this constraint is satisfied automatically, while for others each solution of Lagrangian optimization must be checked for positivity. To find the maximum of (22), the directional (Gâteaux's) derivative of the Lagrangian L with respect to  $g_0$  in the direction of a function  $\psi$  is taken:

$$\int_{\Omega} \left[ \delta - \mu_0 + \frac{\lambda_0}{1 - \alpha} \left( \left( \frac{g_0}{f_0} \right)^{\alpha - 1} - 1 \right) \right] \psi \mathrm{d}\mu.$$
 (23)

$$\int_{\mathcal{I}_2} \left( \frac{k(l_l, l_u)^{\alpha - 1}(l_l^{\alpha - 1} - l_u^{\alpha - 1})}{l_l^{\alpha - 1} - (k(l_l, l_u)l_u)^{\alpha - 1} + (k(l_l, l_u)^{\alpha - 1} - 1)(l/\rho)^{\alpha - 1}} \right)^{\frac{1}{\alpha - 1}}$$
(13)

 $z(l_l, l_u; \alpha, \rho) = \int_{\mathcal{T}_1} f_1 \mathrm{d}\mu + k(l_l, l_u) \int_{\mathcal{T}_2} f_1 \mathrm{d}\mu +$ 

where

$$\mathcal{I}_{1} := \{ y : l(y) < \rho l_{l} \} \equiv \{ y : \hat{l}(y) < \rho \} 
\mathcal{I}_{2} := \{ y : \rho l_{l} \le l(y) \le \rho l_{u} \} \equiv \{ y : \hat{l}(y) = \rho \} 
\mathcal{I}_{3} := \{ y : l(y) > \rho l_{u} \} \equiv \{ y : \hat{l}(y) > \rho \}$$
(14)

and

$$\begin{split} \Phi_1(l, l_l, l_u; \alpha, \rho) &= \\ \frac{1}{z(l_l, l_u; \alpha, \rho)} \cdot \\ & \left( \frac{k(l_l, l_u)^{\alpha - 1}(l_l^{\alpha - 1} - l_u^{\alpha - 1})}{l_l^{\alpha - 1} - (k(l_l, l_u)l_u)^{\alpha - 1} + (k(l_l, l_u)^{\alpha - 1} - 1)(l/\rho)^{\alpha - 1}} \right)^{\frac{1}{\alpha - 1}} \\ \text{with } \Phi_0 &= \Phi_1 l \rho^{-1}. \end{split}$$

#### Theorem III.1. The least favorable densities

$$\hat{g}_{0} = \begin{cases} \frac{l_{l}}{z(l_{l}, l_{u}; \alpha, \rho)} f_{0}, & l < \rho l_{l} \\ \Phi_{0}(l, l_{l}, l_{u}; \alpha, \rho) f_{0}, & \rho l_{l} \le l \le \rho l_{u} \\ \frac{k(l_{l}, l_{u}) l_{u}}{z(l_{l}, l_{u}; \alpha, \rho)} f_{0}, & l > \rho l_{u} \end{cases}$$
(15)

$$\hat{g}_{1} = \begin{cases} \frac{1}{z(l_{l}, l_{u}; \alpha, \rho)} f_{1}, & l < \rho l_{l} \\ \Phi_{1}(l, l_{l}, l_{u}; \alpha, \rho) f_{1}, & \rho l_{l} \le l \le \rho l_{u} \\ \frac{k(l_{l}, l_{u})}{z(l_{l}, l_{u}; \alpha, \rho)} f_{1}, & l > \rho l_{u} \end{cases}$$
(16)

and the robust decision rule

$$\hat{\delta} = \begin{cases} 0, & l < \rho l_l \\ \frac{l_l^{\alpha - 1} (l/\rho)^{1 - \alpha} - 1}{(l_l^{\alpha - 1} - (k(l_l, l_u)l_u)^{\alpha - 1})(l/\rho)^{1 - \alpha} + k(l_l, l_u)^{\alpha - 1} - 1}, & \rho l_l \le l \le \rho l_u \\ 1, & l > \rho l_u \end{cases}$$
(17)

<sup>2</sup>In general arg sup may not always be achieved since  $\mathcal{G}_0$  and  $\mathcal{G}_1$  are noncompact sets in the topologies induced by the  $\alpha$ -divergence distance. In this paper, existence of  $\hat{g}_0$  and  $\hat{g}_1$  is due to the KKT solution of the minimax optimization problem, which is introduced in Section III. Since  $\psi$  is arbitrary, L is maximized whenever

$$\delta - \mu_0 + \frac{\lambda_0}{1 - \alpha} \left( \left( \frac{g_0}{f_0} \right)^{\alpha - 1} - 1 \right) = 0.$$
 (24)

Solving (24) the density function of the LFD  $\hat{G}_0$ ,

$$\hat{g}_0 = \left(\frac{1-\alpha}{\lambda_0} \left(\mu_0 - \delta\right) + 1\right)^{\frac{1}{\alpha-1}} f_0$$
 (25)

is obtained. Writing the Lagrangian for  $P_M$ , in a similar way, with the KKT multipliers  $\mu_0 := \mu_1$  and  $\lambda_0 := \lambda_1$  it follows that

$$\hat{g}_1 = \left(\frac{1-\alpha}{\lambda_1}\left(\mu_1 - 1 + \delta\right) + 1\right)^{\frac{1}{\alpha-1}} f_1.$$
 (26)

Accordingly, the robust likelihood ratio function can be obtained as

$$\hat{l} = \frac{\hat{g}_1}{\hat{g}_0} = \left[\frac{\frac{1-\alpha}{\lambda_1} \left(\mu_1 - 1 + \delta\right) + 1}{\frac{1-\alpha}{\lambda_0} \left(\mu_0 - \delta\right) + 1}\right]^{\frac{1}{\alpha - 1}} l.$$
 (27)

2) *Minimization step:* The minimizing decision function is known to be of type (7) with l to be replaced by  $\hat{l}$  and  $\kappa$  to be determined from (27) via solving  $\hat{l} = \rho$  for  $\delta := \hat{\delta}$ . For every  $\rho$ , this results in

$$\hat{\delta} = \begin{cases} 0, & \hat{l} < \rho \\ \frac{\lambda_0(-1+\alpha+\lambda_1+\mu_1-\alpha\mu_1)}{(-1+\alpha)(\lambda_0+\lambda_1(l/\rho)^{1-\alpha})} - \frac{\lambda_1(\lambda_0+\mu_0-\alpha\mu_0)(l/\rho)^{1-\alpha}}{(-1+\alpha)(\lambda_0+\lambda_1(l/\rho)^{1-\alpha})}, & \hat{l} = \rho \\ 1, & \hat{l} > \rho \end{cases}$$
(28)

Inserting (28) in (25) and (26), the least favorable density functions can be obtained as

$$\hat{g}_{0} = \begin{cases} c_{1}f_{0}, & \hat{l} < \rho \\ \Phi_{0}f_{0}, & \hat{l} = \rho , \\ c_{2}f_{0}, & \hat{l} > \rho \end{cases} \quad \hat{g}_{1} = \begin{cases} c_{3}f_{1}, & \hat{l} < \rho \\ \Phi_{1}f_{1}, & \hat{l} = \rho , \\ c_{4}f_{1}, & \hat{l} > \rho \end{cases}$$
(29)

where

$$c_{1} = \left(\frac{(1-\alpha)\mu_{0} + \lambda_{0}}{\lambda_{0}}\right)^{\frac{1}{\alpha-1}}, c_{2} = \left(\frac{(1-\alpha)(\mu_{0}-1) + \lambda_{0}}{\lambda_{0}}\right)^{\frac{1}{\alpha-1}}$$
$$c_{3} = \left(\frac{(1-\alpha)(\mu_{1}-1) + \lambda_{1}}{\lambda_{1}}\right)^{\frac{1}{\alpha-1}}, c_{4} = \left(\frac{(1-\alpha)\mu_{1} + \lambda_{1}}{\lambda_{1}}\right)^{\frac{1}{\alpha-1}}$$

and

$$\Phi_{0} = \left(\frac{-1 + \lambda_{0} + \lambda_{1} + \mu_{0} + \mu_{1} - \alpha(-1 + \mu_{0} + \mu_{1})}{\lambda_{0} + \lambda_{1}(l/\rho)^{1-\alpha}}\right)^{\frac{1}{\alpha-1}},$$
(30)

$$\Phi_{1} = \left(\frac{-1 + \lambda_{0} + \lambda_{1} + \mu_{0} + \mu_{1} - \alpha(-1 + \mu_{0} + \mu_{1})}{\lambda_{1} + \lambda_{0}(l/\rho)^{\alpha - 1}}\right)^{\alpha - 1}.$$
(31)

In order to determine the unknown parameters, the constraints in the Lagrangian definition, i.e.,  $D(\hat{g}_i, f_i, \alpha) = \epsilon_i$  and  $\Upsilon(\hat{g}_i) = 1, i \in \{0, 1\}$  are imposed. This leads to *four* non-linear equations:

$$c_{1} \int_{\hat{l}<\rho} f_{0} d\mu + \int_{\hat{l}=\rho} \Phi_{0} f_{0} d\mu + c_{2} \int_{\hat{l}>\rho} f_{0} d\mu = 1,$$

$$c_{3} \int_{\hat{l}<\rho} f_{1} d\mu + \int_{\hat{l}=\rho} \Phi_{1} f_{1} d\mu + c_{4} \int_{\hat{l}>\rho} f_{1} d\mu = 1,$$

$$c_{1}^{\alpha} \int_{\hat{l}<\rho} f_{0} d\mu + \int_{\hat{l}=\rho} \Phi_{0}^{\alpha} f_{0} d\mu + c_{2}^{\alpha} \int_{\hat{l}>\rho} f_{0} d\mu = x(\alpha, \epsilon_{0}),$$

$$c_{3}^{\alpha} \int_{\hat{l}<\rho} f_{1} d\mu + \int_{\hat{l}=\rho} \Phi_{1}^{\alpha} f_{1} d\mu + c_{4}^{\alpha} \int_{\hat{l}>\rho} f_{1} d\mu = x(\alpha, \epsilon_{1}),$$
(32)

in four parameters, where  $x(\alpha, \epsilon) = 1 - \alpha(1 - \alpha)\epsilon$ .

3) Optimization Step: In this section, the number of equations as well as the number of parameters are reduced. This allows the re-definition of  $\hat{l}$ ,  $\hat{\delta}$ ,  $\hat{g}_0$  and  $\hat{g}_1$  in a more compact form. Let  $l_l = c_1/c_3$  and  $l_u = c_2/c_4$ , then  $\hat{l} = \hat{g}_1/\hat{g}_0$  from (29) indicates the equivalence of integration domains,  $\mathcal{I}_1$ ,  $\mathcal{I}_2$  and  $\mathcal{I}_3$  as defined by (14). Applying the following steps in (32):

- Consider new domains  $\mathcal{I}_1, \mathcal{I}_2, \mathcal{I}_3$
- Use the substitutions  $c_1 := c_3 l_l$  and  $c_2 := c_4 l_u$
- Divide both sides of the first two equations by  $c_3$
- Equate the resulting equations to each other via  $1/c_3$

leads to  $c_4 = k(l_l, l_u)c_3$ , where  $k(l_l, l_u)$  is as defined by (12). Next, the goal is to find a functional f s.t.  $\Phi_1$  =  $c_3 f(l, l_l, l_u, \alpha)$ . Since  $\Phi_0 f_0 \rho = \Phi_1 f_1$ , it follows that  $\Phi_0 =$  $c_3 f(l, l_l, l_u, \alpha) l \rho^{-1}$ , therefore it suffices to evaluate only  $\Phi_1$ . A step by step derivation of the functional f is given in Appendix A. Accordingly,  $\Phi_0$  is also fully specified in terms of the desired parameters and functions. Inserting  $\Phi_1$  (which is now a functional of  $c_3, l, l_l, l_u, \alpha$ ), c.f., (51), into the second equation in (32) and noticing that  $c_4 = k(l_l, l_u)c_3$  leads to  $c_3 = 1/z(l_l, l_u; \alpha, \rho)$ , where  $z(l_l, l_u; \alpha, \rho)$  is as defined by (13). Applying a similar procedure, which can be found in Appendix B, to  $\hat{\delta}$ , c.f., (28), for the case  $\hat{l} = \rho$  leads to the robust decision rule  $\delta$  as given by Theorem III.1. The least favorable densities,  $\hat{g}_0$  and  $\hat{g}_1$ , and the robust likelihood ratio function  $\hat{l}$  are obtained similarly, by exploiting the connection <sub>z</sub> between the parameters  $c_1$ ,  $c_2$ ,  $c_3$ ,  $c_4$  and  $l_l$ ,  $l_u$ . The same simplifications eventually let the four equations given by (32) to be rewritten as the two equations stated by Theorem III.1. As it was mentioned earlier, both  $\hat{g}_0$  and  $\hat{g}_1$  are obtained <sup>1</sup> uniquely from the Lagrangian L. Hence,  $l = \hat{g}_1/\hat{g}_0$ , and as a result,  $\hat{\delta}$  are also unique. It follows that the solution found for (11) by the KKT multipliers approach is unique as claimed.

Theorem III.1 can be summarized as illustrated in Figure 1. In other words, for any choice of pair of nominal density functions  $f_0$  and  $f_1$ , the robustness parameters  $\epsilon_0$  and  $\epsilon_1$ , the Bayesian threshold  $\rho$  and the distance parameter  $\alpha$ , the robust design outputs the least favorable density functions  $\hat{g}_0$  and  $\hat{g}_1$ and the robust decision rule  $\hat{\delta}$ . Notice that  $\hat{g}_0$  and  $\hat{g}_1$  are the scaled versions (with different scaling factors) of the nominal distributions on  $l < \rho l_l$  and  $l > \rho l_u$ , and in between, they are a composition of both nominals, since  $\Phi_0$  and  $\Phi_1$  are both functionals of  $f_0$  and  $f_1$ . Interpretation of the decision rule  $\hat{\delta}$  is



Fig. 1. Summary of the robust hypothesis testing scheme given by Theorem III.1.



Fig. 2. Nonlinearity relating the nominal likelihood ratios to the robust likelihood ratios.

similar, i.e. in the same two regions the robust decision rule is almost surely zero or one, and in between it is a randomized decision rule. The robust version of the nominal likelihood ratio test is a non-linearly transformed version of the nominal likelihood ratios as illustrated by Figure 2. It is somewhat surprising that the resulting robust likelihood ratio test is the same for the whole family of distances that are parameterized by  $\alpha$ . In other words, the robust version of the likelihood ratio test, which is given by (18) is not explicitly a function of  $\alpha$ . Theorem III.1 is a generalization of [22] in the sense that as  $\alpha \to 1$  and  $\rho = 1$ , the least favorable densities  $\hat{g}_0$  and  $\hat{g}_1$  as well as the robust decision rule  $\delta$  reduce to the ones found in [22]. The flexibility afforded by the generality of considering a set of distances, called the  $\alpha$ -divergence, over [22] is twofold. First, the designer does not need to search for a suitable distance for modeling errors, and each time test for the applicability to the engineering problem at hand, following tedious steps of derivations. Instead, only the parameter  $\alpha$  is required to be determined, which can be done over a training data set via using a suitable search algorithm. Second, the a priori probabilities are not necessarily to be chosen equal. The proposed design with the  $\alpha$ -divergence covers both cases, in addition to the fact that the choice of the nominal probability distributions also does not require any assumption. Additional constraints on the choice of nominal distributions as well as on the robustness parameters simplify the design as introduced in the next section.

## B. Simplified model with additional constraints

In some cases, evidence that the following assumption holds may be available:

**Assumption III.2.** *The nominal likelihood ratio l is monotone and the nominal density functions are symmetric, i.e.,*  $f_1(y) = f_0(-y) \forall y$ 

If, additionally, the robustness parameters are set to be equal,  $\epsilon = \epsilon_0 = \epsilon_1$ , or in other words  $x(\alpha, \epsilon) = x(\alpha, \epsilon_0) = x(\alpha, \epsilon_1)$ , it follows that

$$\delta(y) = 1 - \delta(-y)$$

$$\uparrow$$

$$\begin{cases} l_u = 1/l_l \\ y_u = -y_l \end{cases} \iff \begin{cases} c_2 = c_3 \\ c_1 = c_4 \end{cases} \iff \begin{cases} \lambda_0 = \lambda_1 \\ \mu_0 = \mu_1 \end{cases}$$

$$\downarrow$$

$$g_1(y) = g_0(-y) \qquad (33)$$

where  $y_l = l^{-1}(l_l)$  and  $y_u = l^{-1}(l_u)$ . These relationships are straightforward and therefore the proofs are omitted. Notice that, due to monotonicity of l, the limits of integrals  $\mathcal{I}_1$ ,  $\mathcal{I}_2$  and  $\mathcal{I}_3$  should be re-arranged e.g.,

$$\begin{aligned} \mathcal{I}_1 &:= \{ y : l(y) < \rho l_l \} \\ &\equiv \{ y : y < l^{-1}(\rho l(y_l)) \} \equiv \{ y : y < l^{-1}(\rho l(-y_u)) \}. \end{aligned}$$

The symmetry assumption implies:

$$\begin{aligned} x(\alpha,\epsilon) &= \int_{\mathbb{R}} \left( \frac{g_1(y)}{f_1(y)} \right)^{\alpha} f_1(y) \mathrm{d}y = \int_{\mathbb{R}} \left( \frac{g_1(y)}{f_0(-y)} \right)^{\alpha} f_0(-y) \mathrm{d}y \\ &= \int_{\mathbb{R}} \left( \frac{g_0(y)}{f_0(y)} \right)^{\alpha} f_0(y) \mathrm{d}y = \int_{\mathbb{R}} \left( \frac{g_0(-y)}{f_0(-y)} \right)^{\alpha} f_0(-y) \mathrm{d}y \\ &= \int_{\mathbb{R}} \left( \frac{g_0(-y)}{f_1(y)} \right)^{\alpha} f_1(y) \mathrm{d}y \end{aligned}$$
(34)

for all  $\alpha$  and  $\epsilon$  and, it also implies l(y) = 1/l(-y) and as a result  $\hat{l}(y) = 1/\hat{l}(-y)$  for all y. Hence,  $g_1(y) = g_0(-y) \forall y$  is a solution and all the simplifications in (33) follow. This reduces the four equations given by (32) to two:

$$c_{4} = \left( l(y_{u}) \int_{-\infty}^{y_{l}^{*}} f_{1}(y) dy + \int_{y_{l}^{*}}^{y_{u}^{*}} \left( \frac{1 + l(y_{u})^{\alpha - 1}}{1 + (l(y)/\rho)^{\alpha - 1}} \right)^{\frac{1}{\alpha - 1}} f_{1}(y) dy + \int_{y_{u}^{*}}^{\infty} f_{1}(y) dy \right)^{-1}$$
(35)

and

$$c_{4}^{\alpha} \left( l(y_{u})^{\alpha} \int_{-\infty}^{y_{l}^{*}} f_{1}(y) dy + \int_{y_{l}^{*}}^{y_{u}^{*}} \left( \frac{1 + l(y_{u})^{\alpha - 1}}{1 + (l(y)/\rho)^{\alpha - 1}} \right)^{\frac{\alpha}{\alpha - 1}} f_{1}(y) dy + \int_{y_{u}^{*}}^{\infty} f_{1}(y) dy = x(\alpha, \epsilon),$$
(36)

where  $y_l^*(y_u) = l^{-1}(\rho l(-y_u))$  and  $y_u^*(y_u) = l^{-1}(\rho l(y_u))$ . These two equations can then be combined into a single equation

$$\begin{split} l(y_{u})^{\alpha} & \int_{-\infty}^{y_{l}^{*}} f_{1}(y) \mathrm{d}y + \int_{y_{l}^{*}}^{y_{u}^{*}} \left(\frac{1 + l(y_{u})^{\alpha - 1}}{1 + (l(y)/\rho)^{\alpha - 1}}\right)^{\frac{\alpha}{\alpha - 1}} f_{1}(y) \mathrm{d}y \\ & + \int_{y_{u}^{*}}^{\infty} f_{1}(y) \mathrm{d}y - x(\alpha, \epsilon) \left(l(y_{u}) \int_{-\infty}^{y_{l}^{*}} f_{1}(y) \mathrm{d}y \\ & + \int_{y_{l}^{*}}^{y_{u}^{*}} \left(\frac{1 + l(y_{u})^{\alpha - 1}}{1 + (l(y)/\rho)^{\alpha - 1}}\right)^{\frac{1}{\alpha - 1}} f_{1}(y) \mathrm{d}y + \int_{y_{u}^{*}}^{\infty} f_{1}(y) \mathrm{d}y \right)^{\alpha} = \end{split}$$
(37)

from where the parameter  $y_u$  can easily be determined. Obviously, the computational complexity is reduced considerably with the aforementioned assumptions, i.e., when (37) is compared to (19) and (20). Note that when  $\rho = 1$ , we have  $y_l^* = -y_u$  and  $y_u^* = y_u$  and if additionally  $\alpha \to 1$ , (37) reduces to [2], cf. [9].

#### C. Limiting Robustness Parameters

The existence of a minimax robust test strictly depends on the pre-condition that the uncertainty sets  $\mathcal{G}_i$  are distinct. To satisfy this condition, Huber suggested  $\epsilon_i$  to be chosen small, see [4, p.3]. Dabak [7] does not mention how to choose the parameters, whereas Levy gives an implicit bound as the relative entropy between the half way density  $f_{1/2} = f_0^{1/2} f_1^{1/2} / z$ and the nominal density  $f_0$ , i.e.,  $\epsilon < D(f_{1/2}, f_0)$ , where z is a normalizing constant. In the sequel, we show explicitly which pairs of parameters ( $\epsilon_0, \epsilon_1$ ) are valid to design a minimax robust test for the  $\alpha$ -divergence distance.

The limiting condition for the uncertainty sets to be disjoint is  $\hat{G}_0 = \hat{G}_1 \ \mu$ -a.e. It is clear from the saddle value condition (30) that for any possible choice of  $(\epsilon_0, \epsilon_1)$ , which results in  $\hat{G}_0 = \hat{G}_1$ , it is true that  $P_E \leq 1/2$  for all  $(g_0 \times g_1) \in \mathcal{G}_0 \times \mathcal{G}_1$ . Since infinitesimally smaller parameters guarantee the strict inequality  $P_E < 1/2$ , it is sufficient to determine all possible pairs which result in  $\hat{G}_0 = \hat{G}_1$ . A careful inspection suggests that the LFDs are identical whenever  $l_l \rightarrow \inf l$ and  $l_u \rightarrow \sup l$ . For this choice  $\mathcal{I}_1$  and  $\mathcal{I}_3$  are empty sets and the density functions under each hypothesis are defined only on  $\mathcal{I}_2$ . Without loss of generality, assume that  $\alpha < 1$ ,  $\inf l = 0$  and  $\sup l = \infty$ . For this choice  $l_l \rightarrow 0$  implies  $\mu_1 = \lambda_1/(\alpha - 1) + 1$  and  $l_u \rightarrow \infty$  implies  $\mu_0 = \lambda_0/(\alpha - 1) + 1$ . Inserting these into one of the first two equations in (32), gives

$$\int_{\Omega} \left( \lambda_0 f_0(y)^{1-\alpha} + \lambda_1 \rho^{\alpha - 1} f_1(y)^{1-\alpha} \right)^{\frac{1}{1-\alpha}} \mathrm{d}y = (1-\alpha)^{\frac{1}{1-\alpha}} \,.$$
(38)

Similarly, from the third and fourth equations it follows that

$$\int_{\mathbb{R}} \left( \lambda_0 f_0(y)^{\frac{1-\alpha}{\alpha}} + \lambda_1 \rho^{\alpha-1} f_1(y)^{1-\alpha} f_0(y)^{\frac{(\alpha-1)^2}{\alpha}} \right)^{\frac{\alpha}{1-\alpha}} dy$$
$$= (1-\alpha)^{\frac{\alpha}{1-\alpha}} x(\alpha, \epsilon_0)$$
(39)

and

$$\int_{\mathbb{R}} \left( \lambda_1 f_1(y)^{\frac{1-\alpha}{\alpha}} + \lambda_0 \rho^{1-\alpha} f_0(y)^{1-\alpha} f_1(y)^{\frac{(\alpha-1)^2}{\alpha}} \right)^{\frac{\alpha}{1-\alpha}} dy$$
$$= (1-\alpha)^{\frac{\alpha}{1-\alpha}} x(\alpha, \epsilon_1).$$
(40)

Given  $\rho$  and  $\alpha$ , (38), (39), and (40) can jointly be solved to determine the space of maximum robustness parameters. As an example, consider  $\Omega = \mathbb{R}$ ,  $\rho = 1$  and  $\alpha = 1/2$ . This choice of  $\alpha$  corresponds to the squared Hellinger distance with an additional scaling factor of  $1/\alpha(1 - \alpha) = 4$ . Let  $a = \int_{-\infty}^{\infty} \sqrt{f_0(y)f_1(y)} dy$ . Then, the Equations (38)-(40) reduce to the polynomials in the Lagrangian multipliers  $\lambda_0$  and  $\lambda_1$ ,

$$\lambda_0^2 + \lambda_1^2 + 2\lambda_0\lambda_1 a - \frac{1}{4} = 0, \tag{41}$$

$$4 - 8\lambda_0 - 8\lambda_1 a - \epsilon_0 = 0, \tag{42}$$

$$4 - 8\lambda_1 - 8\lambda_0 a - \epsilon_1 = 0, \tag{43}$$

respectively. Solving (42) and (43) for  $\lambda_0$  and  $\lambda_1$ , respectively, and inserting the results into Equation (41) we get

0,

$$2\epsilon_1(a(\epsilon_0 - 4) + 4) - (4a + \epsilon_0 - 4)^2 - \epsilon_1^2 = 0.$$
 (44)

Equation (44) is quadratic in a and has two roots. One of the roots results in a = 1 for all  $\epsilon_0 = \epsilon_1$ , which is not plausible. Therefore, the correct root is,

$$a = \frac{1}{16} \left( 16 - 4\epsilon_1 + \epsilon_0(\epsilon_1 - 4) - \sqrt{(\epsilon_0 - 8)\epsilon_0(\epsilon_1 - 8)\epsilon_1} \right).$$
(45)

Notice that (45) is symmetric in  $\epsilon_0$  and  $\epsilon_1$ , i.e.,  $a(\epsilon_0, \epsilon_1) = a(\epsilon_1, \epsilon_0)$  for all  $(\epsilon_0, \epsilon_1)$ , as expected. Since  $0 \le a \le 1$  is known a priori, given a choice of  $\epsilon_i$ , the corresponding  $\epsilon_{1-i}$  can be determined from (45) easily, c.f., Section IV. A special case occurs whenever  $\epsilon = \epsilon_0 = \epsilon_1$ , which simplifies (45) to

$$\epsilon_{\max} = 4 - 2\sqrt{2(1+a)}.$$
 (46)

Maximum robustness parameters given by (45) and (46) are in agreement with the ones found in [20]. The case  $\alpha > 1$ , which implies  $\mu_0 = \lambda_0/(\alpha - 1)$  and  $\mu_1 = \lambda_1/(\alpha - 1)$ , can be examined similarly.

#### **IV. SIMULATIONS**

In this section, some simulations are performed to illustrate the theoretical derivations. Consider a simple hypothesis testing problem

$$\mathcal{H}_0^s : Y = W$$
  
$$\mathcal{H}_1^s : Y = W + A \tag{47}$$

where A > 0 is a known DC signal, W is a random variable which follows a symmetric Gaussian mixture distribution

$$W \sim \frac{1}{2} \left( \mathcal{N}(-\mu, \sigma^2) + \mathcal{N}(\mu, \sigma^2) \right), \tag{48}$$

where  $\mathcal{N}(\mu, \sigma^2)$  is a Gaussian distribution with mean  $\mu$  and variance  $\sigma^2$  and Y is a random variable on  $\Omega = \mathbb{R}$ , which is consistent with the data sample y. To account for uncertainties on Y under both hypotheses, let

$$F_0(y) := P(Y < y | \mathcal{H}_0^s) \quad \forall y$$
  

$$F_1(y) := P(Y < y | \mathcal{H}_1^s) \quad \forall y$$
(49)

be the nominal distributions, having the density functions  $f_0$ and  $f_1$  for the binary composite hypothesis testing problem given by (1) and (2). Note that the symmetry condition,  $f_1(y) = f_0(-y)$  for all y, does not hold, and  $l = f_1/f_0$  is not monotone. Assume  $\mu = 2$ ,  $\sigma = 1$  and A = 1 and let the robustness parameters be  $\epsilon_0 = 0.02$  and  $\epsilon_1 = 0.03$  for the  $(\alpha = 4)$ -divergence distance. This example demonstrates an extreme case, for which no straightforward simplification to the equations (19) and (20) exists, both in terms of reducing the number of equations as well as for the domain of integrals. Figure 3 illustrates the nominal density functions  $f_0$  and  $f_1$  along with the density functions of the corresponding least favorable densities (LFD)s  $g_0$  and  $g_1$ , for an equal a priori probability  $\rho = 1$ . It can be observed that LFDs intersect in three distinct intervals, each at the neighborhood of y = -1.5 + j for  $j \in \{0, 2, 4\}$ . In Fig. 4, the same simulation is repeated for  $\rho = 1.2$ . In Fig. 5 the nominal and least favorable likelihood ratios for the same example are shown. As it was given by (18), robustification of the simple hypothesis test corresponds to a non-linear transformation of the nominal likelihood ratios.

In the next simulation, all the parameters are fixed as before, except for  $\alpha$ . We are especially interested in the change in the lower and upper thresholds,  $l_l$  and  $l_u$ , for varying  $\alpha$ . Figure 6 illustrates the outcome of this simulation for  $\rho = 1$ . We can see that  $l_l$  and  $l_u$  tend to 1 for  $\alpha \to \infty$ . It is not straightforward to derive this from (19) and (20) for any  $f_0$ and  $f_1$ . However, if there exists a solution, which is true and unique by the KKT multipliers approach, it should satisfy  $D(f, g; \alpha) = \epsilon_i$  for any  $\alpha > 0$  and for all allowable  $\epsilon_i$ , cf. Section III-C. Assume that g is fixed and it does not depend on  $\alpha$ . Then, the integral  $\int_{\mathbb{R}} g^{\alpha} f^{1-\alpha} d\mu$  is 1 at  $\alpha = 0$  and  $\alpha = 1$ , convex in  $\alpha$ , and it is positive for all  $\alpha > 0$ , f and g. Hence,  $\lim_{\alpha\to\infty} \int_{\mathbb{R}} g^{\alpha} f^{1-\alpha} d\mu = \infty$  and  $\lim_{\alpha\to\infty} D(f, g; \alpha)$ is indeterminate. Using L'Hospital's rule twice we obtain

$$K = \lim_{\alpha \to \infty} D(g, f; \alpha) = \lim_{\alpha \to \infty} \frac{\int_{\mathbb{R}} \log^2(g/f) (g/f)^{\alpha} f d\mu}{2}.$$
 (50)

The integral  $\int_{\mathbb{R}} \log^2(g/f)(g/f)^{\alpha} f d\mu$  is also positive and convex in  $\alpha$ . This implies  $K \to \infty$  for  $\alpha \to \infty$ . Now, assume that g depends on  $\alpha$  and tends to a limiting distribution  $g^*$  for  $||g^* - f|| > 0$ , when  $\alpha \to \infty$ . Then, our conclusion does not change, i.e.,  $K \to \infty$  for  $\alpha \to \infty$ . Since  $D(f, g; \alpha)$  is finite, we require that  $\alpha \to \infty \Longrightarrow g^* \to f$ . Consequently, from (15) and (16),  $\hat{g}_i \to f_i$  whenever  $l_l \to 1$  and  $l_u \to 1$  explains the asymptotic of Figure 6 for any pair  $(f_0, f_1)$ .

Based on simulation results the following are conjectured:

- For a fixed ε<sub>0</sub> and ε<sub>1</sub>, increasing α leads to a monotone decrease in l<sub>u</sub> and monotone increase in l<sub>l</sub> on ℝ<sup>+</sup>\{0,1}.
- For a fixed  $\alpha$ , increasing  $\epsilon_0$ ,  $\epsilon_1$  or both introduces a nondecrease in  $l_u$ , non-increase in  $l_l$ , or both, given that  $\epsilon_0$  and  $\epsilon_1$  are less than their allowable maximum, cf. Section III-C.

The proof of these conjectures is an open problem.

From (19) and (20), it is clear that given a pair  $(\epsilon_0, \epsilon_1)$ , a slight change in  $\alpha$  changes the equations completely and in general  $l_l$  and  $l_u$  are functions of  $\alpha$ . In Figure 7, the robust decision rule  $\hat{\delta}$  for various  $\alpha$  values is plotted, without considering the dependency of  $l_l$  and  $l_u$  on  $\alpha$ . To do this,  $l_l \approx 0.605$  and  $l_u \approx 1.618$ , that are found for  $\rho = 1$ ,  $\alpha = 4$ ,  $\epsilon_0 = 0.02$  and  $\epsilon_0 = 0.03$ , are fixed constants in (17). Then, for  $\alpha = \{0.01, 10, 100\}$ , (17) is plotted. The decision rule  $\hat{\delta}$  tends to a step like function for an increasing  $\alpha$ , whereas for a smaller  $\alpha$ , i.e.,  $\alpha = 0.01$ , the decision rule is almost linear at the domain of the likelihood ratio for which  $\hat{l} = 1$ . This result is also in agreement with the previous findings;  $\hat{\delta}$  tends to a non-randomized likelihood ratio test for  $\alpha \to \infty$ , for which we obtained  $\hat{g}_i \to f_i$  and for  $(f_0, f_1)$  optimum decision rule is known to be a non-randomized likelihood ratio test.

In the following simulation, the simplified model  $(f_0(y) = f_1(-y))$  is tested for mean shifted Gaussian distributions;  $F_0 \sim \mathcal{N}(\mu_0, \sigma^2)$  and  $F_1 \sim \mathcal{N}(\mu_1, \sigma^2)$  with means  $\mu_0 = -1$ ,  $\mu_1 = 1$  and variance  $\sigma^2 = 1$ . The parameters of the composite test are chosen to be  $\rho = 1$ ,  $\epsilon_0 = 0.1$  and  $\epsilon_1 = 0.1$ . Here, our main interest is to observe the change in overlapping regions of least favorable density pairs for various  $\alpha$ . Figure 8 illustrates the outcome of this simulation. It can be seen that the overlapping region is convex for a negative  $\alpha$ , ( $\alpha = -10$ ) almost constant for  $\alpha = 0.01$  and concave for a positive  $\alpha$ , ( $\alpha = 10$ ). For the sake of clarity only three examples of  $\alpha$  are plotted.

In Figure 9, the false alarm and miss detection probabilities of the likelihood ratio test  $\delta$  for  $(f_0, f_1)$  are graphed and compared with the robust test  $\hat{\delta}$  for  $(\hat{g}_0, \hat{g}_1)$ . Two different robust parameter pairs and various signal to noise ratios (SNR)s, i.e., SNR =  $20 \log(A/\sigma)$  are considered. It can be seen that increasing the robustness parameters increases the false alarm and miss detection probabilities for all SNRs, as expected. The difference between false alarm and miss detection probabilities for the same robust test is small and it is more pronounced for low SNRs. For high SNRs the performance of two robust tests are close to each other. The reason is that for high SNRs maximum allowable robustness parameters become relatively high compared to the parameters of both robust settings. Although the nominal test has the lowest error rates, its performance can degrade considerably under uncertainties in the nominal model. The robust tests, on the other hand, have slightly higher error rates, but guaranteed power of the test, which indicates the trade-off between performance and robustness. Finally, in the last simulation, the 3D boundary surface of the maximum robustness parameters is determined for  $\alpha = 0.5$  (45) and is shown in Figure 10. This surface has a cropped rotated cone like shape, which is symmetric about its main diagonal, i.e., with respect to the plane  $\epsilon_0 = \epsilon_1$  on the space  $(\epsilon_0, \epsilon_1, a)$ . Notice that except for the points on the cone like shape that intersect with the  $(\epsilon_0, \epsilon_1, a = 0)$  plane, all other points on  $(\epsilon_0, \epsilon_1, a = 0)$  that are plotted in blue color are un-defined (rather than being valid points with a = 0), implying that for those points no minimax robust test exists.

## V. CONCLUSION

A robust version of the likelihood ratio test considering  $\alpha$ -divergence as the distance to characterize the uncertainty sets has been proposed. The existence of a saddle value to the minimax optimization problem was shown by adopting Sion's minimax theorem. The least favorable distributions, the robust decision rule as well as the robust version of the likelihood



Fig. 3. Nominal densities and the corresponding least favorable densities for  $\rho = 1$ ,  $\alpha = 4$ ,  $\epsilon_0 = 0.02$  and  $\epsilon_1 = 0.03$ .



ratio test were derived in two parameters and in three distinct regions on the co-domain of the nominal likelihood ratio. Two equations from where the parameters can be determined were also derived. It was found that the robust likelihood ratio doesn't depend on the parameter  $\alpha$  that characterizes the distance between the probability measures. When the nominal density functions satisfy a symmetry constraint, the two nonlinear equations were combined into a single equation. Finally, the upper bounds on the parameters that control the degree of robustness were derived. Open problems include proving the monotonicity of the parameters  $l_l$  and  $l_u$  for increasing  $(\epsilon_0, \epsilon_1)$ , or  $\alpha$ . It was shown that simulation results illustrate the theoretical results.

## APPENDIX A SIMPLIFICATION OF $\Phi_1$

From (31) consider the following steps for

$$\Phi_1 = \left(\frac{-1 + \lambda_0 + \lambda_1 + \mu_0 + \mu_1 - \alpha(-1 + \mu_0 + \mu_1)}{\lambda_1 + \lambda_0 (l/\rho)^{\alpha - 1}}\right)^{\frac{1}{\alpha - 1}}$$

• Dividing the numerator and the denominator by  $\lambda_0$  and replacing the term  $1 + \mu_0/\lambda_0 - \alpha \mu_0/\lambda_0$  by  $c_1^{\alpha-1}$  results



Fig. 5. Nominal and least favorable likelihood ratios  $(\hat{g}_1/\hat{g}_0 \text{ for } \rho = 1 \text{ and } \hat{g}_1^*/\hat{g}_0^*$  for  $\rho = 1.2$ ) for  $\alpha = 4$ ,  $\epsilon_0 = 0.02$  and  $\epsilon_1 = 0.03$ .



Fig. 6. Lower and upper thresholds,  $l_l$  and  $l_u$ , for a variable  $\alpha$ ,  $\rho = 1$ ,  $\epsilon_0 = 0.02$  and  $\epsilon_1 = 0.03$ .

in

$$\Phi_1 = \left(\frac{c_1^{\alpha-1} + (\lambda_1 - 1 + \mu_1 + \alpha - \alpha\mu_1)/\lambda_0}{(\lambda_1/\lambda_0) + (l/\rho)^{\alpha-1}}\right)^{\frac{1}{\alpha-1}}.$$

• Multiplying the numerator and the denominator of the result of the previous step by  $\lambda_0/\lambda_1$ , replacing the term  $1 - 1/\lambda_1 + \mu_1/\lambda_1 + \alpha/\lambda_1 - \alpha\mu_1/\lambda_1$  by  $c_3^{\alpha-1}$  and again multiplying both the numerator and the denominator by  $\lambda_1$  gives

$$\Phi_1 = \left(\frac{\lambda_0 c_1^{\alpha-1} + \lambda_1 c_3^{\alpha-1}}{\lambda_1 + \lambda_0 (l/\rho)^{\alpha-1}}\right)^{\frac{1}{\alpha-1}}$$

• The result of the previous step is free of parameters  $\mu_0$ and  $\mu_1$ , but still parameterized by  $\lambda_0$  and  $\lambda_1$ . To eliminate them, using the identities  $\lambda_0 = (1 - \alpha)/(c_1^{\alpha - 1} - c_2^{\alpha - 1})$ and  $\lambda_1 = (1 - \alpha)/(c_4^{\alpha - 1} - c_3^{\alpha - 1})$  leads to

$$\Phi_1 = \left(\frac{(c_1c_4)^{\alpha-1} + (c_2c_3)^{\alpha-1}}{c_1^{\alpha-1} - c_2^{\alpha-1} + (c_4^{\alpha-1} - c_3^{\alpha-1})(l/\rho)^{\alpha-1}}\right)^{\frac{1}{\alpha-1}}$$



Fig. 7. The decision rule  $\hat{\delta}$  for  $\alpha = \{0.01, 10, 100\}$ ,  $\rho = 1$ ,  $\epsilon_0 = 0.02$  and  $\epsilon_1 = 0.03$ .



Fig. 8. Nominal densities and the corresponding least favorable densities for  $\rho = 1, \, \epsilon_0 = 0.1$  and  $\epsilon_1 = 0.1$ .

• The result from the previous step depends only on  $c_1$ ,  $c_2$ ,  $c_3$ ,  $c_4$  and  $\alpha$ . Using the substitutions  $c_1 = c_3 l_l$ ,  $c_2 = c_4 l_u$ 



Fig. 9. False alarm and miss detection probabilities of  $\delta$ , (2), ( $\rho = 1$ ) for  $(f_0, f_1)$  compared to that of the robust decision rule  $\hat{\delta}$  for  $(\hat{g}_0, \hat{g}_1)$  when SNR is varied.



Fig. 10. All allowable pairs of maximum robustness parameters,  $(\epsilon_0, \epsilon_1)$ , w.r.t. all distances  $a \in [0, 1]$  for  $\alpha = 0.5$ .

and 
$$c_4 = k(l_l, l_u)c_3$$
 yields  

$$\Phi_1(l, l_l, l_u, c_3; \alpha, \rho) = c_3 \left( \frac{k(l_l, l_u)^{\alpha - 1}(l_l^{\alpha - 1} - l_u^{\alpha - 1})}{l_l^{\alpha - 1} - (k(l_l, l_u)l_u)^{\alpha - 1} + (k(l_l, l_u)^{\alpha - 1} - 1)(l/\rho)^{\alpha - 1}} \right)^{\frac{1}{\alpha - 1}}$$
(51)

## APPENDIX B SIMPLIFICATION OF $\hat{\delta}$

Since the equivalence of integration domains are given by (14), only

$$\hat{\delta} = \frac{\lambda_0 (-1 + \alpha + \lambda_1 + \mu_1 - \alpha \mu_1)}{(-1 + \alpha)(\lambda_0 + \lambda_1 (l/\rho)^{1-\alpha})} - \frac{\lambda_1 (\lambda_0 + \mu_0 - \alpha \mu_0) (l/\rho)^{1-\alpha}}{(-1 + \alpha)(\lambda_0 + \lambda_1 (l/\rho)^{1-\alpha})}, \quad \hat{l} = \rho$$

is required to be simplified. In the following, the simplification is performed in three steps and the domain term  $\hat{l} = \rho$  is omitted for the sake of simplicity:

Dividing the numerator and the denominator of the first term by λ<sub>1</sub> and the second term by λ<sub>0</sub>, and replacing the related terms by c<sub>1</sub><sup>α-1</sup> and c<sub>3</sub><sup>α-1</sup> results in

$$\hat{\delta} = \frac{\lambda_0}{-1+\alpha} \cdot \frac{c_3^{\alpha-1}}{\frac{\lambda_0}{\lambda_1} + (l/\rho)^{1-\alpha}} - \frac{\lambda_1}{-1+\alpha} \cdot \frac{c_1^{\alpha-1}(l/\rho)^{1-\alpha}}{1 + \frac{\lambda_1}{\lambda_0}(l/\rho)^{1-\alpha}} \\ = \frac{c_3^{\alpha-1} - c_1^{\alpha-1}(l/\rho)^{1-\alpha}}{(-1+\alpha)\left(\frac{1}{\lambda_1} + \frac{1}{\lambda_0}(l/\rho)^{1-\alpha}\right)}.$$

• The result of the previous step is free of parameters  $\mu_0$ and  $\mu_1$ , but still parameterized by  $\lambda_0$  and  $\lambda_1$ . To eliminate them, using the identities  $\lambda_0 = (1 - \alpha)/(c_1^{\alpha - 1} - c_2^{\alpha - 1})$ and  $\lambda_1 = (1 - \alpha)/(c_4^{\alpha - 1} - c_3^{\alpha - 1})$  leads to

$$\hat{\delta} = \frac{(l/\rho)^{1-\alpha}c_1^{\alpha-1} - c_3^{\alpha-1}}{c_4^{\alpha-1} - c_3^{\alpha-1} + (c_1^{\alpha-1} - c_2^{\alpha-1})(l/\rho)^{1-\alpha}}.$$

• The result from the previous step depends only on  $c_1$ ,  $c_2$ ,  $c_3$ ,  $c_4$  and  $\alpha$ . Using the substitutions  $c_1 = c_3 l_l$ ,  $c_2 = c_4 l_u$ 

and  $c_4 = k(l_l, l_u)c_3$  yields

$$\hat{\delta} = \frac{l_l^{\alpha-1}(l/\rho)^{1-\alpha} - 1}{(l_l^{\alpha-1} - (k(l_l, l_u)l_u)^{\alpha-1})(l/\rho)^{1-\alpha} + k(l_l, l_u)^{\alpha-1} - 1},$$
as wanted.

is wanted.

#### ACKNOWLEDGMENT

The authors would like to sincerely thank the anonymous reviewers for their valuable comments and suggestions to improve the quality of the paper. This work was supported by the LOEWE Priority Program Cocoon (http://www.cocoon.tu-darmstadt.de).

#### REFERENCES

- B. C. Levy, Principles of Signal Detection and Parameter Estimation, 1st ed. Springer Publishing Company, Incorporated, 2008.
- [2] —, "Robust hypothesis testing with a relative entropy tolerance," *IEEE Transactions on Information Theory*, vol. 55, no. 1, pp. 413–421, 2009.
- [3] G. Gül and A. M. Zoubir, "Robust hypothesis testing for modeling errors," in *IEEE Int. Conf. on Acoustics, Speech, and Signal Processing* (ICASSP), Vancouver, Canada, May 2013, pp. 5514–5518.
- [4] P. J. Huber, "A robust version of the probability ratio test," Ann. Math. Statist., vol. 36, pp. 1753–1758, 1965.
- [5] P. J. Huber and V. Strassen, "Robust confidence limits," Z. Wahrcheinlichkeitstheorie verw. Gebiete, vol. 10, pp. 269–278, 1968.
- [6] —, "Minimax tests and the Neyman-Pearson lemma for capacities," Ann. Statistics, vol. 1, pp. 251–263, 1973.
- [7] A. G. Dabak and D. H. Johnson, "Geometrically based robust detection," in *Proceedings of the Conference on Information Sciences and Systems*, Johns Hopkins University, Baltimore, MD, May 1994, pp. 73–77.
- [8] A. B. Chan, A. B. Chan, N. Vasconcelos, N. Vasconcelos, P. J. Moreno, and P. J. Moreno, "A family of probabilistic kernels based on information divergence," Tech. Rep., 2004.
- [9] A. Cichocki and S. ichi Amari, "Families of alpha- beta- and gammadivergences: Flexible and robust measures of similarities." *Entropy*, no. 6, pp. 1532–1568.
- [10] T. van Erven and P. Harremoës, "Rényi divergence and kullback-leibler divergence," *IEEE Transactions on Information Theory*, vol. 60, no. 7, pp. 3797–3820, 2014.
- [11] E. Tuncel, "On error exponents in hypothesis testing." *IEEE Transac*tions on Information Theory, no. 8, pp. 2945–2950.
- [12] A. Cichocki, H. Lee, Y.-D. Kim, and S. Choi, "Non-negative matrix factorization with alpha-divergence." *Pattern Recognition Letters*, no. 9, pp. 1433–1440.
- [13] A. O. Hero, B. Ma, O. Michel, and J. Gorman, "Alpha-divergence for classification, indexing and retrieval," University of Michigan, Tech. Rep., 2001.
- [14] L. Meziou, A. Histace, and F. Precioso, "Alpha-divergence maximization for statistical region-based active contour segmentation with nonparametric pdf estimations," in Acoustics, Speech and Signal Processing (ICASSP), 2012 IEEE International Conference on, March 2012, pp. 861–864.
- [15] Z. Bian, J. Ma, L. Tian, J. Huang, H. Zhang, Y. Zhang, W. Chen, and Z. Liang, "Penalized weighted alpha-divergence approach to sinogram restoration for low-dose x-ray computed tomography," in *Nuclear Science Symposium and Medical Imaging Conference (NSS/MIC)*, 2012 *IEEE*, Oct 2012, pp. 3675–3678.
- [16] J. Pluim, J. Maintz, and M. Viergever, "f-information measures in medical image registration," *Medical Imaging, IEEE Transactions on*, vol. 23, no. 12, pp. 1508–1516, Dec 2004.
- [17] D. Morales, L. Pardo, and I. Vajda, "Some new statistics for testing hypotheses in parametric models," *Journal of Multivariate Analysis*, vol. 62, no. 1, pp. 137 – 168, 1997.
- [18] T. Hobza, I. Molina, and D. Morales, "Multi-sample Rényi test statistics," Braz. J. Probab. Stat., vol. 23, no. 2, pp. 196–215, 12 2009.
- [19] M. Pardo and J. Pardo, "Use of Rényi's divergence to test for the equality of the coefficients of variation," *Journal of Computational and Applied Mathematics*, vol. 116, no. 1, pp. 93 – 104, 2000.
- [20] G. Gül and A. M. Zoubir, "Robust hypothesis testing with squared Hellinger distance," in *Proceedings of the 22nd European Signal Pro*cessing Conference (EUSIPCO), Lisbon, Portugal, Sep. 2014.

- [21] —, "Robust detection under communication constraints," in Proceedings of the IEEE 14th International Workshop on Advances in Wireless Communications (SPAWC), Darmstadt, Germany, Jun. 2013, pp. 109– 112.
- [22] —, "Minimax robust hypothesis testing," 2015, submitted for publication at *IEEE Transactions on Information Theory*.
  [23] S. Zheng, P.-Y. Kam, Y.-C. Liang, and Y. Zeng, "Spectrum sensing
- [23] S. Zheng, P.-Y. Kam, Y.-C. Liang, and Y. Zeng, "Spectrum sensing for digital primary signals in cognitive radio: A bayesian approach for maximizing spectrum utilization," *Wireless Communications, IEEE Transactions on*, vol. 12, no. 4, pp. 1774–1782, April 2013.
- [24] J.-P. Aubin and I. Ekeland, *Applied Nonlinear Analysis*. New York: J. Wiley, 1984.
- [25] M. Sion, "On general minimax theorems." Pacific Journal of Mathematics, vol. 8, no. 1, pp. 171–176, 1958.
- [26] "Tychonoff's theorem without the axiom of choice," Fundamenta Mathematicae, vol. 113, no. 1, pp. 21–35, 0 1981.
- [27] A. Tychonoff, "Über die topologische erweiterung von rumen," Mathematische Annalen, vol. 102, no. 1, pp. 544–561, 1930.
- [28] W. Rudin, *Principles of Mathematical Analysis*, ser. International series in pure and applied mathematics. Paris: McGraw-Hill, 1976.
- [29] D. Bertsekas, A. Nedić, and A. Ozdaglar, *Convex Analysis and Optimization*, ser. Athena Scientific optimization and computation series. Athena Scientific, 2003.