

A precise bare simulation approach to the minimization of some distances. I. Foundations

Michel Broniatowski and Wolfgang Stummer

Abstract

In information theory — as well as in the adjacent fields of statistics, machine learning, artificial intelligence, signal processing and pattern recognition — many *flexibilizations* of the omnipresent Kullback-Leibler information distance (relative entropy) and of the closely related Shannon entropy have become frequently used tools. To tackle corresponding constrained minimization (respectively maximization) problems by a newly developed *dimension-free bare (pure) simulation* method, is the main goal of this paper. Almost no assumptions (like convexity) on the set of constraints are needed, within our discrete setup of arbitrary dimension, and our method is precise (i.e., converges in the limit). As a side effect, we also derive an innovative way of constructing new useful distances/divergences. To illustrate the core of our approach, we present numerous solved cases. The potential for widespread applicability is indicated, too; in particular, we deliver many recent references for uses of the involved distances/divergences and entropies in various different research fields (which may also serve as an interdisciplinary interface).

Index Terms

f-divergences of Csiszar-Ali-Silvey-Morimoto type, power divergences, Kullback-Leibler information distance, relative entropy, Renyi divergences, Bhattacharyya distance, Jensen-Shannon divergence/distance, alpha-divergences, Shannon entropy, Renyi entropies, Bhattacharyya coefficient, Tsallis (cross) entropies, Cressie-Read measures, Hellinger distance, Euclidean norms, generalized maximum entropy method, importance sampling.

2nd November 2022

I. INTRODUCTION

THIS paper develops a new approach to perform the nonlinear constrained optimization of directed distances — and connected quantities — based on a random simulation method. Given a set $\Omega \subset \mathbb{R}^K$ with mild regularity, the problem is to find the $\inf \{\phi(\mathbf{Q}), \mathbf{Q} \in \Omega\}$ or the $\sup \{\phi(\mathbf{Q}), \mathbf{Q} \in \Omega\}$ (depending on what the optimization entails), where ϕ is a general objective function. Such ϕ that satisfy the assumptions which allow the proposed method to work, are named as “bare-simulation optimizable”.

In particular, this paper motivates the approach by considering $\phi(\mathbf{Q}) = D_\varphi(\mathbf{Q}, \mathbf{P})$ where $\mathbf{P} \in \mathbb{R}^K$ with positive entries and $D_\varphi(\cdot, \cdot)$ is a Csiszar-Ali-Silvey-Morimoto-type ([1],[2],[3]) φ -divergence of specific form for which a basic link between φ and an instrumental probability distribution can be stated; simulation of random variables under this distribution is the cornerstone for the optimization procedure. Most commonly used φ -divergences can enter this scheme, and the wide range of applicability of our proposal is presented through numerous solved cases.

Let us present an outlook of the core steps of the present approach. The *first step* to perform the distance (divergence) minimization is to normalize the vector¹ \mathbf{P} into a probability vector \mathbf{P}^2 (e.g. the φ -entropy triggering case $\mathbf{P} = (1, \dots, 1)$) is converted into the uniform-probability vector $\mathbf{P} = (1/K, \dots, 1/K)$). The *second step* follows from expressing the function φ in form of the Fenchel-Legendre transform of the cumulant (i.e., log moment) generating function of some random variable W ; a probabilistic construction based on i.i.d. copies of W allows to interpret $\inf \{D_\varphi(\mathbf{Q}, \mathbf{P}), \mathbf{Q} \in \Omega\}$ as an asymptotic characteristic for some explicitly constructable scheme involving both \mathbf{P} and the W_i 's. The *third and final step* consists in the construction of this probabilistic scheme, and it differs for the specific problem context.

In general, for a deterministic setup where the (transformed) probability vector $\mathbf{P} = (p_1, \dots, p_K)$ is completely known and Ω has non-void interior, one can construct the integer part $n_i := \lfloor np_i \rfloor$, partition the index set $\{1, \dots, n\}$ into K sets of size n_1, \dots, n_K and build a K -component vector; each component of this vector is an ad hoc weighted empirical mean of the W_i 's; up to standard transformations the empirical count of the visits of this vector in Ω approximates the solution of the optimization problem $\inf \{D_\varphi(\mathbf{Q}, \mathbf{P}), \mathbf{Q} \in \Omega\}$. Therefore, the resulting approximation can be performed straightforwardly: the (typically) very complicated minimization task is replaced by a much more comfortable — nevertheless convergent — random count procedure which can be based on a fast and accurate — pseudo, true, natural, quantum — random number generator.

In case of the statistical problem one has instead of a *known* probability vector \mathbf{P} a data-describing *sample* X_1, \dots, X_n of n i.i.d. (and even more general) copies of a discrete random variable X with *unknown* distribution (described by an unknown probability vector) \mathbf{P} , and Ω is now a subset of the probability simplex \mathbb{S}^K in \mathbb{R}^K , where \mathbb{S}^K obviously has void interior but nevertheless a very useful special structure. For such contexts, we can adapt the above-described bare-simulation method by

M. Broniatowski is with the LPSM, Sorbonne Université, 4 place Jussieu, 75252 Paris, France. ORCID 0000-0001-6301-5531.

W. Stummer is with the Department of Mathematics, University of Erlangen–Nürnberg (FAU), Cauerstrasse 11, 91058 Erlangen, Germany; e-mail: stummer@math.fau.de. ORCID 0000-0002-7831-4558. Corresponding author.

¹in this paper, vectors are taken to be row vectors

²with a slight abuse of notation; see the main text for a more comprehensive notation

basically using the corresponding (vectorized) sample-based *empirical probability mass function* as \mathbf{P} and accordingly calculate a specific scheme which also makes (in a slightly different way) use of the random variables W_i associated with the function φ , and estimate $\inf \{D_\varphi(\mathbf{Q}, \mathbf{P}), \mathbf{Q} \in \Omega\}$.

To work out the above-sketched program in detail, requires many known and newly developed results from a number of different topics from various fields of information theory, applied probability, computer simulation and analysis. This explains the need for a considerable number of different sections. We also found it necessary to present *explicit* solutions for the optimization procedure in a number of cases which are of common use in the corresponding domains; the range of applications of the outcoming directed distances covers many areas in engineering and in natural sciences, for which we selected recent prominently placed contributions with specific focus in relation with the aims of this paper. All this results in a paper of very substantial length.

Let us first review the class of directed (i.e. not necessarily symmetric) distances (also called divergences) $D(\mathbf{P}, \mathbf{Q})$ between two finite discrete³ (probability) distributions \mathbf{P}, \mathbf{Q} or between two general Euclidean vectors \mathbf{P}, \mathbf{Q} which are proved to be bare simulation optimizable; those serve as important (dis)similarity measures, proximity measures and discrepancy measures in various different research areas such as information theory, statistics, artificial intelligence, machine learning, signal processing, pattern recognition, physics, finance, etc.⁴. A major class are the above-mentioned φ -divergences $D_\varphi(\mathbf{P}, \mathbf{Q})$ of *Csiszar-Ali-Silvey-Morimoto* (CASM); this covers — with corresponding choices of φ — e.g. the omnipresent *Kullback-Leibler information distance/divergence* [6] (also known as relative entropy), the *Jensen-Shannon distance/divergence*, as well as the *power divergences* (also known as alpha-divergences, Cressie-Read measures, and Tsallis cross-entropies). For some comprehensive overviews on CASM φ -divergences, the reader is referred to the insightful books [7]–[13], the survey articles [4],[14]–[16] and the references therein; an imbedding of CASM φ -divergences to more general frameworks can be found e.g. in [17]–[21].

Frequently used special cases of the above-mentioned power divergences are e.g. the (squared) *Hellinger distance*, the *Pearson chi-square divergence*, and the *Neyman chi-square divergence*. Moreover, several deterministic transformations of power divergences are also prominently used in research, most notably the *Bhattacharyya distance* [22] and the more general *Renyi divergences* [23] (also known as Renyi cross-entropies); a comprehensive exposition of the latter is given e.g. in [24]. Some other important deterministic transformations of power divergences include the *Bhattacharyya coefficient* (cf. [22],[25],[26]) — which is also called *affinity* (cf. [27]) and *fidelity similarity* (cf. e.g. [28]) — as well as the *Bhattacharyya arccos distance* (cf. [26]) and the *Fisher distance* (also known as *Rao distance*, *geodesic distance*, cf. e.g. [28]). As shown below, by further explicit transformations we can also recover Sundaresan’s divergence [29][30].

Let us mention that from CASM φ -divergences one can also derive the widely used φ -entropies $\mathcal{E}_\varphi(\mathbf{Q})$ of a distribution \mathbf{Q} (and non-probability versions thereof) in the sense of [31] (see also [32]–[39]); these entropies can be constructed from $D_\varphi(\mathbf{Q}, \mathbf{P}^{unif})$ where \mathbf{P}^{unif} denotes the uniform distribution. Moreover, by use of certain deterministic transformations h one can also deduce the more general (h, φ) -entropies (and non-probability versions thereof) in the sense of [40] (see also e.g. [12]). As will be worked out in detail below, from this one can deduce as special cases a variety of prominently used quantities in research, such as for instance:

- the omnipresent *Shannon entropy* [41], the γ -order *Renyi entropy* [23], the γ -order *entropy of Havrda-Charvat* [42] (also called non-additive γ -order *Tsallis entropy* [43] in statistical physics), the $\tilde{\gamma}$ -order *entropy of Arimoto* [44], Vajda’s quadratic entropy [9], *Sharma-Mittal entropies* [45],
- the Euclidean γ -norms, as well as
- measures of diversity, heterogeneity and unevenness, like the *Gini-Simpson diversity index*, the *diversity index of Hill* [46], the *Simpson-Herfindahl index* (which is also known as *index of coincidence*, cf. [47] and its generalization in [48]), the *diversity index of Patil & Taillie* [49], the γ -*mean heterogeneity index* (see e.g. [50]); see also [51] and [52] for some interrelations with the above-mentioned entropies.

Given that the constraint set Ω reflects some *incomplete/partial* information about a system (e.g. moment constraints), the maximization over $\mathbf{Q} \in \Omega$ of the above-mentioned entropies, norms and diversity indices (and the more general (h, φ) -entropies) is important for many research topics, most notably manifested in Jaynes’s [53],[54] omnipresent, “universally applicable” *maximum entropy principle* (which employs the Shannon entropy), and its generalizations (see e.g. the books [55]–[58] for comprehensive surveys).

In the *statistical* context, the *minimization* $\inf_{\mathbf{Q} \in \Omega} D(\mathbf{Q}, \mathbf{P})$ of divergences from one distribution (respectively, its equivalent vector of frequencies) \mathbf{P} to an appropriate set Ω of distributions (frequency vectors) appears in a natural way, as indicated in the following. For instance, let $\mathbf{P} = \mathbf{P}_{true}$ be the *true* distribution of a mechanism which generates non-deterministic data and Ω be a pregiven *model* in the sense of a (parametric or non-parametric) family of distributions which serves as an “approximation”

³for reasons of technicality, in this paper we only deal with such kind of distributions; for instance, these can be also achieved from more involved systems by quantizations of observations represented by finite partitions of the observation/data space, or by making use of the dual representation for CASM φ -divergences (cf. [4], [5]).

⁴since there exists a vast literature on divergences and connected entropies in these fields, for the sake of brevity we will give in this introduction only some basic references; many corresponding concrete applications will be mentioned in the following sections.

(in fact, a collection of approximations) of the “truth” \mathbf{P}_{true} . If $\mathbf{P}_{true} \notin \Omega$ — e.g. since Ω reflects some simplifications of \mathbf{P}_{true} which is in line with the general scientific procedure — then the positive quantity $\Phi_{\mathbf{P}_{true}}(\Omega) := \inf_{\mathbf{Q} \in \Omega} D(\mathbf{Q}, \mathbf{P}_{true})$ can be used as an *index of model adequacy* in the sense of a degree of departure between the model and the truth (cf. [59], see also e.g. [60]–[62]); small index values should indicate high adequacy. If $\mathbf{P}_{true} \in \Omega$, then $\Phi_{\mathbf{P}_{true}}(\Omega) = 0$ which corresponds to full adequacy. This index of model adequacy $\Phi_{\mathbf{P}_{true}}(\Omega)$ can also be seen as *index of goodness/quality of approximation to the truth* or as *model misspecification error*, and it can be used for model assessment as well as for model search (model selection, model hunting) by comparing the indices $\Phi_{\mathbf{P}_{true}}(\Omega_1), \Phi_{\mathbf{P}_{true}}(\Omega_2), \dots$ of competing models $\Omega_1, \Omega_2, \dots$ and choosing the one with the smallest index; this idea can be also used for classification (e.g. analogously to [63] who deal with continuous (rather than discrete) distributions) where the Ω_i are interpreted as (possibly data-derived but fixed) classes which are disjoint and non-exhaustive.

Typically, in statistical analyses the true distribution \mathbf{P}_{true} is unknown and is either replaced by a hypothesis-distribution \mathbf{P}_{hyp} or by a distribution \mathbf{P}_{data} derived from data (generated by \mathbf{P}_{true}) which converges to \mathbf{P}_{true} as the data/sample size tends to infinity (e.g. \mathbf{P}_{data} may be the well-known empirical distribution or a conditional distribution). Correspondingly, $\Phi_{\mathbf{P}_{data}}(\Omega) = \min_{\mathbf{Q} \in \Omega} D(\mathbf{Q}, \mathbf{P}_{data})$ reflects a data-derived approximation (estimate) of the index of model adequacy (resp. of the model misspecification error) from which one can cast corresponding model-adequacy tests and related goodness-of-fit tests. Accordingly, our new above-described procedure is well fitted for model choice. (For reasons of efficiency, especially in high dimension K). After choosing the most adequate model (say Ω_{i_0} , one can then tackle the problem of finding the corresponding (not necessarily existent or unique) best-model-member/element choice (i.e., the *minimizer*) $\arg \min_{\mathbf{Q} \in \Omega_{i_0}} D(\mathbf{Q}, \mathbf{P}_{data})$ which amounts to the well-known corresponding *minimum distance estimator* (for comprehensive surveys on divergence-based statistical testing and estimation, the reader is referred to e.g. the references in [21]); for the sake of brevity, this will be treated in a follow-up paper.

Besides the above-mentioned principal overview, let us now briefly discuss some existing *technical issues* for the minimization of CASM φ -divergences $\Phi_{\mathbf{P}}(\Omega) := \inf_{\mathbf{Q} \in \Omega} D_{\varphi}(\mathbf{Q}, \mathbf{P})$. For (not necessarily discrete) probability distributions/measures \mathbf{P} and sets Ω of probability distributions/measures satisfying a finite set of linear equality constraints, $\Phi_{\mathbf{P}}(\Omega)$ has been characterized in [64] and more recently in [65]–[68] among others, in various contexts; those results extend to inequality constraints. Minimizations of γ -order Rényi divergences on γ -convex sets Ω are studied e.g. in [69], whereas [70] [71] investigate minimizations of Sundaresan’s divergence on certain convex sets Ω .

To our knowledge, no general representation for $\Phi_{\mathbf{P}}(\Omega)$ for a positive distribution/measure \mathbf{P} (respectively, for a Euclidean vector with positive components) and a general set Ω of signed measures (respectively, of Euclidean vectors with components of arbitrary sign) exists. At the contrary, many *algorithmic* approaches for such minimization problems have been proposed; they mostly aim at finding minimizers more than at the evaluation of the minimum divergence itself, which is obtained as a by-product. Moreover, it is well-known that *such kind* of CASM φ -divergence minimization approaches may be hard to tackle or even intractable via usual methods such as the omnipresent gradient descent method and versions thereof, especially for non-parametric or semi-parametric Ω in sufficiently high-dimensional situations. For instance, Ω may consist (only) of constraints on moments or on L-moments (see e.g. [72]); alternatively, Ω may be e.g. a tubular neighborhood of a parametric model (see e.g. [73],[74]). The same intractability problem holds for the above-mentioned (h, φ) -entropy maximization problems. In the light of this, the goals of this paper are:

- to solve constrained minimization problems of a large range of CASM φ -divergences and deterministic transformations thereof (respectively constrained maximization problems of (h, φ) -entropies including Euclidean norms and diversity indices), by means of a newly developed *dimension-free bare (pure) simulation* method which is precise (i.e., converges in the limit) and which needs almost no assumptions (like convexity) on the set Ω of constraints; in doing so, for the sake of brevity we concentrate on finding/computing the minimum divergences themselves rather than the corresponding minimizers (to achieve the latter, e.g. dichotomous search could be used in a subsequent step, however);
- to derive a method of constructing new useful distances/divergences;
- to present numerous solved cases in order to illuminate our method and its potential for wide-spread applicability; as we go along, we also deliver many recent references for uses of the outcoming distances/divergences and entropies (covering in particular all the above-mentioned ones).

This agenda is achieved in the following way. In the next Section II, we briefly introduce the principal idea of our new bare-simulation optimization paradigm. After manifesting the fundamentally employed class of CASM φ -divergences in Section III, the correspondence between the function φ and the distribution of the random variable W is stated in Section IV. The random simulation scheme aiming at bare-simulation minimizability for deterministic problems is presented in Section V, and Section VI adapts this scheme to the statistical context, together with some deterministic simplex-constraints variant, also providing estimations or approximations for bounds of the solutions on the minimization/maximization problem at hand; the asymptotics which justify the simulation scheme as providing approximations for the optimization problem is discussed. Sections VII, VIII and IX present specific optimization schemes for the important Rényi family of divergences, for the constrained optimization of entropies, as well as a number of prominent deterministic optimization problems which can benefit from the wide applicability of this approach (high dimension and highly disconnected constraint sets, linear assignment problem with side constraints, etc).

Section X presents the construction of estimators together with some importance-sampling procedures and explicit algorithms. Section XI provides explicit constructions for the distribution of the instrumental random variable W for wide classes of functions φ ; this section is based on a new theorem which allows for an easy bridge between moment generating functions and their Fenchel-Legendre transform. Finally, Section XII provides explicitly solved cases which cover both deterministic and statistical important bare-simulation amenable problems. Two important proofs are given in the Appendices A and B of the paper; the remaining technical proofs and discussions are presented in the Supplementary Material.

Let us finally mention that a first simulation-based algorithm in vein with the present proposal has been developed by [75] in the restricted setup of risk estimation for power divergences. The present paper extends this very considerably by (amongst other things) treating *general* CASM φ -divergences *as well as* related entropies, by dealing with corresponding *general* optimization problems of *both* deterministic and stochastic type, respectively, and by developing new *types* of more sophisticated simulation algorithms.

II. A NEW MINIMIZATION PARADIGM

We concern with minimization problems of the following type, where \mathcal{M} is a topological space and \mathcal{T} is the Borel σ -field over a given base on \mathcal{M} ; e.g. take $\mathcal{M} = \mathbb{R}^K$ to be the K -dimensional Euclidean space equipped with the Borel σ -field \mathcal{T} .

Definition 1: A measurable function $\Phi : \mathcal{M} \mapsto \mathbb{R} \cup \{-\infty, \infty\}$ and measurable set $\Omega \subset \mathcal{M}$ ⁵ are called “bare-simulation minimizable” (BS-minimizable) respectively “bare-simulation maximizable” (BS-maximizable) if for

$$\Phi(\Omega) := \inf_{Q \in \Omega} \{\Phi(Q)\} \in]-\infty, \infty[\quad \text{respectively} \quad \Phi(\Omega) := \sup_{Q \in \Omega} \{\Phi(Q)\} \in]-\infty, \infty[\quad (1)$$

there exists a measurable function $G : [0, \infty[\mapsto \mathbb{R}$ as well as a sequence $((\mathcal{X}_n, \mathcal{A}_n, \mathbb{P}_n))_{n \in \mathbb{N}}$ of probability spaces and on them a sequence $(\xi_n)_{n \in \mathbb{N}}$ ⁶ of \mathcal{M} -valued random variables such that

$$G\left(-\lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_n[\xi_n \in \Omega]\right) = \inf_{Q \in \Omega} \Phi(Q) = \Phi(\Omega) \quad (2)$$

$$\text{respectively} \quad G\left(-\lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_n[\xi_n \in \Omega]\right) = \sup_{Q \in \Omega} \Phi(Q) = \Phi(\Omega); \quad (3)$$

in situations where Φ is fixed and different Ω 's are considered, we say that “ Φ is bare-simulation minimizable (BS-minimizable) on Ω ” respectively “ Φ is bare-simulation maximizable (BS-maximizable) on Ω ”.

The basic idea/incentive of this new approach is: if a minimization problem (1) has no explicit solution and is computationally intractable (or unfeasible) but can be shown to be BS-minimizable with concretely constructable $(\xi_n)_{n \in \mathbb{N}}$ and $(\mathbb{P}_n)_{n \in \mathbb{N}}$, then one can basically simulate the log-probabilities $-\frac{1}{n} \log \mathbb{P}_n[\xi_n \in \Omega]$ for large enough integer $n \in \mathbb{N}$ to obtain an approximation of (1) without having to evaluate the corresponding (not necessarily unique) minimizer, where the latter is typically time-costly. Finding minimizers can be performed e.g. through dichotomic search, once an algorithm leading to the minimal value of the divergence on adequate families of sets Ω is at hand; for the sake of brevity, this is omitted in the current paper.

For reasons of transparency, we *start* to demonstrate this approach for the following important/prominent class of *deterministic* constrained minimization problems with the following components:

- (i) \mathcal{M} is the K -dimensional Euclidean space \mathbb{R}^K , i.e. Ω is a set of vectors Q with a number of K components (where K may be huge, as it is e.g. the case in big data contexts);
- (ii) $\Phi(\cdot) := \Phi_P(\cdot)$ depends on some known vector P in \mathbb{R}^K with K nonnegative components;
- (iii) $\Phi_P(\cdot)$ is a “directed distance” (divergence) from P into Ω in the sense of $\Omega \ni Q \mapsto \Phi_P(Q) := D(Q, P)$, where $D(\cdot, \cdot)$ has the two properties “ $D(Q, P) \geq 0$ ” and “ $D(Q, P) = 0$ if and only if $Q = P$ ”. In particular, $D(\cdot, \cdot)$ needs neither satisfy the symmetry $D(Q, P) = D(P, Q)$ nor the triangular inequality.

In other words, (1) together with (i)-(iii) constitutes a *deterministic* constrained distance/divergence-minimization problem; we design a “universal” method to solve such problems by constructing appropriate (cf.(2)) sequences $(\xi_n)_{n \in \mathbb{N}}$ of \mathbb{R}^K -valued random variables, for all directed distances $D(\cdot, \cdot)$ from a large subclass of the important omnipresent Csiszar-Ali-Silvey-Morimoto φ -divergences (also called f -divergences) given in Definition 2 below.

⁵i.e. $\Omega \in \mathcal{T}$

⁶in order to emphasize the dependence on Φ , one should use the notations $(\xi_{\Phi, n})_{n \in \mathbb{N}}$, $\mathbb{P}_{\Phi, n}$, etc.; this is avoided for the sake of a better readability.

III. DIRECTED DISTANCES

To begin with, concerning the above-mentioned point (i) we take the K –dimensional Euclidean space $\mathcal{M} = \mathbb{R}^K$, denote from now on — as usual — its elements (i.e. vectors) in boldface letters, and also employ the subsets

$$\begin{aligned}\mathbb{R}_{\neq 0}^K &:= \{\mathbf{Q} := (q_1, \dots, q_K) \in \mathbb{R}^K : q_i \neq 0 \text{ for all } i = 1, \dots, K\}, \\ \mathbb{R}_{> 0}^K &:= \{\mathbf{Q} := (q_1, \dots, q_K) \in \mathbb{R}^K : q_i > 0 \text{ for all } i = 1, \dots, K\}, \\ \mathbb{R}_{\geq 0}^K &:= \{\mathbf{Q} := (q_1, \dots, q_K) \in \mathbb{R}^K : q_i \geq 0 \text{ for all } i = 1, \dots, K\}, \\ \mathbb{R}_{\neq 0}^{\geq 0} &:= \mathbb{R}_{\geq 0}^K \setminus \{\mathbf{0}\} := \{\mathbf{Q} := (q_1, \dots, q_K) \in \mathbb{R}_{\geq 0}^K : q_i \neq 0 \text{ for some } i = 1, \dots, K\}, \\ \mathbb{S}^K &:= \{\mathbf{Q} := (q_1, \dots, q_K) \in \mathbb{R}_{\geq 0}^K : \sum_{i=1}^K q_i = 1\} \quad (\text{simplex of probability vectors, probability simplex}), \\ \mathbb{S}_{> 0}^K &:= \{\mathbf{Q} := (q_1, \dots, q_K) \in \mathbb{R}_{> 0}^K : \sum_{i=1}^K q_i = 1\}.\end{aligned}$$

Concerning the directed distances $D(\cdot, \cdot)$ in (ii) and (iii), we deal with the following

Definition 2: (a) Let the “divergence generator” be a lower semicontinuous convex function $\varphi :]-\infty, \infty[\rightarrow [0, \infty]$ satisfying $\varphi(1) = 0$. Furthermore, for the effective domain $\text{dom}(\varphi) := \{t \in \mathbb{R} : \varphi(t) < \infty\}$ we assume that its interior $\text{int}(\text{dom}(\varphi))$ is non-empty which implies that $\text{int}(\text{dom}(\varphi)) =]a, b[$ for some $-\infty \leq a < 1 < b \leq \infty$. Moreover, we suppose that φ is strictly convex in a non-empty neighborhood $]t_-^{sc}, t_+^{sc}[\subseteq]a, b[$ of one ($t_-^{sc} < 1 < t_+^{sc}$). Also, we set $\varphi(a) := \lim_{t \downarrow a} \varphi(t)$ and $\varphi(b) := \lim_{t \uparrow b} \varphi(t)$ (these limits always exist). The class of all such functions φ will be denoted by $\tilde{\Upsilon}(]a, b[)$. A frequent choice is e.g. $]a, b[=]0, \infty[$ or $]a, b[=]-\infty, \infty[$.

(b) For $\varphi \in \tilde{\Upsilon}(]a, b[)$, $\mathbf{P} := (p_1, \dots, p_K) \in \mathbb{R}_{\neq 0}^K$ and $\mathbf{Q} := (q_1, \dots, q_K) \in \Omega \subset \mathbb{R}^K$, we define the Csiszar-Ali-Silvey-Morimoto (CASM) φ –divergence

$$\Phi_{\mathbf{P}}(\mathbf{Q}) := D_{\varphi}(\mathbf{Q}, \mathbf{P}) := \sum_{k=1}^K p_k \cdot \varphi\left(\frac{q_k}{p_k}\right) \geq 0. \quad (4)$$

As usual, in (4) we employ the three conventions that $p \cdot \varphi\left(\frac{0}{p}\right) = p \cdot \varphi(0) > 0$ for all $p > 0$, and $0 \cdot \varphi\left(\frac{q}{0}\right) = q \cdot \lim_{x \rightarrow \infty} \frac{\varphi(x \cdot \text{sgn}(q))}{x \cdot \text{sgn}(q)} > 0$ for $q \neq 0$ (employing the sign of q), and $0 \cdot \varphi\left(\frac{0}{0}\right) := 0$. Throughout the paper, we only consider constellations $(\varphi, \mathbf{P}, \Omega)$ for which the very mild condition $\Phi_{\mathbf{P}}(\Omega) := \inf_{\mathbf{Q} \in \Omega} D_{\varphi}(\mathbf{Q}, \mathbf{P}) \neq \infty$ holds.

For probability vectors \mathbb{P} and \mathbf{Q} in \mathbb{S}^K , the φ –divergences $D_{\varphi}(\mathbf{Q}, \mathbb{P})$ were introduced by Csiszar [1], Ali & Silvey [2] and Morimoto [3] (where the first two references even deal with more general probability distributions); for some comprehensive overviews — including statistical applications to goodness-of-fit testing and minimum distance estimation — the reader is referred to the insightful books [7]–[13], the survey articles [4], [14]–[16], and the references therein. Some exemplary recent studies and applications of CASM φ –divergences appear e.g. in [76]–[90]. For the setup of $D_{\varphi}(\mathbf{Q}, \mathbf{P})$ for vectors \mathbf{P}, \mathbf{Q} with non-negative components the reader is referred to e.g. [91] (who deal with even more general nonnegative measures and give some statistical as well as information-theoretic applications) and [92] (including applications to iterative proportional fitting). The case of φ –divergences for vectors with arbitrary components can be extracted from e.g. [66] who actually deal with finite signed measures. For a comprehensive technical treatment, see also [20].

Clearly, from (4) it is obvious that in general $D_{\varphi}(\mathbf{Q}, \mathbf{P}) \neq D_{\varphi}(\mathbf{P}, \mathbf{Q})$ (non-symmetry). Moreover, it is straightforward to deduce that $D_{\varphi}(\mathbf{Q}, \mathbf{P}) = 0$ if and only if $\mathbf{Q} = \mathbf{P}$ (reflexivity). By appropriate choice of φ , one can get as special cases many very prominent divergences which are frequently used in information theory and its applications to e.g. statistics, artificial intelligence, and machine learning. We shall address them later on as we go along.

Remark 3: Since, in general, our methods work also for *non-probability* vectors \mathbf{Q} and \mathbf{P} , we can also deal with — plain versions and transformations of — *weighted φ –divergences* of the form

$$D_{\varphi}^{wei}(\mathbf{Q}, \mathbf{P}) := \sum_{k=1}^K c_k \cdot p_k \cdot \varphi\left(\frac{q_k}{p_k}\right) \geq 0 \quad (5)$$

where $c_k > 0$ ($k = 1, \dots, K$) are weights which not necessarily add up to one. Indeed, we can formally rewrite

$$\inf_{\mathbf{Q} \in \Omega} D_{\varphi}^{wei}(\mathbf{Q}, \mathbf{P}) = \inf_{\mathbf{Q}^{wei} \in \Omega^{wei}} D_{\varphi}(\mathbf{Q}^{wei}, \mathbf{P}^{wei})$$

where $\mathbf{P}^{wei} := (c_1 \cdot p_1, \dots, c_K \cdot p_K)$, $\mathbf{Q}^{wei} := (c_1 \cdot q_1, \dots, c_K \cdot q_K)$ and Ω^{wei} is the corresponding rescaling of Ω . Of course, all the necessary technicalities for the φ –divergences (see below) have to be adapted to the weighted φ –divergences; for the sake of brevity, this will not be discussed in detail. Notice that $\mathbf{P}^{wei}, \mathbf{Q}^{wei}$ are generally not probability vectors anymore, even if \mathbf{Q}, \mathbf{P} are probability vectors. In the latter case, and under the assumption $\sum_{k=1}^K c_k = 1$, the divergences (5) coincide with the discrete versions of the (*c*–)local divergences of Avlogiaris et al. [93], [94] who also give absolutely-continuous versions and beyond (see also [20] for an imbedding in a general divergence framework).

IV. CONSTRUCTION PRINCIPLES: THE CORNERSTONE

For the divergence-minimization $\Phi_{\mathbf{P}}(\Omega) := \inf_{\mathbf{Q} \in \Omega} D_{\varphi}(\mathbf{Q}, \mathbf{P})$ (and its variants), in order to obtain the (2)-conform construction of the desired sequences $(\xi_n)_{n \in \mathbb{N}}$ of \mathbb{R}^K -valued random variables and $(\mathbb{P}_n)_{n \in \mathbb{N}}$ of probability distributions, we will assume (directly or after some multiplication) that the divergence generator $\varphi \in \tilde{\Upsilon}(]a, b[)$ has the additional property that it can be represented as

$$\varphi(t) = \sup_{z \in \mathbb{R}} \left(z \cdot t - \log \int_{\mathbb{R}} e^{z \cdot y} d\zeta(y) \right), \quad t \in \mathbb{R}, \quad (6)$$

for some probability distribution/measure ζ on the real line \mathbb{R} such that the function $z \mapsto MGF_{\zeta}(z) := \int_{\mathbb{R}} e^{z \cdot y} d\zeta(y)$ is finite on some open interval containing zero⁷. From this, we shall construct a sequence $(W_n)_{n \in \mathbb{N}}$ of i.i.d. copies of a random variable W whose distribution is ζ (i.e. $\mathbb{P}[W \in \cdot] = \zeta[\cdot]$ under some \mathbb{P}), from which the desired $(\xi_n)_{n \in \mathbb{N}}$ and $(\mathbb{P}_n)_{n \in \mathbb{N}}$ will be constructed. The class of functions $\varphi \in \tilde{\Upsilon}(]a, b[)$ satisfying the representability (6) will be denoted by $\Upsilon(]a, b[)$. Notice that $\Upsilon(]a, b[)$ contains many divergence generators; this and φ -construction principles will be developed in Section XI below.

The representability (6) is *the* cornerstone for our approach, and opens the gate to make use of simulation methods in appropriate contexts. We first develop this approach for *deterministic* minimization problems in the following Section V (where we retransform the generator φ into $\tilde{c} \cdot \varphi$ for strictly positive scales \tilde{c} and where $\mathbb{P}_n \equiv \mathbb{P}$). Thereafter, in Section VI, we deal with the setup where \mathbf{P} is identified with an unknown probability vector in the simplex \mathbb{S}^K which is supposed to be the limit (as n tends to infinity) of the data-based empirical distribution pertaining to a collection of observations $\mathbf{X}_n := (X_1, \dots, X_n)$; this amounts to the estimation of $\Phi_{\mathbf{P}}(\Omega)$ based on \mathbf{X}_n , leading to the important “minimization-distance estimation problem” in statistics, artificial intelligence and machine learning.

V. BS-MINIMIZABILITY/AMENABILITY: DETERMINISTIC MINIMIZATION PROBLEMS

Problem 4: For pregiven $\varphi \in \tilde{\Upsilon}(]a, b[)$, positive-entries vector $\mathbf{P} := (p_1, \dots, p_K) \in \mathbb{R}_{>0}^K$ (or from some subset thereof), and subset $\Omega \subset \mathbb{R}^K$ (also denoted in boldface letters, with a slight abuse of notation) with regularity properties

$$cl(\Omega) = cl(int(\Omega)), \quad int(\Omega) \neq \emptyset, \quad (7)$$

find

$$\Phi_{\mathbf{P}}(\Omega) := \inf_{\mathbf{Q} \in \Omega} D_{\varphi}(\mathbf{Q}, \mathbf{P}), \quad (8)$$

provided that

$$\inf_{\mathbf{Q} \in \Omega} D_{\varphi}(\mathbf{Q}, \mathbf{P}) < \infty. \quad (9)$$

An immediate consequence thereof is — for pregiven $\varphi \in \tilde{\Upsilon}(]a, b[)$ — the treatment of the more flexible problem

$$\Phi_{\mathbf{P}, h}(\Omega) := \inf_{\mathbf{Q} \in \Omega} h(D_{\varphi}(\mathbf{Q}, \mathbf{P})) = h\left(\inf_{\mathbf{Q} \in \Omega} D_{\varphi}(\mathbf{Q}, \mathbf{P})\right) \quad (10)$$

for any continuous strictly increasing function $h : \mathcal{H} \mapsto \mathbb{R}$ with $\mathcal{H} := [0, \infty[$ and extension $h(\infty) := \sup_{y \in \mathcal{H}} h(y)$ (depending on the problem, a sufficiently large $\mathcal{H} \subset [0, \infty[$ may be enough), respectively of

$$\sup_{\mathbf{Q} \in \Omega} h(D_{\varphi}(\mathbf{Q}, \mathbf{P})) = h\left(\inf_{\mathbf{Q} \in \Omega} D_{\varphi}(\mathbf{Q}, \mathbf{P})\right) \quad (11)$$

for any continuous strictly decreasing function $h : \mathcal{H} \mapsto \mathbb{R}$ and extension $h(\infty) := \inf_{y \in \mathcal{H}} h(y)$.

Remark 5: (a) By the basic properties of φ , it follows that for all $c > 0$ the level sets $\varphi_c := \{x \in \mathbb{R} : \varphi(x) \leq c\}$ are compact and so are the level sets $\Gamma_c := \{\mathbf{Q} \in \mathbb{R}^K : D_{\varphi}(\mathbf{Q}, \mathbf{P}) \leq c\}$ for all $c > 0$.

(b) When Ω is not closed but merely satisfies (7), then the infimum in (8) may not be reached in Ω although being finite; however we aim for finding the *infimum/minimum* in (8). Finding the *minimizers* in (8) is another question. For instance, this can be solved whenever, additionally, Ω is a closed set which implies the existence of minimizers in Ω . In this case, and when the number of such minimizers is finite, those can be e.g. approximated by dichotomic search. For the sake of brevity, this will not be addressed in this paper.

(c) The purpose of condition (7) is to get rid of the lim sup type and lim inf type results in our below-mentioned “bare-simulation” approach and to obtain *limit*-statements which motivate our construction. In practice, it is enough to verify $\Omega \subseteq cl(int(\Omega))$, which is equivalent to the left-hand part of (7). Clearly, any open set $\Omega \subset \mathbb{R}^K$ satisfies the left-hand part of (7). In the subsetup where Ω is a closed convex set and $int(\Omega) \neq \emptyset$, (7) is satisfied and the minimizer $\mathbf{Q}_{min} \in \Omega$ in (8) is attained and even unique. When Ω is open and satisfies (7), then the infimum in (8) exists but is reached at some generalized projection of \mathbf{P} on Ω (see [95] for the Kullback-Leibler divergence case of probability measures, which extends to any φ -divergence in

⁷in particular, this implies that ζ has light tails;

our framework).

(d) Without further mentioning, the regularity condition (7) is supposed to hold in the *full* topology. Of course, $\text{int}(\mathbb{S}^K) = \emptyset$ and thus, for the important probability-vector setup $\Omega \subset \mathbb{S}^K$ the condition (7) is violated which requires extra refinements (cf. Section VI below). The same is needed for $\Omega \subset A \cdot \mathbb{S}^K$ for some $A \neq 1$, since obviously $\text{int}(A \cdot \mathbb{S}^K) = \emptyset$; such a context appears naturally e.g. in connection with mass transportation problems (cf. (104) below) and with distributed energy management (cf. the paragraph after (107)).

(e) Often, Ω will present a (discrete) model⁸. Since Ω is assumed to have a non-void interior (cf. the right-hand part of (7)), this will exclude parametric models $\Omega := \{\mathbf{Q}_\theta : \theta \in \Theta\}$ for some $\Theta \subset \mathbb{R}^d$ ($d < K - 1$), for which $\theta \mapsto \mathbf{Q}_\theta$ constitutes a curve/surface in \mathbb{R}^K ; however, for such a situation, one can employ standard minimization principles. Our approach is predestined for *non- or semiparametric* models, instead. For instance, (7) is valid for appropriate *tubular neighborhoods* of parametric models or for more general non-parametric settings such as e.g. shape constraints.

Let us now present our *new bare-simulation approach* (cf. Definition 1) for solving the distance-optimization Problem 4:

(BS1) Step 1: equivalently rewrite (8) such that the vector \mathbf{P} “turns into” a probability vector $\tilde{\mathbb{P}}$. More exactly, define $M_{\mathbf{P}} := \sum_{i=1}^K p_i > 0$ and let $\tilde{\mathbb{P}} := \mathbf{P}/M_{\mathbf{P}}$, and for \mathbf{Q} in Ω , let $\tilde{\mathbf{Q}} := \mathbf{Q}/M_{\mathbf{P}}$ (notice that $\tilde{\mathbf{Q}}$ may be a non-probability vector). With the function $\tilde{\varphi} \in \tilde{\Upsilon}(]a, b[)$ defined through $\tilde{\varphi} := M_{\mathbf{P}} \cdot \varphi$, we obtain

$$D_\varphi(\mathbf{Q}, \mathbf{P}) = \sum_{k=1}^K p_k \cdot \varphi\left(\frac{q_k}{p_k}\right) = \sum_{k=1}^K M_{\mathbf{P}} \cdot \tilde{p}_k \cdot \varphi\left(\frac{M_{\mathbf{P}} \cdot \tilde{q}_k}{M_{\mathbf{P}} \cdot \tilde{p}_k}\right) = D_{\tilde{\varphi}}(\tilde{\mathbf{Q}}, \tilde{\mathbb{P}}). \quad (12)$$

It follows that the solution of (8) coincides with the one of the problem of finding

$$\tilde{\Phi}_{\tilde{\mathbb{P}}}(\tilde{\Omega}) := \inf_{\tilde{\mathbf{Q}} \in \tilde{\Omega}} D_{\tilde{\varphi}}(\tilde{\mathbf{Q}}, \tilde{\mathbb{P}}), \quad \text{with } \tilde{\Omega} := \Omega/M_{\mathbf{P}}; \quad (13)$$

as a side remark, one can see that in such a situation the rescaling of the divergence generator φ is important, which is one incentive that we incorporate multiples of φ below.

As an important special case we get for the choice $\mathbf{P} := (1, \dots, 1) := \mathbf{1}$ that the “prominent/frequent” separable nonlinear optimization problem of finding the optimal value $\inf_{\mathbf{Q} \in \Omega} \sum_{k=1}^K \varphi(q_k)$ — with objective (e.g. cost, energy, purpose) function $\varphi \in \tilde{\Upsilon}(]a, b[)$ and constraint set (choice set, search space) Ω — can be imbedded into our BS-approach by

$$\inf_{\mathbf{Q} \in \Omega} \sum_{k=1}^K \varphi(q_k) = \inf_{\mathbf{Q} \in \Omega} D_\varphi(\mathbf{Q}, \mathbf{1}) = \tilde{\inf}_{\tilde{\mathbf{Q}} \in \Omega/K} D_{K \cdot \varphi}(\tilde{\mathbf{Q}}, \mathbb{P}^{unif}), \quad (14)$$

with $\mathbb{P}^{unif} := (\frac{1}{K}, \dots, \frac{1}{K})$ being the probability vector of frequencies of the uniform distribution on $\{1, \dots, K\}$.

Remark 6: (a) Since $\mathbf{1}$ can be seen as a reference vector with (normalized) equal components, $\inf_{\mathbf{Q} \in \Omega} D_\varphi(\mathbf{Q}, \mathbf{1})$ in (14) can be interpreted as an “index/degree of (in)equality of the set Ω ”, respectively as an “index/degree of diversity of the set Ω ”.

(b) The quantity $\sum_{k=1}^K \varphi(q_k)$ in (14) can be interpreted as (non-probability extension of a) φ -entropy in the sense of Burbea & Rao [31] (see also [32]–[39]); for applications to scalar quantization for lossy coding of information sources see e.g. [96].

Returning to the original distance-minimizing Problem 4, after the first step (12) and (13), we proceed as follows:

(BS2) Step 2: construct an appropriate sequence $(\xi_n)_{n \in \mathbb{N}}$ of \mathbb{R}^K -valued random variables (cf. (2) in Definition 1):

The following condition transposes the minimization problem (13) (and thus the equivalent problem (8)) into a *BS minimizable/amenable problem* in the sense of Definition 1. The connection of this condition with (6) will be discussed in Proposition 27 and its surroundings, see Section XI below.

Condition 7: With $M_{\mathbf{P}} = \sum_{i=1}^K p_i > 0$, the divergence generator φ in (8) (cf. also (12)) satisfies $\tilde{\varphi} := M_{\mathbf{P}} \cdot \varphi \in \Upsilon(]a, b[)$, i.e. $\tilde{\varphi} \in \tilde{\Upsilon}(]a, b[)$ (which is equivalent to $\varphi \in \tilde{\Upsilon}(]a, b[)$) and there holds the representation

$$\tilde{\varphi}(t) = \sup_{z \in \mathbb{R}} \left(z \cdot t - \log \int_{\mathbb{R}} e^{zy} d\tilde{\zeta}(y) \right), \quad t \in \mathbb{R}, \quad (15)$$

for some probability measure $\tilde{\zeta}$ on the real line \mathbb{R} such that the function $z \mapsto MGF_{\tilde{\zeta}}(z) := \int_{\mathbb{R}} e^{zy} d\tilde{\zeta}(y)$ is finite on some open interval containing zero⁹. Notice that $\tilde{\zeta}$ may depend on $M_{\mathbf{P}}$ in a highly non-trivial way (see e.g. Section XII below).

⁸recall that an alternative naming also used in literature is to call Ω a model class (rather than model), and each $\mathbf{Q} \in \Omega$ a model (rather than model element)

⁹in particular, this implies that $\int_{\mathbb{R}} y d\tilde{\zeta}(y) = 1$ and that $\tilde{\zeta}$ has light tails.

Next, we explain the above-mentioned Step 2 in detail: for any $n \in \mathbb{N}$ and any $k \in \{1, \dots, K\}$, let $n_k := \lfloor n \cdot \tilde{p}_k \rfloor$ where $\lfloor x \rfloor$ denotes the integer part of x . We assume $\mathbf{P} \in \mathbb{R}_{>0}^K$, and since thus none of the \tilde{p}_k 's is zero, one has

$$\lim_{n \rightarrow \infty} \frac{n_k}{n} = \tilde{p}_k. \quad (16)$$

Moreover, we assume that $n \in \mathbb{N}$ is large enough, namely $n \geq \max_{k \in \{1, \dots, K\}} \frac{1}{\tilde{p}_k}$, and decompose the set $\{1, \dots, n\}$ of all integers from 1 to n into the following disjoint blocks: $I_1^{(n)} := \{1, \dots, n_1\}$, $I_2^{(n)} := \{n_1 + 1, \dots, n_1 + n_2\}$, and so on until the last block $I_K^{(n)} := \{\sum_{k=1}^{K-1} n_k + 1, \dots, n\}$ which therefore contains all integers from $n_1 + \dots + n_{K-1} + 1$ to n . Clearly, $I_k^{(n)}$ has $n_k \geq 1$ elements (i.e. $\text{card}(I_k^{(n)}) = n_k$ where $\text{card}(A)$ denotes the number of elements in a set A) for all $k \in \{1, \dots, K-1\}$, and the last block $I_K^{(n)}$ has $n - \sum_{k=1}^{K-1} n_k \geq 1$ elements which anyhow satisfies $\lim_{n \rightarrow \infty} \text{card}(I_K^{(n)})/n = \tilde{p}_K$ ¹⁰. Furthermore, consider a vector $\widetilde{\mathbf{W}} := (\widetilde{W}_1, \dots, \widetilde{W}_n)$ where the \widetilde{W}_i 's are i.i.d. copies of the random variable \widetilde{W} whose distribution is associated with the divergence-generator $\tilde{\varphi} := M_{\mathbf{P}} \cdot \varphi$ through (15), in the sense that $\mathbb{P}[\widetilde{W} \in \cdot] = \tilde{\zeta}[\cdot]$. We group the \widetilde{W}_i 's according to the above-mentioned blocks and sum them up blockwise, in order to build the following K -component random vector

$$\boldsymbol{\xi}_n^{\widetilde{\mathbf{W}}} := \left(\frac{1}{n} \sum_{i \in I_1^{(n)}} \widetilde{W}_i, \dots, \frac{1}{n} \sum_{i \in I_K^{(n)}} \widetilde{W}_i \right); \quad (17)$$

notice that the signs of its components may be negative, depending on the nature of the \widetilde{W}_i 's; moreover, the expectation of its k -th component converges to \tilde{p}_k as n tends to infinity (since the expectation of \widetilde{W}_1 is 1), whereas the n -fold of the corresponding variance converges to \tilde{p}_k times the variance of \widetilde{W}_1 .

For such a context, we obtain the following assertion on BS-minimizability (which will be proved in in Appendix A):

Theorem 8: Let $\mathbf{P} \in \mathbb{R}_{>0}^K$, $M_{\mathbf{P}} := \sum_{i=1}^K p_i > 0$, and suppose that the divergence generator φ satisfies Condition 7 above, with $\tilde{\zeta}$ (cf. (15)). Additionally, let $\widetilde{W} := (\widetilde{W}_i)_{i \in \mathbb{N}}$ be a sequence of random variables where the \widetilde{W}_i 's are i.i.d. copies of the random variable \widetilde{W} whose distribution is $\mathbb{P}[\widetilde{W} \in \cdot] = \tilde{\zeta}[\cdot]$ ¹¹. Then, in terms of the random vectors $\boldsymbol{\xi}_n^{\widetilde{\mathbf{W}}}$ (cf. (17)) there holds

$$- \lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P} \left[\boldsymbol{\xi}_n^{\widetilde{\mathbf{W}}} \in \Omega / M_{\mathbf{P}} \right] = \inf_{\mathbf{Q} \in \Omega} D_{\varphi}(\mathbf{Q}, \mathbf{P}) \quad (18)$$

for any $\Omega \subset \mathbb{R}^K$ with regularity properties (7) and finiteness property (9). In particular, for each $\mathbf{P} \in \mathbb{R}_{>0}^K$ the function $\Phi_{\mathbf{P}}(\cdot) := D_{\varphi}(\cdot, \mathbf{P})$ (cf. (4)) is bare-simulation minimizable (BS-minimizable, cf. (2)) on any such $\Omega \subset \mathbb{R}^K$.

Remark 9: (i) For some contexts, one can *explicitly* give the distribution of each of the independent (non-deterministic parts of the) components $\left(\sum_{i \in I_k^{(n)}} \widetilde{W}_i \right)_{k=1, \dots, K}$ of the vector $\boldsymbol{\xi}_n^{\widetilde{\mathbf{W}}}$; this will ease the corresponding concrete simulations. For instance, we shall give those in the solved-cases Section XII below.

(ii) Let us emphasize that we have assumed $\mathbf{P} \in \mathbb{R}_{>0}^K$ in Theorem 8 which excludes \mathbf{P} from having zero components. However, in cases where $\lim_{x \rightarrow \infty} \left| \frac{\varphi(x \cdot \text{sgn}(q))}{x \cdot \text{sgn}(q)} \right| = +\infty$ for $q \neq 0$, then if $p_{k_0} = 0$ for some k_0 it follows that $q_{k_0} = 0$, which proves that $\mathbf{P} \in \mathbb{R}_{>0}^K$ imposes no restriction in Theorem 8, since the projection of \mathbf{P} in Ω then belongs to the subspace of \mathbb{R}^K generated by the non-null components of \mathbf{P} ; such a situation appears e.g. for power divergence generators φ_{γ} with $\gamma > 2$. So there is no loss of generality assuming $\mathbf{P} \in \mathbb{R}_{>0}^K$ in this case.

As examples for the applicability of Theorem 8, one can e.g. combine *each* of the divergence generators φ of Section XII (except for the one in Subsection XII-I) with *any* of the optimization problems (8),(10),(11),(14) as well as the below-mentioned (21),(22); the needed distributions $\mathbb{P}[\widetilde{W} \in \cdot] = \tilde{\zeta}[\cdot]$ correspond to the entry of the corresponding Subsection XII- with the choice $\tilde{c} \cdot M_{\mathbf{P}}$ instead of \tilde{c} . By taking $\zeta := -\varphi$ instead, one can solve the below-mentioned problems (23) and (24).

Returning to the general context, the limit statement (18) provides the principle for the approximation of the solution of Problem (8). Indeed, by replacing the left-hand side in (18) by its finite counterpart, we deduce for given large n

$$- \frac{1}{n} \log \mathbb{P} \left[\boldsymbol{\xi}_n^{\widetilde{\mathbf{W}}} \in \Omega / M_{\mathbf{P}} \right] \approx \inf_{\mathbf{Q} \in \Omega} D_{\varphi}(\mathbf{Q}, \mathbf{P}); \quad (19)$$

it remains to estimate the left-hand side of (19). The latter can be performed either by a *naive estimator* of the frequency of those replications of $\boldsymbol{\xi}_n^{\widetilde{\mathbf{W}}}$ which hit $\Omega / M_{\mathbf{P}}$, or more efficiently by some improved estimator; for details, see Section X below.

¹⁰if all \tilde{p}_k ($k = 1, \dots, K$) are rational numbers in $]0, 1[$ with $\sum_{k=1}^K \tilde{p}_k = 1$ and N is the (always existing) smallest integer such that all $N \cdot \tilde{p}_k$ ($k = 1, \dots, K$) are integers (i.e. $\in \mathbb{N}$), then for any multiple $n = \ell \cdot N$ ($\ell \in \mathbb{N}$) one gets that all $n_k = n \cdot \tilde{p}_k$ are integers and that $\text{card}(I_K^{(n)}) = n_K$.

¹¹and thus, $E_{\mathbb{P}}[\widetilde{W}_i] = 1$

Remark 10: According to (18) of Theorem 8 as well as (19), we can principally tackle the (approximative) computation of the minimum value $\inf_{\mathbf{Q} \in \Omega} D_{\varphi}(\mathbf{Q}, \mathbf{P}) = \inf_{\mathbf{Q} \in \Omega} \sum_{k=1}^K p_k \cdot \varphi\left(\frac{q_k}{p_k}\right)$ and in particular of $\inf_{\mathbf{Q} \in \Omega} \sum_{k=1}^K \varphi(q_k) = \inf_{\mathbf{Q} \in \Omega} D_{\varphi}(\mathbf{Q}, \mathbf{1})$ (cf. (14)) by basically *only employing a fast and accurate — pseudo, true, natural, quantum — random number generator*¹², provided that the constraint set Ω satisfies the mild assumptions (7) and (9). Notice that (7) also covers (e.g. high-dimensional) constraint sets Ω which are *non-convex* and even *highly disconnected*, and for which other minimization methods (e.g. pure enumeration, gradient or steepest descent methods, etc.¹³) may be problematic or intractable. For instance, (7) covers kind of “ K -dimensional (not necessarily regular) polka dot (leopard skin) pattern type” relaxations $\Omega := \bigcup_{i=1}^N \mathcal{U}_i(\mathbf{Q}_i^{dis})$ of finite discrete constraint sets $\Omega^{dis} := \{\mathbf{Q}_1^{dis}, \dots, \mathbf{Q}_N^{dis}\}$ of high cardinality N (e.g. being exponential or factorial in a large K), where each K -dimensional vector \mathbf{Q}_i^{dis} (e.g. having pure integer components only) is surrounded by some small (in particular, non-overlapping/disjoint) neighborhood $\mathcal{U}_i(\mathbf{Q}_i^{dis})$; in such a context, e.g. $\inf_{\mathbf{Q} \in \Omega} \sum_{k=1}^K \varphi(q_k)$ can be regarded as a “*BS-tractable*” relaxation of the nonlinear discrete (e.g. integer, combinatorial¹⁴) optimization program $\inf_{\mathbf{Q} \in \Omega^{dis}} \sum_{k=1}^K \varphi(q_k)$.

Returning to the general context, notice that Theorem 8 does not cover cases where Ω consists of \mathbf{Q} satisfying the additional constraint $\sum_{i=1}^K q_i = A$ for some fixed $A > 0$ (and thus $\text{int}(\Omega) = \emptyset$ violating (7)). However, such situations can be still handled with an adaption of the above-described BS method, see Remark 13(v),(vi), Lemma 14 and Section XII below.

A. Generalizations

Recalling (14), let us point out that with our new BS approach one may even tackle more general optimization problems of the form $\inf_{\check{\mathbf{Q}} \in \check{\Omega}} \sum_{k=1}^K \check{\varphi}(\check{q}_k)$ where basically $\check{\varphi}$ is some function which is finite and convex in a non-empty neighborhood (say, $]t_0 + a - 1, t_0 + b - 1[$ with $a < 1 < b$) of some point $t_0 \in \mathbb{R}$ as well as differentiable and strictly convex in a non-empty sub-neighborhood of t_0 ; for this, the function $\varphi(t) := \check{\varphi}(t + t_0 - 1) - \check{\varphi}'(t_0) \cdot ((t + t_0 - 1) - t_0) - \check{\varphi}(t_0)$, $t \in]a, b[$, (which corresponds to argument-shifting and adding an affine-linear function) should be such that $K \cdot \varphi \in \Upsilon(]a, b[)$, and from the corresponding minimization problem

$$\begin{aligned} \inf_{\check{\mathbf{Q}} \in \check{\Omega}/K} D_{K \cdot \varphi}(\check{\mathbf{Q}}, \mathbb{P}^{unif}) &= \inf_{\mathbf{Q} \in \Omega} \sum_{k=1}^K \varphi(q_k) = \inf_{\mathbf{Q} \in \Omega} \sum_{k=1}^K \left(\check{\varphi}(q_k + t_0 - 1) - \check{\varphi}'(t_0) \cdot ((q_k + t_0 - 1) - t_0) - \check{\varphi}(t_0) \right) \\ &= \inf_{\check{\mathbf{Q}} \in \check{\Omega} + t_0 - 1} \sum_{k=1}^K \left(\check{\varphi}(\check{q}_k) - \check{\varphi}'(t_0) \cdot (\check{q}_k - t_0) - \check{\varphi}(t_0) \right) \\ &= K \cdot \left(t_0 \cdot \check{\varphi}'(t_0) - \check{\varphi}(t_0) \right) + \inf_{\check{\mathbf{Q}} \in \check{\Omega}} \left(\sum_{k=1}^K \check{\varphi}(\check{q}_k) - \check{\varphi}'(t_0) \cdot \sum_{k=1}^K \check{q}_k \right), \quad \text{with } \check{\Omega} := \Omega + t_0 - 1, \end{aligned} \quad (20)$$

the term $\inf_{\check{\mathbf{Q}} \in \check{\Omega}} \sum_{k=1}^K \check{\varphi}(\check{q}_k)$ should be recoverable; for instance, later on we shall employ constraints sets $\check{\Omega}$ which particularly include $\sum_{k=1}^K \check{q}_k = A > 0$, whereas another possibility would be to use a $\check{\varphi}$ which satisfies $\check{\varphi}'(t_0) = 0$. As a different line of flexibilization of (14), we can also deal with the problem $\inf_{\mathbf{Q} \in \Omega} h\left(\sum_{k=1}^K \varphi(q_k)\right)$ through

$$\inf_{\mathbf{Q} \in \Omega} h\left(\sum_{k=1}^K \varphi(q_k)\right) = h\left(\inf_{\check{\mathbf{Q}} \in \check{\Omega}/K} D_{K \cdot \varphi}(\check{\mathbf{Q}}, \mathbb{P}^{unif})\right) \quad (21)$$

for any φ with $K \cdot \varphi \in \Upsilon(]a, b[)$ and any continuous strictly increasing function $h : \mathcal{H} \mapsto \mathbb{R}$ with $\mathcal{H} := [0, \infty[$ (or a sufficiently large subset thereof), and with the problem $\sup_{\mathbf{Q} \in \Omega} h\left(\sum_{k=1}^K \varphi(q_k)\right)$ through

$$\sup_{\mathbf{Q} \in \Omega} h\left(\sum_{k=1}^K \varphi(q_k)\right) = h\left(\inf_{\check{\mathbf{Q}} \in \check{\Omega}/K} D_{K \cdot \varphi}(\check{\mathbf{Q}}, \mathbb{P}^{unif})\right) \quad (22)$$

for any φ with $K \cdot \varphi \in \Upsilon(]a, b[)$ and any continuous strictly decreasing function $h : \mathcal{H} \mapsto \mathbb{R}$. As a continuation of Remark 6(b), the quantity $h\left(\sum_{k=1}^K \varphi(q_k)\right)$ in (21) can be seen as (non-probability extension of a) (h, φ) -entropy in the sense of Salicru et al. [40] (see also e.g. [12],[117] as well as [118],[119] for exemplary applications); important special cases will be discussed in more detail, below. Combining (20) with (21) (respectively, with (22)) leads to a further flexibilization. Of course, we can also apply our BS method to the maximization $\sup_{\mathbf{Q} \in \Omega} h\left(\sum_{k=1}^K \zeta(q_k)\right)$ for any concave function ζ with $-K \cdot \zeta \in \Upsilon(]a, b[)$ and any continuous strictly increasing function $h : \mathcal{H} \mapsto \mathbb{R}$ with $\mathcal{H} := -[\infty, 0]$ (or a sufficiently large subset thereof), via

$$\sup_{\mathbf{Q} \in \Omega} h\left(\sum_{k=1}^K \zeta(q_k)\right) = h\left(-\inf_{\check{\mathbf{Q}} \in \check{\Omega}/K} D_{-K \cdot \zeta}(\check{\mathbf{Q}}, \mathbb{P}^{unif})\right), \quad (23)$$

¹²see e.g. [97]–[110]

¹³a detailed discussion and comparisons are beyond the scope of this paper, given its current length

¹⁴see e.g. [111]–[116] for comprehensive books on discrete, integer and combinatorial programming and their vast applications

and to $\inf_{\mathbf{Q} \in \Omega} h\left(\sum_{k=1}^K \zeta(q_k)\right)$ for any concave ζ with $-K \cdot \zeta \in \Upsilon(]a, b[)$ and any continuous strictly decreasing $h : \mathcal{H} \mapsto \mathbb{R}$, via

$$\inf_{\mathbf{Q} \in \Omega} h\left(\sum_{k=1}^K \zeta(q_k)\right) = h\left(-\inf_{\tilde{\mathbf{Q}} \in \Omega/K} D_{-K \cdot \zeta}(\tilde{\mathbf{Q}}, \mathbb{P}^{unif})\right). \quad (24)$$

Moreover, we can tackle $\sup_{\tilde{\mathbf{Q}} \in \tilde{\Omega}} \sum_{k=1}^K \check{\zeta}(\check{q}_k)$ where $\check{\zeta}$ is some function which is finite and concave in a non-empty neighborhood $]t_0 + a - 1, t_0 + b - 1[$ (with $a < 1 < b$) of some point $t_0 \in \mathbb{R}$ as well as differentiable and strictly concave in a non-empty sub-neighborhood of t_0 ; for this, the function

$$-\zeta(t) := -\check{\zeta}(t + t_0 - 1) + \check{\zeta}'(t_0) \cdot \left((t + t_0 - 1) - t_0\right) + \check{\zeta}(t_0), \quad t \in]a, b[,$$

should be such that $-K \cdot \zeta \in \Upsilon(]a, b[)$, and from the corresponding minimization problem

$$\begin{aligned} -\inf_{\tilde{\mathbf{Q}} \in \Omega/K} D_{-K \cdot \zeta}(\tilde{\mathbf{Q}}, \mathbb{P}^{unif}) &= \sup_{\mathbf{Q} \in \Omega} \sum_{k=1}^K \zeta(q_k) = \sup_{\mathbf{Q} \in \Omega} \sum_{k=1}^K \left(\check{\zeta}(q_k + t_0 - 1) - \check{\zeta}'(t_0) \cdot ((q_k + t_0 - 1) - t_0) - \check{\zeta}(t_0)\right) \\ &= \sup_{\tilde{\mathbf{Q}} \in \Omega + t_0 - 1} \sum_{k=1}^K \left(\check{\zeta}(\check{q}_k) - \check{\zeta}'(t_0) \cdot (\check{q}_k - t_0) - \check{\zeta}(t_0)\right) \\ &= K \cdot \left(t_0 \cdot \check{\zeta}'(t_0) - \check{\zeta}(t_0)\right) + \sup_{\tilde{\mathbf{Q}} \in \tilde{\Omega}} \left(\sum_{k=1}^K \check{\zeta}(\check{q}_k) - \check{\zeta}'(t_0) \cdot \sum_{k=1}^K \check{q}_k\right), \quad \text{with } \tilde{\Omega} := \Omega + t_0 - 1, \end{aligned} \quad (25)$$

the term $\sup_{\tilde{\mathbf{Q}} \in \tilde{\Omega}} \sum_{k=1}^K \check{\zeta}(\check{q}_k)$ should be recoverable; the left-hand side of (25) corresponds to the special case $h(x) := x$ of the BS-minimizable (23). A combination of (25) with (23) (respectively, with (24)) leads to a further flexibilization.

VI. STOCHASTIC MINIMUM DISTANCE/RISK ESTIMATION AND DETERMINISTIC SIMPLEX CASES

A. General stochastic construction

In contrast to the previous Section V, we now work out our BS method for the important setup where basically P is a *random* (unknown) element of the simplex \mathbb{S}^K of K -component probability (frequency) vectors and $\Omega \subset \mathbb{S}^K$ (which violates (7) since $\text{int}(\Omega) = \emptyset$, cf. Remark 5(d), and thus requires a different treatment). To begin with, in the statistics of discrete data — and in the adjacent research fields of information theory, artificial intelligence and machine learning — one often encounters the following *minimum distance estimation (MDE) problem* which is often also named as *estimation of the empirical risk*:

(MDE1) for index $i \in \mathbb{N}$, let the generation of the i -th (uncertainty-prone) data point be represented by the random variable X_i which takes values in the discrete set $\mathcal{Y} := \{d_1, \dots, d_K\}$ of K distinct values “of any kind”. It is assumed that there exists a probability measure $\mathbb{P}[\cdot]$ on \mathcal{Y} which is the a.s. limit (as n tends to infinity) of the empirical measures \mathbb{P}_n^{emp} defined by the collection (X_1, \dots, X_n) , in formula

$$\lim_{n \rightarrow \infty} \mathbb{P}_n^{emp} := \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \delta_{X_i} = \mathbb{P} \quad \text{a.s.} \quad (26)$$

where δ_y denotes the one-point distribution (Dirac mass) at point y ¹⁵. We will assume that none of the entries of \mathbb{P} bears zero mass so that \mathbb{P} is identified with a point in the interior of \mathbb{S}^K (see below). The underlying probability space (say, $(\mathcal{X}, \mathcal{A}, \mathbb{P})$) where the above a.s. convergence holds, pertains to the random generation of the sequence $(X_i)_{i \in \mathbb{N}}$, of which we do not need to know but for (26). Examples include the i.i.d. case (where the X_i 's are independent and have common distribution \mathbb{P}), ergodic Markov chains on \mathcal{Y} with stationary distribution \mathbb{P} , more globally autoregressive chains with stationary measure \mathbb{P} , etc.

Let us briefly discuss our assumption (26) (resp. its vector form (30) below) on the limit behavior of the empirical distribution of the observed sample $\mathbf{X}_n := (X_1, \dots, X_n)$ as n tends to infinity. In the “basic” statistical context, the sample \mathbf{X}_n consists of i.i.d. replications of a generic random variable X with probability distribution \mathbb{P} . However, our approach captures many other sampling schemes, where the distribution \mathbb{P} is defined implicitly through (26) for which we aim at some estimate of $\Phi_{\mathbb{P}}(\Omega)$ of a family Ω of probability distributions on \mathcal{Y} . Sometimes the sequence of samples $(\mathbf{X}_n)_{n \in \mathbb{N}}$ may stem from a triangular array so that $\mathbf{X}_n = (X_{1,n}, \dots, X_{k_n,n})$ with $k_n \rightarrow \infty$ and (26) is substituted by

$$\lim_{n \rightarrow \infty} \frac{1}{k_n} \sum_{i=1}^{k_n} \delta_{X_{i,n}} = \mathbb{P} \quad \text{a.s.}$$

which does not alter the results of this paper by any means.

¹⁵notice that \mathbb{P}_n^{emp} a probability measure (on the data space \mathcal{Y}), which is random due to its dependence on the X_i 's

(MDE2) given a *model* Ω , i.e. a family Ω of probability distributions \mathbb{Q} on \mathcal{Y} each of which serves as a potential description of the underlying (unknown) data-generating mechanism \mathbb{P} , one would like to find

$$\Phi_{\mathbb{P}}(\Omega) := \inf_{\mathbb{Q} \in \Omega} D_{\varphi}(\mathbb{Q}, \mathbb{P}) \quad (27)$$

which quantifies the *adequacy* of the model Ω for modeling \mathbb{P} , *via* the minimal distance/dissimilarity of Ω to \mathbb{P} ; a lower $\Phi_{\mathbb{P}}$ -value means a better adequacy (in the sense of a lower departure between the model and the truth, cf. [59]–[62]). Hence, especially in the context of *model selection* within complex big-data contexts, for the *search of appropriate models* Ω and model elements/members therein, the (fast and efficient) computation of $\Phi_{\mathbb{P}}(\Omega)$ constitutes a decisive first step, since if the latter is “too large” (respectively “much larger than” $\Phi_{\mathbb{P}}(\bar{\Omega})$ for some competing model $\bar{\Omega}$), then the model Ω is “not adequate enough” (respectively “much less adequate than” $\bar{\Omega}$); in such a situation, the effort of computing the (not necessarily unique) best model element/member $\arg \inf_{\mathbb{Q} \in \Omega} D_{\varphi}(\mathbb{Q}, \mathbb{P})$ within the model Ω is “not very useful” and is thus a “waste of computational time”. Because of such considerations, we concentrate on finding the infimum (27) rather than finding the corresponding minimizer(s). Variants of (27) are of interest, too.

Since $\text{int}(\Omega)$ is required to be a non-empty set (in the relative topology) in the space of probability distributions on \mathcal{Y} , the present procedure is fitted for semi-parametric models Ω , e.g. defined through moment conditions (as extensions of the Empirical Likelihood paradigm, see e.g. [120]), or through L-moment conditions (i.e. moment conditions pertaining to quantile measures, see [72]), or even more involved non-parametric models where the geometry of Ω does not allow for ad-hoc procedures. In such setups, there is typically no closed form of the divergence with respect to any probability distribution available.

The measurement or the estimation of $\Phi_{\mathbb{P}}(\Omega)$ is a tool for the choice of pertinent putative models Ω among a class of specifications. The case when $\Phi_{\mathbb{P}}(\Omega) > 0$ is interesting in its own, since it is quite common in engineering modelling to argue in favor of misspecified models (or (non-void) neighborhoods of such models for sake of robustness issues), due to quest for conservatism; the choice between them is a widely open field e.g. in the practice of reliability. This also opens the question of the choice of the divergence generator φ ; although this will not be discussed in this paper, as a motivating running example the reader may keep in mind the generator $\varphi_2(x) := (x - 1)^2/2$ which induces the divergence $D_{\varphi_2}(\mathbb{Q}, \mathbb{P})$ (see (41) below for details) which quantifies the expected square relative error when substituting the true distribution \mathbb{P} by the model \mathbb{Q} .

An estimate of $\Phi_{\mathbb{P}}(\Omega)$ can be used as a statistics for some test of fit, and indeed the likelihood ratio test adapted to some semi-parametric models has been generalized to the divergence setting (see [120]). The statement of the limit distributions of our estimate, under the model and under misspecification, is postponed to future work.

In the following, we compute/approximate (27) — and some variants thereof — by our *bare simulation (BS)* method, by “mimicking” the deterministic minimization problem (8) respectively (13). Let us first remark that, as usual, each probability distribution (probability measure) \mathbb{P} on $\mathcal{Y} = \{d_1, \dots, d_K\}$ can be uniquely identified with the (row) vector $\mathbb{P} := (p_1, \dots, p_K) \in \mathbb{S}^K$ of the corresponding probability masses (frequencies) $p_k = \mathbb{P}[\{d_k\}]$ via $\mathbb{P}[A] = \sum_{k=1}^K p_k \cdot 1_A(d_k)$ for each $A \subset \mathcal{Y}$, where $1_A(\cdot)$ denotes the indicator function on the set A . In particular, the probability distribution \mathbb{P} in (MDE1) can be identified with (p_1, \dots, p_K) in terms of $p_k = \mathbb{P}[\{d_k\}]$ (which in the i.i.d. case turns into $p_k = \mathbb{P}[X_1 = d_k]$). Along this line, the family Ω of probability distributions in (MDE2) can be identified with a subset $\Omega \subset \mathbb{S}^K$ of probability vectors (viz. of vectors of probability masses). Analogously, each finite nonnegative measure Q on \mathcal{Y} can be uniquely identified with a vector $\mathbf{Q} := (q_1, \dots, q_K) \in \mathbb{R}_{\geq 0}^K$, and each finite signed measure Q with a vector $\mathbf{Q} := (q_1, \dots, q_K) \in \mathbb{R}^K$. The corresponding divergences between distributions/measures are then, as usual, defined through the divergences between their respective masses/frequencies:

$$D_{\varphi}(Q, \mathbb{P}) := D_{\varphi}(\mathbf{Q}, \mathbb{P}). \quad (28)$$

In particular, \mathbb{P}_n^{emp} can be identified with the vector $\mathbb{P}_n^{emp} := (p_{n,1}^{emp}, \dots, p_{n,K}^{emp})$ where

$$p_{n,k}^{emp} := \frac{1}{n} \cdot n_k := \frac{1}{n} \cdot \text{card}(\{i \in \{1, \dots, n\} : X_i = d_k\}) =: \frac{1}{n} \cdot \text{card}(I_k^{(n)}), \quad k \in \{1, \dots, K\}, \quad (29)$$

and accordingly the required limit behaviour (26) is equivalent to the vector-convergence

$$\lim_{n \rightarrow \infty} \left(\frac{n_1}{n}, \dots, \frac{n_K}{n} \right) = (p_1, \dots, p_K) \quad \text{a.s.} \quad (30)$$

Notice that, in contrast to the above Section V, the sets $I_k^{(n)}$ of indexes introduced in (29) and their numbers $n_k = \text{card}(I_k^{(n)})$ of elements are now *random* (due to their dependence on the X_i 's) and $M_{\mathbb{P}_n^{emp}} = 1$. In a *batch procedure*, when $D_{\varphi}(\Omega, \mathbb{P}_n^{emp}) := \inf_{\mathbb{Q} \in \Omega} D_{\varphi}(\mathbf{Q}, \mathbb{P}_n^{emp})$ is estimated once the sample (X_1, \dots, X_n) is observed, we may reorder this sample by putting the n_1 sample points X_i which are equal to d_1 in the first places, and so on; accordingly one ends up with index sets $I_k^{(n)}$ as defined in Section V. When the *online acquisition* of the data X_i 's is required, then we usually do not reorder the sample, and the $I_k^{(n)}$'s do not consist in consecutive indexes, which does not make any change with respect to the resulting construction nor to the estimator.

The above considerations open the gate to our desired “mimicking” of (8) and (13) to achieve (27) (and some variants thereof) by our bare simulation (BS) method. To proceed, we employ a family of random variables $(W_i)_{i \in \mathbb{N}}$ of independent and identically distributed \mathbb{R} -valued random variables with probability distribution $\zeta[\cdot] := \mathbb{P}[W_1 \in \cdot]$ — being connected with the divergence generator $\varphi \in \Upsilon(a, b]$ via the representability (6) — such that $(W_i)_{i \in \mathbb{N}}$ is independent of $(X_i)_{i \in \mathbb{N}}$ ¹⁶.

As a next step, notice that the “natural candidate”

$$\xi_{n, \mathbf{X}}^{\mathbf{W}} := \frac{1}{n} \cdot \sum_{k=1}^K \left(\sum_{i \in I_k^{(n)}} W_i \right) \cdot \delta_{d_k} = \frac{1}{n} \sum_{i=1}^n W_i \cdot \delta_{X_i}$$

is not a probability measure since its total mass is not 1 in general, since in terms of its equivalent vector version

$$\xi_{n, \mathbf{X}}^{\mathbf{W}} := \left(\frac{1}{n} \sum_{i \in I_1^{(n)}} W_i, \dots, \frac{1}{n} \sum_{i \in I_K^{(n)}} W_i \right) \quad (31)$$

the sum $\sum_{k=1}^K \frac{1}{n} \sum_{i \in I_k^{(n)}} W_i = \frac{1}{n} \sum_{j=1}^n W_j$ of the K vector components of (31) is typically not equal to 1; this implies that no limit result of the form (18) with finite limit can hold, since $\xi_{n, \mathbf{X}}^{\mathbf{W}}$ takes values in \mathbb{R}^K and Ω is a subset in the probability simplex \mathbb{S}^K which has *void* interior in \mathbb{R}^K causing a violation of condition (7) (cf. Remark 5(d)); moreover, depending on the concrete form of the generator φ , the corresponding weights may take *negative values*. Therefore, we need some “rescaling”. Indeed, let us introduce the *normalized weighted empirical measure*

$$\xi_{n, \mathbf{X}}^{w\mathbf{W}} := \begin{cases} \frac{1}{\sum_{k=1}^K \sum_{i \in I_k^{(n)}} W_i} \cdot \sum_{k=1}^K \left(\sum_{i \in I_k^{(n)}} W_i \right) \cdot \delta_{d_k} = \sum_{i=1}^n \frac{W_i}{\sum_{j=1}^n W_j} \cdot \delta_{X_i}, & \text{if } \sum_{j=1}^n W_j \neq 0, \\ \infty \cdot \sum_{k=1}^K \delta_{d_k} =: \infty, & \text{if } \sum_{j=1}^n W_j = 0, \end{cases} \quad (32)$$

which will substitute $\xi_{n, \mathbf{X}}^{\mathbf{W}}$ and which may belong to Ω with positive probability. The equivalent vector version of $\xi_{n, \mathbf{X}}^{w\mathbf{W}}$ is

$$\xi_{n, \mathbf{X}}^{w\mathbf{W}} := \begin{cases} \left(\frac{\sum_{i \in I_1^{(n)}} W_i}{\sum_{k=1}^K \sum_{i \in I_k^{(n)}} W_i}, \dots, \frac{\sum_{i \in I_K^{(n)}} W_i}{\sum_{k=1}^K \sum_{i \in I_k^{(n)}} W_i} \right), & \text{if } \sum_{j=1}^n W_j \neq 0, \\ (\infty, \dots, \infty) =: \infty, & \text{if } \sum_{j=1}^n W_j = 0, \end{cases} \quad (33)$$

a point in the linear subset of \mathbb{R}^K spanned by \mathbb{S}^K at infinity.

Remark 11: (i) (Concerning e.g. computer-program command availability) In case of $\sum_{j=1}^n W_j = 0$, in (32) we may equivalently assign to $\xi_{n, \mathbf{X}}^{w\mathbf{W}}$ instead of ∞ any measure (e.g. probability distribution) which does not belong to Ω , respectively, in (33) we may equivalently choose for $\xi_{n, \mathbf{X}}^{w\mathbf{W}}$ any vector outside of Ω instead of ∞ .

(ii) By construction, in case of $\sum_{j=1}^n W_j \neq 0$, the sum of the random K vector components of (33) is now automatically equal to 1, but — as (depending on φ) the W_i ’s may take both positive and negative values¹⁷ — these random components may be negative with probability strictly greater than zero (respectively nonnegative with probability strictly less than 1); in the framework of (32) this means that $\xi_{n, \mathbf{X}}^{w\mathbf{W}}$ is in general a random *signed* measure with total mass 1, in case of $\sum_{j=1}^n W_j \neq 0$. However, $\mathbb{P}[\xi_{n, \mathbf{X}}^{w\mathbf{W}} \in \mathbb{S}_{>0}^K]$ converges to 1 as n tends to infinity, since all the (identically distributed) random variables W_i have expectation 1 (as a consequence of the assumed representability (6)); in case of $\mathbb{P}[W_1 > 0] = 1$ one has even $\mathbb{P}[\xi_{n, \mathbf{X}}^{w\mathbf{W}} \in \mathbb{S}_{>0}^K] = 1$ for all $n \in \mathbb{N}$.

(iii) By generalizing the terminology of e.g. [121], through the right-hand side of (32) one can interpret (for $\sum_{j=1}^n W_j \neq 0$) the normalized weighted empirical measure $\xi_{n, \mathbf{X}}^{w\mathbf{W}}$ as response of an output neuron in a random perceptron consisting of random inputs \mathbf{X} , a layer with n units having one-point-distribution-valued responses $\delta_{X_1}, \dots, \delta_{X_n}$, and independent random synaptic weights $\left(\frac{W_1}{\sum_{j=1}^n W_j}, \dots, \frac{W_n}{\sum_{j=1}^n W_j} \right)$. By our below-mentioned methods, we can approximate $\mathbb{P}[\xi_{n, \mathbf{X}}^{w\mathbf{W}} \in \Omega]$ for nearly any model Ω , and therefore propose proxies of Bayesian rules associated with hidden layers in neural networks, as e.g. suggested in [121].

With the above-mentioned ingredients, we are now in the position to tackle a variant of the distance minimization problem (27), by our bare simulation method through “mimicking” the deterministic minimization problem (8) respectively (13). For this, we also employ the *conditional* distributions $\mathbb{P}_n[\cdot] := \mathbb{P}_{X_1^n}[\cdot] := \mathbb{P}[\cdot | X_1, \dots, X_n]$ and obtain the following

Theorem 12: Suppose that $(X_i)_{i \in \mathbb{N}}$ is a sequence of random variables with values in $\mathcal{Y} := \{d_1, \dots, d_K\}$ such that (26) holds for some probability measure $\mathbb{P}[\cdot]$ on \mathcal{Y} having no zero-mass frequencies (or equivalently, (30) holds for some probability vector $\mathbb{P} \in \mathbb{S}_{>0}^K$). Moreover, let $(W_i)_{i \in \mathbb{N}}$ be a family of independent and identically distributed \mathbb{R} -valued random variables with

¹⁶on the common underlying probability space $(\mathcal{X}, \mathcal{A}, \mathbb{P})$

¹⁷see e.g. the below-mentioned solved Case 4 of Subsection XII-D

probability distribution $\zeta[\cdot] := \mathbb{P}[W_1 \in \cdot]$ being connected with the divergence generator $\varphi \in \Upsilon([a, b])$ via the representability (6), such that $(W_i)_{i \in \mathbb{N}}$ is independent of $(X_i)_{i \in \mathbb{N}}$. Then there holds

$$- \lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_{X_1^n} [\xi_{n, \mathbf{X}}^{w, \mathbf{W}} \in \Omega] = \inf_{\mathbb{Q} \in \Omega} \inf_{m \neq 0} D_\varphi(m \cdot \mathbb{Q}, \mathbb{P}) \quad (34)$$

$$\begin{aligned} &= \inf_{m \neq 0} \inf_{\mathbb{Q} \in \Omega} D_\varphi(m \cdot \mathbb{Q}, \mathbb{P}) \\ &= \inf_{m \neq 0} \inf_{\mathbb{Q} \in \Omega} D_\varphi(m \cdot \mathbb{Q}, \mathbb{P}) \end{aligned} \quad (35)$$

$$= \inf_{\mathbb{Q} \in \Omega} \inf_{m \neq 0} D_\varphi(m \cdot \mathbb{Q}, \mathbb{P}) = - \lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_{X_1^n} [\xi_{n, \mathbf{X}}^{w, \mathbf{W}} \in \Omega] \quad (36)$$

for all sets Ω of probability distributions such that their equivalent probability-vector form \mathfrak{Q} satisfies the regularity properties (7) *in the relative topology* and the finiteness property (9); notice that for the equality (35) we have used the ‘‘divergence link’’ (28). In particular, for each $\mathbb{P} \in \mathbb{S}_{>0}^K$ (respectively, its equivalent probability-distribution \mathbb{P}) the function $\mathbb{Q} \mapsto \inf_{m \neq 0} D_\varphi(m \cdot \mathbb{Q}, \mathbb{P})$ (respectively, the function $\mathbb{Q} \mapsto \inf_{m \neq 0} D_\varphi(m \cdot \mathbb{Q}, \mathbb{P})$) is BS-minimizable (cf. (2)) on all sets $\Omega \subset \mathbb{S}^K$ satisfying (7) *in the relative topology* and (9) (respectively, on their probability-distribution-equivalent Ω).

The proof of Theorem 12 will be given in Appendix B. Analogous to Remark 9(ii), let us emphasize that we have assumed $\mathbb{P} \in \mathbb{S}_{>0}^K$ in Theorem 12. Henceforth, for sets $\Omega \subset \mathbb{S}^K$ of probability vectors we deal with (7) only in the relative topology; thus, the latter will be unmentioned for the sake of brevity. Remark 5(a),(b),(c),(e) applies accordingly.

Remark 13: (i) In strong contrast to Theorem 8, the above result does not provide a direct tool for the solution of Problem (27) since the limit in (34) bears no *direct* information on the minimum divergence $D_\varphi(\Omega, \mathbb{P}) := \inf_{\mathbb{Q} \in \Omega} D_\varphi(\mathbb{Q}, \mathbb{P})$; the link between the corresponding quantities can be emphasized and exploited e.g. in the case of power type divergences, which leads to explicit minimization procedures as shown in Subsection VI-B below. For general divergences, Theorem 12 allows for the estimation of upper and lower bounds of $D_\varphi(\Omega, \mathbb{P})$, as developed in Subsection VI-C below.

(ii) Notice that $\check{D}_\varphi(\mathbb{Q}, \mathbb{P}) := \inf_{m \neq 0} D_\varphi(m \cdot \mathbb{Q}, \mathbb{P})$ satisfies the axioms of a divergence, that is, $\check{D}_\varphi(\mathbb{Q}, \mathbb{P}) \geq 0$, as well as $\check{D}_\varphi(\mathbb{Q}, \mathbb{P}) = 0$ if and only if $\mathbb{Q} = \mathbb{P}$ (reflexivity). Hence, in Theorem 12 we are still within our framework of bare simulation of a divergence minimum w.r.t. its first component (however, notice the difference to (i)).

(iii) Viewed from a ‘‘reverse’’ angle, Theorem 12 gives a crude approximation for the probability for $\xi_{n, \mathbf{X}}^{w, \mathbf{W}}$ to belong to Ω , conditionally upon $\mathbf{X} = (X_1, \dots, X_n)$.

(iv) In the same spirit as Remark 9(i), for some contexts one can *explicitly* give the distribution of each of the independent components $\left(\sum_{i \in I_k^{(n)}} W_i \right)_{k=1, \dots, K}$ of the vector $\xi_{n, \mathbf{X}}^{w, \mathbf{W}}$ given $\mathbf{X} = \mathbf{x}$; this will ease the corresponding concrete simulations in a batch procedure. For instance, we shall give some of those in the solved-cases Section XII below.

(v) Consider the special ‘‘degenerate’’ case where all the data observations are *certain* and thus $(X_i)_{i \in \mathbb{N}}$ is nothing but a *purely deterministic* sequence, say $(\tilde{x}_i)_{i \in \mathbb{N}}$, of elements \tilde{x}_i from the arbitrary set $\mathcal{Y} := \{d_1, \dots, d_K\}$ of K distinct values ‘‘of any kind’’ (e.g., \mathcal{Y} may consist of K distinct numbers); then the corresponding empirical distribution \mathbb{P}_n^{emp} can be identified with the vector $\mathbb{P}_n^{emp} := (p_{n,1}^{emp}, \dots, p_{n,K}^{emp})$ where

$$p_{n,k}^{emp} := \frac{1}{n} \cdot n_k := \frac{1}{n} \cdot \text{card}(\{i \in \{1, \dots, n\} : \tilde{x}_i = d_k\}) =: \frac{1}{n} \cdot \text{card}(I_k^{(n)}), \quad k \in \{1, \dots, K\},$$

and accordingly the required limit behaviour (26) is equivalent to the vector-convergence

$$\lim_{n \rightarrow \infty} \left(\frac{n_1}{n}, \dots, \frac{n_K}{n} \right) = (p_1, \dots, p_K) \quad \text{for some } p_1 > 0, \dots, p_K > 0 \text{ such that } \sum_{k=1}^K p_k = 1.$$

Correspondingly, with the notations $\mathbb{P} := (p_1, \dots, p_K)$ and $\tilde{\mathbf{x}} := (\tilde{x}_1, \dots, \tilde{x}_n)$, the vector-form part of the assertion (34) of Theorem 12 becomes

$$- \lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P} \left[\xi_{n, \tilde{\mathbf{x}}}^{w, \mathbf{W}} \in \Omega \right] = \inf_{\mathbb{Q} \in \Omega} \inf_{m \neq 0} D_\varphi(m \cdot \mathbb{Q}, \mathbb{P}) = \inf_{m \neq 0} \inf_{\mathbb{Q} \in \Omega} D_\varphi(m \cdot \mathbb{Q}, \mathbb{P})$$

for all subsets $\Omega \subset \mathbb{S}^K$ satisfying the regularity properties (7) and the finiteness property (9); notice that the conditional probability $\mathbb{P}_{X_1^n}[\cdot]$ has degenerated to the ordinary probability $\mathbb{P}[\cdot]$.

(vi) In a similar fashion to the proof of (the special degenerate case (v) of) Theorem 12, one can show

$$- \lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P} \left[\xi_n^{w, \mathbf{W}} \in \Omega \right] = \inf_{\mathbb{Q} \in \Omega} \inf_{m \neq 0} D_\varphi(m \cdot \mathbb{Q}, \mathbb{P}) = \inf_{m \neq 0} \inf_{\mathbb{Q} \in \Omega} D_\varphi(m \cdot \mathbb{Q}, \mathbb{P}) \quad (37)$$

for all subsets $\Omega \subset \mathbb{S}^K$ with regularity properties (7) and the finiteness property (9), where

$$\xi_n^{w\mathbf{W}} := \begin{cases} \left(\frac{\sum_{i \in I_1^{(n)}} W_i}{\sum_{k=1}^K \sum_{i \in I_k^{(n)}} W_i}, \dots, \frac{\sum_{i \in I_K^{(n)}} W_i}{\sum_{k=1}^K \sum_{i \in I_k^{(n)}} W_i} \right) = \frac{n \cdot \xi_n^{\mathbf{W}}}{\sum_{i=1}^n W_i}, & \text{if } \sum_{j=1}^n W_j \neq 0, \\ (\infty, \dots, \infty) =: \infty, & \text{if } \sum_{j=1}^n W_j = 0, \end{cases} \quad (38)$$

with $I_1^{(n)} := \{1, \dots, n_1\}$, $I_2^{(n)} := \{n_1 + 1, \dots, n_1 + n_2\}$, \dots , $I_K^{(n)} := \{\sum_{k=1}^{K-1} n_k + 1, \dots, n\}$ and $n_k := \lfloor n \cdot p_k \rfloor$ ($k \in \{1, \dots, K\}$) for some pre-given *known* probability vector $\mathbb{P} := (p_1, \dots, p_K)$. Recall the definition of $\xi_n^{\mathbf{W}}$ in (17) (with \mathbf{W} instead of $\widetilde{\mathbf{W}}$). The limit behaviour (37) contrasts to the one of Theorem 8, where

$$- \lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P} \left[\xi_n^{\widetilde{\mathbf{W}}} \in \Omega / M_{\mathbf{P}} \right] = \inf_{\mathbf{Q} \in \Omega} D_{\varphi}(\mathbf{Q}, \mathbf{P}) \quad (\text{cf. (18)})$$

for any $\Omega \subset \mathbb{R}^K$ with regularity properties (7) and the finiteness property (9); recall that $(\widetilde{W}_i)_{i \in \mathbb{N}}$ are i.i.d. random variables with probability distribution ζ (being connected with the divergence generator $\tilde{\varphi} := M_{\mathbf{P}} \cdot \varphi$ via the representability (15)), whereas $(W_i)_{i \in \mathbb{N}}$ are i.i.d. random variables with probability distribution ζ (being connected with the divergence generator φ via the representability (6)). Indeed, the construction leading to Theorem 8 does not hold any longer when $\Omega \subset \mathbb{S}^K$ is a set of vectors within the probability simplex \mathbb{S}^K and $\mathbf{P} \in \mathbb{S}_{>0}^K$ is a known vector in this simplex with no zero entries. In such a case, one has to use (37) and (38) instead. Notice that for each constant $A > 0$, (37) can be rewritten as

$$- \lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P} \left[\xi_n^{w\mathbf{W}} \in \Omega \right] = \inf_{\mathbf{Q} \in A \cdot \Omega} \inf_{m \neq 0} D_{\varphi} \left(\frac{m}{A} \cdot \mathbf{Q}, \mathbb{P} \right) = \inf_{\mathbf{Q} \in A \cdot \Omega} \inf_{\tilde{m} \neq 0} D_{\varphi}(\tilde{m} \cdot \mathbf{Q}, \mathbb{P}) = \inf_{\tilde{m} \neq 0} \inf_{\mathbf{Q} \in A \cdot \Omega} D_{\varphi}(\tilde{m} \cdot \mathbf{Q}, \mathbb{P}); \quad (39)$$

therein, the constraint $\mathbf{Q} \in A \cdot \Omega$ means geometrically that the vector \mathbf{Q} lives in a subset of a simplex which is parallel to the simplex \mathbb{S}^K of probability vectors and which is cut off at the edges of the first/positive orthant; in view of Remark 5(d) and (39), we can also handle such a situation. Namely, in the light of the third expression in (39) in combination with (12) to (14) for the special case of $\Omega := \mathbf{Q}$ lying in the probability simplex, it makes sense to study e.g. functional relationships between $\inf_{\tilde{m} \neq 0} D_{\tilde{c} \cdot \varphi}(\tilde{m} \cdot \mathbf{Q}, \mathbb{P})$ and $D_{\tilde{c} \cdot \varphi}(\mathbf{Q}, \mathbb{P})$ ($\tilde{c} > 0$) for $\mathbf{Q} \in A \cdot \mathbb{S}^K$ with arbitrary $A > 0$ not necessarily being equal to 1 (i.e. $\mathbf{Q} = A \cdot \mathbf{Q}$ for some probability vector \mathbf{Q}). Indeed, such a context appears naturally e.g. in connection with mass transportation problems (cf. (104) below) and with distributed energy management (cf. the paragraph after (107)); the special case $A = 1/K$ of (39) will also be used below for the application of our BS method to solving (*generalized*) *minimum/maximum entropy problems* for probability vectors (and even for sub-/super-probability vectors) \mathbf{Q} with constraints.

Let us proceed with the main context. As indicated in Remark 13(i), in a number of important cases the limit in the above Theorem 12 can be stated in terms of an invertible function G^{-1} (cf. (2)) of $\inf_{\mathbf{Q} \in \Omega} D_{\varphi}(\mathbf{Q}, \mathbb{P})$ by elimination of m ; as explained above, for the degenerate case (cf. Remark 13(v),(vi)) the search for G^{-1} is even interesting for the more general infimum over non-probability vectors. This m -elimination is the scope of the development in the following Subsection VI-B. For cases where m can not be (yet) explicitly eliminated, we deliver bounds in the second next Subsection VI-C.

B. Construction principle for the estimation of the minimum divergence, the power-type case

Within the context of Theorem 12 respectively Remark 13(v) and (vi), we obtain an explicit solution for the inner (i.e. m -concerning) minimization in (36) for the important case of power-divergence generators $\varphi_{\gamma} : \mathbb{R} \mapsto [0, \infty]$ defined by

$$\varphi_{\gamma}(t) := \begin{cases} \frac{t^{\gamma} - \gamma \cdot t + \gamma - 1}{\gamma \cdot (\gamma - 1)}, & \text{if } \gamma \in] - \infty, 0[\text{ and } t \in] 0, \infty[, \\ -\log t + t - 1, & \text{if } \gamma = 0 \text{ and } t \in] 0, \infty[, \\ \frac{t^{\gamma} - \gamma \cdot t + \gamma - 1}{\gamma \cdot (\gamma - 1)}, & \text{if } \gamma \in] 0, 1[\text{ and } t \in [0, \infty[, \\ t \cdot \log t + 1 - t, & \text{if } \gamma = 1 \text{ and } t \in [0, \infty[, \\ \frac{t^{\gamma} - \gamma \cdot t + \gamma - 1}{\gamma \cdot (\gamma - 1)} \cdot \mathbb{1}_{] 0, \infty[}(t) + \left(\frac{1}{\gamma} - \frac{t}{\gamma - 1} \right) \cdot \mathbb{1}_{]- \infty, 0]}(t), & \text{if } \gamma \in] 1, 2[\text{ and } t \in] - \infty, \infty[, \\ \frac{(t-1)^2}{2}, & \text{if } \gamma = 2 \text{ and } t \in] - \infty, \infty[, \\ \frac{t^{\gamma} - \gamma \cdot t + \gamma - 1}{\gamma \cdot (\gamma - 1)} \cdot \mathbb{1}_{] 0, \infty[}(t) + \left(\frac{1}{\gamma} - \frac{t}{\gamma - 1} \right) \cdot \mathbb{1}_{]- \infty, 0]}(t), & \text{if } \gamma \in] 2, \infty[\text{ and } t \in] - \infty, \infty[, \\ \infty, & \text{else,} \end{cases} \quad (40)$$

which for arbitrary multiplier $\tilde{c} > 0$ generate (the vector-valued form of) the *generalized power divergences* given by

$$D_{\tilde{c} \cdot \varphi_\gamma}(\mathbf{Q}, \mathbf{P}) := \begin{cases} \tilde{c} \cdot \left\{ \frac{\sum_{k=1}^K (q_k)^\gamma \cdot (p_k)^{1-\gamma}}{\gamma \cdot (\gamma-1)} - \frac{1}{\gamma-1} \cdot \sum_{k=1}^K q_k + \frac{1}{\gamma} \cdot \sum_{k=1}^K p_k \right\}, & \text{if } \gamma \in]-\infty, 0[, \mathbf{P} \in \mathbb{R}_{\neq 0}^K \text{ and } \mathbf{Q} \in \mathbb{R}_{>0}^K, \\ \tilde{c} \cdot \left\{ \sum_{k=1}^K p_k \cdot \log\left(\frac{p_k}{q_k}\right) + \sum_{k=1}^K q_k - \sum_{k=1}^K p_k \right\}, & \text{if } \gamma = 0, \mathbf{P} \in \mathbb{R}_{\neq 0}^K \text{ and } \mathbf{Q} \in \mathbb{R}_{>0}^K, \\ \tilde{c} \cdot \left\{ \frac{\sum_{k=1}^K (q_k)^\gamma \cdot (p_k)^{1-\gamma}}{\gamma \cdot (\gamma-1)} - \frac{1}{\gamma-1} \cdot \sum_{k=1}^K q_k + \frac{1}{\gamma} \cdot \sum_{k=1}^K p_k \right\}, & \text{if } \gamma \in]0, 1[, \mathbf{P} \in \mathbb{R}_{\neq 0}^K \text{ and } \mathbf{Q} \in \mathbb{R}_{\neq 0}^K, \\ \tilde{c} \cdot \left\{ \sum_{k=1}^K q_k \cdot \log\left(\frac{q_k}{p_k}\right) - \sum_{k=1}^K q_k + \sum_{k=1}^K p_k \right\}, & \text{if } \gamma = 1, \mathbf{P} \in \mathbb{R}_{>0}^K \text{ and } \mathbf{Q} \in \mathbb{R}_{\geq 0}^K, \\ \tilde{c} \cdot \left\{ \sum_{k=1}^K \frac{(q_k)^\gamma \cdot (p_k)^{1-\gamma}}{\gamma \cdot (\gamma-1)} \cdot 1_{[0, \infty[}(q_k) - \frac{1}{\gamma-1} \cdot \sum_{k=1}^K q_k + \frac{1}{\gamma} \cdot \sum_{k=1}^K p_k \right\}, & \text{if } \gamma \in]1, 2[, \mathbf{P} \in \mathbb{R}_{>0}^K \text{ and } \mathbf{Q} \in \mathbb{R}^K, \\ \tilde{c} \cdot \sum_{k=1}^K \frac{(q_k - p_k)^2}{2 \cdot p_k}, & \text{if } \gamma = 2, \mathbf{P} \in \mathbb{R}_{>0}^K \text{ and } \mathbf{Q} \in \mathbb{R}^K, \\ \tilde{c} \cdot \left\{ \sum_{k=1}^K \frac{(q_k)^\gamma \cdot (p_k)^{1-\gamma}}{\gamma \cdot (\gamma-1)} \cdot 1_{[0, \infty[}(q_k) - \frac{1}{\gamma-1} \cdot \sum_{k=1}^K q_k + \frac{1}{\gamma} \cdot \sum_{k=1}^K p_k \right\}, & \text{if } \gamma \in]2, \infty[, \mathbf{P} \in \mathbb{R}_{>0}^K \text{ and } \mathbf{Q} \in \mathbb{R}^K, \\ \infty, & \text{else;} \end{cases} \quad (41)$$

notice that one has the straightforward relationship $D_{\tilde{c} \cdot \varphi_\gamma}(\cdot, \cdot) = \tilde{c} \cdot D_{\varphi_\gamma}(\cdot, \cdot)$; however, as a motivation for the introduction of $\tilde{c} > 0$, we shall show in the solved-cases Section XII below that the corresponding probability distribution ζ (cf. (6)) of the W_i 's depends on \tilde{c} in a non-straightforward way (see also Remark 13(vi) for another motivation for \tilde{c}). In the course of this, it turns out that $\tilde{c} \cdot \varphi_\gamma \in \Upsilon(]a_\gamma, \infty[)$ with $a_\gamma = 0$ for $\gamma \in]-\infty, 1[$ and $a_\gamma = -\infty$ for $\gamma \in [2, \infty[$.

For $\tilde{c} = 1$ and probability vectors \mathbf{Q}, \mathbb{P} in \mathbb{S}^K respectively $\mathbb{S}_{>0}^K$, the divergences (41) simplify considerably, namely to the well-known *power divergences* $D_{\varphi_\gamma}(\mathbf{Q}, \mathbb{P})$ in the scaling of e.g. Liese & Vajda [7] (in other scalings they are also called *Rathie & Kannapan's non-additive directed divergences of order γ* [122], *Cressie-Read divergences* [123] [8], *relative Tsallis entropies or Tsallis cross-entropies* [124] (see also [125]), *Amari's alpha-divergences* [126]); for some comprehensive overviews on power divergences $D_{\varphi_\gamma}(\mathbf{Q}, \mathbb{P})$ — including statistical applications to goodness-of-fit testing and minimum distance estimation — the reader is referred to the insightful books [7]–[13], the survey articles [4],[14], and the references therein. Prominent and widely used special cases of $D_{\varphi_\gamma}(\mathbf{Q}, \mathbb{P})$ are the omnipresent *Kullback-Leibler information divergence (relative entropy)* where $\gamma = 1$, the equally important *reverse Kullback-Leibler information divergence (reverse relative entropy)* where $\gamma = 0$, the *Pearson chi-square divergence* ($\gamma = 2$), the *Neyman chi-square divergence* ($\gamma = -1$), the *Hellinger divergence* ($\gamma = \frac{1}{2}$, also called squared Hellinger distance, squared Matusita distance [27] or squared Hellinger-Kakutani metric, see e.g. [28] ¹⁸). Some exemplary (relatively) recent studies and applications of power divergences $D_{\varphi_\gamma}(\mathbf{Q}, \mathbb{P})$ — aside from the vast statistical literature (including in particular maximum likelihood estimation and Pearson's chi-square test) — appear e.g. in [88],[127]–[148]; in [149] with $\gamma = 2$; in [150] with $\gamma = 1$; in [151] with $\gamma = -1$ and $\gamma = \frac{1}{2}$; in [152]–[155] with $\gamma = \frac{1}{2}$.

For $\tilde{c} = 1$ and nonnegative-component vectors \mathbf{Q}, \mathbf{P} in $\mathbb{R}_{\geq 0}^K$ respectively $\mathbb{R}_{>0}^K$ respectively $\mathbb{R}_{\neq 0}^K$, the generalized power divergences $D_{\varphi_\gamma}(\mathbf{Q}, \mathbf{P})$ of (41) also (partially) simplify, and were treated by [91] (for even more general probability measures, deriving e.g. also generalized Pinsker's inequalities); for a more general comprehensive technical treatment see also e.g. [20].

Returning to the general context, in Theorem 12 we stated that for each $\mathbb{P} \in \mathbb{S}_{>0}^K$ the function $\mathbf{Q} \mapsto \inf_{m \neq 0} D_\varphi(m \cdot \mathbf{Q}, \mathbb{P})$ is BS-minimizable (cf. (2)) on all sets $\mathfrak{Q} \subset \mathbb{S}^K$ satisfying (7) and (9). The (corresponding subset of the) following Lemma 14 is the *cornerstone* leading from this statement to BS-minimizability of the function $\mathbf{Q} \mapsto D_\varphi(\mathbf{Q}, \mathbb{P})$ on those same sets, for the special divergences in (41). To formulate this in a transparent way, we employ the following three fundamental quantities $H_\gamma(\mathbf{Q}, \mathbb{P})$, $I(\mathbf{Q}, \mathbb{P})$, $\tilde{I}(\mathbf{Q}, \mathbb{P})$ and the arbitrary constant $A > 0$ (where for $A = 1$ all the following vectors \mathbf{Q} will turn into probability vectors \mathbf{Q}). Indeed — for any constellation $(\gamma, \mathbb{P}, \mathbf{Q}) \in \tilde{\Gamma} \times \tilde{\mathcal{M}}_1 \times \tilde{\mathcal{M}}_2$, where $\tilde{\Gamma} \times \tilde{\mathcal{M}}_1 \times \tilde{\mathcal{M}}_2 :=]0, 1[\times \mathbb{S}^K \times A \cdot \mathbb{S}^K$ or $\tilde{\Gamma} \times \tilde{\mathcal{M}}_1 \times \tilde{\mathcal{M}}_2 :=]-\infty, 0[\times \mathbb{S}^K \times A \cdot \mathbb{S}_{>0}^K$ or $\tilde{\Gamma} \times \tilde{\mathcal{M}}_1 \times \tilde{\mathcal{M}}_2 :=]1, \infty[\times \mathbb{S}_{>0}^K \times A \cdot \mathbb{S}^K$ — let

$$0 < H_\gamma(\mathbf{Q}, \mathbb{P}) := \sum_{k=1}^K (q_k)^\gamma \cdot (p_k)^{1-\gamma} = 1 + \gamma \cdot (A - 1) + \gamma \cdot (\gamma - 1) \cdot D_{\varphi_\gamma}(\mathbf{Q}, \mathbb{P}), \quad \gamma \in \mathbb{R} \setminus \{0, 1\}, \quad (42)$$

be the *modified γ -order Hellinger integral of \mathbf{Q} and \mathbb{P}* . Furthermore, for any $\mathbb{P} \in \mathbb{S}_{>0}^K$, $\mathbf{Q} \in A \cdot \mathbb{S}^K$, let

$$-1 < I(\mathbf{Q}, \mathbb{P}) := \sum_{k=1}^K q_k \cdot \log\left(\frac{q_k}{p_k}\right) = D_{\varphi_1}(\mathbf{Q}, \mathbb{P}) + A - 1 \quad (43)$$

¹⁸in some literature, the (square root of the) Hellinger divergence (HD) is misleadingly called Bhattacharyya distance; however, the latter is *basically* some rescaled logarithm of HD, namely $R_{1/2}(\mathbf{Q}, \mathbb{P})$ (cf. (73) with $\gamma = 1/2$)

be the *modified Kullback-Leibler information (modified relative entropy)*. Finally, for any $\mathbb{P} \in \mathbb{S}^K$, $\mathbf{Q} \in A \cdot \mathbb{S}_{>0}^K$, let

$$1 - A \leq \tilde{I}(\mathbf{Q}, \mathbb{P}) := \sum_{k=1}^K p_k \cdot \log \left(\frac{p_k}{q_k} \right) = D_{\varphi_0}(\mathbf{Q}, \mathbb{P}) + 1 - A \quad (44)$$

be the *modified reverse Kullback-Leibler information (modified reverse relative entropy)*. In terms of (42), (43) and (44) we obtain the following assertions which will be proved in Appendix C.

Lemma 14: Let $A > 0$ be an arbitrary constant.

(a) Let $\tilde{c} > 0$ be arbitrary and $(\gamma, \mathbb{P}, \mathbf{Q}) \in \tilde{\Gamma} \times \tilde{\mathcal{M}}_1 \times \tilde{\mathcal{M}}_2$ as above. Then one has

$$\begin{aligned} \inf_{m \neq 0} D_{\tilde{c} \cdot \varphi_\gamma}(m \cdot \mathbf{Q}, \mathbb{P}) &= \inf_{m > 0} D_{\tilde{c} \cdot \varphi_\gamma}(m \cdot \mathbf{Q}, \mathbb{P}) = \frac{\tilde{c}}{\gamma} \cdot \left[1 - A^{\gamma/(\gamma-1)} \cdot \left[1 + \gamma \cdot (A-1) + \frac{\gamma \cdot (\gamma-1)}{\tilde{c}} \cdot D_{\tilde{c} \cdot \varphi_\gamma}(\mathbf{Q}, \mathbb{P}) \right]^{-1/(\gamma-1)} \right] \\ &= \frac{\tilde{c}}{\gamma} \cdot \left[1 - A^{\gamma/(\gamma-1)} \cdot H_\gamma(\mathbf{Q}, \mathbb{P})^{-1/(\gamma-1)} \right] \end{aligned} \quad (45)$$

and consequently for any subset $A \cdot \mathfrak{Q} \subset \tilde{\mathcal{M}}_2$

$$\inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} \inf_{m \neq 0} D_{\tilde{c} \cdot \varphi_\gamma}(m \cdot \mathbf{Q}, \mathbb{P}) = \frac{\tilde{c}}{\gamma} \cdot \left[1 - A^{\gamma/(\gamma-1)} \cdot \left[1 + \gamma \cdot (A-1) + \frac{\gamma \cdot (\gamma-1)}{\tilde{c}} \cdot \inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} D_{\tilde{c} \cdot \varphi_\gamma}(\mathbf{Q}, \mathbb{P}) \right]^{-1/(\gamma-1)} \right], \quad (46)$$

$$\arg \inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} \inf_{m \neq 0} D_{\tilde{c} \cdot \varphi_\gamma}(m \cdot \mathbf{Q}, \mathbb{P}) = \arg \inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} D_{\tilde{c} \cdot \varphi_\gamma}(\mathbf{Q}, \mathbb{P}), \quad (47)$$

$$\inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} \inf_{m \neq 0} D_{\varphi_\gamma}(m \cdot \mathbf{Q}, \mathbb{P}) = \frac{1}{\gamma} \cdot \left[1 - A^{\gamma/(\gamma-1)} \cdot \left[\inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} H_\gamma(\mathbf{Q}, \mathbb{P}) \right]^{-1/(\gamma-1)} \right], \quad \text{for } \gamma < 0 \text{ and } \gamma > 1, \quad (48)$$

$$\arg \inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} \inf_{m \neq 0} D_{\varphi_\gamma}(m \cdot \mathbf{Q}, \mathbb{P}) = \arg \inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} H_\gamma(\mathbf{Q}, \mathbb{P}), \quad \text{for } \gamma < 0 \text{ and } \gamma > 1, \quad (49)$$

$$\inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} \inf_{m \neq 0} D_{\varphi_\gamma}(m \cdot \mathbf{Q}, \mathbb{P}) = \frac{1}{\gamma} \cdot \left[1 - A^{\gamma/(\gamma-1)} \cdot \left[\sup_{\mathbf{Q} \in A \cdot \mathfrak{Q}} H_\gamma(\mathbf{Q}, \mathbb{P}) \right]^{-1/(\gamma-1)} \right], \quad \text{for } \gamma \in]0, 1[, \quad (50)$$

$$\arg \inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} \inf_{m \neq 0} D_{\varphi_\gamma}(m \cdot \mathbf{Q}, \mathbb{P}) = \arg \sup_{\mathbf{Q} \in A \cdot \mathfrak{Q}} H_\gamma(\mathbf{Q}, \mathbb{P}), \quad \text{for } \gamma \in]0, 1[, \quad (51)$$

provided that the infimum on the right-hand side of (46) exists.

(b) For any $\mathbb{P} \in \mathbb{S}_{>0}^K$, $\mathbf{Q} \in A \cdot \mathbb{S}^K$, $\tilde{c} > 0$ one gets

$$\begin{aligned} \inf_{m \neq 0} D_{\tilde{c} \cdot \varphi_1}(m \cdot \mathbf{Q}, \mathbb{P}) &= \inf_{m > 0} D_{\tilde{c} \cdot \varphi_1}(m \cdot \mathbf{Q}, \mathbb{P}) = \tilde{c} \cdot \left[1 - A \cdot \exp \left(- \frac{1}{A \cdot \tilde{c}} \cdot D_{\tilde{c} \cdot \varphi_1}(\mathbf{Q}, \mathbb{P}) + \frac{1}{A} - 1 \right) \right] \\ &= \tilde{c} \cdot \left[1 - A \cdot \exp \left(- \frac{1}{A} \cdot I(\mathbf{Q}, \mathbb{P}) \right) \right] \end{aligned} \quad (52)$$

and consequently for any subset $A \cdot \mathfrak{Q} \subset A \cdot \mathbb{S}^K$

$$\inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} \inf_{m \neq 0} D_{\tilde{c} \cdot \varphi_1}(m \cdot \mathbf{Q}, \mathbb{P}) = \tilde{c} \cdot \left[1 - A \cdot \exp \left(- \frac{1}{A \cdot \tilde{c}} \cdot \inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} D_{\tilde{c} \cdot \varphi_1}(\mathbf{Q}, \mathbb{P}) + \frac{1}{A} - 1 \right) \right], \quad (53)$$

$$\arg \inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} \inf_{m \neq 0} D_{\tilde{c} \cdot \varphi_1}(m \cdot \mathbf{Q}, \mathbb{P}) = \arg \inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} D_{\tilde{c} \cdot \varphi_1}(\mathbf{Q}, \mathbb{P}), \quad (54)$$

$$\inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} \inf_{m \neq 0} D_{\varphi_1}(m \cdot \mathbf{Q}, \mathbb{P}) = \left[1 - A \cdot \exp \left(- \frac{1}{A} \cdot \inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} I(\mathbf{Q}, \mathbb{P}) \right) \right], \quad (55)$$

$$\arg \inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} \inf_{m \neq 0} D_{\varphi_1}(m \cdot \mathbf{Q}, \mathbb{P}) = \arg \inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} I(\mathbf{Q}, \mathbb{P}), \quad (56)$$

provided that the infimum on the right-hand side of (53) exists.

(c) For any $\mathbb{P} \in \mathbb{S}^K$, $\mathbf{Q} \in A \cdot \mathbb{S}_{>0}^K$, $\tilde{c} > 0$ we obtain

$$\inf_{m \neq 0} D_{\tilde{c} \cdot \varphi_0}(m \cdot \mathbf{Q}, \mathbb{P}) = \inf_{m > 0} D_{\tilde{c} \cdot \varphi_0}(m \cdot \mathbf{Q}, \mathbb{P}) = D_{\tilde{c} \cdot \varphi_0}(\mathbf{Q}, \mathbb{P}) + \tilde{c} \cdot (1 - A + \log A) = \tilde{c} \cdot \left(\tilde{I}(\mathbf{Q}, \mathbb{P}) + \log A \right) \quad (57)$$

and consequently for any set subset $A \cdot \mathfrak{Q} \subset A \cdot \mathbb{S}_{>0}^K$

$$\inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} \inf_{m \neq 0} D_{\tilde{c} \cdot \varphi_0}(m \cdot \mathbf{Q}, \mathbb{P}) = \tilde{c} \cdot (1 - A + \log A) + \inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} D_{\tilde{c} \cdot \varphi_0}(\mathbf{Q}, \mathbb{P}), \quad (58)$$

$$\arg \inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} \inf_{m \neq 0} D_{\tilde{c} \cdot \varphi_0}(m \cdot \mathbf{Q}, \mathbb{P}) = \arg \inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} D_{\tilde{c} \cdot \varphi_0}(\mathbf{Q}, \mathbb{P}), \quad (59)$$

$$\inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} \inf_{m \neq 0} D_{\varphi_0}(m \cdot \mathbf{Q}, \mathbb{P}) = \log A + \inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} \tilde{I}(\mathbf{Q}, \mathbb{P}), \quad (60)$$

$$\arg \inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} \inf_{m \neq 0} D_{\varphi_0}(m \cdot \mathbf{Q}, \mathbb{P}) = \arg \inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} \tilde{I}(\mathbf{Q}, \mathbb{P}), \quad (61)$$

provided that the infimum on the right-hand side of (58) exists.

Remark 15: Notice that for $\mathbb{P} \in \mathbb{S}_{>0}^K$ and $\mathbf{Q} \in A \cdot \mathbb{S}^K$, the modified Kullback-Leibler information has the property $I(\mathbf{Q}, \mathbb{P}) \geq 0$ if $A \geq 1$ (cf. (43)); otherwise, $I(\mathbf{Q}, \mathbb{P})$ may become negative, as can be easily seen from the case where $\mathbb{P} := \mathbb{P}^{unif} := (\frac{1}{K}, \dots, \frac{1}{K})$ is the probability vector of frequencies of the uniform distribution on $\{1, \dots, K\}$, and $\mathbf{Q} := (\frac{1}{K+1}, 0, \dots, 0)$. Analogously, for $\mathbb{P} \in \mathbb{S}^K$ and $\mathbf{Q} \in A \cdot \mathbb{S}_{>0}^K$ one gets $\tilde{I}(\mathbf{Q}, \mathbb{P}) \geq 0$ if $A \leq 1$ (cf. (44)); otherwise, $\tilde{I}(\mathbf{Q}, \mathbb{P})$ may become negative (take e.g. $\mathbf{Q} = (\frac{K+1}{K}, \dots, \frac{K+1}{K})$ and $\mathbb{P} := (1, 0, \dots, 0)$).

Remark 16: (a) In the context of Remark 13(vi), according to (39) applied to $\varphi := \tilde{c} \cdot \varphi_\gamma$, for all cases $\gamma \in]-\infty, 0[\cup]0, 1[\cup]2, \infty[$ the left-hand side of each of (46), (48), (50) is independent of $A > 0$ and equal to $-\lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}[\xi_n^{w\mathbf{W}} \in \mathfrak{Q}]$ where — as will be shown in the solved-cases Section XII below — the corresponding \mathbf{W} 's have probability distribution $\zeta[\cdot] = \mathbb{P}[W_1 \in \cdot]$ (cf. (6)) which varies “quite drastically” with γ (and the case $\gamma \in]1, 2[$ has to be even excluded for analytical difficulties¹⁹). Analogously, each of the left-hand sides of (53),(55),(58),(60) is also independent of $A > 0$ and equal to $-\lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}[\xi_n^{w\mathbf{W}} \in \mathfrak{Q}]$ for some \mathbf{W} of respective distribution. Hence, by inversion, all the extremum-describing target quantities $\inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} D_{\tilde{c} \cdot \varphi_\gamma}(\mathbf{Q}, \mathbb{P})$ ($\gamma \in \mathbb{R} \setminus]1, 2[$), $\inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} H_\gamma(\mathbf{Q}, \mathbb{P})$ ($\gamma \in]-\infty, 0[\cup]2, \infty[$), $\sup_{\mathbf{Q} \in A \cdot \mathfrak{Q}} H_\gamma(\mathbf{Q}, \mathbb{P})$ ($\gamma \in]0, 1[$), $\inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} I(\mathbf{Q}, \mathbb{P})$ and $\inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} \tilde{I}(\mathbf{Q}, \mathbb{P})$ can be expressed as $G\left(-\lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}[\xi_n^{w\mathbf{W}} \in \mathfrak{Q}]\right)$ for some explicitly known (A -dependent) function G . This means that — in the sense of Definition 1 — all the corresponding four “cornerstone quantities” $D_{\tilde{c} \cdot \varphi_\gamma}(\mathbf{Q}, \mathbb{P})$, $H_\gamma(\mathbf{Q}, \mathbb{P})$, $I(\mathbf{Q}, \mathbb{P})$, $\tilde{I}(\mathbf{Q}, \mathbb{P})$ are BS-minimizable, respectively BS-maximizable, on $\mathfrak{Q} = A \cdot \mathfrak{Q}$. The above-mentioned inversions (i.e. constructions of $G(\cdot)$) will be concretely carried out in Section XII below — namely in the Propositions 29 to 34.

(b) The case $\varphi := \tilde{c} \cdot \varphi_\gamma$ ($\gamma \in]-\infty, 0[\cup]0, 1[\cup]2, \infty[$) of Theorem 12 works analogously to (a), with the differences that we employ $A = 1$ (instead of arbitrary $A > 0$), (36) (instead of (39)), $\mathbb{P}_{X_1^n}[\cdot]$ (instead of $\mathbb{P}[\cdot]$), and $\xi_{n,\mathbf{X}}^{w\mathbf{W}}$ (instead of $\xi_n^{w\mathbf{W}}$).

(c) As can be seen in the proof of Lemma 14, for the important case $\gamma = 2$ the formulas (45) to (49) also hold for $A < 0$.

For the rest of the paper, in connection with the three points (a),(b),(c) of Remark 16, we will always interpret (without explicit mentioning, for the sake of brevity) the expression “BS minimizable/maximizable” accordingly in terms of $-\lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}[\xi_n^{w\mathbf{W}} \in \cdot]$ respectively of $-\lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_{X_1^n}[\xi_{n,\mathbf{X}}^{w\mathbf{W}} \in \cdot]$.

Let us end this subsection with a comparison: suppose that we have a (sufficiently large) number n of *concrete* data observations $X_i = x_i$ ($i = 1, \dots, n$) from the unknown probability distribution \mathbb{P} (in vector form), and from these we want to approximate/estimate the unknown distance $\inf_{\mathbf{Q} \in \mathfrak{Q}} D_{\varphi_\gamma}(\mathbf{Q}, \mathbb{P})$ from a family of probability models (in vector form) \mathfrak{Q} (e.g. for model-adequacy evaluations, for goodness-of-fit testing purposes): by the elaborations in Section XII below, for the approximation of $\inf_{\mathbf{Q} \in \mathfrak{Q}} D_{\varphi_\gamma}(\mathbf{Q}, \mathbb{P})$ we can use

$$G\left(-\frac{1}{n} \cdot \log \mathbb{P}_{x_1^n}[\xi_{n,\mathbf{X}}^{w\mathbf{W}} \in \mathfrak{Q}]\right) \quad (62)$$

where $\mathbb{P}_{x_1^n}[\cdot] := \mathbb{P}[\cdot | X_1 = x_1, \dots, X_n = x_n]$, $\mathbf{x} := (x_1, \dots, x_n)$, and G (cf. (2)) is e.g. chosen as follows:

$G(z) := -\frac{\tilde{c}}{\gamma \cdot (\gamma - 1)} \cdot \left\{1 - \left(1 - \frac{\tilde{c}}{\gamma} \cdot z\right)^{1-\gamma}\right\}$ for the three cases $\gamma < 0$, $\gamma \in]0, 1[$ and $\gamma \geq 2$, $G(z) := z$ for $\gamma = 0$ (reversed Kullback-Leibler divergence), and $G(z) := -\tilde{c} \cdot \log\left(1 - \frac{1}{\tilde{c}} \cdot z\right)$ for $\gamma = 1$ (Kullback-Leibler divergence). Notice that (62) contrasts to the alternative approximation (of $\inf_{\mathbf{Q} \in \mathfrak{Q}} D_{\varphi_\gamma}(\mathbf{Q}, \mathbb{P})$) given by

$$\inf_{\mathbf{Q} \in \mathfrak{Q}} D_{\varphi_\gamma}(\mathbf{Q}, \mathbb{P}_n^{emp,co}) \quad (63)$$

which is used in the context of “classical” statistical minimum distance estimation (MDE) with power divergences; in (63), we have employed $\mathbb{P}_n^{emp,co} = \frac{1}{n} \cdot \sum_{i=1}^n \delta_{x_i}$ to be the realization of the empirical distribution $\mathbb{P}_n^{emp} = \frac{1}{n} \cdot \sum_{i=1}^n \delta_{X_i}$. Indeed, especially in complicated high-dimensional non-parametric or semi-parametric big-data contexts, we have substituted a quite difficult *optimization problem* (63) by a much easier solvable *counting problem* (62). The same holds analogously for Renyi distances/divergences, etc.

C. Construction principle for bounds of the minimum divergence in the general case

Returning to Theorem 12, we now consider the general case when the divergence generator $\varphi \in \mathcal{T}(]a, b[)$ is *not* of the power type (40). Recall from (35),(36) the crucial terms (with $\mathbb{P} \in \mathbb{S}_{>0}$)

$$\inf_{m \neq 0} D_\varphi(m \cdot \mathfrak{Q}, \mathbb{P}) := \inf_{m \neq 0} \inf_{\mathbf{Q} \in \mathfrak{Q}} D_\varphi(m \cdot \mathbf{Q}, \mathbb{P}) = \inf_{\mathbf{Q} \in \mathfrak{Q}} \inf_{m \neq 0} D_\varphi(m \cdot \mathbf{Q}, \mathbb{P}) < \infty \quad (64)$$

for all sets \mathfrak{Q} satisfying the regularity properties (7) and the convenient, more restrictive finiteness property

¹⁹because in this case there are some indications that the representation (6) only holds for some *signed* probability distribution ζ (e.g. having a density with positive and negative values).

$$\inf_{\mathbf{Q} \in \mathfrak{Q}} \inf_{k=1, \dots, K} \frac{q_k}{p_k} \in \text{dom}(\varphi), \quad \sup_{\mathbf{Q} \in \mathfrak{Q}} \sup_{k=1, \dots, K} \frac{q_k}{p_k} \in \text{dom}(\varphi) \quad (65)$$

which implies (9); notice that $\inf_{k=1, \dots, K} \frac{q_k}{p_k} \leq 1$, $\sup_{k=1, \dots, K} \frac{q_k}{p_k} \geq 1$ with equalities if and only if $\mathbf{Q} = \mathbb{P}$. Since $\mathfrak{Q} \neq \{\mathbb{P}\}$ (cf. the right-hand side of (7)), the double infimum (supremum) in (65) is strictly smaller (larger) than 1. In general, the inner minimization $\inf_{m \neq 0} D_\varphi(m \cdot \mathbf{Q}, \mathbb{P})$ in (64) can not be performed in explicit closed form, but e.g. in the specific case of power divergences (cf. (40), (41)) the optimization $\inf_{m \neq 0} D_{\tilde{c} \cdot \varphi_\gamma}(m \cdot \mathbf{Q}, \mathbb{P})$ produces an explicit form, which in turn leads to a straightforward one-to-one correspondence between $D_{\tilde{c} \cdot \varphi_\gamma}(\mathfrak{Q}, \mathbb{P})$ and $\inf_{m \neq 0} D_{\tilde{c} \cdot \varphi_\gamma}(m \cdot \mathfrak{Q}, \mathbb{P})$ (cf. Lemma 14). Under (7) and (65) one has

$$\inf_{m \neq 0} D_\varphi(m \cdot \mathfrak{Q}, \mathbb{P}) \leq D_\varphi(\mathfrak{Q}, \mathbb{P}) \leq D_\varphi(\mathbf{Q}, \mathbb{P}).$$

For transparency, we first investigate the (widely useable) subset where $\text{dom}(\varphi) =]0, \infty[$ (and thus, $\text{int}(\text{dom}(\varphi)) =]a, b[=]0, \infty[$) and $\mathfrak{Q} \subset \mathbb{S}_{>0}^K$. Let us start with the lower bound $\inf_{m \neq 0} D_\varphi(m \cdot \mathfrak{Q}, \mathbb{P})$. We prove that the minimizer in m is a well defined constant, which belongs to a compact set in $\mathbb{R}_{>0}$. To see this, notice first that from (65) one can obtain

$$\left[\inf_{\mathbf{Q} \in \mathfrak{Q}} \inf_{k=1, \dots, K} \frac{m \cdot q_k}{p_k} \in \text{dom}(\varphi) \quad \text{and} \quad \sup_{\mathbf{Q} \in \mathfrak{Q}} \sup_{k=1, \dots, K} \frac{m \cdot q_k}{p_k} \in \text{dom}(\varphi) \right] \iff m \in]0, \infty[.$$

Moreover, for any fixed \mathbf{Q} in \mathfrak{Q} there is a unique number $m = m(\mathbf{Q}) > 0$ which satisfies the first-order optimality condition

$$\psi_{\mathbf{Q}}(m) := \frac{d}{dm} D_\varphi(m \cdot \mathbf{Q}, \mathbb{P}) = \sum_{k=1}^K q_k \cdot \varphi' \left(\frac{m \cdot q_k}{p_k} \right) = 0 \quad \text{for } m \in]0, \infty[\quad (66)$$

$$\text{and thus} \quad D_\varphi(m(\mathbf{Q}) \cdot \mathbf{Q}, \mathbb{P}) = \inf_{m \neq 0} D_\varphi(m \cdot \mathbf{Q}, \mathbb{P}); \quad (67)$$

indeed, the mapping $]0, \infty[\ni m \rightarrow D_\varphi(m \cdot \mathbf{Q}, \mathbb{P})$ is strictly convex and infinitely differentiable, and the strictly increasing function $\psi_{\mathbf{Q}}$ is such that $\psi_{\mathbf{Q}}(m)$ is strictly negative for all $m \in]0, 1[$ for which $\sup_{k=1, \dots, K} \frac{m \cdot q_k}{p_k} < 1$ whereas $\psi_{\mathbf{Q}}(m)$ is strictly positive for all $m > 1$ for which $\inf_{k=1, \dots, K} \frac{m \cdot q_k}{p_k} > 1$ (recall the note right after (65) and $\varphi'(1) = 0$). Hence, for any $\mathbf{Q} \in \mathfrak{Q}$ the unique zero $m(\mathbf{Q})$ of (66) (and hence, the unique minimizer in (67)) is in the compact interval

$$\left[\frac{1}{\sup_{k=1, \dots, K} \frac{q_k}{p_k}}, \frac{1}{\inf_{k=1, \dots, K} \frac{q_k}{p_k}} \right] \subseteq \left[\frac{1}{\sup_{\mathbf{Q} \in \mathfrak{Q}} \sup_{k=1, \dots, K} \frac{q_k}{p_k}}, \frac{1}{\inf_{\mathbf{Q} \in \mathfrak{Q}} \inf_{k=1, \dots, K} \frac{q_k}{p_k}} \right] \subset \left] \frac{1}{b}, \frac{1}{a} \right[=]0, \infty[.$$

When \mathfrak{Q} is closed in \mathbb{S}^K , then by continuity of the function $\mathbf{Q} \mapsto D_\varphi(m(\mathbf{Q}) \cdot \mathbf{Q}, \mathbb{P})$ there exists a \mathbf{Q}^* in \mathfrak{Q} which achieves the infimum on \mathfrak{Q} . When \mathfrak{Q} is not closed but satisfies (7), then the infimum exists anyway, possibly on the boundary $\partial \mathfrak{Q}$. Anyhow, for such \mathbf{Q}^* there holds

$$D_\varphi(m(\mathbf{Q}^*) \cdot \mathbf{Q}^*, \mathbb{P}) \leq D_\varphi(\mathfrak{Q}, \mathbb{P}) \leq D_\varphi(\mathbf{Q}^*, \mathbb{P}), \quad (68)$$

where we use the continuity of $\mathbf{Q} \mapsto D_\varphi(\mathbf{Q}, \mathbb{P})$ and (7) to obtain the last inequality above, even when $\mathbf{Q}^* \in \partial \mathfrak{Q}$ and $\mathbf{Q}^* \notin \mathfrak{Q}$. That (68) provides sharp bounds can be seen through the case of power divergences. Indeed, for the latter one basically gets (cf. Appendix C) $m(\mathbf{Q}) = (1 + \frac{\gamma(\gamma-1)}{\tilde{c}} \cdot D_{\tilde{c} \cdot \varphi_\gamma}(\mathbf{Q}, \mathbb{P}))^{1/(1-\gamma)}$ and $D_{\tilde{c} \cdot \varphi_\gamma}(m(\mathbf{Q}) \cdot \mathbf{Q}, \mathbb{P}) = \frac{\tilde{c}}{\gamma} (1 - m(\mathbf{Q}))$ for the case $\gamma \in \mathbb{R} \setminus \{0, 1\}$, respectively, $m(\mathbf{Q}) = \exp(-\frac{1}{\tilde{c}} \cdot D_{\tilde{c} \cdot \varphi_1}(\mathbf{Q}, \mathbb{P}))$ and $D_{\tilde{c} \cdot \varphi_1}(m(\mathbf{Q}) \cdot \mathbf{Q}, \mathbb{P}) = \tilde{c} \cdot (1 - m(\mathbf{Q}))$ for the case $\gamma = 1$, respectively, $m(\mathbf{Q}) = 1$ and $D_{\tilde{c} \cdot \varphi_0}(m(\mathbf{Q}) \cdot \mathbf{Q}, \mathbb{P}) = D_{\tilde{c} \cdot \varphi_0}(\mathbf{Q}, \mathbb{P})$ for the remaining case $\gamma = 0$. In all cases, $D_{\tilde{c} \cdot \varphi_\gamma}(m(\mathbf{Q}) \cdot \mathbf{Q}, \mathbb{P})$ is an increasing function of $D_{\tilde{c} \cdot \varphi_\gamma}(\mathbf{Q}, \mathbb{P})$ and therefore, $\mathbf{Q}^* \in \arg \inf_{\mathbf{Q} \in \mathfrak{Q}} D_{\tilde{c} \cdot \varphi_\gamma}(m(\mathbf{Q}) \cdot \mathbf{Q}, \mathbb{P})$ also satisfies $\mathbf{Q}^* \in \arg \inf_{\mathbf{Q} \in \mathfrak{Q}} D_{\tilde{c} \cdot \varphi_\gamma}(\mathbf{Q}, \mathbb{P})$. Hence, the right-hand side and the left-hand side of (68) coincide. Now due to (6), the LHS of (68) can be estimated since by Theorem 12 for each $\mathbb{P} \in \mathbb{S}_{>0}^K$ the divergence $\inf_{m \neq 0} D_\varphi(m \cdot \mathbf{Q}, \mathbb{P})$ is BS-minimizable on sets $\mathfrak{Q} \subset \mathbb{S}^K$. We shall propose in the below-treated Subsection X-B3 an algorithm to handle the estimation of the RHS of (68), whenever \mathbb{P} is known (as in Remark 13(v)) or when \mathbb{P} is approximated by the empirical distribution of the data set (X_1, \dots, X_n) . Also note that (68) holds also for \mathfrak{Q} substituted by $A \cdot \mathfrak{Q}$ for any $A \neq 0$.

Other cases of interest include when $\text{dom}(\varphi)$ is not $]0, \infty[$. We list two cases which extend the above discussion. Firstly, consider φ with $\text{dom}(\varphi) = [0, \infty[$. Then we may extend (68) to cases when $\mathfrak{Q} \subset \mathbb{S}^K$ instead of $\mathfrak{Q} \subset \mathbb{S}_{>0}^K$, hence allowing for possible null entries in \mathfrak{Q} . When $\text{dom}(\varphi) =]a, b[$ for some $a < 0$, then clearly the same argument leading to (68) holds; this case is of interest, for instance, when extending a statistical model to signed measures (see e.g. [156] for the important task of testing the number of components in a parametric probability mixture model).

D. On the difference between minimization problems of deterministic nature and risk minimization

In the context of minimization of the functional $\Phi_{\mathbf{P}}(\mathbf{Q})$ over $\Omega \subset \mathbb{R}^K$ for *known* vector \mathbf{P} , due to Theorem 8 our bare simulation approach allows for the *approximate solution* for any divergence D_φ satisfying the basic representation (15). Indeed, any proxy of $\mathbb{P}[\{\xi_n^{\widetilde{\mathbf{W}}} \in \Omega/M_{\mathbf{P}}\}]$ yields a corresponding proxy for $\Phi_{\mathbf{P}}[\Omega]$. This paves the way to the solution of numerous optimization problems, where the divergence D_φ is specifically suited to the problem at hand. In the statistical context, when the *probability distribution* (in its vector-form) \mathbb{P} is *unknown* up to some indirect information provided by sampling or by any mean providing a sequence $(X_i)_{i \in \mathbb{N}}$ satisfying condition (26) (resp. (30)), Theorem 12 adds a complementary step of complexity; indeed, the *estimation* of $\Phi_{\mathbf{P}}(\Omega)$ over $\Omega \subset \mathbb{S}^K$ results as its subproduct through the optimization upon m which can be performed explicitly only in a number of specific divergences D_φ , e.g. the power divergences D_{φ_γ} , and which carries over also to their monotone transformations such as e.g. the Renyi divergences. It is of relevance to mention that — as already indicated above — these divergences cover a *very broad range* of statistical criteria, indeed most of them, from the (maximum-likelihood estimation connected) likelihood divergence ($\gamma = 0$) to the Kullback-Leibler one ($\gamma = 1$), the two standard Chi-square distances ($\gamma = 2, \gamma = -1$), the Hellinger distance ($\gamma = 1/2$), etc.; in contrast with deterministic minimization problems, the choice of a statistical criterion (or risk function) is not imposed by the modelling of the problem at hand, but is dictated by the need for sharp measures of fit. Other divergences are more difficult to handle and our general results in Section VI-C still prove some usefulness, since estimation of upper and lower bounds for risk is of common use.

As a “preparatory” remark, recall first that each probability distribution (probability measure) \mathbb{P} on $\mathcal{Y} = \{d_1, \dots, d_K\}$ has been uniquely identified with the vector $\mathbb{P} := (p_1, \dots, p_K) \in \mathbb{S}^K$ of the corresponding probability masses (frequencies) $p_k = \mathbb{P}[\{d_k\}]$ via $\mathbb{P}[A] = \sum_{k=1}^K p_k \cdot \mathbb{1}_A(d_k)$ for each $A \subset \mathcal{Y}$; from this, we have measured the distance/divergence between two probability distributions \mathbb{P}, \mathbb{Q} through the distance/divergence between their frequency vectors \mathbb{P}, \mathbb{Q} :

$$D_\varphi(\mathbb{Q}, \mathbb{P}) := D_\varphi(\mathbb{Q}, \mathbb{P}) \quad (\text{cf. (28)}).$$

However, it has been noted in [157] in a context of even more general divergences $D(\mathbf{Q}, \mathbf{P})$ between vectors \mathbf{P}, \mathbf{Q} that — alternatively — the latter two may consist of components $p_k = \mathbb{P}[\{E_k\}]$, $q_k = \mathbb{Q}[\{E_k\}]$ which are probabilities of only some *selected* (e.g. increasing) events $(E_k)_{k \in \{1, \dots, M\}}$ of *application-based concrete* interest (within *not necessarily discrete* probability models), e.g. related to cumulative distribution functions (see [21]). Of course, we can apply our BS method to such a vector context. As other alternatives, we can also deal with divergences between *non-probabilistic* uncertainty quantifications, such as fuzzy sets and basic belief assignments (see our paper’s full arXiv-version [158]).

VII. RENYI DIVERGENCES AND FRIENDS

It is well known that Renyi divergences are widely used important tools in information theory as well as in the adjacent fields of statistics, machine learning, artificial intelligence, signal processing and pattern recognition. As a consequence of our considerations in the above Subsection VI-B in combination with the solved cases in Section XII below, we can also apply our BS method to the constrained optimization of Renyi divergences and closely related important quantities. To start with, let us fix $\tilde{c} = 1$ and an arbitrary triple $(\gamma, \mathbb{P}, \mathbb{Q})$ which satisfies the assumptions of Lemma 14(a) with $A := \sum_{k=1}^K q_k > 0$. For such a setup, we have obtained in (42) the γ -order Hellinger integral (of \mathbb{Q} and \mathbb{P}) $H_\gamma(\mathbb{Q}, \mathbb{P}) > 0$, which is not a divergence; as a terminology-concerning side remark, let us mention that $H_\gamma(\mathbb{Q}, \mathbb{P})$ ($\gamma \geq 1$) is called *relative information generating function* in [159], see e.g. also [160]; moreover, $H_\gamma(\mathbb{Q}, \mathbb{P})$ is sometimes termed (γ -order) *Chernoff coefficient* being a component of the Chernoff distances/informations [161]. In [162] the name (γ -order) *Hellinger transform* is used. Notice that the special case $\gamma = \frac{1}{2}$ is nothing but (a multiple of) the well-known important *Bhattacharyya coefficient* (cf. [22],[25],[26])

$$BC(\mathbb{Q}, \mathbb{P}) := H_{1/2}(\mathbb{Q}, \mathbb{P}) = \sum_{k=1}^K \sqrt{q_k \cdot p_k} = 1 + \frac{1}{2} \cdot (A - 1) + \frac{1}{2} \cdot \left(\frac{1}{2} - 1\right) \cdot D_{\varphi_{\frac{1}{2}}}(\mathbb{Q}, \mathbb{P})$$

which is also known as *affinity* (cf. [27], see e.g. also [163]) and (*classic, non-quantum*) *fidelity similarity* (cf. e.g. [28]); for non-probability vectors $\mathbf{P} \in \mathbb{R}_{\geq 0}^K$ with $M_{\mathbf{P}} > 0$ one can simply retransform $\mathbb{P} := \frac{\mathbf{P}}{M_{\mathbf{P}}}$ and thus imbed $BC(\mathbf{Q}, \mathbf{P}) = \sqrt{M_{\mathbf{P}}} \cdot BC(\mathbb{Q}, \mathbb{P})$ into our BS context. There is a vast literature on very recent applications of the Bhattacharyya coefficient, for instance it appears exemplarily in [164]–[195]. To proceed with the general context, for any $\gamma \in]-\infty, 0[\cup]0, 1[\cup]1, \infty[$ let the function $h_\gamma :]0, \infty[\mapsto]-\infty, \infty[$ be such that $x \mapsto h_\gamma(1 + \gamma \cdot (A - 1) + \gamma \cdot (\gamma - 1) \cdot x)$ is continuous and strictly increasing (respectively, strictly decreasing) for all $x \geq 0$ with $1 + \gamma \cdot (A - 1) + \gamma \cdot (\gamma - 1) \cdot x > 0$; since $D_{\varphi_\gamma}(\mathbf{Q}, \mathbf{P})$ is BS-minimizable on $\Omega = A \cdot \Omega$, then also the — not necessarily nonnegative — quantity $h_\gamma\left(1 + \gamma \cdot (A - 1) + \gamma \cdot (\gamma - 1) \cdot D_{\varphi_\gamma}(\mathbf{Q}, \mathbf{P})\right) = h_\gamma\left(H_\gamma(\mathbf{Q}, \mathbf{P})\right)$ is BS-minimizable (respectively, BS-maximizable) on $\Omega = A \cdot \Omega$. If h_γ satisfies additionally $h_\gamma(1) = 0$ as well as $h_\gamma(1 + \gamma \cdot (A - 1) + \gamma \cdot (\gamma - 1) \cdot x) \geq 0$ for all $x \geq 0$ with $1 + \gamma \cdot (A - 1) + \gamma \cdot (\gamma - 1) \cdot x > 0$, then $D_{h_\gamma}(\mathbf{Q}, \mathbf{P}) := h_\gamma\left(1 + \gamma \cdot (A - 1) + \gamma \cdot (\gamma - 1) \cdot D_{\varphi_\gamma}(\mathbf{Q}, \mathbf{P})\right) = h_\gamma\left(H_\gamma(\mathbf{Q}, \mathbf{P})\right)$ constitutes a divergence²⁰ which is BS-minimizable on $\Omega = A \cdot \Omega$ (respectively, BS-maximizable on $\Omega = A \cdot \Omega$).

²⁰in the usual sense that $D_{h_\gamma}(\mathbf{Q}, \mathbf{P}) \geq 0$ with equality iff $\mathbf{Q} = \mathbf{P}$.

Let us consider some important examples. For the identity mapping $h_\gamma^{Id}(y) := y$ ($y > 0$) the function $x \mapsto 1 + \gamma \cdot (A - 1) + \gamma \cdot (\gamma - 1) \cdot x$ is strictly increasing for $\gamma < 0$ and $\gamma > 1$ (on the required domain of x), and strictly decreasing for $\gamma \in]0, 1[$. Accordingly, $H_\gamma(\mathbf{Q}, \mathbb{P})$ is BS-minimizable on $\Omega = A \cdot \mathfrak{Q}$ for $\gamma < 0$ and $\gamma \geq 2$ and BS-maximizable on $\Omega = A \cdot \mathfrak{Q}$ for $\gamma \in]0, 1[$ (this is consistent with (48),(50)); in particular, the Bhattacharyya coefficient $BC(\mathbf{Q}, \mathbb{P})$ is BS-maximizable on $\Omega = A \cdot \mathfrak{Q}$. Some other important choices are

$$h_\gamma(y) := h_{c_1, c_2, c_3}(y) := c_1 \cdot (y^{c_2} - c_3), \quad y > 0, c_1, c_2 \in \mathbb{R} \setminus \{0\}, c_3 \in \mathbb{R}, \quad (69)$$

$$h_\gamma(y) := h_{c_4, f}^R(y) := \lim_{c_2 \rightarrow 0} h_{c_4/f(c_2), c_2, 1}(y) = \frac{c_4}{f'(0)} \cdot \log(y), \quad y > 0, c_4 \in \mathbb{R} \setminus \{0\}, \quad (70)$$

$$h_\gamma(y) := h_{c_5, c_6}^{GB2}(y) := c_5 \cdot (\arccos(y))^{c_6}, \quad \gamma \in]0, 1[, y \in]0, 1[, c_5 > 0, c_6 > 0, \quad (71)$$

$$h_\gamma(y) := h_{\nu, c_7}^{BB}(y) := c_7 \cdot \frac{\log(1 - \frac{1-y}{\nu})}{\log(1 - \frac{1}{\nu})}, \quad \gamma \in]0, 1[, y \in]0, 1[, c_7 > 0, \nu \in]-\infty, 0[\cup]1, \infty[, \quad (72)$$

where the constants c_1 to c_7 may depend on γ , and f is some (maybe γ -dependent) function which is differentiable in a neighborhood of 0 and satisfies $f(0) = 0, f'(0) \neq 0$ (e.g. $f(z) = c_8 \cdot z$ for some non-zero constant c_8). Clearly, $h_{c_1, c_2, c_3}(\cdot)$ is strictly increasing (respectively, strictly decreasing) if and only if $c_1 \cdot c_2 > 0$ (respectively, $c_1 \cdot c_2 < 0$). Moreover, $h_{c_4, f}^R(\cdot)$ is strictly increasing (respectively, strictly decreasing) if and only if $\frac{c_4}{f'(0)} > 0$ (respectively, $\frac{c_4}{f'(0)} < 0$). Furthermore, both $h_{c_5, c_6}^{GB2}(\cdot)$ and $h_{\nu, c_7}^{BB}(\cdot)$ are strictly decreasing. For instance, the special case $h_\gamma(y) = h_{c_4, Id}^R(y)$ with $c_4 := \frac{1}{\gamma \cdot (\gamma - 1)}$ (recall that $\gamma \in]-\infty, 0[\cup]0, 1[\cup]2, \infty[$) and identity function $f := Id$ leads to the quantities

$$\begin{aligned} R_\gamma(\mathbf{Q}, \mathbb{P}) &:= D_{h_{c_4, Id}^R}(\mathbf{Q}, \mathbb{P}) = \frac{\log\left(1 + \gamma \cdot (A - 1) + \gamma \cdot (\gamma - 1) \cdot D_{\varphi_\gamma}(\mathbf{Q}, \mathbb{P})\right)}{\gamma \cdot (\gamma - 1)} = \frac{\log\left(H_\gamma(\mathbf{Q}, \mathbb{P})\right)}{\gamma \cdot (\gamma - 1)} \\ &= \frac{\log\left(\sum_{k=1}^K (q_k)^\gamma \cdot (p_k)^{1-\gamma}\right)}{\gamma \cdot (\gamma - 1)}, \quad \gamma \in]-\infty, 0[\cup]0, 1[\cup]2, \infty[, \end{aligned} \quad (73)$$

(provided that all involved power divergences are finite), which are thus BS-minimizable on $\Omega = A \cdot \mathfrak{Q}$; notice that $R_\gamma(\mathbf{Q}, \mathbb{P}) \geq 0$ if $\gamma \in [2, \infty[$ together with $A \in [1, \infty[$, and if $\gamma \in]-\infty, 0[\cup]0, 1[$ together with $A \in]0, 1[$. The special subcase $A = 1$ in (73) (and thus, \mathbf{Q} is a probability vector \mathfrak{Q}) corresponds to the prominent *Renyi divergences/distances* [23] (in the scaling of e.g. Liese & Vajda [7] and in probability-vector form), see e.g. [24] for a comprehensive study of their properties; as a side remark, $\gamma \cdot (\gamma - 1) \cdot R_\gamma(\mathbf{Q}, \mathbb{P})$ is also employed in the Chernoff distances/informations [161]. The special subcase $R_{1/2}(\mathbf{Q}, \mathbb{P})$ (i.e. $\gamma = 1/2$ and $A = 1$ in (73)) corresponds to (a multiple of) the widely used *Bhattacharyya distance* (of type 1) between \mathbf{Q} and \mathbb{P} , cf. [22] (see e.g. also [196]). Sometimes, $\exp(R_\gamma(\mathbf{Q}, \mathbb{P}))$ is also called *Renyi divergence/distance*. Some exemplary (relatively) recent studies and applications of Renyi divergences $R_\gamma(\mathbf{Q}, \mathbb{P})$ (respectively, their multiple or exponential) — aside from the substantial statistical literature — appear e.g. in [30],[197]–[215]. There is vast literature on recent applications of the above-mentioned special case $R_{1/2}(\mathbf{Q}, \mathbb{P})$ — that is, the Bhattacharyya distance (of type 1); for instance, it appears in [216]–[229]. As a further example, consider (cf. (71))

$$\begin{aligned} \mathcal{B}_{\gamma, c_5, c_6}(\mathbf{Q}, \mathbb{P}) &:= D_{h_{c_5, c_6}^{GB2}}(\mathbf{Q}, \mathbb{P}) = c_5 \cdot \left(\arccos\left(1 + \gamma \cdot (\gamma - 1) \cdot D_{\varphi_\gamma}(\mathbf{Q}, \mathbb{P})\right)\right)^{c_6} = c_5 \cdot \left(\arccos\left(H_{\varphi_\gamma}(\mathbf{Q}, \mathbb{P})\right)\right)^{c_6} \\ &= c_5 \cdot \left(\arccos\left(\sum_{k=1}^K (q_k)^\gamma \cdot (p_k)^{1-\gamma}\right)\right)^{c_6} \geq 0, \quad \gamma \in]0, 1[, c_5 > 0, c_6 > 0, \end{aligned}$$

which is BS-minimizable on \mathfrak{Q} . The case $\mathcal{B}_{1/2, 1, 1}(\mathbf{Q}, \mathbb{P})$ corresponds to the well-known *Bhattacharyya arccos distance* (*Bhattacharyya distance of type 2*) in [26] (which is also called Wootters distance [230]), and $\mathcal{B}_{1/2, 1, 2}(\mathbf{Q}, \mathbb{P})$ to its variant in [25]; the case $\mathcal{B}_{1/2, 2, 1}(\mathbf{Q}, \mathbb{P})$ is known as *Fisher distance* or *Rao distance* or *geodesic distance* (see e.g. [28]); a nice graphical illustration of the geometric connection between the Fisher distance $\mathcal{B}_{1/2, 2, 1}(\mathbf{Q}, \mathbb{P})$ and the Hellinger distance/metric $\sqrt{\frac{1}{2} \cdot D_{\varphi_{1/2}}(\mathbf{Q}, \mathbb{P})}$ can be found e.g. on p.35 in [231]. Some exemplary applications of the Bhattacharyya arccos distance $\mathcal{B}_{1/2, 1, 1}(\mathbf{Q}, \mathbb{P})$ can be found e.g. in [232]–[235],[155]. Let us give another example, namely (cf. (72))

$$\begin{aligned} \tilde{\mathcal{B}}_{\gamma, \nu, c_7}(\mathbf{Q}, \mathbb{P}) &:= D_{h_{\nu, c_7}^{BB}}(\mathbf{Q}, \mathbb{P}) = \frac{c_7}{\log(1 - \frac{1}{\nu})} \cdot \log\left(1 - \frac{1 - \left(1 + \gamma \cdot (\gamma - 1) \cdot D_{\varphi_\gamma}(\mathbf{Q}, \mathbb{P})\right)}{\nu}\right) \\ &= \frac{c_7}{\log(1 - \frac{1}{\nu})} \cdot \log\left(1 - \frac{1 - H_{\varphi_\gamma}(\mathbf{Q}, \mathbb{P})}{\nu}\right) = \frac{c_7}{\log(1 - \frac{1}{\nu})} \cdot \log\left(1 - \frac{1 - \sum_{k=1}^K (q_k)^\gamma \cdot (p_k)^{1-\gamma}}{\nu}\right) \in [0, c_7[, \\ &\quad \gamma \in]0, 1[, c_7 > 0, \nu \in]-\infty, 0[\cup]1, \infty[, \end{aligned}$$

which is BS-minimizable on \mathfrak{Q} . The case $\tilde{\mathcal{B}}_{1/2,\nu,1}(\mathfrak{Q}, \mathbb{P})$ corresponds to the *Bounded Bhattacharyya Distance Measures* of [236]. We can also employ divergences of the form $\check{R}_\gamma(\mathfrak{Q}, \mathbb{P}) := R_\gamma(T_1(\mathfrak{Q}), T_2(\mathbb{P}))$ ²¹ where $T_1 : \mathcal{D}_1 \mapsto \mathcal{R}_1, T_2 : \mathcal{D}_2 \mapsto \mathcal{R}_2$ are (say) invertible functions on appropriately chosen subsets $\mathcal{D}_1, \mathcal{D}_2, \mathcal{R}_1, \mathcal{R}_2$ of the probability-vector simplex \mathbb{S}^K . For instance, consider the following special case (with a slight abuse of notation):

$$\check{R}_\gamma(\mathfrak{Q}, \mathbb{P}) := R_\gamma(\tilde{\mathfrak{Q}}, \tilde{\mathbb{P}}) = \frac{1}{\gamma \cdot (\gamma - 1)} \cdot \log \left(\sum_{k=1}^K \left(\frac{(q_k)^{\nu_1}}{\sum_{j=1}^K (q_j)^{\nu_1}} \right)^\gamma \cdot \left(\frac{(p_k)^{\nu_2}}{\sum_{j=1}^K (p_j)^{\nu_2}} \right)^{1-\gamma} \right) \quad (74)$$

where (i) $\tilde{\mathfrak{Q}} := (\tilde{q}_k)_{k=1}^K$ with $\tilde{q}_k := \frac{(q_k)^{\nu_1}}{\sum_{j=1}^K (q_j)^{\nu_1}}$ is the *escort probability distribution (in vector form) associated with the probability distribution (in vector form) $\mathfrak{Q} := (q_k)_{k=1}^K \in \mathbb{S}_{>0}^K$* , and (ii) $\tilde{\mathbb{P}} := (\tilde{p}_k)_{k=1}^K$ with $\tilde{p}_k := \frac{(p_k)^{\nu_2}}{\sum_{j=1}^K (p_j)^{\nu_2}}$ is the *escort probability distribution associated with the probability distribution $\mathbb{P} := (p_k)_{k=1}^K \in \mathbb{S}_{>0}^K$* , in terms of some fixed escort parameters $\nu_1 > 0, \nu_2 > 0$. For the special choice $\nu_1 = \nu_2 > 0$ and $\gamma := \frac{\nu}{\nu_1}$ with $\nu \in]0, \nu_1[\cup]2\nu_1, \infty[$ we obtain from (74)

$$\begin{aligned} 0 &\leq \frac{\nu}{\nu_1} \cdot R_{\nu/\nu_1}(\tilde{\mathfrak{Q}}, \tilde{\mathbb{P}}) = \frac{\log \left(\sum_{k=1}^K (\tilde{q}_k)^{\nu/\nu_1} \cdot (\tilde{p}_k)^{1-(\nu/\nu_1)} \right)}{\frac{\nu}{\nu_1} - 1} \\ &= \frac{\nu_1}{\nu - \nu_1} \cdot \log \left(\sum_{k=1}^K (q_k)^\nu \cdot (p_k)^{\nu_1 - \nu} \right) - \frac{\nu}{\nu - \nu_1} \cdot \log \left(\sum_{k=1}^K (q_k)^{\nu_1} \right) + \log \left(\sum_{k=1}^K (p_k)^{\nu_1} \right) =: \check{R}_{\nu/\nu_1}(\mathfrak{Q}, \mathbb{P}) \end{aligned} \quad (75)$$

which is BS-minimizable (in $\tilde{\mathfrak{Q}}$) on \mathfrak{Q} . Our divergence $\check{R}_{\nu/\nu_1}(\mathfrak{Q}, \mathbb{P})$ in (75) is basically a multiple of a divergence which has been very recently used in [237]. Moreover, $\check{R}_{1/\nu_1}(\mathfrak{Q}, \mathbb{P})$ (i.e. the special case $\nu = 1$ in (75)) is equal to *Sundaresan's divergence* [29] [30] (see also [238], [70], [71], [239]); for our BS-approach, we need the restriction $\nu_1 \in]0, \frac{1}{2}[\cup]1, \infty[$. Notice that Sundaresan's divergence can be employed in mismatch-cases of (i) Campbell's coding problem, (ii) Arian's guessing problem, (iii) memoryless guessing, and (iv) task partitioning problems; see e.g. [30], [198], [200].

Returning to the general context, functions of the modified Kullback-Leibler information $I(\mathfrak{Q}, \mathbb{P})$ and the modified reverse Kullback-Leibler information $\tilde{I}(\mathfrak{Q}, \mathbb{P})$ can be treated analogously. For the sake of brevity, we only deal with the former and fix arbitrary $\mathbb{P} \in \mathbb{S}_{>0}^K$ and $\mathfrak{Q} \in A \cdot \mathbb{S}^K$ with $A := \sum_{k=1}^K q_k > 0$. For this, in (43) we have obtained $I(\mathfrak{Q}, \mathbb{P})$ which is generally not a divergence (cf. Remark 15). In the following, let the function $h_1 :]-1, \infty[\mapsto]-\infty, \infty[$ be continuous and strictly increasing (respectively, strictly decreasing); since $D_{\varphi_1}(\mathfrak{Q}, \mathbb{P})$ is BS-minimizable on $\Omega = A \cdot \mathfrak{Q}$, also the quantity $h_1(A - 1 + D_{\varphi_1}(\mathfrak{Q}, \mathbb{P})) = h_1(I(\mathfrak{Q}, \mathbb{P}))$ is BS-minimizable on $\Omega = A \cdot \mathfrak{Q}$ (respectively, BS-maximizable on $\Omega = A \cdot \mathfrak{Q}$). In particular, by using the negative identity mapping $h_\gamma^{-Id}(y) := -y$ ($y > -1$) we get that $-I(\mathfrak{Q}, \mathbb{P})$ is BS-maximizable. Another exemplary choice for h_1 is (cf. [240] in the scaling of e.g. [241])

$$h_1(y) := h_s^{SM}(y) := \frac{e^{(s-1) \cdot y} - 1}{s - 1}, \quad y \in \mathbb{R}, s \in]0, 1[\cup]1, \infty[, \quad (76)$$

which is strictly increasing; hence, $h_s^{SM}(I(\mathfrak{Q}, \mathbb{P}))$ (and also $h_s^{SM}(D_{\varphi_1}(\mathfrak{Q}, \mathbb{P}))$) is BS-minimizable on $\Omega = A \cdot \mathfrak{Q}$.

VIII. CONSTRAINED ENTROPY MAXIMIZATIONS

Of course, entropies are extremely important tools in information theory, as well as in the adjacent fields of statistics, machine learning, artificial intelligence, signal processing and pattern recognition. As a consequence of our considerations in the above Subsection VI-B in combination with the solved cases in Section XII, we can also apply our BS method to the constrained optimization of a wide range of entropies and closely related diversity indices. To begin with, let us fix any $(\gamma, \mathfrak{Q}) \in (\tilde{\Gamma} \setminus]1, 2[) \times \mathcal{M}_2$ (cf. Lemma 14(a)) with $A := \sum_{k=1}^K q_k > 0$. Moreover, we take $\mathbb{P} := \mathbb{P}^{unif} := (\frac{1}{K}, \dots, \frac{1}{K})$ to be the probability vector of frequencies of the uniform distribution on $\{1, \dots, K\}$. Then, for $\gamma \in]-\infty, 0[\cup]0, 1[\cup]2, \infty[$ one gets $H_\gamma(\mathfrak{Q}, \mathbb{P}^{unif}) = K^{\gamma-1} \cdot \sum_{k=1}^K q_k^\gamma$. One can rewrite $K^{1-\gamma} \cdot H_\gamma(\mathfrak{Q}, \mathbb{P}^{unif}) = \sum_{k=1}^K q_k^\gamma$; the latter is sometimes called *heterogeneity index of type γ* , see e.g. [50], with $\gamma = 2$ being the *Simpson-Herfindahl index* which is also known as *index of coincidence* (cf. [47] and its generalization in [48]). Alternatively, $\sum_{k=1}^K q_k^\gamma$ is also called *Onicescu's information energy* in case of $\gamma = 2$ (cf. [242], see also [243] for comprehensive investigations) and in general *information energy of order γ* (cf. [244], see also e.g. [245]); for exemplary applications the reader may take (discretized versions of) e.g. [246]–[248]. In some other literature (see e.g. [160]), $\sum_{k=1}^K q_k^\gamma$ is alternatively called *Golomb's [249] information generating function (of a probability distribution \mathfrak{Q})*; yet another name is *generalized information potential* and for $\gamma = 2$ *information potential* (cf. e.g. [250], [251]). From the above

²¹and analogously power divergences $\check{D}_{\tilde{c}, \varphi_\gamma}(\mathfrak{Q}, \mathbb{P}) := D_{\tilde{c}, \varphi_\gamma}(T_1(\mathfrak{Q}), T_2(\mathbb{P}))$ etc.

investigations, we obtain that $\sum_{k=1}^K q_k^\gamma$ is BS-minimizable on $\Omega = A \cdot \mathfrak{Q}$ for $\gamma < 0$ and $\gamma \geq 2$, and BS-maximizable on $\Omega = A \cdot \mathfrak{Q}$ for $\gamma \in]0, 1[$. More generally, by employing (69) and (70), for the class of entropies (diversity indices)

$$\begin{aligned} \mathcal{E}_{\gamma, c_1, c_2, c_3}(\mathbf{Q}) &:= h_{c_1, c_2, c_3} \left(\sum_{k=1}^K q_k^\gamma \right) = c_1 \cdot \left(\left(\sum_{k=1}^K q_k^\gamma \right)^{c_2} - c_3 \right) \\ &= c_1 \cdot \left(K^{c_2 \cdot (1-\gamma)} \cdot H_\gamma(\mathbf{Q}, \mathbb{P}^{unif})^{c_2} - c_3 \right), \quad c_1, c_2 \in \mathbb{R} \setminus \{0\}, c_3 \in \mathbb{R}, \end{aligned} \quad (77)$$

$$\begin{aligned} \mathcal{E}_{c_4, f}^R(\mathbf{Q}) &:= h_{c_4, f}^R \left(\sum_{k=1}^K q_k^\gamma \right) = \frac{c_4}{f'(0)} \cdot \log \left(\sum_{k=1}^K q_k^\gamma \right), \\ &= \frac{c_4}{f'(0)} \cdot \left(\log(H_\gamma(\mathbf{Q}, \mathbb{P}^{unif})) + (1-\gamma) \cdot \log(K) \right), \quad c_4 \in \mathbb{R} \setminus \{0\}, \end{aligned} \quad (78)$$

(which is similar to the entropy-class of Morales et al. [252] who use a different, more restrictive parametrization and probability distributions \mathbf{Q}), one gets the following extremum-behaviour:

- $\mathcal{E}_{\gamma, c_1, c_2, c_3}(\mathbf{Q})$ is BS-minimizable if $\gamma < 0$ and $c_1 \cdot c_2 > 0$;
- $\mathcal{E}_{\gamma, c_1, c_2, c_3}(\mathbf{Q})$ is BS-minimizable if $\gamma \geq 2$ and $c_1 \cdot c_2 > 0$;
- $\mathcal{E}_{\gamma, c_1, c_2, c_3}(\mathbf{Q})$ is BS-minimizable if $\gamma \in]0, 1[$ and $c_1 \cdot c_2 < 0$;
- $\mathcal{E}_{\gamma, c_1, c_2, c_3}(\mathbf{Q})$ is BS-maximizable if $\gamma < 0$ and $c_1 \cdot c_2 < 0$;
- $\mathcal{E}_{\gamma, c_1, c_2, c_3}(\mathbf{Q})$ is BS-maximizable if $\gamma \geq 2$ and $c_1 \cdot c_2 < 0$;
- $\mathcal{E}_{\gamma, c_1, c_2, c_3}(\mathbf{Q})$ is BS-maximizable if $\gamma \in]0, 1[$ and $c_1 \cdot c_2 > 0$;
- $\mathcal{E}_{c_4, f}^R(\mathbf{Q})$ is BS-minimizable if $\gamma < 0$ and $\frac{c_4}{f'(0)} > 0$;
- $\mathcal{E}_{c_4, f}^R(\mathbf{Q})$ is BS-minimizable if $\gamma \geq 2$ and $\frac{c_4}{f'(0)} > 0$;
- $\mathcal{E}_{c_4, f}^R(\mathbf{Q})$ is BS-minimizable if $\gamma \in]0, 1[$ and $\frac{c_4}{f'(0)} < 0$;
- $\mathcal{E}_{c_4, f}^R(\mathbf{Q})$ is BS-maximizable if $\gamma < 0$ and $\frac{c_4}{f'(0)} < 0$;
- $\mathcal{E}_{c_4, f}^R(\mathbf{Q})$ is BS-maximizable if $\gamma \geq 2$ and $\frac{c_4}{f'(0)} < 0$;
- $\mathcal{E}_{c_4, f}^R(\mathbf{Q})$ is BS-maximizable if $\gamma \in]0, 1[$ and $\frac{c_4}{f'(0)} > 0$.

From this, one can deduce that our new BS method works for the constrained minimization/maximization of the following well-known, prominently used measures of entropy respectively measures of diversity, and beyond:

(E1) $c_1 = 1, c_2 = \frac{1}{\gamma}, c_3 = 0$: the Euclidean γ -norm (also known as γ -norm heterogeneity index, see e.g. [50]) $\|\mathbf{Q}\|_\gamma := \left(\sum_{k=1}^K q_k^\gamma \right)^{1/\gamma} = K^{(1-\gamma)/\gamma} \cdot \left(H_\gamma(\mathbf{Q}, \mathbb{P}^{unif}) \right)^{1/\gamma}$ is BS-minimizable on $\Omega = A \cdot \mathfrak{Q}$ for $\gamma \geq 2$, and BS-maximizable on $\Omega = A \cdot \mathfrak{Q}$ for $\gamma < 0$ and $\gamma \in]0, 1[$ (note that $\|\mathbf{Q}\|_1 = A$);

similarly, the γ -mean heterogeneity index (see e.g. [50], as well as [52] for its interpretation as “effective number of species” respectively as “numbers equivalent”) given by $\mathcal{E}^{HI}(\mathbf{Q}) := \left(\sum_{k=1}^K q_k^\gamma \right)^{1/(\gamma-1)} = \frac{1}{K} \cdot \left(H_\gamma(\mathbf{Q}, \mathbb{P}^{unif}) \right)^{1/(\gamma-1)}$ is BS-minimizable on $\Omega = A \cdot \mathfrak{Q}$ for $\gamma \in]0, 1[$ and $\gamma \geq 2$, and BS-maximizable on $\Omega = A \cdot \mathfrak{Q}$ for $\gamma < 0$. Alternatively, $\mathcal{E}^{HI}(\mathbf{Q})$ is also called (γ -order) Hill diversity index or Hill number [46], respectively (γ -order) Hannah-Kay index [253], respectively (γ -order) Renyi heterogeneity (cf. [254]), respectively (γ -order) exponential Renyi entropy or exponential entropy (cf. [255]) since it is equal to $\exp(\mathcal{E}^{gR}(\mathbf{Q}))$ (cf. (E6) below). The γ -mean heterogeneity index (under one of the above-mentioned namings) was recently employed e.g. by [256]–[260].

(E2) $c_1 = \frac{1}{2^{1-\gamma}-1}, c_2 = 1, c_3 = 1$: the entropy

$$\mathcal{E}^{gHC}(\mathbf{Q}) := \frac{1}{2^{1-\gamma}-1} \cdot \left(\sum_{k=1}^K q_k^\gamma - 1 \right) = \frac{1}{2^{1-\gamma}-1} \cdot \left(K^{1-\gamma} \cdot H_\gamma(\mathbf{Q}, \mathbb{P}^{unif}) - 1 \right) \quad (79)$$

is BS-minimizable on $\Omega = A \cdot \mathfrak{Q}$ for $\gamma < 0$, and BS-maximizable on $\Omega = A \cdot \mathfrak{Q}$ for $\gamma \in]0, 1[$ and $\gamma \geq 2$; the special subcase $A = 1$ in (79) (and thus, $\mathbf{Q} = \mathbf{Q}$ is a probability vector) corresponds to the γ -order entropy of Havrda-Charvat [42] (also called non-additive γ -order Tsallis entropy [43] in statistical physics) where the special case $\gamma = 2$ is (a multiple of) Vajda’s quadratic entropy [9] and Ahlswede’s identification entropy [261] (see also [262]). Some exemplary (relatively) recent studies and applications of $\mathcal{E}^{gHC}(\mathbf{Q})$ appear e.g. in [212],[246],[248],[263]–[272]. For $\gamma = 2$, a directly connected quantity is the measure of concentration (cf. e.g. [273]) $\mathcal{E}^{gMC}(\mathbf{Q}) := 1 - \frac{1}{K} - \mathcal{E}^{gHC}(\mathbf{Q}) = \sum_{k=1}^K \left(q_k - \frac{1}{K} \right)^2$ which (up to a multiple) was introduced by [274] as an appropriate measure of information for quantum experiments.

(E3) $\gamma := \frac{1}{\tilde{\gamma}}, c_1 = \frac{1}{\tilde{\gamma}-1}, c_2 = \tilde{\gamma}, c_3 = 1$: the entropy

$$\mathcal{E}^{gA}(\mathbf{Q}) := \frac{1}{\tilde{\gamma}-1} \cdot \left(\left(\sum_{k=1}^K q_k^{1/\tilde{\gamma}} \right)^{\tilde{\gamma}} - 1 \right) = \frac{1}{\tilde{\gamma}-1} \cdot \left(K^{\tilde{\gamma}-1} \cdot H_{1/\tilde{\gamma}}(\mathbf{Q}, \mathbb{P}^{unif})^{\tilde{\gamma}} - 1 \right) \quad (80)$$

is BS-minimizable on $\Omega = A \cdot \mathfrak{Q}$ for $\tilde{\gamma} < 0$, and BS-maximizable on $\Omega = A \cdot \mathfrak{Q}$ for $\tilde{\gamma} \in]0, \frac{1}{2}]$ and $\tilde{\gamma} > 1$; the special subcase $A = 1$ in (80) (and thus, $\mathbf{Q} = \mathfrak{Q}$ is a probability vector) corresponds to the $\tilde{\gamma}$ -order entropy of Arimoto [44].

(E4) $s \in \mathbb{R} \setminus \{1\}$, $c_1 = \frac{1}{1-s}$, $c_2 = \frac{1-s}{1-\gamma}$, $c_3 = 1$: the entropy

$$\mathcal{E}^{gSM1}(\mathbf{Q}) := \frac{1}{1-s} \cdot \left(\left(\sum_{k=1}^K q_k^\gamma \right)^{(1-s)/(1-\gamma)} - 1 \right) = \frac{1}{1-s} \cdot \left(K^{1-s} \cdot H_\gamma(\mathbf{Q}, \mathbb{P}^{unif})^{(1-s)/(1-\gamma)} - 1 \right) \quad (81)$$

is BS-minimizable on $\Omega = A \cdot \mathfrak{Q}$ for $\gamma < 0$ and BS-maximizable on $\Omega = A \cdot \mathfrak{Q}$ for $\gamma \in]0, 1[$ and $\gamma \geq 2$; the special subcase $A = 1$ in (81) (and thus, $\mathbf{Q} = \mathfrak{Q}$ is a probability vector) corresponds to the entropy of order γ and degree s of Sharma & Mittal [45] in the scaling of e.g. [40].

(E5) $s \in \mathbb{R} \setminus \{0\}$, $\gamma = s + 1$, $c_1 = -\frac{1}{s}$, $c_2 = 1$, $c_3 = 1$: the diversity index

$$\mathcal{E}^{gPT}(\mathbf{Q}) := -\frac{1}{s} \cdot \left(\sum_{k=1}^K q_k^{s+1} - 1 \right) = -\frac{1}{s} \cdot \left(K^{-s} \cdot H_{s+1}(\mathbf{Q}, \mathbb{P}^{unif}) - 1 \right) \quad (82)$$

is BS-minimizable on $\Omega = A \cdot \mathfrak{Q}$ for $s < -1$ and BS-maximizable on $\Omega = A \cdot \mathfrak{Q}$ for $s \in]-1, 0[$ and $s \geq 1$; the special subcase $A = 1$ in (82) (and thus, $\mathbf{Q} = \mathfrak{Q}$ is a probability vector) corresponds to the diversity index of degree s of Patil & Taillie [49]; the case $s = 1$ for probability measures $\mathbf{Q} = \mathfrak{Q}$ gives the well-known Gini-Simpson diversity index.

(E6) $c_4 = \frac{1}{1-\gamma}$, $f(z) = z$: the entropy

$$\mathcal{E}^{gR}(\mathbf{Q}) := \frac{1}{1-\gamma} \cdot \log \left(\sum_{k=1}^K q_k^\gamma \right) = \frac{1}{1-\gamma} \cdot \left(\log(H_\gamma(\mathbf{Q}, \mathbb{P}^{unif})) + (1-\gamma) \cdot \log(K) \right) = \frac{\log 2}{1-\gamma} \cdot \log_2 \left(\sum_{k=1}^K q_k^\gamma \right) \quad (83)$$

is BS-minimizable on $\Omega = A \cdot \mathfrak{Q}$ for $\gamma < 0$, and BS-maximizable on $\Omega = A \cdot \mathfrak{Q}$ for $\gamma \in]0, 1[$ and $\gamma \geq 2$; the special subcase $A = 1$ in (83) (and thus, $\mathbf{Q} = \mathfrak{Q}$ is a probability vector) corresponds to the prominent (additive) γ -order Renyi entropy [23]. As well known, there is a vast literature on Renyi entropies $\mathcal{E}^{gR}(\mathbf{Q})$. Some exemplary (mostly recent) studies and applications appear e.g. in [30],[198],[200],[212],[246],[248],[267],[275]–[292].

Remark 17: (i) For Renyi entropies there are also matrix versions $\mathcal{E}^{gR}(X) := \frac{1}{1-\gamma} \cdot \log \left(\sum_{i=1}^{K_1} \sum_{j=1}^{K_2} x_{ij}^\gamma \right)$ where $X := (x_{ij})_{i=1, \dots, K_1}^{j=1, \dots, K_2}$ is a $K_1 \times K_2$ -matrix whose elements x_{ij} are (say) strictly positive and sum up to A . Such a setup with $A = 1$ is e.g. used in *time-frequency analyses of signals* where the i 's correspond to discrete time points, the j 's to discrete frequencies, and x_{ij} to the probability that (i, j) occurs; see e.g. [293]. Another line of application is to use as X the normalized communicability matrix of a directed network (respectively the upper triangular part of X in case of an unweighted and undirected network). Of course, the matrix version $\mathcal{E}^{gR}(X)$ can be easily and equivalently rewritten in our vector version $\mathcal{E}^{gR}(\mathbf{Q})$ by setting $\mathbf{Q} := (q_1, \dots, q_{K_1 \cdot K_2})$ such that $x_{ij} = q_{(i-1) \cdot K_2 + j}$ ($i = 1, \dots, K_1, j = 1, \dots, K_2$) and hence $K := K_1 \cdot K_2$; accordingly, we can apply our BS method.

(ii) The latter conversion works analogously also for matrix versions of all the other entropies, divergences, etc. of this paper; more flexible versions where $i \in \{1, \dots, K_1\}$, $j \in J_i$ for some $J_i \subseteq \{1, \dots, K_2\}$ as well as multidimensional-array/tensor versions can be transformed in a similar book-keeping manner, too. For instance, within the above-mentioned framework of unweighted and undirected networks, [294] and [295] employ communicability matrix versions of the Shannon entropy and the Jensen-Shannon divergence (JSD); see also [296] for similar network applications of the JSD. Moreover, [297] use “3D versions” of Tsallis entropies for brain magnetic resonance (MR) image segmentation.

Remark 18: All the above cases which are BS-maximizable can be interpreted as bare-simulation approach to the solution of *generalized maximum entropy problems on $\Omega = A \cdot \mathfrak{Q}$* .

Remark 19: (i) If (all) the above- and below-mentioned entropies are used for probability vectors $\mathbf{Q} \in \mathbb{S}^K$ — i.e. one employs $\mathcal{E}(\mathbf{Q})$ — then typically the components q_k of \mathbf{Q} represent a genuine probability mass (frequency) $q_k = \mathbb{P}[\{d_k\}]$ of some data point (state) d_k . However, $\mathbf{Q} \in \mathbb{S}^K$ may alternatively be *artificially* generated. For instance, for the purpose of fault detections of mechanical drives, [298] use Renyi entropies where the q_k 's are normalized squared energy-describing coefficients of the wavelet packet transform of measured vibration records. Another exemplary “artificial” operation is concatenation.

(ii) An analogous statement holds for the employment of (all) the above- and below-mentioned divergences $D(\mathbf{Q}, \mathbb{P})$ — and their transformations — between genuine respectively artificially generated probability vectors $\mathbf{Q}, \mathbb{P} \in \mathbb{S}^K$.

To proceed with our general investigations, let us mention that the remaining parameter cases $\gamma = 0$ and $\gamma = 1$ can be treated analogously. For the sake of brevity, we only deal with the latter. For this, let $\mathbf{Q} \in A \cdot \mathbb{S}^K$ with $A := \sum_{k=1}^K q_k > 0$ and $\mathbb{P} := \mathbb{P}^{unif}$. Clearly, $I(\mathbf{Q}, \mathbb{P}^{unif}) - A \cdot \log K = \sum_{k=1}^K q_k \cdot \log(q_k)$; thus the latter is BS-minimizable on $\Omega = A \cdot \mathfrak{Q}$. More generally, for any continuous strictly increasing (respectively strictly decreasing) function $h_1 :]-1, \infty[\mapsto \mathbb{R}$, the quantity $h_1 \left(\sum_{k=1}^K q_k \cdot \log(q_k) \right)$ is BS-minimizable on $\Omega = A \cdot \mathfrak{Q}$ (respectively BS-maximizable on $\Omega = A \cdot \mathfrak{Q}$). Important special cases are:

(E7) $h_1(y) := h_1^{-Id}(y) = -y$: the entropy

$$\mathcal{E}^{Sh}(\mathbf{Q}) := h_1^{-Id}\left(\sum_{k=1}^K q_k \cdot \log(q_k)\right) = -\sum_{k=1}^K q_k \cdot \log(q_k) \quad (84)$$

is BS-maximizable on $\Omega = A \cdot \mathfrak{Q}$; the special subcase $A = 1$ in (84) (and thus, $\mathbf{Q} = \mathbf{Q}$ is a probability vector) corresponds to the omnipresent *Shannon entropy* [41]; hence, by our bare-simulation approach we can particularly tackle *maximum entropy problems* on almost arbitrary sets \mathfrak{Q} of probability vectors. Analogously, we can treat $\frac{1}{\log(K)} \cdot \mathcal{E}^{Sh}(\mathbf{Q})$ which is called *Pielou's evenness index* [299], and $1 - \frac{1}{\log(K)} \cdot \mathcal{E}^{Sh}(\mathbf{Q}) \in [0, 1]$ which is sometimes used as *clonality (clonotype diversity) index* (see e.g. [300] and [301] (with supplementary private communication)). As a further example for Remark 19, [302] uses q_k 's which are normalized squared coefficients of an orthogonal wavelet decomposition, and accordingly, $\frac{1}{\log(K)} \cdot \mathcal{E}^{Sh}(\mathbf{Q})$ can be interpreted as the entropy of the distribution of energy of oscillations at various frequency and time scales. Some further exemplary studies and applications of the maximization of $\mathcal{E}^{Sh}(\mathbf{Q})$ — aside from the vast physics literature — appear e.g. in [303]–[314].

(E8) $s \in]0, 1[\cup]1, \infty[$, $h_1(y) := h_s^{SM2}(y) := \frac{e^{(s-1)y} - 1}{1-s}$ (cf. (76)) with $y \in \mathbb{R}$: the entropy

$$\mathcal{E}^{SM2}(\mathbf{Q}) := h_s^{SM2}\left(\sum_{k=1}^K q_k \cdot \log(q_k)\right) = \frac{1}{1-s} \cdot \left(\exp\left\{(s-1) \cdot \sum_{k=1}^K q_k \cdot \log(q_k)\right\} - 1\right) \quad (85)$$

is BS-maximizable on $\Omega = A \cdot \mathfrak{Q}$; the special subcase $A = 1$ in (85) (and thus, $\mathbf{Q} = \mathbf{Q}$ is a probability vector) corresponds to the (second type) *entropy of Sharma & Mittal* [45] in the scaling of e.g. [12] (p.20).

IX. FURTHER IMPORTANT DETERMINISTIC OPTIMIZATION PROBLEMS

In order to support the importance and universality of information-theoretic methods for other research fields, let us show how our BS method can be used to tackle — in a new way — important classes of (say) deterministic optimization problems, which are not directly connected to divergences *at first sight*.

For instance, as a consequence of the above Subsection VI-B in combination with the solved cases in Section XII, by retransformation we can even apply our BS method to optimizations of nonnegative *linear* objective functions with constraint sets on Euclidean γ -norm spheres. Indeed, for nonnegative $\check{\mathbf{Q}} := (\check{q}_1, \dots, \check{q}_K)$ and $\check{\mathbf{P}} := (\check{p}_1, \dots, \check{p}_K)$ one can rewrite their scalar product as γ -order Hellinger integrals

$$\sum_{k=1}^K \check{q}_k \cdot \check{p}_k = c_1 \cdot \sum_{k=1}^K q_k^\gamma \cdot p_k^{1-\gamma} = c_1 \cdot H_\gamma(\mathbf{Q}, \mathbb{P}) \quad \text{where} \quad (86)$$

$$\gamma \in]0, 1[\cup]2, \infty[\quad \text{if } \check{\mathbf{Q}} \in [0, \infty[^K, \check{\mathbf{P}} \in]0, \infty[^K \quad \text{respectively} \quad \gamma \in]-\infty, 0[\quad \text{if } \check{\mathbf{Q}} \in]0, \infty[^K, \check{\mathbf{P}} \in]0, \infty[^K,$$

$$q_k := \check{q}_k^{1/\gamma}, \quad p_k := \frac{\check{p}_k^{1/(1-\gamma)}}{\sum_{i=1}^K \check{p}_i^{1/(1-\gamma)}}, \quad c_1 := \left(\sum_{i=1}^K \check{p}_i^{1/(1-\gamma)}\right)^{1-\gamma} =: \|\check{\mathbf{P}}\|_{1/(1-\gamma)}. \quad (87)$$

The required constraint $\sum_{k=1}^K q_k = A > 0$ retransforms to $\|\check{\mathbf{Q}}\|_{1/\gamma} = A^\gamma$ and thus, $\check{\mathbf{Q}}$ must lie on (the positive/nonnegative part of) the $\frac{1}{\gamma}$ -norm-sphere $\partial B_{1/\gamma}(0, A^\gamma)$ around the origin with radius A^γ . Accordingly, for $\gamma \in [2, \infty[$ we have

$$\inf_{\check{\mathbf{Q}} \in \check{\Omega}} \sum_{k=1}^K \check{q}_k \cdot \check{p}_k = c_1 \cdot \inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} H_\gamma(\mathbf{Q}, \mathbb{P}) \quad (88)$$

and we can apply (141) of Proposition 31(a) respectively Proposition 32(a) below ²², as long as the original constraint set $\check{\Omega} \in \partial B_{1/\gamma}(0, A^\gamma) \cap [0, \infty[^K$ transforms (via $q_k = \check{q}_k^{1/\gamma}$) into a constraint set $A \cdot \mathfrak{Q}$ which satisfies the regularity assumption (7) in the relative topology (as a side remark, notice that $\text{int}(\partial B_{1/\gamma}(0, A^\gamma)) = \emptyset$ in the full topology). For the case $\gamma \in]-\infty, 0[$ we also have (88) and apply (141) of Proposition 29(a) for any original constraint set $\check{\Omega} \in \partial B_{1/\gamma}(0, A^\gamma) \cap]0, \infty[^K$ which transforms into $A \cdot \mathfrak{Q}$ satisfying (7) in the relative topology. In contrast, for the case $\gamma \in]0, 1[$ we get

$$\sup_{\check{\mathbf{Q}} \in \check{\Omega}} \sum_{k=1}^K \check{q}_k \cdot \check{p}_k = c_1 \cdot \sup_{\mathbf{Q} \in A \cdot \mathfrak{Q}} H_\gamma(\mathbf{Q}, \mathbb{P})$$

and apply (149) of Proposition 30(a) for any original constraint set $\check{\Omega} \in \partial B_{1/\gamma}(0, A^\gamma) \cap [0, \infty[^K$ which transforms into $A \cdot \mathfrak{Q}$ satisfying (7) in the relative topology.

²²here and analogously henceforth, by this we mean the condition (141) as it appears in the Proposition 31(a) respectively Proposition 32(a)

Analogously to Remark 10 in the previous Section V (where we have dealt with constraints sets Ω of considerably different topological nature than here) we can also principally tackle all the optimization problems of Subsection VI-B, the Sections VII,VIII and the upper part of this Section IX by basically *only employing a fast and accurate — pseudo, true, natural, quantum — random number generator*, provided that the constraint set $A \cdot \Omega$ satisfies the mild assumptions (7) (in the relative topology) and (9). Recall that $A > 0$ (and for φ_2 even $A \in \mathbb{R} \setminus \{0\}$) and that $\mathbf{Q} \in A \cdot \Omega$ implies in particular the constraint $\sum_{k=1}^K q_k = A$. The regularity assumption (7) allows for e.g. high-dimensional constraint sets $A \cdot \Omega$ which are *non-convex* and even *highly disconnected*, and for which other minimization methods (e.g. pure enumeration, gradient or steepest descent methods, etc.) may be problematic or intractable. For example, (7) covers kind of “ K –dimensional (not necessarily regular) polka dot pattern type” relaxations $A \cdot \Omega := \dot{\bigcup}_{i=1}^N \mathcal{U}_i(\mathbf{Q}_i^{dis})$ of finite discrete constraint sets $A \cdot \Omega^{dis} := \{\mathbf{Q}_1^{dis}, \dots, \mathbf{Q}_N^{dis}\}$ of high cardinality N (e.g. being exponential or factorial in a large K), where each K –dimensional vector \mathbf{Q}_i^{dis} has total-sum-of-components equal to A and is surrounded by some small (“flat”, i.e. in the relative topology) neighborhood $\mathcal{U}_i(\mathbf{Q}_i^{dis})$. For the sake of brevity, in the following discussion we confine ourselves to the deterministic setup (e.g. Proposition 32(a) rather than (b)) which particularly involves $\mathbb{I}[\cdot]$ (rather than $\mathbb{I}_{X_1^+}[\cdot]$) and $\xi_n^{w\mathbf{W}}$ (rather than $\xi_{n,X}^{w\mathbf{W}}$). In such a context, all the optimization problems of Subsection VI-B, the Sections VII,VIII and the upper part of this Section IX — subsumed as (cf. (1) to (3))

$$\inf_{\mathbf{Q} \in A \cdot \Omega} \Phi(\mathbf{Q}) \quad \text{respectively} \quad \sup_{\mathbf{Q} \in A \cdot \Omega} \Phi(\mathbf{Q}) \quad \text{—}$$

can be regarded as a “BS-tractable” *relaxations* of the nonlinear discrete (e.g. integer, combinatorial) programming problems

$$\inf_{\mathbf{Q} \in A \cdot \Omega^{dis}} \Phi(\mathbf{Q}) \quad \text{respectively} \quad \sup_{\mathbf{Q} \in A \cdot \Omega^{dis}} \Phi(\mathbf{Q});$$

as examples take e.g. $\Phi(\mathbf{Q}) = c_1 \cdot \left(\left(\sum_{k=1}^K q_k^\gamma \right)^{c_2} - c_3 \right)$ (with $\gamma \in \mathbb{R} \setminus \{0, 1\}$) or $\Phi(\mathbf{Q}) = \Phi_{\mathbb{P}}(\mathbf{Q}) = D_{\tilde{\mathbf{C}}, \varphi_\gamma}(\mathbf{Q}, \mathbb{P})$. For instance, $A \cdot \Omega^{dis}$ may contain only K –dimensional vectors \mathbf{Q}_i^{dis} ($i = 1, \dots, N$) whose components stem from a finite set \mathcal{B} of nonnegative integers and add up to A . If $\mathcal{B} = \{0, 1\}$, then we can even deal with nonnegative *linear* objective functions $\Phi(\mathbf{Q}) = \sum_{k=1}^K \check{p}_k \cdot q_k$ where $\mathbf{Q} := (q_1, \dots, q_K)$ with $q_k \in \{0, 1\}$ and $\check{\mathbf{P}} := (\check{p}_1, \dots, \check{p}_K)$ has components $\check{p}_k > 0$ which reflect e.g. the cost associated with the k –th state. Indeed, by noticing that $q_k^\gamma = q_k$ for $\gamma \in]0, 1[\cup]2, \infty[$, we can employ (86) and (87) to end up with

$$\begin{aligned} \inf_{\mathbf{Q} \in A \cdot \Omega^{dis}} \sum_{k=1}^K q_k \cdot \check{p}_k &= \|\check{\mathbf{P}}\|_{1/(1-\gamma)} \cdot \inf_{\mathbf{Q} \in A \cdot \Omega^{dis}} \sum_{k=1}^K q_k^\gamma \cdot \left(\frac{\check{p}_k^{1/(1-\gamma)}}{\sum_{i=1}^K \check{p}_i^{1/(1-\gamma)}} \right)^{1-\gamma} = \|\check{\mathbf{P}}\|_{1/(1-\gamma)} \cdot \inf_{\mathbf{Q} \in A \cdot \Omega^{dis}} H_\gamma(\mathbf{Q}, \mathbb{P}), \quad \gamma \in [2, \infty[, \quad (89) \\ \sup_{\mathbf{Q} \in A \cdot \Omega^{dis}} \sum_{k=1}^K q_k \cdot \check{p}_k &= \|\check{\mathbf{P}}\|_{1/(1-\gamma)} \cdot \sup_{\mathbf{Q} \in A \cdot \Omega^{dis}} H_\gamma(\mathbf{Q}, \mathbb{P}), \quad \gamma \in]0, 1[. \quad (90) \end{aligned}$$

The corresponding relaxations are

$$\inf_{\mathbf{Q} \in A \cdot \Omega} \sum_{k=1}^K q_k \cdot \check{p}_k = \|\check{\mathbf{P}}\|_{1/(1-\gamma)} \cdot \inf_{\mathbf{Q} \in A \cdot \Omega} H_\gamma(\mathbf{Q}, \mathbb{P}), \quad \gamma \in [2, \infty[, \quad (91)$$

$$\sup_{\mathbf{Q} \in A \cdot \Omega} \sum_{k=1}^K q_k \cdot \check{p}_k = \|\check{\mathbf{P}}\|_{1/(1-\gamma)} \cdot \sup_{\mathbf{Q} \in A \cdot \Omega} H_\gamma(\mathbf{Q}, \mathbb{P}), \quad \gamma \in]0, 1[; \quad (92)$$

for (91) we can apply (141) of Proposition 31(a) respectively Proposition 32(a), whereas for (92) we apply (149) of Proposition 30(a) — as long as the relaxation constraint set $A \cdot \Omega$ satisfies (7) in the relative topology. For the sake of illustration, let us consider a sum-minimization-type *linear assignment problem* with side constraints (for a comprehensive book on assignment problems see e.g. [315]). Suppose that there are K individuals (people, machines, etc.) to carry out K tasks (jobs, etc.). Each individual is assigned to carry out exactly one task. There is cost $c_{ij} > 0$ if individual i is assigned to (i.e., carries out) task j . We want to find the minimum total cost amongst all assignments. There may be side constraints, e.g. each assignment has a value $v_{ij} > 0$ and the total value of the assignment should be above a pregiven threshold. As usual, the problem can be formulated with the help of binary variables x_{ij} where $x_{ij} = 1$ if individual i is assigned to task j , and $x_{ij} = 0$ otherwise.

Accordingly, we want to compute

$$\inf_{K \times K\text{-matrices } x=(x_{ij})} \sum_{i=1}^K \sum_{j=1}^K c_{ij} \cdot x_{ij} \quad (93)$$

subject to

$$\sum_{j=1}^K x_{ij} = 1 \quad \text{for all } i \in \{1, \dots, K\}, \quad (\text{i.e. each individual } i \text{ does one task}), \quad (94)$$

$$\sum_{i=1}^K x_{ij} = 1 \quad \text{for all } j \in \{1, \dots, K\}, \quad (\text{i.e. each task } j \text{ is done by one individual}), \quad (95)$$

$$x_{ij} \in \{0, 1\} \quad \text{for all } i \in \{1, \dots, K\}, j \in \{1, \dots, K\}, \quad (96)$$

$$\text{side (i.e. additional) constraints on } x = (x_{ij})_{i,j=1,\dots,K}. \quad (97)$$

Of course, this can be equivalently rewritten in terms of K^2 -dimensional vectors as follows: let $\mathbf{Q} := (q_1, \dots, q_{K^2})$ and $\check{\mathbf{P}} := (\check{p}_1, \dots, \check{p}_{K^2})$ be such that $c_{ij} = \check{p}_{(i-1) \cdot K + j}$ and $x_{ij} = q_{(i-1) \cdot K + j}$ for $i, j \in \{1, \dots, K\}$ and compute

$$\inf_{\mathbf{Q} \in K \cdot \mathfrak{Q}^{dis}} \sum_{k=1}^{K^2} q_k \cdot \check{p}_k \quad (98)$$

where $K \cdot \mathfrak{Q}^{dis} \subset \mathbb{R}^{K^2}$ is the set of all vectors $\mathbf{Q} = (q_1, \dots, q_{K^2})$ which satisfy the constraints

$$\sum_{j=1}^K q_{(i-1) \cdot K + j} = 1 \quad \text{for all } i \in \{1, \dots, K\}, \quad (99)$$

$$\sum_{i=1}^K q_{(i-1) \cdot K + j} = 1 \quad \text{for all } j \in \{1, \dots, K\}, \quad (100)$$

$$q_k \in \{0, 1\} \quad \text{for all } k \in \{1, \dots, K^2\}, \quad (101)$$

$$\text{side constraints on } \mathbf{Q}. \quad (102)$$

As seen above, this can be rewritten as γ -order Hellinger-integral minimization problem (89), with $\gamma \geq 2$. We can obtain a *highly disconnected* “non-void-interior-type” relaxation of the binary integer programming problem (98) to (102) by replacing (101) with

$$q_k \in [0, \varepsilon_1] \cup [1 - \varepsilon_2, 1] \quad \text{for all } k \in \{1, \dots, K^2\}, \quad (103)$$

for some (possibly arbitrarily) small $\varepsilon_1, \varepsilon_2 > 0$ with $\varepsilon_1 + \varepsilon_2 < 1$. We denote by $K \cdot \mathfrak{Q}$ the outcoming set manifested by the constraints (99), (100), (102) and (103), and accordingly we end up with a minimization problem of type (91), which we can tackle by (141) of Proposition 31(a) respectively Proposition 32(a), as long as (7) (in the relative topology) is satisfied. For instance, we can take $\gamma = 2$ and basically solve the corresponding optimization problem by basically simulating K^2 -dimensional Gaussian random variables (even though the cardinality of $K \cdot \mathfrak{Q}^{dis}$ may be high). As a side remark, let us mention that our relaxation (103) contrasts considerably to the frequently used continuous *linear programming (LP) relaxation*

$$q_k \in [0, 1] \quad \text{for all } k \in \{1, \dots, K^2\}.$$

Let us finally mention that an important special case of a minimization problem (93) to (97) is — the integer programming formulation of — the omnipresent (asymmetric) *traveling salesman problem (TSP)* with possible side constraints²³. There, one has K cities and the cost of traveling from city i to city $j \neq i$ is given by $c_{ij} > 0$. Moreover, one sets $x_{ij} = 1$ if the traveler goes directly from city i to city j (in that order), and $x_{ij} = 0$ otherwise. For technical reasons, for $i = j$ we attribute a cost $c_{ii} > 0$ (e.g. hotel costs), but we require that always $x_{ii} = 0$ which we subsume as the first part of the constraints (97). Then, the constraint (94) means that the traveler leaves from city i exactly once, whereas (95) reflects that the traveler arrives at city j exactly once. The goal is to find a directed tour — i.e. a directed cycle/circuit that visits all K cities once — of minimum cost. Within this context, the second part of the constraints (97) should basically exclude solutions which consist of disconnected subtours (subtour elimination constraints (of e.g. the seminal Dantzig et al. [319]), connectivity constraints, cut-set constraints). Here, we also allow for additional/side constraints which we subsume as the third part (97) of the constraints. Hence, our above-mentioned considerations open the gate to *principally tackle* such kind of TSP problems with our BS method.

For sum-*maximization*-type linear assignment problems with side constraints, where e.g. c_{ij} is a profit (rather than a cost) and the ultimate goal is total profit maximization, we can proceed analogously, by employing (90),(92) (instead of (89),(91)).

²³see e.g. [316]–[318] for comprehensive books on TSP, its variations and its applications to logistics, machine scheduling, printed circuit board drilling, communication-network frequencying, genome sequencing, data clustering, and many others.

Another line of applications of our method to deterministic optimization problems is that the BS minimizability of (77) with $\gamma = 2$ and $c_1 \cdot c_2 > 0$ (see also (145) of Proposition 32(a) below) can be employed to solve the following discrete Monge-Kantorovich-type optimal mass transportation problem (optimal coupling problem) with *side (i.e. additional) constraints*: given two nonnegative-entries vectors $\boldsymbol{\mu} := (\mu_1, \dots, \mu_{K_1}) \in [0, \infty]^{K_1}$ and $\boldsymbol{\nu} := (\nu_1, \dots, \nu_{K_2}) \in [0, \infty]^{K_2}$ with equal total “mass” $\sum_{k=1}^{K_1} \mu_k = \sum_{k=1}^{K_2} \nu_k = A > 0$, compute

$$\begin{aligned} & \inf_{K_1 \times K_2\text{-matrices } \pi} K_1 \cdot K_2 \cdot \sum_{u=1}^{K_1} \sum_{v=1}^{K_2} \left(\pi_{u,v} - \frac{1}{K_1 \cdot K_2} \right)^2 \quad \text{subject to} \\ & \sum_{v=1}^{K_2} \pi_{u,v} = \mu_u \quad \text{for all } u \in \{1, \dots, K_1\}, \quad \sum_{u=1}^{K_1} \pi_{u,v} = \nu_v \quad \text{for all } v \in \{1, \dots, K_2\}, \\ & \pi_{u,v} \in [0, A] \quad \text{for all } u \in \{1, \dots, K_1\}, v \in \{1, \dots, K_2\}, \quad \text{side constraints on } \pi, \boldsymbol{\mu}, \boldsymbol{\nu}. \end{aligned}$$

Indeed, this problem can be equivalently rewritten in terms $K_1 \cdot K_2$ -dimensional vectors as follows: given two nonnegative-entries vectors $\boldsymbol{\mu}, \boldsymbol{\nu}$ as above, compute

$$\inf_{\mathbf{Q} \in \Omega} K_1 \cdot K_2 \cdot \sum_{k=1}^{K_1 \cdot K_2} \left(q_k - \frac{1}{K_1 \cdot K_2} \right)^2 = \inf_{\mathbf{Q} \in \Omega} K_1 \cdot K_2 \cdot \sum_{k=1}^{K_1 \cdot K_2} q_k^2 + 1 - 2A \quad (104)$$

where $\Omega \subset \mathbb{R}^{K_1 \cdot K_2}$ is the set of all vectors $\mathbf{Q} = (q_1, \dots, q_{K_1 \cdot K_2})$ which satisfy the constraints

$$\begin{aligned} & \sum_{j=1}^{K_2} q_{(i-1) \cdot K_2 + j} = \mu_i \quad \text{for all } i \in \{1, \dots, K_1\}, \quad \sum_{i=1}^{K_1} q_{(i-1) \cdot K_2 + j} = \nu_j \quad \text{for all } j \in \{1, \dots, K_2\}, \\ & q_k \in [0, A] \quad \text{for all } k \in \{1, \dots, K_1 \cdot K_2\}, \\ & \text{side constraints on } \mathbf{Q}, \boldsymbol{\mu}, \boldsymbol{\nu}. \end{aligned} \quad (105)$$

Clearly, via divisions by A , one can equivalently rewrite $\Omega = A \cdot \tilde{\Omega}$ for some $\tilde{\Omega} \subset \mathbb{S}^{K_1 \cdot K_2}$ in the $K_1 \cdot K_2$ -dimensional probability simplex. Hence, we can employ (77) (see also (145)) with $K = c_1 = K_1 \cdot K_2$, $\gamma = 2$, $c_2 = 1$ and $c_3 = \frac{2A-1}{K_1 \cdot K_2}$, provided that the side constraints (105) are such that $\tilde{\Omega}$ satisfies the regularity property (7) and the finiteness property (9). Notice that (104) is equal to $\inf_{\mathbf{Q} \in \Omega} D_{2, \varphi_2}(\mathbf{Q}, \mathbb{P}^{unif})$ where $\mathbb{P}^{unif} := (\frac{1}{K_1 \cdot K_2}, \dots, \frac{1}{K_1 \cdot K_2})$ is the probability vector of frequencies of the uniform distribution on $\{1, \dots, K_1 \cdot K_2\}$, and $\tilde{c} = 2$. The special case $A = 1$ with side constraint (105) of the form $K_1 \cdot \min_{i \in \{1, \dots, K_1\}} \mu_i + K_2 \cdot \min_{j \in \{1, \dots, K_2\}} \nu_j \geq 1$ was explicitly solved by e.g. [320] [321], who also give applications to cryptographic guessing problems (spy problems), task partitioning and graph clustering.

The importance of the case $\gamma = 2$ stems also from the fact that one can equivalently rewrite *separable quadratic minimization problems* as minimization problems of Pearson chi-square divergences. Indeed, one can straightforwardly derive that

$$\inf_{\mathbf{Q} \in \tilde{\Omega}} \sum_{k=1}^K (c_{1,k} + c_{2,k} \cdot \check{q}_k + c_{3,k} \cdot \check{q}_k^2), \quad c_{1,k} \in \mathbb{R}, c_{2,k} \in \mathbb{R} \setminus \{0\}, c_{3,k} \in]0, \infty[, \quad (106)$$

is equal to (recall that $\varphi_2(t) := \frac{(t-1)^2}{2}$, cf. (40))

$$c_4 + \inf_{\mathbf{Q} \in \tilde{\Omega}} D_{\varphi_2}(\mathbf{Q}, \mathbf{P}), \quad (107)$$

where $\mathbf{Q} := (q_1, \dots, q_K)$ with $q_k := -c_{2,k} \cdot \check{q}_k$, $\mathbf{P} := (p_1, \dots, p_K)$ with $p_k := \frac{c_{2,k}^2}{2 \cdot c_{3,k}} > 0$, $c_4 := \sum_{k=1}^K (c_{1,k} - \frac{c_{2,k}^2}{4 \cdot c_{3,k}})$, and $\tilde{\Omega}$ is the corresponding reformulation of the constraint set $\tilde{\Omega}$. To achieve the applicability of our BS method, we further transform (107) into its equal form (cf. (12),(13))

$$c_4 + \inf_{\tilde{\mathbf{Q}} \in \tilde{\Omega}/M_{\mathbf{P}}} D_{M_{\mathbf{P}} \cdot \varphi_2}(\tilde{\mathbf{Q}}, \tilde{\mathbf{P}}) \quad (108)$$

with $M_{\mathbf{P}} := \sum_{k=1}^K p_k > 0$ and $\tilde{\mathbf{P}} := \mathbf{P}/M_{\mathbf{P}}$. If $\tilde{\Omega}/M_{\mathbf{P}}$ satisfies (7) and (9) (e.g. it may be highly disconnected), then we can apply Theorem 8. In contrast, if $\tilde{\Omega}/M_{\mathbf{P}} = A \cdot \tilde{\Omega}$ for some $A \in \mathbb{R} \setminus \{0\}$ and some $\tilde{\Omega} \subset \mathbb{S}_{\geq 0}^K$ satisfying (7), then we can apply Theorem 12, Remark 13(vi), Lemma 14(a) (see also Proposition 32(a) below) together with Remark 16(c); for instance, this may appear if $\tilde{\Omega}$ contains (amongst others) the original constraint $\sum_{k=1}^K \check{q}_k = C$ for some constant $C > 0$, and $c_{2,k} = c_2$ does not depend on k , which leads to the choice $A = -\frac{c_2 \cdot C}{M_{\mathbf{P}}}$. Notice that $A < 0$ if $c_2 > 0$. For example, optimization problems (106) with $c_{1,k} > 0$, $c_{2,k} > 0$, $c_{3,k} > 0$ and constraints $\sum_{k=1}^K \check{q}_k = C$, $\check{q}_k \in [\check{q}_k^-, \check{q}_k^+]$ appear in distributed energy management as *economic dispatch problems* in smart grids of power generators, where \check{q}_k is the active power generation of the k -th generator, C is the total power demand, \check{q}_k^- resp. \check{q}_k^+ represent the lower resp. upper bound of the k -th generator’s output, and the cost of power generation is $c_{1,k} + c_{2,k} \cdot \check{q}_k + c_{3,k} \cdot \check{q}_k^2$ (cf. e.g. [322]–[325]). Another important special case of (106) to (108) is the omnipresent L_2 -minimization; indeed, with the choices $c_{3,k} = 1$, $c_{2,k} = -2v_k$, and $c_{1,k} = v_k^2$ for some $\mathbf{V} = (v_1, \dots, v_K)$, the minimization problem (106) is nothing but $\inf_{\tilde{\mathbf{Q}} \in \tilde{\Omega}} \|\tilde{\mathbf{Q}} - \mathbf{V}\|_2^2$; if $\tilde{\Omega}$ depends on a pregiven L -dimensional vector \mathbf{x} (with $L < K$), this can be regarded as a *non-parametric regression problem* in a wide sense.

X. ESTIMATORS

We demonstrate how one can *principally* implement our BS approach; further, deeper analyses are given in a follow-up paper.

A. Estimators for the deterministic minimization problem

We address the minimization problem

$$D_\varphi(\Omega, \mathbf{P}) := \inf_{\mathbf{Q} \in \Omega} D_\varphi(\mathbf{Q}, \mathbf{P}) = \inf_{\tilde{\mathbf{Q}} \in \tilde{\Omega}} D_{\tilde{\varphi}}(\tilde{\mathbf{Q}}, \tilde{\mathbb{P}}) =: D_{\tilde{\varphi}}(\tilde{\Omega}, \tilde{\mathbb{P}}) \quad \text{with } \tilde{\Omega} := \Omega/M_{\mathbf{P}} \quad (\text{cf. (8) and (13)}), \quad (109)$$

whose numerical solution is based on Theorem 8 which basically states that for large integer $n \in \mathbb{N}$ one has

$$\inf_{\mathbf{Q} \in \Omega} D_\varphi(\mathbf{Q}, \mathbf{P}) \approx -\frac{1}{n} \log \mathbb{P}[\boldsymbol{\xi}_n^{\tilde{\mathbf{W}}} \in \tilde{\Omega}] \quad (110)$$

in terms of $\tilde{\varphi} := M_{\mathbf{P}} \cdot \varphi$ and the random vectors $\boldsymbol{\xi}_n^{\tilde{\mathbf{W}}} = \left(\frac{1}{n} \sum_{i \in I_1^{(n)}} \tilde{W}_i, \dots, \frac{1}{n} \sum_{i \in I_K^{(n)}} \tilde{W}_i \right)$ (cf. (17)) with $n_k := \lfloor n \cdot \tilde{p}_k \rfloor$ leading to the disjoint index blocks $I_1^{(n)} := \{1, \dots, n_1\}$, $I_2^{(n)} := \{n_1 + 1, \dots, n_1 + n_2\}$, \dots , $I_K^{(n)} := \{\sum_{k=1}^{K-1} n_k + 1, \dots, n\}$. Recall that $\tilde{\mathbf{W}} := (\tilde{W}_1, \dots, \tilde{W}_n)$ is a random vector consisting of components \tilde{W}_i which are i.i.d. copies of the random variable \tilde{W} whose distribution is $\mathbb{P}[\tilde{W} \in \cdot] = \tilde{\zeta}[\cdot]$ obeying the representation $\tilde{\varphi}(t) = \sup_{z \in \mathbb{R}} \left(z \cdot t - \log \int_{\mathbb{R}} e^{zy} d\tilde{\zeta}(y) \right)$, $t \in \mathbb{R}$, (cf. (15)). Hence, due to (110), the estimation of $D_\varphi(\Omega, \mathbf{P})$ amounts to the estimation of $\mathbb{P}[\boldsymbol{\xi}_n^{\tilde{\mathbf{W}}} \in \tilde{\Omega}]$. For the rest of this subsection, we assume that $\tilde{\mathbb{P}} \in \mathbb{S}_{>0}^K$, that n is chosen such that all $n \cdot \tilde{p}_k$ are integers (and hence, $n = \sum_{k=1}^K n_k$), and that $\tilde{\Omega} \subset \mathbb{R}^K$ satisfies the regularity property $cl(\tilde{\Omega}) = cl(int(\tilde{\Omega}))$, $int(\tilde{\Omega}) \neq \emptyset$ which implies that the same condition holds for Ω ; moreover, we suppose that $D_{\tilde{\varphi}}(\tilde{\Omega}, \tilde{\mathbb{P}}) \in]0, \infty[$ (and thus $\mathbb{P} \notin cl(\Omega)$). For the ease of the following discussions, we introduce the notations

$$T(\mathbf{x}) := \left(\frac{1}{n_1} \sum_{i \in I_1^{(n)}} x_i, \dots, \frac{1}{n_K} \sum_{i \in I_K^{(n)}} x_i \right) \quad \text{for any } \mathbf{x} := (x_1, \dots, x_n) \in \mathbb{R}^n,$$

as well as \mathfrak{D} for the diagonal matrix with diagonal entries $1/\tilde{p}_1, \dots, 1/\tilde{p}_K$ and null entries off the diagonal. Accordingly, the probability in (110) becomes $\mathbb{P}[\boldsymbol{\xi}_n^{\tilde{\mathbf{W}}} \in \tilde{\Omega}] = \mathbb{P}[T(\tilde{\mathbf{W}}) \in \Lambda]$ where $\Lambda := \tilde{\Omega} \cdot \mathfrak{D}$ is a set of (row) vectors in \mathbb{R}^K which is known/derived from the concrete context. The *naive estimator* $\hat{\Pi}_L^{naive}$ of $\mathbb{P}[\boldsymbol{\xi}_n^{\tilde{\mathbf{W}}} \in \tilde{\Omega}]$ is constructed through the following procedure: simulate independently L copies $\tilde{\mathbf{W}}^{(1)}, \dots, \tilde{\mathbf{W}}^{(L)}$ of the vector $\tilde{\mathbf{W}} := (\tilde{W}_1, \dots, \tilde{W}_n)$, with independent entries under $\tilde{\zeta}$, and define (with a slight abuse of notation) $\hat{\Pi}_L^{naive} := \frac{1}{L} \sum_{\ell=1}^L \mathbf{1}_\Lambda(T(\tilde{\mathbf{W}}^{(\ell)}))$; however this procedure is time costly, since this estimate has a very bad hit rate. Thus, in the following, a so-called “efficient Importance Sampling (IS)” scheme — in the sense of Sadowsky & Bucklew [326] (denoted [SB] hereunder) — is adapted for the sophisticated (i.e. non-naive) estimation of $\mathbb{P}[\boldsymbol{\xi}_n^{\tilde{\mathbf{W}}} \in \tilde{\Omega}]$. For this, we employ the following additional Assumption (OM) on the set $\tilde{\Omega}$:

(OM) For any $\tilde{\omega} \in cl(\tilde{\Omega})$ there exists a vector $\mathbf{x} = (x_1, \dots, x_n) \in]t_-^{sc}, t_+^{sc}[^n$ such that $\tilde{\omega} = \left(\frac{1}{n} \sum_{i \in I_1^{(n)}} x_i, \dots, \frac{1}{n} \sum_{i \in I_K^{(n)}} x_i \right)$, or equivalently, for any $\boldsymbol{\lambda} \in cl(\Lambda)$ there exists a vector $\mathbf{x} = (x_1, \dots, x_n) \in]t_-^{sc}, t_+^{sc}[^n$ such that $\boldsymbol{\lambda} = T(\mathbf{x})$.

In the current setup, Assumption (OM) is for instance feasible if for all $\tilde{\omega} \in cl(\tilde{\Omega})$ there holds for all its components $\tilde{\omega}_k \in]\tilde{p}_k \cdot t_-^{sc}, \tilde{p}_k \cdot t_+^{sc}[=]\frac{n_k}{n} \cdot t_-^{sc}, \frac{n_k}{n} \cdot t_+^{sc}[$ ($k = 1, \dots, K$)²⁴ — which e.g. is always satisfied in the common case $dom(\tilde{\varphi}) = dom(\varphi) =]a, b[=]t_-^{sc}, t_+^{sc}[=]0, \infty[$ (e.g. for the power-divergence generators $\tilde{\varphi} = \tilde{c} \cdot \varphi_\gamma$, $\gamma \leq 0$, cf. Subsections XII-A, XII-E below) together with $cl(\tilde{\Omega}) \in \mathbb{R}_{>0}^K$.

To proceed, for any distribution \tilde{S} on \mathbb{R}^n with support included in the support of the product measure $\tilde{\zeta}^{\otimes n}$ it holds (with a slight abuse of notation)

$$\mathbb{P}[\boldsymbol{\xi}_n^{\tilde{\mathbf{W}}} \in \tilde{\Omega}] = E_{\tilde{\zeta}^{\otimes n}}[\mathbf{1}_\Lambda(T(\tilde{\mathbf{W}}))] = E_{\tilde{S}}\left[\mathbf{1}_\Lambda(T(\tilde{\mathbf{V}})) \cdot \frac{d\tilde{\zeta}^{\otimes n}}{d\tilde{S}}(\tilde{\mathbf{V}})\right]$$

from where the *improved IS estimator* of $\mathbb{P}[\boldsymbol{\xi}_n^{\tilde{\mathbf{W}}} \in \tilde{\Omega}]$ is obtained by sampling L i.i.d. replications $\tilde{\mathbf{V}}^{(1)}, \dots, \tilde{\mathbf{V}}^{(L)}$ of the random vector $\tilde{\mathbf{V}}$ with distribution \tilde{S} and by defining

$$\hat{\Pi}_L^{improved} := \frac{1}{L} \sum_{\ell=1}^L \mathbf{1}_\Lambda(T(\tilde{\mathbf{V}}^{(\ell)})) \cdot \frac{d\tilde{\zeta}^{\otimes n}}{d\tilde{S}}(\tilde{\mathbf{V}}^{(\ell)}) \quad (111)$$

The precise form of the efficient IS distribution \tilde{S}^{opt} relies on the definition of a “dominating point” of Λ , which we characterize in the present context. For $\mathbf{x} := (x_1, \dots, x_n)$ in \mathbb{R}^n we define $I_{\tilde{\mathbf{W}}}(\mathbf{x}) := \sup_{\mathbf{z} \in \mathbb{R}^n} \left(\langle \mathbf{z}, \mathbf{x} \rangle - \log E_{\tilde{\zeta}^{\otimes n}}[\exp(\langle \mathbf{z}, \tilde{\mathbf{W}} \rangle)] \right)$, and for each $\boldsymbol{\lambda} \in cl(\Lambda)$ we let $I(\boldsymbol{\lambda}) := \inf \{ I_{\tilde{\mathbf{W}}}(\mathbf{x}) : T(\mathbf{x}) = \boldsymbol{\lambda} \}$ (notice that the set is non-empty because of (OM)). We call a point $\underline{\boldsymbol{\lambda}} := (\underline{\lambda}_1, \dots, \underline{\lambda}_K)$ a *minimal rate point (mrp)* of Λ if (a) $\underline{\boldsymbol{\lambda}} \in \partial\Lambda$ and (b) $I(\underline{\boldsymbol{\lambda}}) \leq I(\boldsymbol{\lambda})$ for all $\boldsymbol{\lambda} \in \Lambda$. A point $\underline{\boldsymbol{\lambda}}$ is

²⁴since then e.g. we can uniformly take $x_i := \frac{n}{n_k} \cdot \tilde{\omega}_k = \lambda_k \in]\frac{n}{n_k} \cdot \frac{n_k}{n} \cdot t_-^{sc}, \frac{n}{n_k} \cdot \frac{n_k}{n} \cdot t_+^{sc}[=]t_-^{sc}, t_+^{sc}[$ for all $i \in I_k^{(n)}$ ($k = 1, \dots, K$)

called a *dominating point* of Λ if (a) $\underline{\lambda} \in \partial\Lambda$, and (b) $I(\underline{\lambda}) \leq I(\lambda)$ for all $\lambda \in \Lambda$ with attainment, namely there exists a vector $\underline{\mathbf{x}} \in]t_-^{sc}, t_+^{sc}[^n$ such that $I_{\widetilde{\mathbf{W}}}(\underline{\mathbf{x}}) = I(\underline{\lambda})$ with $\underline{\lambda} = T(\underline{\mathbf{x}})$. The characterization of a dominating point $\underline{\lambda}$ is settled in the following

Lemma 20: Suppose that Assumption (OM) holds, and let $\underline{\lambda}$ be a (always existent) mrp of Λ . Then, $\underline{\lambda}$ is a dominating point, and $\inf \{I_{\widetilde{\mathbf{W}}}(\mathbf{x}) : T(\mathbf{x}) = \underline{\lambda}\}$ is reached at some vector $\underline{\mathbf{x}}$ in $]t_-^{sc}, t_+^{sc}[^n$ such that for all $k \in \{1, \dots, K\}$ and all $i \in I_k^{(n)}$ there holds $\underline{x}_i = \underline{\lambda}_k$ and $I_{\widetilde{\mathbf{W}}}(\underline{\mathbf{x}}) = n \cdot \sum_{k=1}^K \tilde{p}_k \cdot \tilde{\varphi}(\underline{\lambda}_k)$.

The proof of Lemma 20 is given in Appendix F. Notice that (OM) implies the *existence* of a dominating point $\underline{\lambda}$, but *uniqueness* may not hold. In the latter case, one can try to proceed as in Theorem 2 of [SB] and the discussion thereafter.

However, we assume now uniqueness of $\underline{\lambda}$; this allows for the identification of \tilde{S}^{opt} . By Theorem 1 of [SB] and Theorem 3.1 of [64], the asymptotically optimal IS distribution \tilde{S}^{opt} is obtained as the Kullback-Leibler-divergence projection of $\tilde{\zeta}^{\otimes n}$ on the set of all probability distributions on \mathbb{R}^n centered at point $\underline{\mathbf{x}}$, whose coordinates are — according to Lemma 20 — functions of the coordinates of $\tilde{\mathbf{Q}} := \underline{\lambda} \cdot \mathcal{D}^{-1}$ such that $T(\underline{\mathbf{x}}) = \tilde{\mathbf{Q}} \cdot \mathcal{D}$.

The above definition of \tilde{S}^{opt} presumes the knowledge of $\underline{\lambda}$, which cannot be assumed (otherwise the minimization problem is solved in advance). The aim of the following construction is to provide a proxy \tilde{S} to \tilde{S}^{opt} , where \tilde{S} is the Kullback-Leibler-divergence projection of $\tilde{\zeta}^{\otimes n}$ on the set of all probability distributions on \mathbb{R}^n centered at some point \mathbf{x}^* which is close to $\underline{\mathbf{x}}$. For this sake, we need to have at hand a *proxy* of $\underline{\lambda}$ or, equivalently, a *preliminary guess* $\tilde{\mathbf{Q}}^*$ of $\tilde{\mathbf{Q}} := \arg \inf_{\tilde{\mathbf{Q}} \in \tilde{\Omega}} \sum_{k=1}^K \tilde{p}_k \cdot \tilde{\varphi}(\tilde{q}_k / \tilde{p}_k)$. This guess is by no means produced in order to provide a direct estimate of $D_{\tilde{\varphi}}(\tilde{\Omega}, \tilde{\mathbb{P}})$ but merely to provide the IS distribution \tilde{S} which in turn leads to a sharp estimate of $D_{\tilde{\varphi}}(\tilde{\Omega}, \tilde{\mathbb{P}})$.

Proxy method 1: in some cases we might have at hand some particular point $\tilde{\mathbf{Q}}^* := (\tilde{q}_1^*, \dots, \tilde{q}_K^*)$ in $\tilde{\Omega}$; the resulting IS distribution \tilde{S} with $\tilde{\mathbf{Q}}$ substituted by $\tilde{\mathbf{Q}}^*$ is not optimal in the sense of [SB], but anyhow produces an estimator with good hitting rate, possibly with a loss in the variance. One possible way to obtain such a point $\tilde{\mathbf{Q}}^*$ in $\tilde{\Omega}$ is to simulate runs of (say) M -variate i.i.d. vectors $\tilde{\mathbf{W}}$ under $\tilde{\zeta}^{\otimes M}$ until the first time where $\xi_M^{\tilde{\mathbf{W}}}$ belongs to $\tilde{\Omega}$; then we set $\tilde{\mathbf{Q}}^* := \xi_M^{\tilde{\mathbf{W}}}$ for the succeeding realization $\tilde{\mathbf{W}}$. Before we proceed, it is useful to mention that the need for a drastic fall in the number of simulation runs pertains for cases when $D_{\tilde{\varphi}}(\tilde{\Omega}, \tilde{\mathbb{P}})$ is large. The following construction is suited to this case, which is of relevance in applications both in optimization and in statistics when choosing between competing models none of which is assumed to represent the true one, but merely less inadequate ones.

Proxy method 2: when $D_{\tilde{\varphi}}(\tilde{\Omega}, \tilde{\mathbb{P}})$ is presumably large, we make use of asymptotic approximation to get a proxy of $\tilde{\mathbf{Q}}$. For this, we define a sampling distribution on \mathbb{R}^K fitted to the divergence through

$$f(\tilde{\mathbf{Q}}) := C \cdot \exp\left(-\sum_{k=1}^K \tilde{p}_k \cdot \tilde{\varphi}(\tilde{q}_k / \tilde{p}_k)\right) = C \cdot \exp\left(-D_{\tilde{\varphi}}(\tilde{\mathbf{Q}}, \tilde{\mathbb{P}})\right) \quad (112)$$

where C is a normalizing constant. Let \mathbf{T} be a K -variate random variable with density f . The distribution of \mathbf{T} given $(\mathbf{T} \in \tilde{\Omega})$ concentrates around $\arg \inf_{\tilde{\mathbf{Q}} \in \tilde{\Omega}} D_{\tilde{\varphi}}(\tilde{\mathbf{Q}}, \tilde{\mathbb{P}})$ when $D_{\tilde{\varphi}}(\tilde{\Omega}, \tilde{\mathbb{P}})$ is large. Indeed, for any $\tilde{\mathbf{Q}} \in \tilde{\Omega}$ denote by $\mathbf{V}_\varepsilon(\tilde{\mathbf{Q}})$ a small neighborhood of $\tilde{\mathbf{Q}}$ in \mathbb{R}^K with radius ε ; clearly, the probability of the event $(\mathbf{T} \in \mathbf{V}_\varepsilon(\tilde{\mathbf{Q}}))$ when restricted to $\tilde{\mathbf{Q}} \in \tilde{\Omega}$ is maximum when $\tilde{\mathbf{Q}} = \tilde{\mathbf{Q}}^*$, where $\tilde{\mathbf{Q}}^*$ is the “dominating point of $\tilde{\Omega}$ ” in the sense that $\tilde{\mathbf{Q}}^* := \underline{\lambda} \cdot \mathcal{D}^{-1} \in \partial\tilde{\Omega}$ is the above-defined transform of the dominating point $\underline{\lambda}$ (assuming uniqueness); a precise argumentation under adequate conditions is postponed to Appendix F. Accordingly, we obtain a proxy $\tilde{\mathbf{Q}}^*$ of $\tilde{\mathbf{Q}}$ by simulating a sequence of independent K -variate random variables \mathbf{T}_1, \dots with distribution (112) until (say) \mathbf{T}_m belongs to $\tilde{\Omega}$ and set $\tilde{\mathbf{Q}}^* := \mathbf{T}_m$.

To proceed with the derivation of the IS sampling distribution \tilde{S} on \mathbb{R}^n , we fix $\tilde{\mathbf{Q}}^* := (\tilde{q}_1^*, \dots, \tilde{q}_K^*)$ to be a proxy of $\tilde{\mathbf{Q}}$ or an initial guess in $\tilde{\Omega}$. As an intermediate step, we construct the probability distribution \tilde{U}_k on \mathbb{R} given by

$$d\tilde{U}_k(v) := \exp\left(\tau_k \cdot v - \Lambda_{\tilde{\zeta}}(\tau_k)\right) d\tilde{\zeta}(v) = \frac{\exp(\tau_k \cdot v)}{MGF_{\tilde{\zeta}}(\tau_k)} d\tilde{\zeta}(v) \quad (113)$$

where $\tau_k \in \text{int}(\text{dom}(MGF_{\tilde{\zeta}}))$ is the unique solution of the equation $\Lambda_{\tilde{\zeta}}'(\tau_k) = \frac{\tilde{q}_k}{\tilde{p}_k} \in]t_-^{sc}, t_+^{sc}[$ and thus — by relation (184) of Appendix E — we can compute explicitly $\tau_k = \tilde{\varphi}'\left(\frac{\tilde{q}_k}{\tilde{p}_k}\right)$. Therefore, \tilde{U}_k is the Kullback-Leibler-divergence projection of $\tilde{\zeta}$ on the class of all probability distributions on \mathbb{R} whose expectation is \tilde{q}_k^* . As a side remark, notice that one possible way of obtaining an explicit form of the probability distribution \tilde{U}_k may be by identification through its moment generating function

$$\text{dom}(MGF_{\tilde{\zeta}}) - \tau_k \ni z \mapsto MGF_{\tilde{U}_k}(z) = \frac{MGF_{\tilde{\zeta}}(z + \tau_k)}{MGF_{\tilde{\zeta}}(\tau_k)}$$

of which all ingredients are principally available; for instance, this will be used in the solved-cases Section XII below. From (113), we define $\tilde{S}_k := \underbrace{\tilde{U}_k \otimes \cdots \otimes \tilde{U}_k}_{n_k \text{ times}}$ for all $k \in \{1, \dots, K\}$, whence

$$d\tilde{S}_k(v_{k,1}, \dots, v_{k,n_k}) = \exp\left(\left(\sum_{i \in I_k^{(n)}} \tau_k \cdot v_{k,i}\right) - n_k \cdot \Lambda_{\tilde{\zeta}}(\tau_k)\right) d\tilde{\zeta}(v_{k,1}) \cdots d\tilde{\zeta}(v_{k,n_k}),$$

which manifests \tilde{S}_k as the Kullback-Leibler-divergence projection of $\underbrace{\tilde{\zeta} \otimes \cdots \otimes \tilde{\zeta}}_{n_k \text{ times}}$ on the class of all probability distributions on \mathbb{R}^k whose expectation vector is $\tilde{\mathbf{Q}}^* = (\tilde{q}_1^*, \dots, \tilde{q}_K^*) \in \mathbb{R}^k$. Let now

$$\tilde{S} := \tilde{S}_1 \otimes \cdots \otimes \tilde{S}_K,$$

which therefore satisfies (recall that $\sum_{k=1}^K n_k = n$)

$$d\tilde{S}(v_{1,1}, \dots, v_{1,n_1}, \dots, v_{K,1}, \dots, v_{K,n_K}) = \exp\left(\sum_{k=1}^K \left(\sum_{i \in I_k^{(n)}} \tau_k \cdot v_{k,i}\right) - n_k \cdot \Lambda_{\tilde{\zeta}}(\tau_k)\right) d\tilde{\zeta}^{\otimes n}(v_{1,1}, \dots, v_{1,n_1}, \dots, v_{K,1}, \dots, v_{K,n_K}). \quad (114)$$

The same procedure with all \tilde{q}_k^* substituted by the coordinates \tilde{q}_k of $\tilde{\mathbf{Q}}$ produces S^{opt} . Therefore, \tilde{S} is a substitute for S^{opt} with the change in the centering from the unknown vector $\tilde{\mathbf{Q}}$ to its proxy $\tilde{\mathbf{Q}}^*$.

As a straightforward consequence of (111) and (114), we obtain the improved IS estimator of $\mathbb{P}[\xi_n^{\tilde{\mathbf{W}}} \in \tilde{\Omega}]$ as

$$\hat{\Pi}_L^{improved} = \frac{1}{L} \sum_{\ell=1}^L \mathbf{1}_{\Lambda}(T(\tilde{\mathbf{V}}^{(\ell)})) \cdot \prod_{k=1}^K IS_k(\tilde{\mathbf{V}}_k^{(\ell)}) \quad (115)$$

where $\tilde{\mathbf{V}}_k^{(\ell)} := \left(\tilde{V}_i^{(\ell)}\right)_{i \in I_k^{(n)}}$ is the k -th block of the ℓ -th replication $\tilde{\mathbf{V}}^{(\ell)}$ of $\tilde{\mathbf{V}}$ under \tilde{S} , and the k -th importance-sampling factor is $\tilde{IS}_k(v_{k,1}, \dots, v_{k,n_k}) := \frac{d\tilde{\zeta}^{\otimes n_k}}{d\tilde{S}_k}(v_{k,1}, \dots, v_{k,n_k}) = \exp\left(n_k \cdot \Lambda_{\tilde{\zeta}}(\tau_k) - \tau_k \cdot \sum_{i=1}^{n_k} v_{k,i}\right)$ with $n_k = \text{card}(I_k^{(n)})$.

Summing up, we arrive at the following algorithm in case that $\tilde{\Omega}$ has a unique dominating point (in the above-defined sense): **Step D1.** Exemplarily, we start with proxy method 2 (the other proxy method 1 works analogously): get a proxy $\tilde{\mathbf{Q}}^*$ of $\tilde{\mathbf{Q}}$ by simulating a sequence of independent K -variate random variables \mathbf{T}_1, \dots with distribution (112) until (say) \mathbf{T}_m belongs to $\tilde{\Omega}$ and set $\tilde{\mathbf{Q}}^* := \mathbf{T}_m$.

Step D2. For all k in $\{1, \dots, K\}$ compute $\tau_k = \tilde{\varphi}'\left(\frac{\tilde{q}_k^*}{p_k}\right)$.

Step D3. For all ℓ in $\{1, \dots, L\}$ perform a run of $\tilde{\mathbf{V}}^{(\ell)}$ under \tilde{S} as follows:

For all k in $\{1, \dots, K\}$ simulate n_k i.i.d. random variables $\tilde{V}_{k_1}^{(\ell)}, \dots, \tilde{V}_{k_{n_k}}^{(\ell)}$ with common distribution \tilde{U}_k defined in (113). Set $\tilde{\mathbf{V}}_k^{(\ell)} := (\tilde{V}_{k_1}^{(\ell)}, \dots, \tilde{V}_{k_{n_k}}^{(\ell)})$ to be the corresponding row vector. Construct $\tilde{\mathbf{V}}^{(\ell)}$ as the row vector obtained by concatenating the

$\tilde{\mathbf{V}}_k^{(\ell)}$, i.e. $\tilde{\mathbf{V}}^{(\ell)} := \left(\tilde{\mathbf{V}}_1^{(\ell)}, \dots, \tilde{\mathbf{V}}_K^{(\ell)}\right)$, and make use of $\hat{\Pi}_L^{improved}$ given in (115) with the τ_k 's obtained in Step D2 above to define (in the light of (109),(110)) the *BS minimum-distance estimator*

$$\widehat{D}_{\varphi}(\Omega, \mathbf{P}) := \widehat{D}_{\tilde{\varphi}}(\tilde{\Omega}, \tilde{\mathbf{P}}) := -\frac{1}{n} \log \hat{\Pi}_L^{improved}. \quad (116)$$

For many cases, the simulation burden needed for the computation of $\hat{\Pi}_L^{improved}$ — and thus of $\widehat{D}_{\varphi}(\Omega, \mathbf{P})$ — can be drastically reduced, especially for high dimensions K and large sample size $n \cdot L$. In fact, in terms of the notations $n_k := \text{card}(I_k^{(n)})$, $\widehat{W}_k^{(\ell)} := \sum_{i \in I_k^{(n)}} \tilde{V}_i^{(\ell)}$ and

$$\widetilde{ISF}_k(x) := \frac{d\tilde{\zeta}^{*n_k}}{d\tilde{U}_k^{*n_k}}(x) = \exp(n_k \cdot \Lambda_{\tilde{\zeta}}(\tau_k) - x \cdot \tau_k) \quad (117)$$

(where $\tilde{\zeta}^{*n_k}$ is the n_k -convolution of the measure $\tilde{\zeta}$), one can rewrite (115) as

$$\hat{\Pi}_L^{improved} = \frac{1}{L} \sum_{\ell=1}^L \mathbf{1}_{\Lambda}\left(\left(\frac{1}{n_1} \widehat{W}_1^{(\ell)}, \dots, \frac{1}{n_K} \widehat{W}_K^{(\ell)}\right)\right) \cdot \prod_{k=1}^K \widetilde{ISF}_k(\widehat{W}_k^{(\ell)}) \quad (118)$$

with K -vector $(\frac{1}{n_1}\widehat{W}_1^{(\ell)}, \dots, \frac{1}{n_K}\widehat{W}_K^{(\ell)})$. Clearly, the random variable $\widehat{W}_k^{(\ell)}$ ($k = 1, \dots, K$) has distribution $\widetilde{U}_k^{*n_k}$. Hence, if $\widetilde{U}_k^{*n_k}$ can be explicitly constructed, then for the computation of $\widehat{\Pi}_L^{improved}$ it suffices to simulate the $K \cdot L$ random variables $\widehat{W}_k^{(\ell)}$ rather than the $n \cdot L$ random variables $\widetilde{V}_i^{(\ell)}$; notice that according to the right-hand side of (117), one can explicitly compute $\widetilde{ISF}_k(\cdot)$ which can be interpreted as *Importance Sampling Factor pertaining to the block k*. In the case that ζ is infinitely divisible, simulation issues may become especially comfortable. The tractability of this reduction effect will be exemplarily demonstrated in the solved-cases Subsections XII-A to XII-F below, for the BS minimization of the important power divergences (for which the infinite divisibility holds).

Let us finally remark that from the above-mentioned Steps D1 to D4 (and analogously S1 to S4 below) one can see that for our BS method we basically need only a fast and accurate — pseudo, true, natural, quantum — random number generator. The corresponding computations can be principally run in parallel, and require relatively moderate computer memory/storage; a detailed discussion is beyond the scope of this paper, given its current length.

B. Estimators for the statistical minimization problem

1) Estimation of $\mathbb{P}_{X_1^n}[\xi_{n,\mathbf{X}}^{w\mathbf{W}} \in \Omega]$:

In relation with Theorem 12 — by making use of the notations in (29),(30),(33) — we start by remarking that the development of the estimator $\widehat{\Pi}_L^{improved}$ works quite analogously to that of $\widehat{\Pi}_L^{improved}$ in the previous Subsection X-A. To make this even more transparent, we label (w.l.o.g.) all random vectors of length n in the same way as above: we sort the already given and thus fixed data X_i 's in such a way that the first n_1 of them share the same value d_1 , and so on, until the last block with length n_K in which the data have common value d_K . With this, analogously to Subsection X-A we apply again an “efficient Importance Sampling (IS)” scheme in the sense of Sadowsky & Buckle [326]. This will involve the simulation of L independent n -tuples $\mathbf{V}^{(\ell)} := (V_n^{(\ell)}, \dots, V_1^{(\ell)})$ with common distribution S on \mathbb{R}^n , such that $\zeta^{\otimes n}$ is (measure-)equivalent with respect to S . In fact, we rewrite $\mathbb{P}_{X_1^n}[\xi_{n,\mathbf{X}}^{w\mathbf{W}} \in \Omega]$ as

$$\mathbb{P}_{X_1^n}[\xi_{n,\mathbf{X}}^{w\mathbf{W}} \in \Omega] = E_S \left[\frac{d\zeta^{\otimes n}}{dS}(V_1, \dots, V_n) \cdot \mathbf{1}_\Omega(\xi_{n,\mathbf{X}}^{w\mathbf{V}}) \right] \quad (119)$$

where S designates any IS distribution of the vector $\mathbf{V} := (V_1, \dots, V_n)$, and $E_S[\cdot]$ denotes the corresponding expectation operation. Notice that S is a *random* probability distribution on \mathbb{R}^n ; in fact, S is a conditional probability distribution given X_1^n , and thus it would be more precise to write $S|X_1^n$ instead of S ; for the sake of brevity, we omit $|X_1^n$. As a consequence of (119), for adequately chosen S , an improved estimator of $\mathbb{P}_{X_1^n}[\xi_{n,\mathbf{X}}^{w\mathbf{W}} \in \Omega]$ is given by

$$\widehat{\Pi}_L^{improved} := \frac{1}{L} \sum_{\ell=1}^L \frac{d\zeta^{\otimes n}}{dS}(V_1^{(\ell)}, \dots, V_n^{(\ell)}) \cdot \mathbf{1}_\Omega(\xi_{n,\mathbf{X}}^{w\mathbf{V}^{(\ell)}}), \quad (120)$$

which by the virtue of (36) also estimates $\inf_{\mathbf{Q} \in \Omega} \inf_{m \neq 0} D_\varphi(m \cdot \mathbf{Q}, \mathbb{P})$ (which we here suppose to be in $]0, \infty[$ as well as $\mathbb{P}_n^{emp} \notin cl(\Omega)$ for large enough n). Let us now deal with the concrete construction of a reasonable S . For this, as above, we need a particular $\mathbf{Q}^* \in int(\Omega)$. This (i) may be given in advance or (ii) it may be achieved by simulation; in the following, we only work with the latter. Indeed, given some (typically) large integer M , we employ a realization $\mathbf{W}^* := (W_1^*, \dots, W_M^*)$ such that $\mathbf{Q}^* := \xi_{M,\mathbf{X}}^{w\mathbf{W}^*} \in int(\Omega)$, by drawing replicates $\mathbf{W} = (W_1, \dots, W_M)$ under $\zeta^{\otimes M}$ until the first time where $\xi_{M,\mathbf{X}}^{w\mathbf{W}}$ belongs to $int(\Omega)$; (only) at this point, for consistency we artificially add $m - 1$ copies of each observed data point and transparently keep the same notations, leading e.g. to $p_{M,k}^{emp} = p_{n,k}^{emp}$ for $M := m \cdot n$. Notice that by the nature of Ω , \mathbf{Q}^* is a probability vector which has the K components

$$q_k^* := \sum_{i=1}^M \frac{W_i^*}{\sum_{j=1}^M W_j^*} \mathbf{1}_{\{d_k\}}(X_i), \quad k = 1, \dots, K. \quad (121)$$

Before we proceed, let us give the substantial remark that changing (V_1, \dots, V_n) drawn under S to $(c \cdot V_1, \dots, c \cdot V_n)$ for any $c \neq 0$ yields $\xi_{n,\mathbf{X}}^{w\mathbf{V}} = \xi_{n,c\mathbf{V}}^{w\mathbf{V}}$ so that the distribution S is not uniquely determined. Amongst all candidates, we choose the — uniquely determined — S which is the Kullback-Leibler-divergence projection of $\zeta^{\otimes n}$ on the set of all probability distributions on \mathbb{R}^n such that the K “non-normalized” moment constraints

$$E_S[\xi_{n,\mathbf{X}}^{\mathbf{V}}] = \xi_{M,\mathbf{X}}^{w\mathbf{W}^*} \quad (122)$$

(rather than the normalized $E_S[\xi_{n,\mathbf{X}}^{w\mathbf{V}}] = \xi_{M,\mathbf{X}}^{w\mathbf{W}^*}$) are satisfied, with the non-normalized vectors

$$\xi_{M,\mathbf{X}}^{w\mathbf{W}^*} := \left(\frac{1}{M} \sum_{j=1}^M W_j^* \right) \cdot \mathbf{Q}^* =: \overline{W}^* \cdot \mathbf{Q}^*, \quad \xi_{n,\mathbf{X}}^{\mathbf{V}} := \left(\frac{1}{n} \sum_{j=1}^n V_j \right) \cdot \xi_{n,\mathbf{X}}^{w\mathbf{V}}.$$

As already indicated above, this projection S is a well-determined unique distribution on \mathbb{R}^n and — as we shall see in Proposition 21 below — it is such that $\xi_{n,\mathbf{X}}^{w\mathbf{V}}$ belongs to Ω with probability bounded away from 0 as n increases, when (V_1, \dots, V_n) are drawn under S . Therefore, this IS distribution produces an estimate of $\mathbb{P}_{X_1^n}[\xi_{n,\mathbf{X}}^{w\mathbf{W}} \in \Omega]$.

In order to justify the above construction of S , we give the following result, which states that the IS sampling distribution S yields a good hitting rate. Its proof will be given in Appendix G.

Proposition 21: With the above definition of S , $\liminf_{n \rightarrow \infty} S \left[\boldsymbol{\xi}_{n, \mathbf{X}}^{w, \mathbf{V}} \in \boldsymbol{\Omega} \right]$ is bounded away from 0.

We now come to the detailed construction of S . The constraints (122) can be written in explicit form as

$$E_S \left[\frac{1}{n_k} \sum_{i \in I_k^{(n)}} V_i \right] = \overline{W}^* \cdot q_k^*, \quad k = 1, \dots, K. \quad (123)$$

The distribution S can be obtained by blocks. Indeed, let us define S^k as the Kullback-Leibler-divergence (KL) projection of $\zeta^{\otimes n_k}$ on the set of all distributions on \mathbb{R}^{n_k} such that (123) holds. We define the resulting S as the product distribution of those S^k 's. To obtain the latter, we start by defining U_k as the KL projection of ζ on the set of all measures Q on \mathbb{R} under (123). Then,

$$dU_k(v) = \exp(\tau_k \cdot v - \Lambda_\zeta(\tau_k)) d\zeta(v),$$

where $\tau_k \in \text{int}(\text{dom}(MGF_\zeta))$ is (under the appropriately adapted Assumption (OM)) the unique solution of the equation $\Lambda'_\zeta(\tau_k) = \overline{W}^* \cdot \frac{q_k^*}{p_{n,k}^*}$ and thus — by relation (184) of Appendix E — we can compute explicitly $\tau_k = \varphi' \left(\frac{\overline{W}^* \cdot q_k^*}{p_{n,k}^*} \right)$. The distribution S^k is then defined by $S^k := \underbrace{U_k \otimes \dots \otimes U_k}_{n_k \text{ times}}$ from which we obtain $S := S^1 \otimes \dots \otimes S^K$. With this construction, it holds

$$\frac{dS}{d\zeta^{\otimes n}}(v_{1,1}, \dots, v_{1,n_1}, \dots, v_{K,1}, \dots, v_{K,n_K}) = \exp \left(\sum_{k=1}^K \sum_{i \in I_k^{(n)}} \left(\tau_k \cdot v_{k,i} - \Lambda_\zeta(\tau_k) \right) \right)$$

which proves that S is indeed the KL projection of $\zeta^{\otimes n}$ we aimed at. Therefore, \mathbf{V} is composed of K independent blocks of length n_k each, and the k -th subvector \mathbf{V}_k consists of all the random variables V_i whose index i satisfies $X_i = d_k$. Within \mathbf{V}_k , all components are i.i.d. with same distribution U_k on \mathbb{R} defined through $\frac{dU_k}{d\zeta}(u) = \exp\{\tau_k \cdot u - \Lambda_\zeta(\tau_k)\} = \frac{\exp\{\tau_k \cdot u\}}{MGF_\zeta(\tau_k)}$, which leads to the moment generating function $\text{dom}(MGF_\zeta) - \tau_k \ni z \mapsto MGF_{U_k}(z) := \int_{\mathbb{R}} e^{zy} dU_k(y) = \frac{MGF_\zeta(z + \tau_k)}{MGF_\zeta(\tau_k)}$. Notice that U_k is a distorted distribution of ζ with the distortion parameter τ_k (in some cases, this distortion even becomes a tilting/dampening). The estimator $\widehat{\Pi}_L^{\text{improved}}$ defined in (120) can be implemented through the following algorithm:

Step S1. Choose some (typically large) M and simulate (in the above-described fashion) repeatedly i.i.d. vectors (W_1, \dots, W_M) — whose independent components have common distribution ζ — until $\boldsymbol{\xi}_{M, \mathbf{X}}^{w, \mathbf{W}}$ belongs to $\boldsymbol{\Omega}$. Call (W_1^*, \dots, W_M^*) the corresponding vector and \overline{W}^* the arithmetic mean of its components. Moreover, denote by $\boldsymbol{\xi}_{M, \mathbf{X}}^{w, \mathbf{W}^*}$ the corresponding normalized weighted empirical measure, identified with the K -component vector $Q^* := (q_1^*, \dots, q_K^*)$ with q_k^* defined in (121).

Step S2. For all $k \in \{1, \dots, K\}$ compute $\tau_k = \varphi' \left(\frac{\overline{W}^* \cdot q_k^*}{p_{n,k}^*} \right)$.

Step S3. For all $\ell \in \{1, \dots, L\}$ simulate independently for all $k \in \{1, \dots, K\}$ a row vector $\mathbf{V}_k^{(\ell)} := (V_{k_1}^{(\ell)}, \dots, V_{k_{n_k}}^{(\ell)})$ with independent components with common distribution U_k defined above. Concatenate these vectors to the row vector $\mathbf{V}^{(\ell)}$.

Step S4. Compute the estimator $\widehat{\Pi}_L^{\text{improved}}$ by making use of the formula (120) which turns into the explicit form

$$\widehat{\Pi}_L^{\text{improved}} = \frac{1}{L} \sum_{\ell=1}^L \exp \left(\sum_{k=1}^K \left(n_k \cdot \Lambda_\zeta(\tau_k) - \tau_k \cdot \sum_{i \in I_k^{(n)}} V_i^{(\ell)} \right) \right) \cdot \mathbf{1}_\Omega \left(\boldsymbol{\xi}_{n, \mathbf{X}}^{w, \mathbf{V}^{(\ell)}} \right). \quad (124)$$

Analogous to the paragraph right after (116), in many cases we may improve the simulation burden needed for the computation of the estimator $\widehat{\Pi}_L^{\text{improved}}$. In fact, in terms of the notations $\widehat{W}_k^{(\ell)} := \sum_{i \in I_k^{(n)}} V_i^{(\ell)}$ we can rewrite (124) as

$$\widehat{\Pi}_L^{\text{improved}} = \frac{1}{L} \sum_{\ell=1}^L \mathbf{1}_\Omega \left(\boldsymbol{\xi}_{n, \mathbf{X}}^{w, \mathbf{V}^{(\ell)}} \right) \cdot \prod_{k=1}^K ISF_k \left(\widehat{W}_k^{(\ell)} \right) \quad (125)$$

$$\text{with } ISF_k(x) := \exp(n_k \cdot \Lambda_\zeta(\tau_k) - x \cdot \tau_k) \quad (126)$$

$$\text{and } \boldsymbol{\xi}_{n, \mathbf{X}}^{w, \mathbf{V}^{(\ell)}} = \begin{cases} \left(\frac{\widehat{W}_1^{(\ell)}}{\sum_{k=1}^K \widehat{W}_k^{(\ell)}}, \dots, \frac{\widehat{W}_K^{(\ell)}}{\sum_{k=1}^K \widehat{W}_k^{(\ell)}} \right), & \text{if } \sum_{k=1}^K \widehat{W}_k^{(\ell)} \neq 0, \\ (\infty, \dots, \infty) =: \infty, & \text{if } \sum_{k=1}^K \widehat{W}_k^{(\ell)} = 0. \end{cases} \quad (127)$$

Clearly, the random variable $\widehat{W}_k^{(\ell)}$ ($k = 1, \dots, K$) has distribution $U_k^{*n_k}$. Hence, if $U_k^{*n_k}$ can be explicitly constructed, then for the computation of $\widehat{\Pi}_L^{\text{improved}}$ it suffices to independently simulate the $K \cdot L$ random variables $\widehat{W}_k^{(\ell)}$ (rather than the $n \cdot L$ random variables $V_i^{(\ell)}$).

2) *Direct BS-estimability of $D_\varphi(\mathfrak{Q}, \mathbb{P})$* :

For the cases $\varphi := \tilde{c} \cdot \varphi_\gamma$ (cf. (40)), by making use of Theorem 12, Lemma 14 (especially (46),(53),(58)) and the results of the previous subsection, we get through inversion the *power-divergence estimators (BS estimators of power divergences)*

$$\begin{aligned} D_{\tilde{c} \cdot \varphi_\gamma}(\widehat{\mathfrak{Q}}, \mathbb{P}) &:= -\frac{\tilde{c}}{\gamma(\gamma-1)} \left\{ 1 - \left(1 + \frac{\gamma}{\tilde{c}} \cdot \frac{1}{n} \cdot \log \widehat{\Pi}_L^{\text{improved}} \right)^{1-\gamma} \right\}, & \gamma \in]-\infty, 0[\cup]0, 1[\cup]2, \infty[, \\ D_{\tilde{c} \cdot \varphi_0}(\widehat{\mathfrak{Q}}, \mathbb{P}) &:= -\frac{1}{n} \log \widehat{\Pi}_L^{\text{improved}}, & \gamma = 0, \\ D_{\tilde{c} \cdot \varphi_1}(\widehat{\mathfrak{Q}}, \mathbb{P}) &:= -\tilde{c} \cdot \log \left(1 + \frac{1}{\tilde{c}} \cdot \frac{1}{n} \cdot \log \widehat{\Pi}_L^{\text{improved}} \right), & \gamma = 1. \end{aligned}$$

For more details including the correspondingly involved simulation distributions ζ , see the solved-cases Section XII below.

3) *BS-estimability of bounds of $D_\varphi(\mathfrak{Q}, \mathbb{P})$* :

For divergence cases which are not directly BS-estimable, we present now the algorithm for the BS evaluation of the bounds

$$\inf_{m \neq 0} D_\varphi(m \cdot \mathfrak{Q}, \mathbb{P}) = \inf_{\mathbb{Q} \in \mathfrak{Q}} D_\varphi(m(\mathbb{Q}) \cdot \mathbb{Q}, \mathbb{P}) \stackrel{(\Delta)}{=} D_\varphi(m(\mathbb{Q}^*) \cdot \mathbb{Q}^*, \mathbb{P}) \leq D_\varphi(\mathfrak{Q}, \mathbb{P}) \leq D_\varphi(\mathbb{Q}^*, \mathbb{P}) \quad (128)$$

obtained in Subsection VI-C, where \mathbb{Q}^* satisfies the above equality (Δ) . The estimator of the lower bound in (128) is

$\widehat{D} := -\frac{1}{n} \log \widehat{\Pi}_L^{\text{improved}}$ defined in (124). We now turn to an estimate of the upper bound. Consider for any fixed $\mathbb{Q} := (q_1, \dots, q_K)$ in $\mathbb{S}_{>0}^K$ the real number $m_n(\mathbb{Q})$ which satisfies $D_\varphi(m_n(\mathbb{Q}) \cdot \mathbb{Q}, \mathbb{P}_n^{\text{emp}}) = \inf_{m \neq 0} D_\varphi(m \cdot \mathfrak{Q}, \mathbb{P}_n^{\text{emp}})$ where $\mathbb{P}_n^{\text{emp}}$ was defined in the course of (29). Such $m_n(\mathbb{Q})$ is well defined for all \mathbb{Q} since it satisfies the equation (in m)

$$\frac{d}{dm} D_\varphi(m \cdot \mathbb{Q}, \mathbb{P}_n^{\text{emp}}) = \sum_{k=1}^K q_k \cdot \varphi' \left(\frac{m \cdot q_k}{p_{n,k}^{\text{emp}}} \right) = 0. \quad (129)$$

Since the mapping $m \rightarrow D_\varphi(m \cdot \mathbb{Q}, \mathbb{P})$ is convex and differentiable, existence and uniqueness of $m_n(\mathbb{Q})$ hold; furthermore, $m_n(\mathbb{Q}) \in \left] \min_k p_{n,k}^{\text{emp}}/q_k, \max_k p_{n,k}^{\text{emp}}/q_k \right]$ since $\frac{d}{dm} D_\varphi(m \cdot \mathbb{Q}, \mathbb{P}_n^{\text{emp}})$ is negative when $m = \min_k p_{n,k}^{\text{emp}}/q_k$ and positive when $m = \max_k p_{n,k}^{\text{emp}}/q_k$. An estimate of the distribution \mathbb{Q}^* is required. This can be achieved as follows:

- Estimate $\inf_{m \neq 0} D_\varphi(m \cdot \mathfrak{Q}, \mathbb{P})$ through $\widehat{D} := -\frac{1}{n} \log \widehat{\Pi}_L^{\text{improved}}$ defined in (124).
- Set $i = 0$.
- Get some $\mathbb{Q}_i := (q_{i,1}, \dots, q_{i,K})$ in \mathfrak{Q} ; this can be obtained by simulating runs of vectors (W_1, \dots) through i.i.d. sampling under ζ . Evaluate $m_n(\mathbb{Q}_i)$ by solving (129) (with $q_{i,k}$ instead of q_k) for m , which is a fast calculation by the bisection method.
- If $D_\varphi(m_n(\mathbb{Q}_i) \cdot \mathbb{Q}_i, \mathbb{P}_n^{\text{emp}}) < \widehat{D} + \eta$ for some small $\eta > 0$, then the proxy of \mathbb{Q}^* is \mathbb{Q}_i , denoted by $\widehat{\mathbb{Q}}^*$.
- Else set $i \leftarrow i + 1$ and get \mathbb{Q}_i in $\mathfrak{Q} \cap \{\mathbb{Q} : D_\varphi(\mathbb{Q}, \mathbb{P}_n^{\text{emp}}) < D_\varphi(\mathbb{Q}_{i-1}, \mathbb{P}_n^{\text{emp}})\}$ and iterate.

That this algorithm converges in the sense that it produces some $\widehat{\mathbb{Q}}^*$ is clear. Since by (128)

$$D_\varphi(m(\mathbb{Q}^*) \cdot \mathbb{Q}^*, \mathbb{P}) \leq D_\varphi(\mathfrak{Q}, \mathbb{P}) \leq D_\varphi(\mathbb{Q}^*, \mathbb{P}),$$

we have obtained both estimated lower and upper bounds for $D_\varphi(\mathfrak{Q}, \mathbb{P})$.

That the upper bound is somehow optimal can be seen from the power-divergence Cases 1 to 6 developed in the solved-cases Section XII below. Indeed, in this context the solution of equation (129) is explicit and produces $m(\mathbb{Q})$ as a function of $D_\varphi(\mathbb{Q}, \mathbb{P})$ through a Hellinger integral, and the mapping $\mathbb{Q} \rightarrow D_\varphi(m(\mathbb{Q}) \cdot \mathbb{Q}, \mathbb{P})$ is increasing with respect to $D(\mathbb{Q}, \mathbb{P})$. Hence, $\mathbb{Q} \rightarrow \inf_{m \neq 0} D_\varphi(m \cdot \mathbb{Q}, \mathbb{P})$ is minimal when $D_\varphi(\mathbb{Q}, \mathbb{P})$ is minimal as $\mathbb{Q} \in \mathfrak{Q}$. Therefore, $\mathbb{Q}^* \in \arg \inf_{\mathbb{Q} \in \mathfrak{Q}} D_\varphi(m(\mathbb{Q}) \cdot \mathbb{Q}, \mathbb{P})$ also satisfies $\mathbb{Q}^* \in \arg \inf_{\mathbb{Q} \in \mathfrak{Q}} D_\varphi(\mathbb{Q}, \mathbb{P})$.

XI. FINDING/CONSTRUCTING THE DISTRIBUTION OF THE WEIGHTS

Recall first that in Theorem 12, one crucial component is the sequence $(W_i)_{i \in \mathbb{N}}$ of weights being i.i.d. copies of a random variable W whose probability distribution is ζ (i.e. $\mathbb{P}[W \in \cdot] = \zeta[\cdot]$), where the latter has to be connected with the divergence generator $\varphi \in \Upsilon(]a, b[)$ through the representation

$$\varphi(t) = \sup_{z \in \mathbb{R}} \left(z \cdot t - \log \int_{\mathbb{R}} e^{zy} d\zeta(y) \right), \quad t \in \mathbb{R}, \quad (\text{cf. (6)})$$

under the additional requirement that the function $z \mapsto MGF_\zeta(z) := \int_{\mathbb{R}} e^{zy} d\zeta(y)$ is finite on some open interval containing zero (“light-tailedness”); for Theorem 8, we need the corresponding variant (15) for $M_{\mathbf{P}} \cdot \varphi \in \Upsilon(]a, b[)$ (rather than φ).

Hence, finding such “BS-associated pairs (φ, ζ) ” is an important — but highly nontrivial — issue. Indeed, one approach is to start from a given divergence generator $\varphi \in \Upsilon(]a, b[)$ having some additional properties, switch to its Fenchel-Legendre

transform φ^* (and some exponentially-linear transforms thereof), and verify some sufficient conditions for the outcome to be a moment-generating function MGF_ζ of a unique probability distribution ζ which has light tails. For finding the concrete ζ , one typically should know the explicit form of φ^* . However, it is well known that it can sometimes be hard to determine the explicit form of the Fenchel-Legendre transform of a convex function. This hardness issue also applies for the reverse direction of starting from a concrete probability distribution ζ with light tails, computing its log-moment-generating function (called cumulant-generating function) $z \mapsto \Lambda_\zeta(z) := \log MGF_\zeta(z)$ and the corresponding Fenchel-Legendre transform Λ_ζ^* which is nothing but the associated divergence generator φ (cf. (6)).

As will be demonstrated via numerous solved cases in the following Section XII, the — “kind of intermediate” — new construction method given in the below-mentioned Theorem 22 can help to ease these two tasks. To formulate this, we employ the class \mathfrak{F} of functions $F :]-\infty, \infty[\mapsto]-\infty, \infty[$ with the following properties:

- (F1) $\text{int}(\text{dom}(F)) =]a_F, b_F[$ for some $-\infty \leq a_F < 1 < b_F \leq \infty$;
- (F2) F is smooth (infinitely continuously differentiable) on $]a_F, b_F[$;
- (F3) F is strictly increasing on $]a_F, b_F[$.

Clearly, for any $F \in \mathfrak{F}$ one gets the existence of $F(a_F) := \lim_{t \downarrow a_F} F(t) \in]-\infty, \infty[$ and $F(b_F) := \lim_{t \uparrow b_F} F(t) \in]-\infty, \infty[$; moreover, its inverse $F^{-1} : \mathcal{R}(F) \mapsto]a_F, b_F[$ exists, where $\mathcal{R}(F) := \{F(t) : t \in \text{dom}(F)\}$. Furthermore, F^{-1} is strictly increasing and smooth (infinitely continuously differentiable) on the open interval $\text{int}(\mathcal{R}(F)) = \{F(t) : t \in]a_F, b_F[\} =]F(a_F), F(b_F)[$, and $F^{-1}(\text{int}(\mathcal{R}(F))) =]a_F, b_F[$. Within such a context, we obtain

Theorem 22: Let $F \in \mathfrak{F}$ and fix an arbitrary point $c \in \text{int}(\mathcal{R}(F))$. Moreover, introduce the notations²⁵ $]\lambda_-, \lambda_+[:= \text{int}(\mathcal{R}(F)) - c$ and $]t_-^{sc}, t_+^{sc}[:=]1 + a_F - F^{-1}(c), 1 + b_F - F^{-1}(c)[$ (which implies $\lambda_- < 0 < \lambda_+$ and $t_-^{sc} < 1 < t_+^{sc}$). Furthermore, define the functions $\Lambda :]-\infty, \infty[\mapsto]-\infty, \infty[$ and $\varphi :]-\infty, \infty[\mapsto [0, \infty[$ by

$$\Lambda(z) := \Lambda^{(c)}(z) := \begin{cases} \int_0^z F^{-1}(u+c) du + z \cdot (1 - F^{-1}(c)) \in]-\infty, \infty[, & \text{if } z \in]\lambda_-, \lambda_+[, \\ \int_0^{\lambda_-} F^{-1}(u+c) du + \lambda_- \cdot (1 - F^{-1}(c)) \in]-\infty, \infty[, & \text{if } z = \lambda_- > -\infty, \\ \int_0^{\lambda_+} F^{-1}(u+c) du + \lambda_+ \cdot (1 - F^{-1}(c)) \in]-\infty, \infty[, & \text{if } z = \lambda_+ < \infty, \\ \infty, & \text{else,} \end{cases} \quad (130)$$

where the second respectively third line are meant as $\lim_{z \downarrow \lambda_-} (\int_0^z F^{-1}(u+c) du + z \cdot (1 - F^{-1}(c)))$ respectively $\lim_{z \uparrow \lambda_+} (\int_0^z F^{-1}(u+c) du + z \cdot (1 - F^{-1}(c)))$, and

$$\varphi(t) := \varphi^{(c)}(t) := \begin{cases} (t + F^{-1}(c) - 1) \cdot [F(t + F^{-1}(c) - 1) - c] - \int_0^{F(t + F^{-1}(c) - 1) - c} F^{-1}(u+c) du \in [0, \infty[, & \text{if } t \in]t_-^{sc}, t_+^{sc}[, \\ (t_-^{sc} + F^{-1}(c) - 1) \cdot [F(t_-^{sc} + F^{-1}(c) - 1) - c] - \int_0^{F(t_-^{sc} + F^{-1}(c) - 1) - c} F^{-1}(u+c) du \in [0, \infty[, & \text{if } t = t_-^{sc} > -\infty, \\ (t_+^{sc} + F^{-1}(c) - 1) \cdot [F(t_+^{sc} + F^{-1}(c) - 1) - c] - \int_0^{F(t_+^{sc} + F^{-1}(c) - 1) - c} F^{-1}(u+c) du \in [0, \infty[, & \text{if } t = t_+^{sc} < \infty, \\ \varphi(t_-^{sc}) + \lambda_- \cdot (t - t_-^{sc}) \in [0, \infty[, & \text{if } t_-^{sc} > -\infty \text{ and } t \in]-\infty, t_-^{sc}[, \\ \varphi(t_+^{sc}) + \lambda_+ \cdot (t - t_+^{sc}) \in [0, \infty[, & \text{if } t_+^{sc} < \infty \text{ and } t \in]t_+^{sc}, \infty[, \\ \infty, & \text{else,} \end{cases} \quad (131)$$

where the second respectively third line are again meant as lower respectively upper limit.

Then, Λ and φ have the following properties:

- (i) On $]\lambda_-, \lambda_+[$, the function Λ is smooth and strictly convex and consequently, $\exp(\Lambda)$ is smooth and strictly log-convex; moreover, there holds $\Lambda(0) = 0$, $\Lambda'(0) = 1$.
- (ii) $\varphi \in \tilde{\Upsilon}(]a, b[)$, where $a := t_-^{sc} \cdot 1_{\{-\infty\}}(\lambda_-) - \infty \cdot 1_{] - \infty, 0[}(\lambda_-)$, $b := t_+^{sc} \cdot 1_{\{\infty\}}(\lambda_+) + \infty \cdot 1_{] 0, \infty[}(\lambda_+)$.
- (iii) $\varphi(t) = \Lambda^*(t) = \sup_{z \in]-\infty, \infty[} (z \cdot t - \Lambda(z)) = \sup_{z \in]\lambda_-, \lambda_+[} (z \cdot t - \Lambda(z))$ for all $t \in \mathbb{R}$.
- (iv) $\Lambda(z) = \varphi^*(z) = \sup_{t \in]-\infty, \infty[} (t \cdot z - \varphi(t)) = \sup_{t \in]a, b[} (t \cdot z - \varphi(t))$ for all $z \in \mathbb{R}$.

The proof of Theorem 22 will be given Appendix E.

Remark 23: (a) Notice that the newly constructed Λ and φ (cf. (130),(131)) depend on the choice of the *anchor point* c ; this is e.g. illustrated in Subsection XII-F below. Hence, as a side effect, by using whole families $(F_\vartheta)_\vartheta$ together with different anchor points c , via Theorem 22 one can generate new classes (and new classifications) of φ -divergence generators — and thus

²⁵for the sake of brevity, we avoid here the more complete notation $\lambda_-^{F,c}, \lambda_+^{F,c}, t_-^{sc,F,c}, t_+^{sc,F,c}$ indicating the dependence on F and c .

of corresponding φ -divergences — which can be of great use, even in other contexts beyond our BS optimization framework. (b) If F satisfies $F(1) = 0$ and thus $F^{-1}(0) = 1$, then the natural choice $c := 0$ induces $] \lambda_-, \lambda_+[= \text{int}(\mathcal{R}(F))$ and $]t_-^{sc}, t_+^{sc}[=]a_F, b_F[$, and consequently (due to $F^{-1}(c) - 1 = 0$) leads to the simplification of “the first lines of” (130),(131) to

$$\Lambda(z) := \Lambda^{(0)}(z) := \int_0^z F^{-1}(u) du, \quad z \in \text{int}(\mathcal{R}(F)), \quad (132)$$

$$\varphi(t) := \varphi^{(0)}(t) := t \cdot F(t) - \int_0^{F(t)} F^{-1}(u) du, \quad t \in]a_F, b_F[; \quad (133)$$

the simplifications of the respective other lines of (130) and (131) are straightforward.

One can even prove many more properties of φ given by (131) which strongly indicate that the F -constructed function $z \mapsto \exp(\Lambda(z)) = \exp(\varphi^*(z))$ is a *good candidate* for a moment generating function of a probability distribution ζ , and hence for the representability (6) (i.e. $\varphi \in \Upsilon(]a, b[)$). However, one still needs to verify one of the conditions (a) to (c) of the following Proposition 24. This may go wrong, as the case of power divergences φ_γ with $\gamma \in]1, 2[$ indicates (cf. the conjecture of Subsection XII-F3 below).

Proposition 24: Let φ be given by (131) of Theorem 22. Then, $\varphi \in \Upsilon(]a, b[)$ if one of the following three conditions holds: (a) $a > -\infty$, $\lambda_- = -\infty$, and the function $z \mapsto M(z) := e^{-a \cdot z + \varphi^*(z)}$ is absolutely monotone on $] -\infty, 0[$ (i.e. all derivatives exist and satisfy $\frac{\partial^k}{\partial z^k} M(z) \geq 0$ for all $k \in \mathbb{N}_0$, $z \in] -\infty, 0[$), (b) $b < \infty$, $\lambda_+ = \infty$, and the function $z \mapsto M(z) := e^{b \cdot z + \varphi^*(-z)}$ is absolutely monotone on $] -\infty, 0[$, (c) $a = -\infty$, $b = -\infty$, and the function $z \mapsto M(z) := e^{\varphi^*(z)}$ is exponentially convex on $] \lambda_-, \lambda_+[$ (i.e. $M(\cdot)$ is continuous and satisfies $\sum_{i=1}^n \sum_{j=1}^n c_i \cdot c_j \cdot M\left(\frac{z_i + z_j}{2}\right) \geq 0$ for all $n \in \mathbb{N}$, $c_i, c_j \in \mathbb{R}$ and $z_i, z_j \in] \lambda_-, \lambda_+[$). If one of the three conditions (a) to (c) holds, then the associated probability distribution ζ (cf. (6)) has expectation $\int_{\mathbb{R}} y d\zeta(y) = 1$ and finite moments of all orders, i.e. $\int_{\mathbb{R}} y^j d\zeta(y) < \infty$ for all $j \in \mathbb{N}_0$; in terms of $\zeta[\cdot] := \mathbb{P}[W \in \cdot]$ this means that $E_{\mathbb{P}}[W] = 1$ and $E_{\mathbb{P}}[W^j] < \infty$.

Proposition 24(a),(b) follow from Theorem 22 and the well-known (probability-version of) *Bernstein's theorem* [327] (see e.g. also [328]) — which needs to be appropriately shifted — whereas (c) follows from Theorem 22 and the well-known (probability-version of) *Widder's theorem* [329]²⁶ (see e.g. also [331]–[335]). As far as applicability is concerned, it is well known that verifying absolute monotonicity is typically more comfortable than verifying exponential convexity. Fortunately, one can often use the former, since for many known divergence generators there holds $a > -\infty$ (often $a = 0$) or/and $b < \infty$.

For the identification of light-tailed *semi-/half-lattice* distributions, we obtain the following two sets of conditions, which even allow for the desired explicit determination of ζ :

Proposition 25: Let φ be given by (131) of Theorem 22, with some $a > -\infty$. Furthermore, assume that there exists some constant $\check{c} > 0$ as well as some function $H : [0, \infty[\mapsto [0, \infty[$ which is continuous on $[0, 1]$ with $H(1) = 1$ and absolutely monotone on $]0, 1[$, such that $e^{\varphi^*(\frac{z}{\check{c}}) - a \cdot \frac{z}{\check{c}}} = H(e^z)$ ($z \in] -\infty, \check{c} \cdot \lambda_+[$). Then one has $\varphi \in \Upsilon(]a, b[)$ and $\zeta = \sum_{n=0}^{\infty} p_n \cdot \delta_{a + \check{c} \cdot n}$ with $p_n := \frac{1}{n!} \cdot \frac{d^n H}{dt^n}(0)$, i.e. $\mathbb{P}[W = a + \check{c} \cdot n] = p_n$ ($n \in \mathbb{N}_0$).

Proposition 26: Let φ be given by (131) of Theorem 22, with some $b < \infty$. Furthermore, assume that there exists some constant $\check{c} > 0$ as well as some function $H : [0, \infty[\mapsto [0, \infty[$ which is continuous on $[0, 1]$ with $H(1) = 1$ and absolutely monotone on $]0, 1[$, such that $e^{\varphi^*(-\frac{z}{\check{c}}) + b \cdot \frac{z}{\check{c}}} = H(e^z)$ ($z \in] -\infty, -\check{c} \cdot \lambda_-[$). Then one has $\varphi \in \Upsilon(]a, b[)$ and $\zeta = \sum_{n=0}^{\infty} p_n \cdot \delta_{b - \check{c} \cdot n}$ with $p_n := \frac{1}{n!} \cdot \frac{d^n H}{dt^n}(0)$, i.e. $\mathbb{P}[W = b - \check{c} \cdot n] = p_n$ ($n \in \mathbb{N}_0$).

The Propositions 25 and 26 follow from some straightforward transformations and a well-known characterization of probability generating functions H (see e.g. in Theorem 1.2.10 of [336]).

As an incentive for the following investigations, let us recall the discussion in the surroundings of Condition 7 pertaining to the minimization problem (13), where we have addressed possible connections between the two representabilities (6) (needed e.g. for Theorem 12) and (15) (needed e.g. for Theorem 8); this strongly relates to the question, for which constants $\tilde{c} > 0$ the validity $\varphi \in \Upsilon(]a, b[)$ triggers the validity of $\tilde{c} \cdot \varphi \in \Upsilon(]a, b[)$. To begin with, it is straightforward to see that $\varphi \in \Upsilon(]a, b[)$ *always implies* $\tilde{c} \cdot \varphi \in \Upsilon(]a, b[)$ for all *integers* $\tilde{c} \in \mathbb{N}$; indeed, if φ satisfies (6) for some $\zeta = \mathbb{P}[W \in \cdot]$, then for each integer $\tilde{c} \in \mathbb{N}$ one gets that $\tilde{c} \cdot \varphi$ satisfies (6) for $\tilde{\zeta} = \mathbb{P}[\sum_{j=1}^{\tilde{c}} \frac{W_j}{\tilde{c}} \in \cdot]$; in the latter, the W_j 's are i.i.d. copies of W . Clearly, $MGF_{\tilde{\zeta}}$ is then finite on some open interval containing zero (differing from the one for MGF_{ζ} only by a scaling with $1/\tilde{c}$).

²⁶for the relevant conversion between the involved Riemann-Stieltjes integral with nondecreasing (but not necessarily right-continuous) integrator into a measure integral, one can apply e.g. the general theory in Chapter 6 of [330].

For the following family of distributions, one can even trigger $\tilde{c} \cdot \varphi \in \Upsilon(]a, b[)$ for all $\tilde{c} > 0$: for the sake of a corresponding precise formulation, recall first the common knowledge that, generally speaking, a probability distribution ζ on \mathbb{R} with light tails — in the sense that its moment generating function $z \mapsto MGF_\zeta(z) := \int_{\mathbb{R}} e^{z \cdot y} d\zeta(y)$ is finite on some open interval $] \lambda_-, \lambda_+[$ containing zero — is (said to be) *infinitely divisible* if there holds

$$\text{for each } n \in \mathbb{N} \text{ there exists a probability distribution } \zeta_n \text{ on } \mathbb{R} \text{ such that } \int_{\mathbb{R}} e^{z \cdot y} d\zeta(y) = \left(\int_{\mathbb{R}} e^{z \cdot y} d\zeta_n(y) \right)^n, \quad z \in] \lambda_-, \lambda_+[; \quad (134)$$

in fact, (134) means that the (light-tailed) moment generating function MGF_ζ is *infinitely divisible* in the sense that each n -th root $(MGF_\zeta)^{1/n}$ must be the moment generating function of some (light-tailed) probability distribution (denoted here by ζ_n). In particular, (134) implies that ζ_n is unique, and that ζ must necessarily have (one-sided or two-sided) unbounded support $\text{supp}(\zeta)$. The latter may differ from $\text{supp}(\zeta_n)$. In our BS context (6), (134) equivalently means that the associated random variable W is *infinitely divisible* (with light-tailed distribution), in the sense that

$$\text{for each } n \in \mathbb{N} \text{ there exists a sequence of i.i.d. random variables } Y_{n,1}, \dots, Y_{n,n} \text{ such that } W \stackrel{d}{=} Y_{n,1} + \dots + Y_{n,n},$$

where $\stackrel{d}{=}$ means “have equal probability distributions” and $\mathbb{P}[W \in \cdot] = \zeta[\cdot]$, $\mathbb{P}[Y_{n,1} \in \cdot] = \zeta_n[\cdot]$.

For the above-mentioned context, we obtain the useful assertion (which will be proved in Appendix D):

Proposition 27: Suppose that $\varphi \in \Upsilon(]a, b[)$ with connected probability distribution ζ from (6). Then there holds:

$$\tilde{c} \cdot \varphi \in \Upsilon(]a, b[) \text{ for all } \tilde{c} > 0 \iff \zeta \text{ is infinitely divisible.}$$

Notice that Proposition 27 covers especially the important prominent *power divergences* (cf. the solved-cases Sections XII-A to XII-F below) for which we provide the corresponding infinitely divisible distributions explicitly. More generally, for the identification of light-tailed infinitely divisible distributions, we obtain the following three sets of sufficient conditions:

Proposition 28: Let φ be given by (131) of Theorem 22. Then, $\varphi \in \Upsilon(]a, b[)$ and the associated probability distribution ζ is infinitely divisible, if one of the following three conditions holds:

- (a) $a > -\infty$, $\lambda_- = -\infty$, and the function $z \mapsto \varphi^{*'}(z) - a = (\varphi')^{-1}(z) - a$ is absolutely monotone on $] -\infty, 0[$,
- (b) $b < \infty$, $\lambda_+ = \infty$, and the function $z \mapsto -\varphi^{*'}(-z) + b = -(\varphi')^{-1}(-z) + b$ is absolutely monotone on $] -\infty, 0[$,
- (c) $a = -\infty$, $b = \infty$, and the function $z \mapsto \frac{\varphi^{*''}(z)}{\varphi^{*''}(0)} = \frac{\varphi''(1)}{\varphi''((\varphi')^{-1}(z))}$ is exponentially convex on $] \lambda_-, \lambda_+[$.

In the first case (a) there automatically follows $b = \infty$, whereas in the second case (b) one automatically gets $a = -\infty$.

The proof of Proposition 28 is given in Appendix D.

Let us end this section by giving some further comments on the task of finding concretely the probability distribution (if existent) $\zeta[\cdot] = \mathbb{P}[W \in \cdot]$ from the Fenchel-Legendre transform $\Lambda = \varphi^*$ of a pregiven divergence generator φ , which should satisfy $MGF(z) := \exp(\Lambda(z)) = \int_{\mathbb{R}} e^{z \cdot y} d\zeta(y) = E_{\mathbb{P}}[\exp(z \cdot W)]$ ($z \in \mathbb{R}$). Recall that this is used for the simulation of the weights $(W_i)_{i \in \mathbb{N}}$ which are i.i.d. copies of W and which are the crucial building ingredients of ξ_n^W in Theorem 8, respectively, of $\xi_n^{w, \mathbf{X}}$ in Theorem 12. The search for ζ can be done e.g. by inversion of a candidate moment generating function MGF, or by search in tables or computer software which list distributions and their MGF. Also notice that ζ needs not necessarily be explicitly known in full detail (e.g. in terms of a computationally tractable density or frequency); for instance, as well known from insurance applications, for — comfortably straightforwardly simulable — doubly-random sums $W := \sum_{i=1}^N A_i$ of nonnegative i.i.d. random variables $(A_i)_{i \in \mathbb{N}}$ with known law $\Pi_A[\cdot] := \mathbb{P}[A \in \cdot]$ being independent of a counting-type random variable N with known law Π_N , one can mostly compute explicitly $MGF_\zeta(z) = PGF_{\Pi_N}(MGF_{\Pi_A}(z))$ with the help of the probability generating function PGF_{Π_N} of Π_N , but the corresponding density/frequency of ζ may not be known explicitly in a tractable form. The below-mentioned solved Case 2 in Subsection XII-B of power divergences with generator φ_γ ($\gamma \in]0, 1[$) manifests such a situation. In the end, if no explicit distribution ζ and no comfortably simulable W —construction are available, one can still try to simulate an i.i.d. sequence $(W_i)_{i \in \mathbb{N}}$ from the pregiven moment generating function (which is $\exp(\Lambda(z))$ here); see e.g. [337] and references therein which also contains saddle point methods approximation techniques.

Let us finally mention that a more complete picture on finding the distribution ζ of the weights W (including necessary and sufficient conditions, boundary behaviours, etc.) can be found in our paper’s full arXiv-version [158].

XII. EXPLICITLY SOLVED CASES

With the help of Theorem 22 we can comfortably generate various important solved cases, which we demonstrate now. Notice especially in the Cases XII-A to XII-F the astonishing effect that the “homogeneous/continuous” class of power-divergence generators $(\varphi_\gamma)_{\gamma \in \mathbb{R}}$ (cf. (40)) are connected to a “very inhomogeneous” family of underlying “cornerstone simulation laws” $(\zeta_\gamma)_{\gamma \in \mathbb{R}}$ of W -distributions: discrete, continuous, mixture of discrete and continuous, as the parameter γ varies.

A. Case 1

For $\gamma < 0$, $\tilde{c} \in]0, \infty[$ and $]a_{F_{\gamma, \tilde{c}}}, b_{F_{\gamma, \tilde{c}}}[:=]0, \infty[$ we define

$$F_{\gamma, \tilde{c}}(t) := \begin{cases} \frac{\tilde{c}}{\gamma-1} \cdot (t^{\gamma-1} - 1), & \text{if } t \in]0, \infty[, \\ -\infty, & \text{if } t \in]-\infty, 0]. \end{cases} \quad (135)$$

Clearly, $\mathcal{R}(F_{\gamma, \tilde{c}}) =]-\infty, \frac{\tilde{c}}{1-\gamma}[$. Furthermore, $F_{\gamma, \tilde{c}}(\cdot)$ is strictly increasing and smooth on $]a_{F_{\gamma, \tilde{c}}}, b_{F_{\gamma, \tilde{c}}}[$, and thus, $F_{\gamma, \tilde{c}} \in \mathfrak{F}$. Since $F_{\gamma, \tilde{c}}(1) = 0$, let us choose the natural anchor point $c := 0 \in \text{int}(\mathcal{R}(F_{\gamma, \tilde{c}}))$, which leads to $]\lambda_-, \lambda_+[:= \text{int}(\mathcal{R}(F_{\gamma, \tilde{c}})) - c =]-\infty, \frac{\tilde{c}}{1-\gamma}[$, $]t_-^{sc}, t_+^{sc}[:=]1 + a_{F_{\gamma, \tilde{c}}} - F_{\gamma, \tilde{c}}^{-1}(c), 1 + b_{F_{\gamma, \tilde{c}}} - F_{\gamma, \tilde{c}}^{-1}(c)[=]0, \infty[$ and $]a, b[=]0, \infty[$ since $a := t_-^{sc} \cdot \mathbf{1}_{\{-\infty\}}(\lambda_-) - \infty \cdot \mathbf{1}_{]-\infty, 0]}(\lambda_-) = 0$, $b := t_+^{sc} \cdot \mathbf{1}_{\{\infty\}}(\lambda_+) + \infty \cdot \mathbf{1}_{]0, \infty]}(\lambda_+) = \infty$. By using $F_{\gamma, \tilde{c}}^{-1}(x) = (1 + \frac{(\gamma-1) \cdot x}{\tilde{c}})^{\frac{1}{\gamma-1}}$ for $x \in \text{int}(\mathcal{R}(F_{\gamma, \tilde{c}}))$, by straightforward calculations we can deduce from formula (131) (see also (133))

$$\varphi_{\gamma, \tilde{c}}(t) := \varphi_{\gamma, \tilde{c}}^{(0)}(t) = \begin{cases} \tilde{c} \cdot \frac{t^\gamma - \gamma \cdot t + \gamma - 1}{\gamma \cdot (\gamma - 1)} \in [0, \infty[, & \text{if } t \in]0, \infty[, \\ \infty, & \text{if } t \in]-\infty, 0], \end{cases}$$

which coincides with $\tilde{c} \cdot \varphi_\gamma(t)$ for $\varphi_\gamma(t)$ from (40) and which generates the γ -corresponding power divergences given in (41). Moreover, we can derive from formula (130) (see also (132))

$$\Lambda_{\gamma, \tilde{c}}(z) := \Lambda_{\gamma, \tilde{c}}^{(0)}(z) = \begin{cases} \frac{\tilde{c}}{\gamma} \cdot \left\{ \left(\frac{\gamma-1}{\tilde{c}} \cdot z + 1 \right)^{\frac{\gamma}{\gamma-1}} - 1 \right\}, & \text{if } z \in]-\infty, \frac{\tilde{c}}{1-\gamma}[, \\ -\frac{\tilde{c}}{\gamma} > 0, & \text{if } z = \frac{\tilde{c}}{1-\gamma}, \\ \infty, & \text{if } z \in]\frac{\tilde{c}}{1-\gamma}, \infty[. \end{cases}$$

The latter is the cumulant generating function of a “tilted (i.e. negatively distorted) stable distribution” $\zeta[\cdot] = \mathbb{P}[W \in \cdot]$ of a random variable W , which can be constructed as follows: let Z be an auxiliary random variable (having density f_Z and support $\text{supp}(Z) = [0, \infty[$) of a stable law with parameter-quadruple $(\frac{-\gamma}{1-\gamma}, 1, 0, -\frac{\tilde{c}^{1/(1-\gamma)} \cdot (1-\gamma)^{-\gamma/(1-\gamma)}}{\gamma})$ in terms of the “form-B notation” on p.12 in [338]; by applying a general Laplace-transform result on p.112 of the same text we can deduce

$$M_Z(z) := E_{\mathbb{P}}[\exp(z \cdot Z)] = \int_0^\infty \exp(z \cdot y) \cdot f_Z(y) dy = \begin{cases} \exp\left(\frac{\tilde{c}^{1/(1-\gamma)} \cdot (1-\gamma)^{-\gamma/(1-\gamma)}}{\gamma} \cdot (-z)^\alpha\right), & \text{if } z \in]-\infty, 0], \\ \infty, & \text{if } z \in]0, \infty[, \end{cases} \quad (136)$$

where $\alpha := -\frac{\gamma}{1-\gamma} \in]0, 1[$. Since $0 \notin \text{int}(\text{dom}(M_Z))$ (and thus, Z does not have light-tails) we have to tilt (dampen) the density in order to extend the effective domain. Accordingly, let W be a random variable having density²⁷

$$f_W(y) := \frac{\exp\{-\frac{y \cdot \tilde{c}}{1-\gamma}\}}{\exp\{\tilde{c}/\gamma\}} \cdot f_Z(y) \cdot \mathbf{1}_{]0, \infty[}(y), \quad y \in \mathbb{R}. \quad (137)$$

Then one can straightforwardly deduce from (136) that $\int_0^\infty f_W(y) dy = 1$ and that

$$M_W(z) := E_{\mathbb{P}}[\exp(z \cdot W)] = \int_0^\infty \exp(z \cdot y) \cdot f_W(y) dy = \begin{cases} \exp\left(\frac{\tilde{c}}{\gamma} \cdot \left\{ \left(\frac{\gamma-1}{\tilde{c}} \cdot z + 1 \right)^{\frac{\gamma}{\gamma-1}} - 1 \right\}\right), & \text{if } z \in]-\infty, \frac{\tilde{c}}{1-\gamma}[, \\ \infty, & \text{if } z \in]\frac{\tilde{c}}{1-\gamma}, \infty[. \end{cases}$$

Notice that ζ is an infinitely divisible (cf. Proposition 27) continuous distribution with density f_W , and that $\zeta[]0, \infty[= \mathbb{P}[W > 0] = 1$. Concerning the important Remark 9(i), for i.i.d. copies $(W_i)_{i \in \mathbb{N}}$ of W we obtain the probability distribution $\zeta^{*n_k}[\cdot] := \mathbb{P}[\check{W} \in \cdot]$ of $\check{W} := \sum_{i \in I_k^{(n)}} W_i$ having the density

$$f_{\check{W}}(y) := \frac{\exp\{-\frac{y \cdot \tilde{c}}{1-\gamma}\}}{\exp\{\tilde{c} \cdot \text{card}(I_k^{(n)})/\gamma\}} \cdot f_Z(y) \cdot \mathbf{1}_{]0, \infty[}(y), \quad y \in \mathbb{R}, \quad (138)$$

where \check{Z} is a random variable with density $f_{\check{Z}}$ of a stable law with parameters $(\frac{-\gamma}{1-\gamma}, 1, 0, -\frac{\tilde{c}^{1/(1-\gamma)} \cdot (1-\gamma)^{-\gamma/(1-\gamma)}}{\gamma} \cdot \text{card}(I_k^{(n)}))$. We are now in the position to state explicitly the bare-simulation-minimizations (respectively maximizations) of the corresponding (γ -order) power divergences, Renyi divergences, Hellinger integrals and measures of entropy (diversity), which can be deduced from Theorem 12, Remark 13(vi), Lemma 14(a), (69), (73), (77), (78):

²⁷in the classical sense, with respect to Lebesgue measure

Proposition 29: (a) Consider $\varphi := \tilde{c} \cdot \varphi_\gamma$ with $\gamma < 0$, and let $\mathbb{P} \in \mathbb{S}_{>0}^K$ as well as $\tilde{c} > 0$ be arbitrary but fixed. Furthermore, let $W := (W_i)_{i \in \mathbb{N}}$ be an i.i.d. sequence of non-negative real-valued random variables having density (137). Then for all $A > 0$ and all $\Omega \subset \mathbb{S}_{>0}^K$ with (7) there holds

$$-\lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P} \left[\xi_n^{w\mathbf{W}} \in \Omega \right] = \inf_{\mathbf{Q} \in A \cdot \Omega} \frac{\tilde{c}}{\gamma} \cdot \left[1 - A^{\gamma/(\gamma-1)} \cdot \left[1 + \gamma \cdot (A-1) + \frac{\gamma \cdot (\gamma-1)}{\tilde{c}} \cdot D_{\tilde{c} \cdot \varphi_\gamma}(\mathbf{Q}, \mathbb{P}) \right]^{-1/(\gamma-1)} \right] \quad (139)$$

as well as the BS minimizabilities/maximizabilites (cf. Definition 1)

$$\inf_{\mathbf{Q} \in A \cdot \Omega} D_{\tilde{c} \cdot \varphi_\gamma}(\mathbf{Q}, \mathbb{P}) = \lim_{n \rightarrow \infty} \frac{\tilde{c}}{\gamma \cdot (\gamma-1)} \cdot \left\{ A^\gamma \cdot \left(1 + \frac{\gamma}{\tilde{c}} \cdot \frac{1}{n} \cdot \log \mathbb{P} \left[\xi_n^{w\mathbf{W}} \in \Omega \right] \right)^{1-\gamma} + \gamma \cdot (1-A) - 1 \right\}, \quad (140)$$

$$\inf_{\mathbf{Q} \in A \cdot \Omega} H_\gamma(\mathbf{Q}, \mathbb{P}) = \lim_{n \rightarrow \infty} A^\gamma \cdot \left(1 + \gamma \cdot \frac{1}{n} \cdot \log \mathbb{P} \left[\check{\xi}_n^{w\mathbf{W}} \in \Omega \right] \right)^{1-\gamma}, \quad (141)$$

$$\inf_{\mathbf{Q} \in A \cdot \Omega} c_1 \cdot \left(H_\gamma(\mathbf{Q}, \mathbb{P})^{c_2} - c_3 \right) = \lim_{n \rightarrow \infty} c_1 \cdot \left\{ A^{c_2 \cdot \gamma} \cdot \left(1 + \frac{\gamma}{n} \cdot \log \mathbb{P} \left[\check{\xi}_n^{w\mathbf{W}} \in \Omega \right] \right)^{c_2 \cdot (1-\gamma)} - c_3 \right\}, \quad \text{if } c_1 \cdot c_2 > 0, c_3 \in \mathbb{R}, \quad (142)$$

$$\sup_{\mathbf{Q} \in A \cdot \Omega} c_1 \cdot \left(H_\gamma(\mathbf{Q}, \mathbb{P})^{c_2} - c_3 \right) = \lim_{n \rightarrow \infty} c_1 \cdot \left\{ A^{c_2 \cdot \gamma} \cdot \left(1 + \frac{\gamma}{n} \cdot \log \mathbb{P} \left[\check{\xi}_n^{w\mathbf{W}} \in \Omega \right] \right)^{c_2 \cdot (1-\gamma)} - c_3 \right\}, \quad \text{if } c_1 \cdot c_2 < 0, c_3 \in \mathbb{R}, \quad (143)$$

$$\inf_{\mathbf{Q} \in A \cdot \Omega} R_\gamma(\mathbf{Q}, \mathbb{P}) = \lim_{n \rightarrow \infty} \frac{1}{\gamma \cdot (\gamma-1)} \cdot \log \left(A^\gamma \cdot \left(1 + \gamma \cdot \frac{1}{n} \cdot \log \mathbb{P} \left[\check{\xi}_n^{w\mathbf{W}} \in \Omega \right] \right)^{1-\gamma} \right), \quad (144)$$

$$\inf_{\mathbf{Q} \in A \cdot \Omega} c_1 \cdot \left(\left(\sum_{k=1}^K q_k^\gamma \right)^{c_2} - c_3 \right) = \lim_{n \rightarrow \infty} c_1 \cdot \left\{ K^{c_2 \cdot (1-\gamma)} \cdot A^{c_2 \cdot \gamma} \cdot \left(1 + \frac{\gamma}{n} \cdot \log \mathbb{P} \left[\check{\xi}_n^{w\mathbf{W}} \in \Omega \right] \right)^{c_2 \cdot (1-\gamma)} - c_3 \right\},$$

if $c_1 \cdot c_2 > 0, c_3 \in \mathbb{R}, \quad (145)$

$$\sup_{\mathbf{Q} \in A \cdot \Omega} c_1 \cdot \left(\left(\sum_{k=1}^K q_k^\gamma \right)^{c_2} - c_3 \right) = \lim_{n \rightarrow \infty} c_1 \cdot \left\{ K^{c_2 \cdot (1-\gamma)} \cdot A^{c_2 \cdot \gamma} \cdot \left(1 + \frac{\gamma}{n} \cdot \log \mathbb{P} \left[\check{\xi}_n^{w\mathbf{W}} \in \Omega \right] \right)^{c_2 \cdot (1-\gamma)} - c_3 \right\},$$

if $c_1 \cdot c_2 < 0, c_3 \in \mathbb{R}, \quad (146)$

$$\inf_{\mathbf{Q} \in A \cdot \Omega} \frac{1}{1-\gamma} \cdot \log \left(\sum_{k=1}^K q_k^\gamma \right) = \lim_{n \rightarrow \infty} \frac{1}{1-\gamma} \cdot \left[\log \left(A^\gamma \cdot \left(1 + \gamma \cdot \frac{1}{n} \cdot \log \mathbb{P} \left[\check{\xi}_n^{w\mathbf{W}} \in \Omega \right] \right)^{1-\gamma} \right) + (1-\gamma) \cdot \log(K) \right], \quad (147)$$

where $\xi_n^{w\mathbf{W}}$ is the normalized randomly weighted empirical measure given in (38), $\check{\xi}_n^{w\mathbf{W}}$ is its special case for $\tilde{c} = 1$, and $\check{\check{\xi}}_n^{w\mathbf{W}}$ is its special case for $\tilde{c} = 1$ together with $\mathbb{P} = \mathbb{P}^{unif}$ ²⁸. From this, the BS-minimizability/maximizability of the important norms/entropies/diversity indices (E1) to (E6) follow immediately as special cases.

(b) The special case $\varphi := \tilde{c} \cdot \varphi_\gamma$ ($\gamma < 0$) of Theorem 12 works analogously to (a), with the differences that we employ (i) additionally a sequence $(X_i)_{i \in \mathbb{N}}$ of random variables being independent of $(W_i)_{i \in \mathbb{N}}$ and satisfying condition (26) (resp. (30)), (ii) $A = 1$ (instead of arbitrary $A > 0$), (iii) $\mathbb{P}_{X_1^n}[\cdot]$ (instead of $\mathbb{P}[\cdot]$), (iv) $\xi_{n,\mathbf{X}}^{w\mathbf{W}}$ (instead of $\xi_n^{w\mathbf{W}}$), (v) $\check{\xi}_{n,\mathbf{X}}^{w\mathbf{W}}$ (instead of $\check{\xi}_n^{w\mathbf{W}}$), and (vi) $\check{\check{\xi}}_{n,\mathbf{X}}^{w\mathbf{W}}$ (instead of $\check{\xi}_n^{w\mathbf{W}}$).

Within the context of Subsection X-A, for the concrete simulative estimation $\widehat{D_{\tilde{c} \cdot \varphi_\gamma}}(\Omega, \mathbf{P})$ via (116) and (118), we obtain — in terms of $M_{\mathbf{P}} := \sum_{i=1}^K p_i > 0$, $n_k = n \cdot \tilde{p}_k \in \mathbb{N}$ and \tilde{q}_k^* from proxy method 1 or 2 — that $\tilde{U}_k^{*n_k}$ has the (Lebesgue-)density

$$f_{\tilde{U}_k^{*n_k}}(x) := \frac{\exp\left(\left(\tau_k - \frac{\tilde{c} \cdot M_{\mathbf{P}}}{1-\gamma}\right) \cdot x\right)}{\exp\left(n_k \cdot \frac{\tilde{c} \cdot M_{\mathbf{P}}}{\gamma} \cdot \left(1 + \frac{\gamma-1}{\tilde{c} \cdot M_{\mathbf{P}}} \cdot \tau_k\right)^{\gamma/(\gamma-1)}\right)} \cdot f_{\check{Z}}(x) \cdot \mathbb{1}_{0,\infty[}(x), \quad x \in \mathbb{R},$$

where $\tau_k = \tilde{c} \cdot M_{\mathbf{P}} \cdot \frac{1 - \left(\frac{\tilde{q}_k^*}{p_k}\right)^{\gamma-1}}{1-\gamma}$ for $\tilde{q}_k^* > 0$, and \check{Z} is a random variable with density $f_{\check{Z}}$ of a stable law with parameter-quadruple $\left(\frac{-\gamma}{1-\gamma}, 1, 0, -n_k \cdot \frac{(\tilde{c} \cdot M_{\mathbf{P}})^{1/(1-\gamma)} \cdot (1-\gamma)^{-\gamma/(1-\gamma)}}{\gamma}\right)$ (analogously to \check{Z} of (138) but with \tilde{c} replaced by $\tilde{c} \cdot M_{\mathbf{P}}$). Also,

$$\widehat{ISF}_k(x) = e^{-\tau_k \cdot x} \cdot \exp\left(\frac{n_k \cdot \tilde{c} \cdot M_{\mathbf{P}}}{\gamma} \cdot \left(\left(1 + \frac{\gamma-1}{\tilde{c} \cdot M_{\mathbf{P}}} \cdot \tau_k\right)^{\frac{\gamma}{\gamma-1}} - 1\right)\right), \quad x > 0.$$

For the above random variables, algorithms for simulation can be obtained by adapting e.g. the works of [339] and [340]. Within the different context of Subsection X-B, the corresponding estimators $\widehat{\Pi}_L^{improved}$ can be obtained as follows:

²⁸the latter two notations will be also used in the following Propositions 30 to 34

- (i) proceed as above but set $M_{\mathbf{P}} = 1$, replace \tilde{q}_k^* by $\overline{W}^* \cdot q_k^*$ as well as \tilde{p}_k by $p_{n,k}^{emp}$; accordingly, $\tilde{U}_k^{*n_k}$ turns into $U_k^{*n_k}$ and \widetilde{ISF}_k into ISF_k ;
- (ii) simulate independently the random variables $\widehat{W}_k^{(\ell)}$ from $U_k^{*n_k}$ ($k \in \{1, \dots, K\}$, $\ell \in \{1, \dots, L\}$);
- (iii) plug in the results of (i),(ii) into (125), (126), and (127) in order to concretely compute $\widehat{\Pi}_L^{improved}$.

From this, we can easily generate improved estimators of the power divergences $\inf_{\mathbf{Q} \in \mathfrak{Q}} D_{\tilde{c}, \varphi_\gamma}(\mathbf{Q}, \mathbb{P})$ — and more generally, improved estimators of all the infimum-quantities (e.g. Renyi divergences) respectively supremum-quantities in the parts (b) of the Proposition 29 with $A = 1$ — by simply replacing $\mathbb{P}_{X_1^n}[\xi_n^{w\mathbf{W}} \in \mathfrak{Q}]$ (respectively, its variants) by the corresponding estimator $\widehat{\Pi}_L^{improved}$. If — in the light of Remark 13(vi) — the $\mathbb{P} = (p_1, \dots, p_K)$ is a pre-given known probability vector²⁹ (rather than the limit of the vector of empirical frequencies/masses of a sequence of random variables X_i , cf. (30)), then we proceed analogously as above by replacing $p_{n,k}^{emp}$ with p_k ; correspondingly, we obtain improved estimators of all the infimum-quantities respectively supremum-quantities (e.g. Renyi entropies, diversity indices) in the parts (a) of Proposition 29 with $A = 1$. For the sake of brevity, here we only present explicitly the outcoming improved estimator for the power divergences (in the “ X_i -context”). Indeed, we simply replace the $\mathbb{P}_{X_1^n}[\xi_n^{w\mathbf{W}} \in \mathfrak{Q}]$ in the “part-(b)-version of” formula (140) (with $A = 1$) by the improved estimator $\widehat{\Pi}_L^{improved}$ obtained through (i) to (iii); for arbitrarily fixed $\tilde{c} > 0$, this leads to the *improved power-divergence estimators (BS estimators of power divergences)*

$$D_{\tilde{c}, \varphi_\gamma}(\widehat{\mathfrak{Q}}, \mathbb{P}) := -\frac{\tilde{c}}{\gamma(\gamma-1)} \left\{ 1 - \left(1 + \frac{\gamma}{\tilde{c}} \cdot \frac{1}{n} \cdot \log \widehat{\Pi}_L^{improved} \right)^{1-\gamma} \right\}. \quad (148)$$

B. Case 2

For $\gamma \in]0, 1[$, $\tilde{c} \in]0, \infty[$ and $a_{F_{\gamma, \tilde{c}}}, b_{F_{\gamma, \tilde{c}}} :=]0, \infty[$ we obtain the same $F_{\gamma, \tilde{c}}(t)$ of (135), $\mathcal{R}(F_{\gamma, \tilde{c}}) =]-\infty, \frac{\tilde{c}}{1-\gamma}[$, $]\lambda_-, \lambda_+ [=]-\infty, \frac{\tilde{c}}{1-\gamma}[$ (with $c := 0$), $]t_-^{sc}, t_+^{sc} [=]0, \infty[$ and $]a, b [=]0, \infty[$. Accordingly, we deduce from (131),(133)

$$\varphi_{\gamma, \tilde{c}}(t) := \varphi_{\gamma, \tilde{c}}^{(0)}(t) = \begin{cases} \tilde{c} \cdot \frac{t^\gamma \cdot t + \gamma - 1}{\gamma \cdot (\gamma - 1)} \in [0, \infty[, & \text{if } t \in]0, \infty[, \\ \frac{\tilde{c}}{\gamma} > 0, & \text{if } t = 0, \\ \infty, & \text{if } t \in]-\infty, 0[, \end{cases}$$

which coincides with $\tilde{c} \cdot \varphi_\gamma(t)$ for $\varphi_\gamma(t)$ from (40) and which generates the γ -corresponding power divergences given in (41). Moreover, we can derive from (130), (132)

$$\Lambda_{\gamma, \tilde{c}}(z) := \Lambda_{\gamma, \tilde{c}}^{(0)}(z) = \begin{cases} \frac{\tilde{c}}{\gamma} \cdot \left\{ \left(\frac{\gamma-1}{\tilde{c}} \cdot z + 1 \right)^{\frac{\gamma}{\gamma-1}} - 1 \right\}, & \text{if } z \in]-\infty, \frac{\tilde{c}}{1-\gamma}[, \\ \infty, & \text{if } z \in \left[\frac{\tilde{c}}{1-\gamma}, \infty \right[, \end{cases}$$

which is the cumulant generating function of the Compound-Poisson-Gamma distribution $\zeta = C(POI(\theta), GAM(\alpha, \beta))$ with $\theta = \frac{\tilde{c}}{\gamma} > 0$, rate parameter (inverse scale parameter) $\alpha = \frac{\tilde{c}}{1-\gamma} > 0$, and shape parameter $\beta = \frac{\gamma}{1-\gamma} > 0$. In other words, W has the comfortably simulable form $W = \sum_{i=1}^N Z_i$ ³⁰ for some i.i.d. sequence $(Z_i)_{i \in \mathbb{N}}$ of Gamma $GAM(\alpha, \beta)$ distributed random variables (with parameter-pair (α, β))³¹ and some independent $POI(\theta)$ -distributed random variable N . Notice that ζ is an infinitely divisible distribution (cf. Proposition 27) which is a mixture of a one-point distribution at zero and a continuous distribution on $]0, \infty[$, with $\zeta[]0, \infty[] = \mathbb{P}[W \geq 0] = 1$, $\zeta[\{0\}] = \mathbb{P}[W = 0] = e^{-\theta}$ and $\zeta[B] = \mathbb{P}[W \in B] = \int_B f_{\tilde{c}, \gamma}(u) du$ for every (measurable) subset of $]0, \infty[$ having density

$$\begin{aligned} f_{C(POI(\theta), GAM(\alpha, \beta))}(y) &:= \frac{\exp(-\alpha \cdot y - \theta)}{y} \cdot \sum_{k=1}^{\infty} \frac{\theta^k \cdot (\alpha y)^{k\beta}}{k! \cdot \Gamma(k\beta)} \cdot \mathbb{1}_{]0, \infty[}(y) \\ &= \frac{1}{y} \cdot \exp\left(-\tilde{c} \cdot \left(\frac{y}{1-\gamma} + \frac{1}{\gamma}\right)\right) \cdot \sum_{k=1}^{\infty} \frac{a_k}{k!} \cdot \tilde{c}^{k/(1-\gamma)} \cdot \gamma^{-k} \cdot (1-\gamma)^{-k\gamma/(1-\gamma)} \cdot y^{k\gamma/(1-\gamma)} \cdot \mathbb{1}_{]0, \infty[}(y) =: f_{\tilde{c}, \gamma}(y), \quad y \in \mathbb{R}, \end{aligned}$$

where $a_k := 1/\Gamma(\frac{k\gamma}{1-\gamma})$ (see e.g. [341] with a different parametrization). Concerning the important Remark 9(i), for i.i.d. copies $(W_i)_{i \in \mathbb{N}}$ of W we obtain the probability distribution $\zeta^{*n_k}[\cdot] := \mathbb{P}[\check{W} \in \cdot]$ of $\check{W} := \sum_{i \in I_k^{(n)}} W_i$ to be $C(POI(\check{\theta}), GAM(\alpha, \beta))$ with $\check{\theta} = \frac{\tilde{c} \cdot \text{card}(I_k^{(n)})}{\gamma} > 0$, $\alpha = \frac{\tilde{c}}{1-\gamma} > 0$, $\beta = \frac{\gamma}{1-\gamma} > 0$. For the bare-simulation-minimizations (respectively maximizations)

²⁹e.g. the uniform distribution \mathbb{P}^{unif} on $\{1, \dots, K\}$

³⁰with the usual convention $\sum_{i=1}^0 Z_i := 0$

³¹here and henceforth, we use the notation that a Gamma distribution $GAM(\alpha, \beta)$ with rate parameter (inverse scale parameter) $\alpha > 0$ and shape parameter $\beta > 0$ has (Lebesgue-)density $f(y) := \frac{\alpha^\beta \cdot y^{\beta-1} \cdot e^{-\alpha \cdot y}}{\Gamma(\beta)} \cdot \mathbb{1}_{]0, \infty[}(y)$, $y \in \mathbb{R}$; its cumulant generating function is $\Lambda(z) = \beta \cdot \log(\frac{\alpha}{\alpha-z})$ for $z \in]-\infty, \alpha[$ (and $\Lambda(z) = \infty$ for $z \geq \alpha$).

of the corresponding (γ -order) power divergences, Renyi divergences, Hellinger integrals and measures of entropy (diversity), we obtain from Theorem 12, Remark 13(vi), Lemma 14(a), (69), (73), (77), (78) the following:

Proposition 30: (a) Consider $\varphi := \tilde{c} \cdot \varphi_\gamma$ with $\gamma \in]0, 1[$, and let $\mathbb{P} \in \mathbb{S}_{>0}^K$ as well as $\tilde{c} > 0$ be arbitrary but fixed. Furthermore, let $W := (W_i)_{i \in \mathbb{N}}$ be an i.i.d. sequence of random variables with Compound-Poisson-Gamma distribution $\zeta = C(POI(\theta), GAM(\alpha, \beta))$ having parameters $\theta = \frac{\tilde{c}}{\gamma} > 0$, $\alpha = \frac{\tilde{c}}{1-\gamma} > 0$, $\beta = \frac{\gamma}{1-\gamma} > 0$. Then for all $A > 0$ and all $\mathfrak{Q} \subset \mathbb{S}^K$ with (7) there hold (139), (140), (144) as well as

$$\sup_{\mathfrak{Q} \in \mathcal{A} \cdot \mathfrak{Q}} H_\gamma(\mathfrak{Q}, \mathbb{P}) = \lim_{n \rightarrow \infty} A^\gamma \cdot \left(1 + \gamma \cdot \frac{1}{n} \cdot \log \mathbb{P} \left[\check{\xi}_n^{w\mathbf{W}} \in \mathfrak{Q} \right] \right)^{1-\gamma}, \quad (149)$$

$$\sup_{\mathfrak{Q} \in \mathcal{A} \cdot \mathfrak{Q}} c_1 \cdot \left(H_\gamma(\mathfrak{Q}, \mathbb{P})^{c_2} - c_3 \right) = \lim_{n \rightarrow \infty} c_1 \cdot \left\{ A^{c_2 \cdot \gamma} \cdot \left(1 + \frac{\gamma}{n} \cdot \log \mathbb{P} \left[\check{\xi}_n^{w\mathbf{W}} \in \mathfrak{Q} \right] \right)^{c_2 \cdot (1-\gamma)} - c_3 \right\}, \quad \text{if } c_1 \cdot c_2 > 0, c_3 \in \mathbb{R},$$

$$\inf_{\mathfrak{Q} \in \mathcal{A} \cdot \mathfrak{Q}} c_1 \cdot \left(H_\gamma(\mathfrak{Q}, \mathbb{P})^{c_2} - c_3 \right) = \lim_{n \rightarrow \infty} c_1 \cdot \left\{ A^{c_2 \cdot \gamma} \cdot \left(1 + \frac{\gamma}{n} \cdot \log \mathbb{P} \left[\check{\xi}_n^{w\mathbf{W}} \in \mathfrak{Q} \right] \right)^{c_2 \cdot (1-\gamma)} - c_3 \right\}, \quad \text{if } c_1 \cdot c_2 < 0, c_3 \in \mathbb{R},$$

$$\sup_{\mathfrak{Q} \in \mathcal{A} \cdot \mathfrak{Q}} c_1 \cdot \left(\left(\sum_{k=1}^K q_k^\gamma \right)^{c_2} - c_3 \right) = \lim_{n \rightarrow \infty} c_1 \cdot \left\{ K^{c_2 \cdot (1-\gamma)} \cdot A^{c_2 \cdot \gamma} \cdot \left(1 + \frac{\gamma}{n} \cdot \log \mathbb{P} \left[\check{\xi}_n^{w\mathbf{W}} \in \mathfrak{Q} \right] \right)^{c_2 \cdot (1-\gamma)} - c_3 \right\},$$

if $c_1 \cdot c_2 > 0, c_3 \in \mathbb{R}$,

$$\inf_{\mathfrak{Q} \in \mathcal{A} \cdot \mathfrak{Q}} c_1 \cdot \left(\left(\sum_{k=1}^K q_k^\gamma \right)^{c_2} - c_3 \right) = \lim_{n \rightarrow \infty} c_1 \cdot \left\{ K^{c_2 \cdot (1-\gamma)} \cdot A^{c_2 \cdot \gamma} \cdot \left(1 + \frac{\gamma}{n} \cdot \log \mathbb{P} \left[\check{\xi}_n^{w\mathbf{W}} \in \mathfrak{Q} \right] \right)^{c_2 \cdot (1-\gamma)} - c_3 \right\},$$

if $c_1 \cdot c_2 < 0, c_3 \in \mathbb{R}$.

$$\sup_{\mathfrak{Q} \in \mathcal{A} \cdot \mathfrak{Q}} \frac{1}{1-\gamma} \cdot \log \left(\sum_{k=1}^K q_k^\gamma \right) = \lim_{n \rightarrow \infty} \frac{1}{1-\gamma} \cdot \left[\log \left(A^\gamma \cdot \left(1 + \gamma \cdot \frac{1}{n} \cdot \log \mathbb{P} \left[\check{\xi}_n^{w\mathbf{W}} \in \mathfrak{Q} \right] \right)^{1-\gamma} \right) + (1-\gamma) \cdot \log(K) \right]. \quad (150)$$

From this, the BS-minimizability/maximizability of the important norms/entropies/diversity indices (E1) to (E6) follows immediately as special cases.

(b) The special case $\varphi := \tilde{c} \cdot \varphi_\gamma$ ($\gamma \in]0, 1[$) of Theorem 12 works analogously to (a), with the differences that we employ (i) additionally a sequence $(X_i)_{i \in \mathbb{N}}$ of random variables being independent of $(W_i)_{i \in \mathbb{N}}$ and satisfying condition (26) (resp. (30)), (ii) $A = 1$ (instead of arbitrary $A > 0$), (iii) $\mathbb{P}_{X_T^\dagger}[\cdot]$ (instead of $\mathbb{P}[\cdot]$), (iv) $\xi_{n,\mathbf{X}}^{w\mathbf{W}}$ (instead of $\xi_n^{w\mathbf{W}}$), (v) $\check{\xi}_{n,\mathbf{X}}^{w\mathbf{W}}$ (instead of $\check{\xi}_n^{w\mathbf{W}}$), and (vi) $\check{\xi}_{n,\mathbf{X}}^{w\mathbf{W}}$ (instead of $\check{\xi}_n^{w\mathbf{W}}$).

Within the context of Subsection X-A, for the concrete simulative estimation $\widehat{D}_{\tilde{c}, \varphi_\gamma}(\mathfrak{Q}, \mathbf{P})$ via (116) and (118), we derive — in terms of $M_{\mathbf{P}} := \sum_{i=1}^K p_i > 0$, $n_k = n \cdot \tilde{p}_k \in \mathbb{N}$ and \tilde{q}_k^* from proxy method 1 or 2 — that

$\tilde{U}_k^{*n_k} = C(POI(n_k \cdot \tilde{\theta}), GAM(\frac{\tilde{c} \cdot M_{\mathbf{P}}}{1-\gamma} - \tau_k, \frac{\gamma}{1-\gamma}))$ with $\tilde{\theta} := \frac{\tilde{c} \cdot M_{\mathbf{P}}}{\gamma} \cdot \left(\frac{(\gamma-1) \cdot \tau_k}{\tilde{c} \cdot M_{\mathbf{P}}} + 1 \right)^{\gamma/(\gamma-1)}$ and $\tau_k = \tilde{c} \cdot M_{\mathbf{P}} \cdot \frac{1 - \left(\frac{\tilde{q}_k^*}{\tilde{p}_k}\right)^{\gamma-1}}{1-\gamma}$ for $\tilde{q}_k^* > 0$. Furthermore,

$$\widetilde{ISF}_k(x) = e^{-\tau_k x} \cdot \exp \left(\frac{n_k \cdot \tilde{c} \cdot M_{\mathbf{P}}}{\gamma} \cdot \left(\left(1 + \frac{\gamma-1}{\tilde{c} \cdot M_{\mathbf{P}}} \cdot \tau_k \right)^{\frac{\gamma}{\gamma-1}} - 1 \right) \right), \quad x \geq 0.$$

Within the different context of Subsection X-B, the corresponding estimators $\widehat{\Pi}_L^{improved}$ can be obtained analogously to the last paragraph of Subsection XII-A, with Proposition 30 instead of Proposition 29 (and (148) remains the same).

C. Case 3

For $\gamma > 2$, $\tilde{c} \in]0, \infty[$ and $]a_{F_{\gamma, \tilde{c}}}, b_{F_{\gamma, \tilde{c}}}[:=]0, \infty[$ we obtain

$$F_{\gamma, \tilde{c}}(t) := \begin{cases} \frac{\tilde{c}}{\gamma-1} \cdot (t^{\gamma-1} - 1), & \text{if } t \in]0, \infty[, \\ -\frac{\tilde{c}}{\gamma-1}, & \text{if } t = 0, \\ -\infty, & \text{if } t \in]-\infty, 0[, \end{cases} \quad (151)$$

$\mathcal{R}(F_{\gamma, \tilde{c}}) = \left[-\frac{\tilde{c}}{\gamma-1}, \infty[\right], \lambda_-, \lambda_+[=] - \frac{\tilde{c}}{\gamma-1}, \infty[$ (with $c := 0$), $]t_-^{sc}, t_+^{sc}[=]0, \infty[$ and $]a, b[=]-\infty, \infty[$. Correspondingly, we deduce from (131),(133)

$$\varphi_{\gamma, \tilde{c}}(t) := \varphi_{\gamma, \tilde{c}}^{(0)}(t) = \begin{cases} \tilde{c} \cdot \frac{t^{\gamma-\gamma \cdot t + \gamma - 1}}{\gamma \cdot (\gamma-1)} \in [0, \infty[, & \text{if } t \in]0, \infty[, \\ \frac{\tilde{c}}{\gamma} > 0, & \text{if } t = 0, \\ \frac{\tilde{c}}{\gamma} - \frac{\tilde{c}}{\gamma-1} \cdot t \in]0, \infty[, & \text{if } t \in]-\infty, 0[, \end{cases} \quad (152)$$

which coincides with $\tilde{c} \cdot \varphi_\gamma(t)$ for $\varphi_\gamma(t)$ from (40) and which generates the γ -corresponding power divergences given in (41). Moreover, we can derive from formula (130) (see also (132))

$$\Lambda_{\gamma, \tilde{c}}(z) := \Lambda_{\gamma, \tilde{c}}^{(0)}(z) = \begin{cases} \frac{\tilde{c}}{\gamma} \cdot \left\{ \left(\frac{\gamma-1}{\tilde{c}} \cdot z + 1 \right)^{\frac{\gamma}{\gamma-1}} - 1 \right\}, & \text{if } z \in] -\frac{\tilde{c}}{\gamma-1}, \infty[, \\ -\frac{\tilde{c}}{\gamma} < 0, & \text{if } z = -\frac{\tilde{c}}{\gamma-1}, \\ \infty, & \text{if } z \in] -\infty, -\frac{\tilde{c}}{\gamma-1}]. \end{cases} \quad (153)$$

The latter is the cumulant generating function of a ‘‘distorted stable distribution’’ $\zeta[\cdot] = \mathbb{P}[W \in \cdot]$ of a random variable W , which can be constructed as follows: let Z be an auxiliary random variable (having density f_Z and support $\text{supp}(Z) =]-\infty, \infty[$) of a stable law with parameter-quadruple $(\frac{\gamma}{\gamma-1}, 1, 0, \frac{\tilde{c}^{1/(1-\gamma)} \cdot (\gamma-1)^{\gamma/(\gamma-1)}}{\gamma})$ in terms of the above-mentioned ‘‘form-B notation’’; by applying a general Laplace-transform result on p. 112 of [338] we can derive

$$M_Z(z) := E_{\mathbb{P}}[\exp(z \cdot Z)] = \int_{-\infty}^{\infty} \exp(z \cdot y) \cdot f_Z(y) dy = \begin{cases} \exp\left(\frac{\tilde{c}^{1/(1-\gamma)} \cdot (\gamma-1)^{\gamma/(\gamma-1)}}{\gamma} \cdot (-z)^\alpha\right), & \text{if } z \in]-\infty, 0], \\ \infty, & \text{if } z \in]0, \infty[, \end{cases} \quad (154)$$

where $\alpha := \frac{\gamma}{\gamma-1} \in]1, 2[$. Since $0 \notin \text{int}(\text{dom}(M_Z))$ (and thus, Z does not have light-tails) we have to distort the density in order to extend the effective domain. Accordingly, let W be a random variable having density

$$f_W(y) := \frac{\exp\{\frac{y \cdot \tilde{c}}{\gamma-1}\}}{\exp\{\tilde{c}/\gamma\}} \cdot f_Z(-y), \quad y \in \mathbb{R}. \quad (155)$$

Then one can straightforwardly deduce from (154) that $\int_{-\infty}^{\infty} f_W(y) dy = 1$ and that

$$M_W(z) := E_{\mathbb{P}}[\exp(z \cdot W)] = \int_{-\infty}^{\infty} \exp(z \cdot y) \cdot f_W(y) dy = \begin{cases} \exp\left(\frac{\tilde{c}}{\gamma} \cdot \left\{ \left(\frac{\gamma-1}{\tilde{c}} \cdot z + 1 \right)^{\frac{\gamma}{\gamma-1}} - 1 \right\}\right), & \text{if } z \in [-\frac{\tilde{c}}{\gamma-1}, \infty[, \\ \infty, & \text{if } z \in]-\infty, -\frac{\tilde{c}}{\gamma-1}]. \end{cases}$$

Notice that ζ is an infinitely divisible (cf. Proposition 27) continuous distribution with density f_W , and that $\zeta[]0, \infty[] = \mathbb{P}[W > 0] = \int_0^{\infty} f_W(u) du \in]0, 1[$, $\zeta[\{0\}] = \mathbb{P}[W = 0] = 0$. Concerning the important Remark 9(i), for i.i.d. copies $(W_i)_{i \in \mathbb{N}}$ of W we obtain the probability distribution $\zeta^{*n_k}[\cdot] := \mathbb{P}[\tilde{W} \in \cdot]$ of $\tilde{W} := \sum_{i \in I_k^{(n)}} W_i$ having the density

$$f_{\tilde{W}}(y) := \frac{\exp\{\frac{y \cdot \tilde{c}}{\gamma-1}\}}{\exp\{\tilde{c} \cdot \text{card}(I_k^{(n)})/\gamma\}} \cdot f_Z(-y), \quad y \in \mathbb{R}, \quad (156)$$

where \tilde{Z} is a random variable with density $f_{\tilde{Z}}$ of a stable law with parameters $(\frac{\gamma}{\gamma-1}, 1, 0, \text{card}(I_k^{(n)}) \cdot \frac{\tilde{c}^{1/(1-\gamma)} \cdot (\gamma-1)^{\gamma/(\gamma-1)}}{\gamma})$. For the bare-simulation-minimizations (respectively maximizations) of the corresponding (γ -order) power divergences, Renyi divergences, Hellinger integrals and measures of entropy (diversity), we obtain from Theorem 12 respectively Remark 13(vi), Lemma 14(a), (69), (73), (77), (78) the following

Proposition 31: (a) Consider $\varphi := \tilde{c} \cdot \varphi_\gamma$ with $\gamma > 2$, and let $\mathbb{P} \in \mathbb{S}_{>0}^K$ as well as $\tilde{c} > 0$ be arbitrary but fixed. Furthermore, let $W := (W_i)_{i \in \mathbb{N}}$ be an i.i.d. sequence of real-valued random variables having density (155). Then for all $A > 0$ and $\mathbf{\Omega} \subset \mathbb{S}^K$ with (7) there hold all the BS-extremizabilities (139) to (146) as well as (150). From this, the BS-minimizability/maximizability of the important norms/entropies/diversity indices (E1) to (E6) follow immediately as special cases.

(b) The special case $\varphi := \tilde{c} \cdot \varphi_\gamma$ ($\gamma > 2$) of Theorem 12 works analogously to (a), with the differences that we employ (i) additionally a sequence $(X_i)_{i \in \mathbb{N}}$ of random variables being independent of $(W_i)_{i \in \mathbb{N}}$ and satisfying condition (26) (resp. (30)), (ii) $A = 1$ (instead of arbitrary $A > 0$), (iii) $\mathbb{P}_{X_{\mathbb{P}}}[\cdot]$ (instead of $\mathbb{P}[\cdot]$), (iv) $\xi_{n, \mathbf{X}}^{w \mathbf{W}}$ (instead of $\xi_n^{w \mathbf{W}}$), (v) $\check{\xi}_{n, \mathbf{X}}^{w \mathbf{W}}$ (instead of $\check{\xi}_n^{w \mathbf{W}}$), and (vi) $\check{\xi}_{n, \mathbf{X}}^{w \mathbf{W}}$ (instead of $\check{\xi}_n^{w \mathbf{W}}$).

Within the context of Subsection X-A, for the concrete simulative estimation $\widehat{D}_{\tilde{c}, \varphi_\gamma}(\mathbf{\Omega}, \mathbf{P})$ via (116) and (118), we derive — in terms of $M_{\mathbf{P}} := \sum_{i=1}^K p_i > 0$, $n_k = n \cdot \tilde{p}_k \in \mathbb{N}$ and \tilde{q}_k^* from proxy method 1 or 2 — that $\tilde{U}_k^{*n_k}$ has the (Lebesgue-)density

$$f_{\tilde{U}_k^{*n_k}}(x) := \frac{\exp((\tau_k + \frac{\tilde{c} \cdot M_{\mathbf{P}}}{\gamma-1}) \cdot x)}{\exp\left(n_k \cdot \frac{\tilde{c} \cdot M_{\mathbf{P}}}{\gamma} \cdot \left(1 + \frac{\gamma-1}{\tilde{c} \cdot M_{\mathbf{P}}} \cdot \tau_k\right)^{\gamma/(\gamma-1)}\right)} \cdot f_{\tilde{Z}}(-x), \quad x \in \mathbb{R},$$

where $\tau_k = -\frac{\tilde{c} \cdot M_{\mathbf{P}}}{\gamma-1} \cdot \left(1 - \left(\frac{\tilde{q}_k^*}{p_k}\right)^{\gamma-1}\right) \cdot \mathbb{1}_{]0, \infty[}(\tilde{q}_k^*)$ for $\tilde{q}_k^* \in \mathbb{R}$, and \check{Z} is a random variable with density $f_{\check{Z}}$ of a stable law with parameter-quadruple $(\frac{\gamma}{\gamma-1}, 1, 0, n_k \cdot \frac{(\tilde{c} \cdot M_{\mathbf{P}})^{1/(1-\gamma)} \cdot (\gamma-1)^{\gamma/(\gamma-1)}}{\gamma})$ (analogously to \check{Z} of (156) but with \tilde{c} replaced by $\tilde{c} \cdot M_{\mathbf{P}}$).

Furthermore, $\widehat{ISF}_k(x) = e^{-\tau_k \cdot x} \cdot \exp\left(\frac{n_k \cdot \tilde{c} \cdot M_{\mathbf{P}}}{\gamma} \cdot \left(\left(1 + \frac{\gamma-1}{\tilde{c} \cdot M_{\mathbf{P}}} \cdot \tau_k\right)^{\frac{\gamma}{\gamma-1}} - 1\right)\right)$, $x \in \mathbb{R}$. For the above random variables, algorithms for simulation can be obtained by adapting e.g. the works of [339] and [340]. Within the different context of Subsection X-B, the corresponding estimators $\widehat{\Pi}_L^{improved}$ can be obtained analogously to the last paragraph of Subsection XII-A, with Proposition 31 instead of Proposition 29 (and (148) remains the same).

D. Case 4

For $\gamma = 2$, $\tilde{c} \in]0, \infty[$ and $]a_{F_{\gamma, \tilde{c}}}, b_{F_{\gamma, \tilde{c}}}[=]-\infty, \infty[$ we define $F_{2, \tilde{c}}(t) := \tilde{c} \cdot (t-1)$ ($t \in]-\infty, \infty[$). Clearly, $\mathcal{R}(F_{2, \tilde{c}}) =]-\infty, \infty[$ and $F_{2, \tilde{c}} \in \mathfrak{F}$. Since $F_{2, \tilde{c}}(1) = 0$, let us choose the natural anchor point $c := 0$, which leads to $] \lambda_-, \lambda_+[=]-\infty, \infty[$, $]t_-^{sc}, t_+^{sc}[=]-\infty, \infty[$ and $]a, b[=]-\infty, \infty[$. By using $F_{2, \tilde{c}}^{-1}(x) = 1 + \frac{x}{\tilde{c}}$ for $x \in \text{int}(\mathcal{R}(F_{2, \tilde{c}}))$, from (131),(133) we deduce

$$\varphi_{2, \tilde{c}}(t) := \varphi_{2, \tilde{c}}^{(0)}(t) = \tilde{c} \cdot \frac{(t-1)^2}{2} \in [0, \infty[, \quad t \in]-\infty, \infty[,$$

which coincides with $\tilde{c} \cdot \varphi_2(t)$ for $\varphi_2(t)$ from (40) which generates the \tilde{c} -fold of the half Pearson-chisquare divergence given in the sixth line of (41). Moreover, from (130),(132) we derive

$$\Lambda_{2, \tilde{c}}(z) := \Lambda_{2, \tilde{c}}^{(0)}(z) = \frac{z^2}{2\tilde{c}} + z = \frac{\tilde{c}}{2} \cdot \left\{ \left(\frac{1}{\tilde{c}} \cdot z + 1 \right)^2 - 1 \right\}, \quad z \in]-\infty, \infty[,$$

which is the well-known cumulant generating function of the Normal distribution (Gaussian distribution) $\zeta = NOR(1, \frac{1}{\tilde{c}})$ with mean 1 and variance $\frac{1}{\tilde{c}}$. Notice that ζ is an infinitely divisible (cf. Proposition 27) continuous distribution with density $f_{NOR(1, \frac{1}{\tilde{c}})}(y) := \sqrt{\frac{\tilde{c}}{2\pi}} \cdot \exp\left(-\frac{\tilde{c}(y-1)^2}{2}\right)$ ($y \in \mathbb{R}$) and that $\zeta[]0, \infty[= \mathbb{P}[W > 0] = \int_0^\infty f_{NOR(1, \frac{1}{\tilde{c}})}(u) du \in]0, 1[$, $\zeta[\{0\}] = \mathbb{P}[W = 0] = 0$. Concerning the important Remark 9(i), for i.i.d. copies $(W_i)_{i \in \mathbb{N}}$ of W the probability distribution $\zeta^{*n_k}[\cdot] := \mathbb{P}[\check{W} \in \cdot]$ of $\check{W} := \sum_{i \in I_k^{(n)}} W_i$ is $NOR(\text{card}(I_k^{(n)}), \frac{\text{card}(I_k^{(n)})}{\tilde{c}})$. For the desired bare-simulation-optimizations we obtain from Theorem 12 respectively Remark 13(vi), Lemma 14(a), (69), (73), (77), (78) the following

Proposition 32: (a) Consider $\varphi := \tilde{c} \cdot \varphi_\gamma$ with $\gamma = 2$, and let $\mathbb{P} \in \mathbb{S}_{>0}^K$ as well as $\tilde{c} > 0$ be arbitrary but fixed. Furthermore, let $W := (W_i)_{i \in \mathbb{N}}$ be an i.i.d. sequence of real-valued random variables with probability distribution $\zeta = NOR(1, \frac{1}{\tilde{c}})$. Then for all $A > 0$ and $\mathfrak{Q} \subset \mathbb{S}^K$ with (7) there hold all the BS-extremizabilites (139) to (146) as well as (150) with plugging-in $\gamma = 2$. From this, the BS-minimizability/maximizability of the important norms/entropies/diversity indices (E1) to (E6) follow as special cases. By Remark 16(c), one can even take $A < 0$ in (139) to (146) and (150) as well as in (E1),(E2),(E4),(E6). (b) The special case $\varphi := \tilde{c} \cdot \varphi_\gamma$ ($\gamma = 2$) of Theorem 12 works analogously to (a), with the differences that we employ (i) additionally a sequence $(X_i)_{i \in \mathbb{N}}$ of random variables being independent of $(W_i)_{i \in \mathbb{N}}$ and satisfying condition (26) (resp. (30)), (ii) $A = 1$ (instead of arbitrary $A > 0$), (iii) $\mathbb{P}_{X_{\mathbf{P}}}[\cdot]$ (instead of $\mathbb{P}[\cdot]$), (iv) $\xi_{n, \mathbf{X}}^{w\mathbf{W}}$ (instead of $\xi_n^{w\mathbf{W}}$), (v) $\check{\xi}_{n, \mathbf{X}}^{w\mathbf{W}}$ (instead of $\check{\xi}_n^{w\mathbf{W}}$), and (vi) $\check{\xi}_{n, \mathbf{X}}^{w\mathbf{W}}$ (instead of $\check{\xi}_n^{w\mathbf{W}}$).

Within the context of Subsection X-A, for the concrete simulative estimation $\widehat{D}_{\tilde{c}, \varphi_2}(\mathfrak{Q}, \mathbf{P})$ via (116) and (118), we derive — in terms of $M_{\mathbf{P}} := \sum_{i=1}^K p_i > 0$, $n_k = n \cdot \tilde{p}_k \in \mathbb{N}$ and \tilde{q}_k^* from proxy method 1 or 2 — that $\tilde{U}_k^{*n_k} = NOR(n_k \cdot (1 + \frac{\tau_k}{\tilde{c} \cdot M_{\mathbf{P}}}), \frac{n_k}{\tilde{c} \cdot M_{\mathbf{P}}})$ with $\tau_k = \tilde{c} \cdot M_{\mathbf{P}} \cdot \left(\frac{\tilde{q}_k^*}{p_k} - 1\right)$ for $\tilde{q}_k^* \in \mathbb{R}$. Moreover, for all $x \in \mathbb{R}$ one obtains $\widehat{ISF}_k(x) = \exp\left(\frac{n_k \cdot \tau_k^2}{2\tilde{c} \cdot M_{\mathbf{P}}} - (x - n_k) \cdot \tau_k\right)$. Within the different context of Subsection X-B, the corresponding estimators $\widehat{\Pi}_L^{improved}$ can be obtained analogously to the last paragraph of Subsection XII-A, with Proposition 32 instead of Proposition 29 (and (148) remains the same).

E. Case 5

For $\gamma = 0$, $\tilde{c} \in]0, \infty[$ and $]a_{F_{\gamma, \tilde{c}}}, b_{F_{\gamma, \tilde{c}}}[:=]0, \infty[$ we obtain the same $F_{\gamma, \tilde{c}}(t)$ of (135), $\mathcal{R}(F_{\gamma, \tilde{c}}) =]-\infty, \frac{\tilde{c}}{1-\tilde{c}}[$, $] \lambda_-, \lambda_+[=]-\infty, \frac{\tilde{c}}{1-\tilde{c}}[$ (with $c := 0$), $]t_-^{sc}, t_+^{sc}[=]0, \infty[$ and $]a, b[=]0, \infty[$. By using $F_{0, \tilde{c}}^{-1}(x) = \frac{1}{1-\frac{x}{\tilde{c}}}$ for $x \in \text{int}(\mathcal{R}(F_{0, \tilde{c}})) =]-\infty, \tilde{c}[$, we can deduce from (131),(133)

$$\varphi_{0, \tilde{c}}(t) := \varphi_{0, \tilde{c}}^{(0)}(t) = \begin{cases} \tilde{c} \cdot (-\log t + t - 1) \in [0, \infty[, & \text{if } t \in]0, \infty[, \\ \infty, & \text{if } t \in]-\infty, 0], \end{cases}$$

which coincides with $\tilde{c} \cdot \varphi_0(t)$ for the generator $\varphi_0(t)$ from (40) which generates the reverse Kullback-Leibler divergence (reverse relative entropy) given in the second line of (41) with $\tilde{c} = 1$. Furthermore, we can derive from (130),(132)

$$\Lambda_{0, \tilde{c}}(z) := \Lambda_{0, \tilde{c}}^{(0)}(z) = \begin{cases} -\tilde{c} \cdot \log\left(1 - \frac{z}{\tilde{c}}\right), & \text{if } z \in]-\infty, \tilde{c}[, \\ \infty, & \text{if } z \in [\tilde{c}, \infty[, \end{cases}$$

which is the cumulant generating function of the Gamma distribution $\zeta = GAM(\tilde{c}, \tilde{c})$ with rate parameter (inverse scale parameter) \tilde{c} and shape parameter \tilde{c} ; the special case $\tilde{c} = 1$ leads to $\zeta = GAM(1, 1) = EXP(1)$ being the exponential distribution with mean 1. Notice that ζ is an infinitely divisible (cf. Proposition 27) continuous distribution with density $f(y) := \frac{\tilde{c}^{\tilde{c}} \cdot y^{\tilde{c}-1} \cdot e^{-\tilde{c} \cdot y}}{\Gamma(\tilde{c})} \cdot \mathbb{1}_{]0, \infty[}(y)$ ($y \in \mathbb{R}$), and that $\zeta[]0, \infty[= \mathbb{P}[W > 0] = 1$. Concerning the important Remark 9(i), for i.i.d. copies $(W_i)_{i \in \mathbb{N}}$ of W the probability distribution $\zeta^{*n_k}[\cdot] := \mathbb{P}[\tilde{W} \in \cdot]$ of $\tilde{W} := \sum_{i \in I_k^{(n)}} W_i$ is $GAM(\tilde{c}, \tilde{c} \cdot card(I_k^{(n)}))$. For the desired bare-simulation-optimizations we obtain from Theorem 12, Remark 13(vi) and Lemma 14(c) the following

Proposition 33: (a) Consider $\varphi := \tilde{c} \cdot \varphi_\gamma$ with $\gamma = 0$, and let $\mathbb{P} \in \mathbb{S}_{>0}^K$ as well as $\tilde{c} > 0$ be arbitrary but fixed. Furthermore, let $W := (W_i)_{i \in \mathbb{N}}$ be an i.i.d. sequence of non-negative real-valued random variables with Gamma distribution $\zeta = GAM(\tilde{c}, \tilde{c})$. Then for all $A > 0$ and all $\Omega \subset \mathbb{S}_{>0}^K$ with (7) there hold the BS minimizabilites (cf. (2))

$$\inf_{\mathbf{Q} \in A \cdot \Omega} D_{\tilde{c}, \varphi_0}(\mathbf{Q}, \mathbb{P}) = - \lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P} \left[\xi_n^{w\mathbf{W}} \in \Omega \right] + \tilde{c} \cdot (A - 1 - \log A),$$

$$\inf_{\mathbf{Q} \in A \cdot \Omega} \tilde{I}(\mathbf{Q}, \mathbb{P}) = \inf_{\mathbf{Q} \in A \cdot \Omega} \sum_{k=1}^K p_k \cdot \log \left(\frac{p_k}{q_k} \right) = - \lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P} \left[\xi_n^{w\mathbf{W}} \in \Omega \right] - \log A.$$

(b) The special case $\varphi := \tilde{c} \cdot \varphi_\gamma$ ($\gamma = 0$) of Theorem 12 works analogously to (a), with the differences that we employ (i) additionally a sequence $(X_i)_{i \in \mathbb{N}}$ of random variables being independent of $(W_i)_{i \in \mathbb{N}}$ and satisfying condition (26) (resp. (30)), (ii) $A = 1$ (instead of arbitrary $A > 0$), (iii) $\mathbb{P}_{X_1^n}[\cdot]$ (instead of $\mathbb{P}[\cdot]$), and (iv) $\xi_{n, \mathbf{X}}^{w\mathbf{W}}$ (instead of $\xi_n^{w\mathbf{W}}$).

Within the context of Subsection X-A, for the concrete simulative estimation $\widehat{D_{\tilde{c}, \varphi_0}}(\Omega, \mathbf{P})$ via (116) and (118), we derive — in terms of $M_{\mathbf{P}} := \sum_{i=1}^K p_i > 0$, $n_k = n \cdot \tilde{p}_k \in \mathbb{N}$ and \tilde{q}_k^* from proxy method 1 or 2 — that $\tilde{U}_k^{*n_k} = GAM(\tilde{c} \cdot M_{\mathbf{P}} - \tau_k, n_k \cdot \tilde{c} \cdot M_{\mathbf{P}})$, with $\tau_k = \tilde{c} \cdot M_{\mathbf{P}} \cdot (1 - \frac{\tilde{p}_k}{\tilde{q}_k^*})$ for $\tilde{q}_k^* > 0$ (the latter is equivalent to $\tau_k < \tilde{c} \cdot M_{\mathbf{P}}$). Moreover, for all $x > 0$ one gets $\widehat{ISF}_k(x) = \left(\frac{\tilde{c} \cdot M_{\mathbf{P}}}{\tilde{c} \cdot M_{\mathbf{P}} - \tau_k} \right)^{n_k \cdot \tilde{c} \cdot M_{\mathbf{P}}} \cdot e^{-\tau_k \cdot x}$. Within the different context of Subsection X-B, the corresponding estimators $\widehat{\Pi}_L^{improved}$ can be obtained analogously to the last paragraph of Subsection XII-A, with Proposition 33 instead of Proposition 29 and with $D_{\tilde{c}, \varphi_0}(\Omega, \mathbb{P}) := -\frac{1}{n} \log \widehat{\Pi}_L^{improved}$ instead of (148).

F. Case 6

1) *Case 6a: anchor point $c = 0$:* for $\gamma = 1$, $\tilde{c} \in]0, \infty[$ and $]a_{F_{1, \tilde{c}}}, b_{F_{1, \tilde{c}}}[=]0, \infty[$ we define

$$F_{1, \tilde{c}}(t) := \begin{cases} \tilde{c} \cdot \log t = \lim_{\gamma \rightarrow 1} F_{\gamma, \tilde{c}}(t), & \text{if } t \in]0, \infty[, \\ -\infty, & \text{if } t \in]-\infty, 0]. \end{cases}$$

Clearly, $\mathcal{R}(F_{1, \tilde{c}}) =]-\infty, \infty[$ and $F_{1, \tilde{c}} \in \mathcal{F}$. Since $F_{1, \tilde{c}}(1) = 0$, let us *first* choose the natural anchor point $c := 0$, which leads to $] \lambda_-, \lambda_+ [= \text{int}(\mathcal{R}(F_{1, \tilde{c}})) =]-\infty, \infty[$, $]t_-^{sc}, t_+^{sc}[=]0, \infty[$ and $]a, b[=]0, \infty[$. By using $F_{1, \tilde{c}}^{-1}(x) = \exp(\frac{x}{\tilde{c}})$ for $x \in \mathcal{R}(F_{1, \tilde{c}})$, we deduce from (131),(133)

$$\varphi_{1, \tilde{c}}(t) := \varphi_{1, \tilde{c}}^{(0)}(t) := \begin{cases} \tilde{c} \cdot (t \cdot \log t + 1 - t) \in [0, \infty[, & \text{if } t \in]0, \infty[, \\ \tilde{c}, & \text{if } t = 0, \\ \infty, & \text{if } t \in]-\infty, 0], \end{cases}$$

which coincides with $\tilde{c} \cdot \varphi_1(t)$ for the generator $\varphi_1(t)$ from (40) which generates the Kullback-Leibler divergence (relative entropy) given in the fourth line of (41) with $\tilde{c} = 1$. Moreover, we derive from (130),(132)

$$\Lambda_{1, \tilde{c}}(z) := \Lambda_{1, \tilde{c}}^{(0)}(z) := \int_0^z F_{1, \tilde{c}}^{-1}(u) du = \tilde{c} \cdot \left(\exp\left(\frac{z}{\tilde{c}}\right) - 1 \right), \quad z \in]-\infty, \infty[,$$

which is the cumulant generating function of $\zeta = \frac{1}{\tilde{c}} \cdot POI(\tilde{c})$ being the “ $\frac{1}{\tilde{c}}$ -fold Poisson distribution with mean \tilde{c} ” which means that $W = \frac{1}{\tilde{c}} \cdot Z$ for a $POI(\tilde{c})$ -distributed random variable Z ; for the special case $\tilde{c} = 1$ one particularly gets $\zeta = POI(1)$ to be the Poisson distribution with mean 1. Notice that ζ is an infinitely divisible (cf. Proposition 27) discrete distribution with the frequencies $\mathbb{P}[W = \ell \cdot \frac{1}{\tilde{c}}] = \exp(-\tilde{c}) \cdot \frac{\tilde{c}^\ell}{\ell!}$ for all nonnegative integers $\ell \in \mathbb{N}_0$ (and zero elsewhere). Hence, $\mathbb{P}[W \geq 0] = 1$, $\mathbb{P}[W = 0] = \exp(-\tilde{c})$. Concerning the important Remark 9(i), for i.i.d. copies $(W_i)_{i \in \mathbb{N}}$ of W the probability distribution of $\tilde{W} := \sum_{i \in I_k^{(n)}} W_i$ is $\frac{1}{\tilde{c}} \cdot POI(\tilde{c} \cdot card(I_k^{(n)}))$.

For the desired bare-simulation-optimizations we obtain from Theorem 12, Remark 13(vi) and Lemma 14(b) the following

Proposition 34: (a) Consider $\varphi := \tilde{c} \cdot \varphi_\gamma$ with $\gamma = 1$, and let $\mathbb{P} \in \mathbb{S}_{>0}^K$ as well as $\tilde{c} > 0$ be arbitrary but fixed. Furthermore, let $W := (W_i)_{i \in \mathbb{N}}$ be an i.i.d. sequence of non-negative real-valued random variables with distribution $\zeta = \frac{1}{\tilde{c}} \cdot \text{POI}(\tilde{c})$. Then for all $A > 0$ and all $\mathfrak{Q} \subset \mathbb{S}^K$ with (7) there holds

$$-\lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P} \left[\xi_n^{w\mathbf{W}} \in \mathfrak{Q} \right] = \inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} \tilde{c} \cdot \left[1 - A \cdot \exp \left(-\frac{1}{A \cdot \tilde{c}} \cdot D_{\tilde{c} \cdot \varphi_1}(\mathbf{Q}, \mathbb{P}) + \frac{1}{A} - 1 \right) \right]$$

and the BS minimizabilities/maximizabilites (cf. Definition 1)

$$\begin{aligned} \inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} D_{\tilde{c} \cdot \varphi_1}(\mathbf{Q}, \mathbb{P}) &= \lim_{n \rightarrow \infty} \tilde{c} \cdot \left\{ 1 - A \cdot \left[1 + \log \left(\frac{1}{A} \cdot \left(1 + \frac{1}{\tilde{c}} \cdot \frac{1}{n} \cdot \log \mathbb{P} \left[\xi_n^{w\mathbf{W}} \in \mathfrak{Q} \right] \right) \right) \right] \right\}, \\ \inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} I(\mathbf{Q}, \mathbb{P}) &= \inf_{\mathbf{Q} \in A \cdot \mathfrak{Q}} \sum_{k=1}^K q_k \cdot \log \left(\frac{q_k}{p_k} \right) = -\lim_{n \rightarrow \infty} A \cdot \log \left(\frac{1}{A} \cdot \left(1 + \frac{1}{n} \cdot \log \mathbb{P} \left[\xi_n^{w\mathbf{W}} \in \mathfrak{Q} \right] \right) \right), \\ \max_{\mathbf{Q} \in A \cdot \mathfrak{Q}} \mathcal{E}^{Sh}(\mathbf{Q}) &= \max_{\mathbf{Q} \in A \cdot \mathfrak{Q}} (-1) \cdot \sum_{k=1}^K q_k \cdot \log(q_k) = \lim_{n \rightarrow \infty} A \cdot \log K + A \cdot \log \left(\frac{1}{A} \cdot \left(1 + \frac{1}{n} \cdot \log \mathbb{P} \left[\xi_n^{w\mathbf{W}} \in \mathfrak{Q} \right] \right) \right), \quad (157) \\ \max_{\mathbf{Q} \in A \cdot \mathfrak{Q}} \mathcal{E}^{gSM2}(\mathbf{Q}) &= \max_{\mathbf{Q} \in A \cdot \mathfrak{Q}} \frac{1}{1-s} \cdot \exp \left\{ (s-1) \cdot \sum_{k=1}^K q_k \cdot \log(q_k) - 1 \right\} \\ &= \lim_{n \rightarrow \infty} \frac{1}{1-s} \cdot \exp \left\{ (1-s) \cdot \left[A \cdot \log K + A \cdot \log \left(\frac{1}{A} \cdot \left(1 + \frac{1}{n} \cdot \log \mathbb{P} \left[\xi_n^{w\mathbf{W}} \in \mathfrak{Q} \right] \right) \right) \right] - 1 \right\}, \quad s \in]0, 1[\cup]1, \infty[. \end{aligned}$$

The special subcase $A = 1$ in (157) (and thus, \mathbf{Q} is a probability vector) corresponds to the *maximum entropy problem* for the Shannon entropy $\mathcal{E}^{Sh}(\cdot)$. This can hence be tackled by our BS approach for almost arbitrary sets \mathfrak{Q} of probability vectors.

(b) The special case $\varphi := \tilde{c} \cdot \varphi_\gamma$ ($\gamma = 1$) of Theorem 12 works analogously to (a), with the differences that we employ (i) additionally a sequence $(X_i)_{i \in \mathbb{N}}$ of random variables being independent of $(W_i)_{i \in \mathbb{N}}$ and satisfying condition (26) (resp. (30)), (ii) $A = 1$ (instead of arbitrary $A > 0$), (iii) $\mathbb{P}_{X_1^n}[\cdot]$ (instead of $\mathbb{P}[\cdot]$), (iv) $\xi_{n,\mathbf{X}}^{w\mathbf{W}}$ (instead of $\xi_n^{w\mathbf{W}}$), (v) $\xi_{n,\mathbf{X}}^{w\mathbf{W}}$ (instead of $\xi_n^{w\mathbf{W}}$), and (vi) $\xi_{n,\mathbf{X}}^{w\mathbf{W}}$ (instead of $\xi_n^{w\mathbf{W}}$).

Within the context of Subsection X-A, for the concrete simulative estimation $\widehat{D_{\tilde{c} \cdot \varphi_1}}(\mathfrak{Q}, \mathbf{P})$ via (116) and (118), we derive — in terms of $M_{\mathbf{P}} := \sum_{i=1}^K p_i > 0$, $n_k = n \cdot \tilde{p}_k \in \mathbb{N}$ and \tilde{q}_k^* from proxy method 1 or 2 — that $\tilde{U}_k^{*n_k}$ is the probability distribution $\frac{1}{\tilde{c} \cdot M_{\mathbf{P}}} \cdot \text{POI} \left(n_k \cdot \tilde{c} \cdot M_{\mathbf{P}} \cdot \exp \left(\frac{\tau_k}{\tilde{c} \cdot M_{\mathbf{P}}} \right) \right)$ with support on the lattice $\left\{ \frac{j}{\tilde{c} \cdot M_{\mathbf{P}}}, j \in \mathbb{N}_0 \right\}$, where $\tau_k = \tilde{c} \cdot \log \left(\frac{\tilde{q}_k^*}{\tilde{p}_k} \right)$ for $\tilde{q}_k^* > 0$. Moreover, for all $j \in \mathbb{N}_0$ we obtain (by setting $x := \frac{j}{\tilde{c} \cdot M_{\mathbf{P}}}$)

$$\widehat{ISF}_k \left(\frac{j}{\tilde{c} \cdot M_{\mathbf{P}}} \right) = \exp \left(n_k \cdot \tilde{c} \cdot M_{\mathbf{P}} \cdot \left(\exp \left(\frac{\tau_k}{\tilde{c} \cdot M_{\mathbf{P}}} \right) - 1 \right) - j \cdot \frac{\tau_k}{\tilde{c} \cdot M_{\mathbf{P}}} \right).$$

Within the different context of Subsection X-B, the corresponding estimators $\widehat{\Pi}_L^{improved}$ can be obtained analogously to the last paragraph of Subsection XII-A, with Proposition 34 instead of Proposition 29 and with

$$D_{\tilde{c} \cdot \varphi_1}(\mathfrak{Q}, \mathbb{P}) := -\tilde{c} \cdot \log \left(1 + \frac{1}{\tilde{c}} \cdot \frac{1}{n} \cdot \log \widehat{\Pi}_L^{improved} \right) \text{ instead of (148).}$$

2) Case 6b: Different anchor point :

For several fields of applications it is important to have φ -divergences which also allow for vectors with negative components. In order to construct corresponding generators φ which are e.g. (i) close to classical ones for nonnegative entries, and (ii) which can be used for our BS method, one can appropriately vary the anchor point c in Theorem 22. This is exemplarily shown in the following. Indeed, let $\gamma = 1$, $]a_{F_{1,\tilde{c}}}, b_{F_{1,\tilde{c}}}[=]0, \infty[$ and $F_{1,\tilde{c}}(t)$ as in Case 6a above; for brevity let us fix (say) $\tilde{c} := 1$. By choosing a *general* anchor point $c \in \mathcal{R}(F_{1,1}) =]-\infty, \infty[$ (instead of $c = 0$), we obtain $] \lambda_-, \lambda_+[= \text{int}(\mathcal{R}(F_{1,1})) - c =]-\infty, \infty[$, $]t_-^{sc}, t_+^{sc}[=]1 + a_{F_{1,1}} - F_{1,1}^{-1}(c), 1 + b_{F_{1,1}} - F_{1,1}^{-1}(c)[=]1 - e^c, \infty[$ and $]a, b[=]1 - e^c, \infty[$. From (131),(133) we deduce

$$\varphi_{1,1}(t) := \varphi_{1,1}^{(c)}(t) := \begin{cases} (t + e^c - 1) \cdot [\log(t + e^c - 1) - c] + 1 - t \in [0, \infty[, & \text{if } t \in]1 - e^c, \infty[, \\ e^c, & \text{if } t = 1 - e^c, \\ \infty, & \text{if } t \in]-\infty, 1 - e^c[. \end{cases}$$

The corresponding divergence is

$$\begin{aligned} D_{\varphi_{1,1}^{(c)}}(\mathbf{Q}, \mathbf{P}) &:= \sum_{k=1}^K \left(q_k + p_k \cdot (e^c - 1) \right) \cdot \left\{ \log \left(\frac{q_k}{p_k} + e^c - 1 \right) - c \right\} - \sum_{k=1}^K q_k + \sum_{k=1}^K p_k, \\ &\text{if } \mathbf{P} \in \mathbb{R}_{>0}^K \text{ and } \mathbf{Q} \in \mathbb{R}^K \text{ with } \mathbf{Q} \in [(1 - e^c) \cdot \mathbf{P}, \infty[\text{ component-wise,} \end{aligned}$$

which for the special anchor-point choice $c = 0$ coincides with the Kullback-Leibler divergence (relative entropy) given in the fourth line of (41). Notice that $D_{\varphi_{1,1}^{(c)}}(\mathbf{Q}, \mathbb{P})$ has been recently used in [156] for the important task of testing mixtures of probability distributions; in fact, in order to get considerable comfort in testing mixture-type hypotheses against corresponding marginal-type alternatives, they employ choices $c > 0$ since then $\varphi_{1,1}^{(c)}(t)$ is finite especially for some range of negative values $t < 0$. Returning to our general considerations, we can employ (130),(132) to derive

$$\Lambda_{1,1}(z) := \Lambda_{1,1}^{(c)}(z) := \int_0^z F_{1,1}^{-1}(u+c) du + z \cdot (1 - F_{1,1}^{-1}(c)) = e^c \cdot (e^z - 1) + z \cdot (1 - e^c), \quad z \in]-\infty, \infty[,$$

which is the cumulant generating function of the “shifted Poisson distribution” $\zeta = POI(e^c) + 1 - e^c$, i.e. $W := Z + 1 - e^c$ with a $POI(e^c)$ -distributed random variable Z . Notice that ζ is a discrete distribution with frequencies $\mathbb{P}[W = \ell + 1 - e^c] = \exp(-e^c) \cdot \frac{e^{c \cdot \ell}}{\ell!}$ for all $\ell \in \mathbb{N}_0$ (and zero elsewhere). Moreover, $\mathbb{P}[W > 0] = 1$ iff $c < 0$, $\mathbb{P}[W < 0] > 0$ iff $c > 0$, $\mathbb{P}[W = 0] \neq 0$ iff “ $c = \log(1+k)$ for some $k \in \mathbb{N}_0$ ”. Concerning the important Remark 9(i), for i.i.d. copies $(W_i)_{i \in \mathbb{N}}$ of W the probability distribution $\zeta^{*n_k}[\cdot] := \mathbb{P}[\check{W} \in \cdot]$ of $\check{W} := \sum_{i \in I_k^{(n)}} W_i$ is $POI(\text{card}(I_k^{(n)}) \cdot e^c) + (1 - e^c) \cdot \text{card}(I_k^{(n)})$.

Within the context of Subsection X-A, for our concrete simulations we derive — in terms of (say) $M_{\mathbf{P}} := \sum_{i=1}^K p_i = 1$, $n_k = n \cdot \tilde{p}_k \in \mathbb{N}$ and \tilde{q}_k^* from proxy method 1 or 2 — that $\tilde{U}_k^{*n_k}$ is the shifted Poisson distribution $POI(n_k \cdot e^{c+\tau_k}) + n_k \cdot (1 - e^c)$ with support on the lattice $\{j + n_k \cdot (1 - e^c), j \in \mathbb{N}_0\}$, where $\tau_k = \log(\frac{\tilde{q}_k^*}{\tilde{p}_k} + e^c - 1) - c$ for $\tilde{q}_k^* > \tilde{p}_k \cdot (1 - e^c)$. Furthermore, for all $j \in \mathbb{N}_0$ we obtain (by setting $x := j + n_k \cdot (1 - e^c)$)

$$\widetilde{ISF}_k(j + n_k \cdot (1 - e^c)) = \exp(n_k \cdot e^c \cdot (e^{\tau_k} - 1) - j \cdot \tau_k).$$

Notice that the mass of $\tilde{U}_k^{*n_k}$ at zero depends crucially on the value of the anchor point c , since $\tilde{U}_k^{*n_k}[\{0\}] > 0$ if and only if $c = \log(1 + \frac{\ell}{n_k})$ for some $\ell \in \mathbb{N}_0$; moreover, $\tilde{U}_k^{*n_k}[]_0, \infty[= 1$ if $c < 0$ and $\tilde{U}_k^{*n_k}[]-\infty, 0[> 0$ if $c > 0$.

3) Case 6c: The remaining γ -constellations :

For $\gamma \in]1, 2[$, $\tilde{c} \in]0, \infty[$, $]a_{F_{\gamma, \tilde{c}}}, b_{F_{\gamma, \tilde{c}}}[:=]0, \infty[$ and anchor point $c := 0$, one can proceed as in Subsection XII-C with $F_{\gamma, \tilde{c}}$ from (151) and deduce from (131),(133) the same $\varphi_{\gamma, \tilde{c}}$ of (152) which coincides with $\tilde{c} \cdot \varphi_{\gamma}(t)$ for $\varphi_{\gamma}(t)$ from (40) and which generates the γ -corresponding power divergences given in (41). Moreover, we derive from formula (130),(132) the same $\Lambda_{\gamma, \tilde{c}}$ of (153); however, in contrast to Subsection XII-C, one gets for the therein involved crucial exponent $\frac{\gamma}{\gamma-1} > 2$. From this, we conjecture that there is no *probability* measure ζ such that $\Lambda_{\gamma, \tilde{c}}$ is the cumulant generating function of ζ . Indeed, we conjecture that ζ becomes a *signed* finite measure with total mass 1, i.e. it has a density (with respect to some dominating measure) with positive and negative values which “integrates to 1”; accordingly, our BS method can not be applied to this situation.

Remark 35: The characterization of the probability distribution ζ in (6) which may result from Theorem 22 — as seen through the above Cases 1 to 6b — considerably improves other approaches which make use of their identification through the concept of power variance functions of Natural Exponential Families, as developed by [342]–[344] and others. The latter approach has been used in [75] in a similar perspective as developed here, but can not be extended outside the range of power divergences, in contrast with the following Cases 7 to 10 which can only be handled as a consequence of Theorem 22.

G. Case 7

Consider the interesting “generalization” of the Kullback-Leibler divergence: for $\tilde{c} > 0$ and $\alpha \in]-1, 0[\cup]0, \infty[$ define

$$F_{gKL, \alpha, \tilde{c}}(t) := \begin{cases} \tilde{c} \cdot \log\left(\frac{(1+\alpha) \cdot t}{1+\alpha t}\right), & \text{if } \{ \alpha \in]0, \infty[\text{ and } t \in]0, \infty[\} \text{ or } \{ \alpha \in]-1, 0[\text{ and } t \in]0, -\frac{1}{\alpha}[\}, \\ -\infty, & \text{if } \alpha \in]-1, 0[\cup]0, \infty[\text{ and } t \in]-\infty, 0[, \\ \infty, & \text{if } \alpha \in]-1, 0[\text{ and } t \in]-\frac{1}{\alpha}, \infty[, \end{cases}$$

(notice that $\lim_{\alpha \rightarrow 0+} F_{gKL, \alpha, \tilde{c}}(t) = F_{1, \tilde{c}}(t)$, cf. Case 6a). Clearly, $]a_{F_{gKL, \alpha, \tilde{c}}}, b_{F_{gKL, \alpha, \tilde{c}}}[:=]0, \infty[$ for $\alpha \in]0, \infty[$ and $]a_{F_{gKL, \alpha, \tilde{c}}}, b_{F_{gKL, \alpha, \tilde{c}}}[:=]0, -\frac{1}{\alpha}[$ for $\alpha \in]-1, 0[$. Moreover, $\mathcal{R}(F_{gKL, \alpha, \tilde{c}}) =]-\infty, \tilde{c} \cdot \log(1 + \frac{1}{\alpha})[$ for $\alpha \in]0, \infty[$ and $\mathcal{R}(F_{gKL, \alpha, \tilde{c}}) =]-\infty, \infty[$ for $\alpha \in]-1, 0[$. Furthermore, $F_{gKL, \alpha, \tilde{c}}(\cdot)$ is strictly increasing and smooth on the respective $]a_{F_{gKL, \alpha, \tilde{c}}}, b_{F_{gKL, \alpha, \tilde{c}}}[$, and thus, $F_{gKL, \alpha, \tilde{c}} \in \mathfrak{F}$. Since $F_{gKL, \alpha, \tilde{c}}(1) = 0$, let us choose the natural anchor point $c := 0$, which leads to $] \lambda_-, \lambda_+ [= \text{int}(\mathcal{R}(F_{gKL, \alpha, \tilde{c}})) =]-\infty, \tilde{c} \cdot \log(1 + \frac{1}{\alpha})[$ and $]t_-^{sc}, t_+^{sc}[=]0, \infty[=]a, b[$ for the case $\alpha \in]0, \infty[$,

respectively, to $]\lambda_-, \lambda_+[= \text{int}(\mathcal{R}(F_{gKL,\alpha,\tilde{c}})) =]-\infty, \infty[$ and $]t_-^{sc}, t_+^{sc}[=]0, -\frac{1}{\alpha}[=]a, b[$ for the case $\alpha \in]-1, 0[$. By employing $F_{gKL,\alpha,\tilde{c}}^{-1}(x) = \frac{1}{(1+\alpha) \cdot e^{-x/\tilde{c}} - \alpha}$ for $x \in]\lambda_-, \lambda_+[$, we derive from formula (131) (see also (133))

$$\varphi_{gKL,\alpha,\tilde{c}}(t) := \varphi_{gKL,\alpha,\tilde{c}}^{(0)}(t)$$

$$:= \begin{cases} \tilde{c} \cdot \left[t \cdot \log t + (t + \frac{1}{\alpha}) \cdot \log \left(\frac{1+\alpha}{1+\alpha \cdot t} \right) \right] \in [0, \infty[, & \text{if } \{ \alpha \in]0, \infty[\text{ and } t \in]0, \infty[\} \text{ or } \{ \alpha \in]-1, 0[\text{ and } t \in]0, -\frac{1}{\alpha}[\}, \\ \frac{\tilde{c}}{\alpha} \cdot \log(1 + \alpha) \in]0, \infty[, & \text{if } \alpha \in]-1, 0[\cup]0, \infty[\text{ and } t = 0, \\ \infty, & \text{if } \alpha \in]-1, 0[\cup]0, \infty[\text{ and } t \in]-\infty, 0[, \\ \infty, & \text{if } \alpha \in]-1, 0[\text{ and } t \in [-\frac{1}{\alpha}, \infty[; \end{cases}$$

notice the *new effect* that this divergence generator is finite only on a *finite* interval, in case of $\alpha \in]-1, 0[$. In any α -case, we build from $\varphi_{gKL,\alpha,\tilde{c}}$ the corresponding divergence (cf. (4))

$$D_{\varphi_{gKL,\alpha,\tilde{c}}}(\mathbf{Q}, \mathbf{P}) = \tilde{c} \cdot \left\{ \sum_{k=1}^K q_k \cdot \log \left(\frac{q_k}{(1 - \frac{1}{1+\alpha}) \cdot q_k + \frac{1}{1+\alpha} \cdot p_k} \right) + \frac{1}{\alpha} \cdot \sum_{k=1}^K p_k \cdot \log \left(\frac{p_k}{(1 - \frac{1}{1+\alpha}) \cdot q_k + \frac{1}{1+\alpha} \cdot p_k} \right) \right\}$$

if $\{ \alpha \in]0, \infty[, \mathbf{P} \in \mathbb{R}_{\geq 0}^K \text{ and } \mathbf{Q} \in \mathbb{R}_{\geq 0}^K \}$ or $\{ \alpha \in]-1, 0[, \mathbf{P} \in \mathbb{R}_{> 0}^K \text{ and } \mathbf{Q} \in \mathbb{R}_{\geq 0}^K \text{ with } \mathbf{Q} \leq -\frac{1}{\alpha} \cdot \mathbf{P} \}$.

Notice that the symmetry $D_{\varphi_{gKL,\alpha,\tilde{c}}}(\mathbf{Q}, \mathbf{P}) = D_{\varphi_{gKL,\alpha,\tilde{c}}}(\mathbf{P}, \mathbf{Q})$ generally holds only if $\alpha = 1$; indeed, this special case leads to

$$\varphi_{snKL,\tilde{c}}(t) := \varphi_{gKL,1,\tilde{c}}(t) := \begin{cases} \tilde{c} \cdot \left[t \cdot \log t + (t + 1) \cdot \log \left(\frac{2}{t+1} \right) \right] \in [0, \infty[, & \text{if } t \in]0, \infty[, \\ \tilde{c} \cdot \log 2, & \text{if } t = 0, \\ \infty, & \text{if } t \in]-\infty, 0[, \end{cases}$$

and

$$D_{\varphi_{snKL,\tilde{c}}}(\mathbf{Q}, \mathbf{P}) := D_{\varphi_{gKL,1,\tilde{c}}}(\mathbf{Q}, \mathbf{P}) = \tilde{c} \cdot \left\{ \sum_{k=1}^K q_k \cdot \log \left(\frac{2q_k}{q_k + p_k} \right) + \sum_{k=1}^K p_k \cdot \log \left(\frac{2p_k}{q_k + p_k} \right) \right\} \text{ if } \mathbf{P}, \mathbf{Q} \in \mathbb{R}_{\geq 0}^K \text{ with } \mathbf{P} + \mathbf{Q} \in \mathbb{R}_{> 0}^K. \quad (158)$$

For the special subcase that $\tilde{c} = 1$ and that $\mathbf{P} = \mathbb{P}$, $\mathbf{Q} = \mathbb{Q}$ are probability vectors, the divergence (158) can be rewritten as a sum of two Kullback-Leibler divergences (cf. (41))

$$D_{\varphi_{snKL,1}}(\mathbb{Q}, \mathbb{P}) = D_{\varphi_1}(\mathbb{Q}, (\mathbb{Q} + \mathbb{P})/2) + D_{\varphi_1}(\mathbb{P}, (\mathbb{Q} + \mathbb{P})/2), \quad \text{if } \mathbb{P}, \mathbb{Q} \in \mathbb{S}^K \text{ with } \frac{\mathbb{P} + \mathbb{Q}}{2} \in \mathbb{S}_{> 0}^K,$$

which is the well-known (cf. [31],[345]–[350]) *Jensen-Shannon divergence* (being also called symmetrized and normalized Kullback-Leibler divergence, symmetrized and normalized relative entropy, capacity discrimination); this is equal to the $(2 \log 2)$ -fold of a special (namely, equally-weighted two-population) case of the Sibson information radius of order 1 (cf. [351]) which has also been addressed e.g. by [232] for genetic cluster analysis. By the way, for $\alpha > 0$ the divergence $D_{\varphi_{gKL,\alpha,\tilde{c}}}(\mathbb{Q}, \mathbb{P})$ can also be interpreted as a multiple of a special non-equally-weighted Sibson information radius of order 1. In a context of comparison of — not necessarily connected — networks where \mathbb{Q}, \mathbb{P} are probability vectors derived from matrices (cf. Remark 17) which are transforms of corresponding graph invariants (e.g. network portraits), the (matrix-equivalent of the) Jensen-Shannon divergence $D_{\varphi_{snKL,1}}(\mathbb{Q}, \mathbb{P})$ is also called the *network portrait divergence*, cf. [296]. There is a vast literature on recent applications of the Jensen-Shannon divergence in different research fields, for instance it appears exemplarily in [155],[211],[290],[352]–[392].

Remark 36: Let us transform $\varphi_{gSH,\alpha}(t) := \frac{1-t}{\alpha} \cdot \log(1 + \alpha) - \varphi_{gKL,\alpha,1}(t) = -t \cdot \log t + \frac{1}{\alpha} \cdot (1 + \alpha \cdot t) \cdot \log(1 + \alpha \cdot t) - \frac{1}{\alpha} \cdot (1 + \alpha) \cdot t \cdot \log(1 + \alpha)$ (for $t \in [0, 1]$). The function $\varphi_{gSH,\alpha}(\cdot)$ is strictly concave on $[0, 1]$ with $\varphi_{gSH,\alpha}(0) = \varphi_{gSH,\alpha}(1) = 0$. Hence, for probability vectors $\mathbb{Q} = (q_k)_{k=1,\dots,K}$, the φ -entropy $\sum_{k=1}^K \varphi_{gSH,\alpha}(q_k)$ is Kapur's [393] generalization of the Shannon entropy (which corresponds to $\alpha = 0$ in the limit) whose maximization has been connected with generalizations of the Bose-Einstein statistics and the Fermi-Dirac statistics e.g. in [56].

To proceed with our general considerations, one can deduce from formula (130) (see also (132))

$$\Lambda_{gKL,\alpha,\tilde{c}}(z) := \Lambda_{gKL,\alpha,\tilde{c}}^{(0)}(z)$$

$$:= \begin{cases} \int_0^z F_{gKL,\alpha,\tilde{c}}^{-1}(u) du = -\frac{\tilde{c}}{\alpha} \cdot \log((1 + \alpha) - \alpha \cdot e^{z/\tilde{c}}), & \text{if } \alpha \in]0, \infty[\text{ and } z \in]-\infty, \tilde{c} \cdot \log(1 + \frac{1}{\alpha})[, \\ \int_0^z F_{gKL,\alpha,\tilde{c}}^{-1}(u) du = -\frac{\tilde{c}}{\alpha} \cdot \log((1 + \alpha) - \alpha \cdot e^{z/\tilde{c}}), & \text{if } \alpha \in]-1, 0[\text{ and } z \in]-\infty, \infty[, \\ \infty, & \text{if } \alpha \in]0, \infty[\text{ and } z \in [\tilde{c} \cdot \log(1 + \frac{1}{\alpha}), \infty[. \end{cases} \quad (159)$$

(a) Subcase $\alpha \in]0, \infty[$: the $\Lambda_{gKL,\alpha,\tilde{c}}$ of (159) is the cumulant generating function of $\zeta = \frac{1}{\tilde{c}} \cdot NB(\frac{\tilde{c}}{\alpha}, \frac{1}{1+\alpha})$ being the " $\frac{1}{\tilde{c}}$ -fold Negative-Binomial distribution with parameters $\frac{\tilde{c}}{\alpha}$ and $\frac{1}{1+\alpha}$ " which means that $W = \frac{1}{\tilde{c}} \cdot Z$ for a $NB(\frac{\tilde{c}}{\alpha}, \frac{1}{1+\alpha})$ -distributed

random variable Z . For the special case $\tilde{c} = 1$, $\alpha = 1$ this reduces to $\zeta = NB(1, \frac{1}{2})$ to be the Negative-Binomial distribution with parameters 1 and $\frac{1}{2}$. Notice that ζ is an infinitely divisible (cf. Proposition 27) discrete distribution with frequencies $\mathbb{P}[W = \ell \cdot \frac{1}{\tilde{c}}] = (-1)^\ell \cdot \binom{-\tilde{c}}{\ell} \cdot \alpha^\ell \cdot (1 + \alpha)^{-\ell - \tilde{c}/\alpha}$ for all nonnegative integers $\ell \in \mathbb{N}_0$ (and zero elsewhere). Moreover, one has $\mathbb{P}[W \geq 0] = 1$, $\mathbb{P}[W = 0] = \frac{1}{(1 + \alpha)^{\tilde{c}/\alpha}}$. Concerning the important Remark 9(i), for i.i.d. copies $(W_i)_{i \in \mathbb{N}}$ of W the probability distribution $\zeta^{*n_k}[\cdot] := \mathbb{P}[\check{W} \in \cdot]$ of $\check{W} := \sum_{i \in I_k^{(n)}} W_i$ is $\frac{1}{\tilde{c}} \cdot NB(\frac{\tilde{c}}{\alpha} \cdot \text{card}(I_k^{(n)}), \frac{1}{1 + \alpha})$. Within the context of Subsection X-A, for the concrete simulative estimation $D_{gKL, \alpha, \tilde{c}}(\mathbf{Q}, \mathbf{P})$ via (116) and (118), we obtain — in terms of $M_{\mathbf{P}} := \sum_{i=1}^K p_i > 0$, $n_k = n \cdot \tilde{p}_k \in \mathbb{N}$ and \tilde{q}_k^* from proxy method 1 or 2 — that $\tilde{U}_k^{*n_k} = \frac{1}{\tilde{c} \cdot M_{\mathbf{P}}} \cdot NB(\frac{\tilde{c} \cdot M_{\mathbf{P}}}{\alpha} \cdot \text{card}(I_k^{(n)}), 1 - \frac{\alpha}{1 + \alpha} \cdot \exp(\frac{\tau_k}{\tilde{c} \cdot M_{\mathbf{P}}}))$ where $\tau_k = F_{gKL, \alpha, \tilde{c}}(\frac{\tilde{q}_k^*}{\tilde{p}_k})$ for $\tilde{q}_k^* > 0$; moreover, \widetilde{ISF}_k can be straightforwardly computed by (117).

(b) Subcase $\alpha \in] -1, 0[$: for any integer $m \in \mathbb{N}$ being strictly larger than \tilde{c} and the choice $\alpha = -\frac{\tilde{c}}{m}$, we obtain $\Lambda_{gKL, -\tilde{c}/m, \tilde{c}}(z) = m \cdot \log((1 - \frac{\tilde{c}}{m}) + \frac{\tilde{c}}{m} \cdot e^{z/\tilde{c}})$ (cf. (159)) which is the cumulant generating function of $\zeta = \frac{1}{\tilde{c}} \cdot BIN(m, \frac{\tilde{c}}{m})$ being the “ $\frac{1}{\tilde{c}}$ -fold Binomial distribution with parameters m and $\frac{\tilde{c}}{m}$ ” which means that $W = \frac{1}{\tilde{c}} \cdot Z$ for a $BIN(m, \frac{\tilde{c}}{m})$ -distributed random variable Z . Notice that ζ is a *non*-infinitely divisible discrete distribution with frequencies $\mathbb{P}[W = \ell \cdot \frac{1}{\tilde{c}}] = \binom{m}{\ell} \cdot (\frac{\tilde{c}}{m})^\ell \cdot (1 - \frac{\tilde{c}}{m})^{m-\ell}$ for $\ell \in \{0, 1, \dots, m\}$ (and zero elsewhere). Furthermore, there holds $\mathbb{P}[W \geq 0] = 1$, $\mathbb{P}[W = 0] = (1 - \frac{\tilde{c}}{m})^m$. Concerning the important Remark 9(i), for i.i.d. copies $(W_i)_{i \in \mathbb{N}}$ of W the probability distribution $\zeta^{*n_k}[\cdot] := \mathbb{P}[\check{W} \in \cdot]$ of $\check{W} := \sum_{i \in I_k^{(n)}} W_i$ is $\frac{1}{\tilde{c}} \cdot BIN(m \cdot \text{card}(I_k^{(n)}), \frac{\tilde{c}}{m})$. Analogously to (a), for $D_{gKL, \alpha, \tilde{c}}(\mathbf{Q}, \mathbf{P})$ we employ that $\tilde{U}_k^{*n_k} = \frac{1}{\tilde{c} \cdot M_{\mathbf{P}}} \cdot BIN(m \cdot \text{card}(I_k^{(n)}), \tilde{p})$ where $m \in \mathbb{N}$ is strictly larger than $\tilde{c} \cdot M_{\mathbf{P}}$, $\tilde{p} := \frac{\tilde{c} \cdot M_{\mathbf{P}} \cdot \exp(\frac{\tau_k}{\tilde{c} \cdot M_{\mathbf{P}}})}{m - \tilde{c} \cdot M_{\mathbf{P}} + \tilde{c} \cdot M_{\mathbf{P}} \cdot \exp(\frac{\tau_k}{\tilde{c} \cdot M_{\mathbf{P}}})}$ and $\tau_k = F_{gKL, \alpha, \tilde{c}}(\frac{\tilde{q}_k^*}{\tilde{p}_k})$ for $\tilde{q}_k^* \in]0, -\frac{\tilde{p}_k}{\alpha}[$; furthermore, for \widetilde{ISF}_k we use (117).

As far as the construction of bounds of $D_{\varphi_{gKL, \alpha, \tilde{c}}}(\mathbf{Q}, \mathbf{P})$ in the light of (68) in Subsection VI-C is concerned, for the sake of brevity we confine ourselves to the above-mentioned important Jensen-Shannon divergence special case $\alpha = 1$ and $\tilde{c} = 1$, and thus $\varphi_{gKL, 1, 1} = \varphi_{snKL, 1}$. Correspondingly, we abbreviate (158) as

$$J(\mathbf{Q}, \mathbf{P}) := D_{\varphi_{gKL, 1, 1}}(\mathbf{Q}, \mathbf{P}) = I(\mathbf{Q}, (\mathbf{Q} + \mathbf{P})/2) + I(\mathbf{P}, (\mathbf{Q} + \mathbf{P})/2) \quad \text{if } \mathbf{P}, \mathbf{Q} \in \mathbb{R}_{\geq 0}^K \text{ with } \mathbf{P} + \mathbf{Q} \in \mathbb{R}_{> 0}^K,$$

where we use (as an extension of (43)) $I(\check{\mathbf{Q}}, \check{\mathbf{P}}) := \sum_{k=1}^K \check{q}_k \cdot \log(\frac{\check{q}_k}{\check{p}_k})$ for $\check{\mathbf{P}} \in \mathbb{R}_{> 0}^K$ and $\check{\mathbf{Q}} \in \mathbb{R}_{\geq 0}^K$. We explore the sharpness of the bounds for $J(\mathbf{Q}, \mathbf{P})$ as defined in (68). For this, we consider a given probability distribution \mathbb{P} on \mathcal{Y} with strictly positive entries; the set \mathfrak{Q} consists of all probability distributions \mathbf{Q} on \mathcal{Y} whose total variation distance $V(\mathbf{Q}, \mathbb{P}) := \sum_{k=1}^K |q_k - p_k|$ ³² to \mathbb{P} lies between v and $v + h$ for $v > 0$ and small h and which also satisfies $\sup\left(\sup_{k=1, \dots, K} \frac{p_k}{q_k}, \sup_{k=1, \dots, K} \frac{q_k}{p_k}\right) \leq L$ for some strictly positive finite L . This set \mathfrak{Q} defines a class of distributions \mathbf{Q} away from \mathbb{P} still keeping some regularity w.r.t. \mathbb{P} . Also, \mathfrak{Q} satisfies (7). We will prove that the bounds in (68) are sharp in this case. Notice that $J(m \cdot \mathbf{Q}, \mathbb{P}) = \infty$ for $m < 0$ and hence $\inf_{\mathbf{Q} \in \mathfrak{Q}} \inf_{m \neq 0} J(m \cdot \mathbf{Q}, \mathbb{P}) = \inf_{\mathbf{Q} \in \mathfrak{Q}} \inf_{m > 0} J(m \cdot \mathbf{Q}, \mathbb{P})$. We first provide a lower bound for the latter. It holds for all $m > 0$ and \mathbf{Q} in \mathfrak{Q} that

$$J(m \cdot \mathbf{Q}, \mathbb{P}) = (m + 1) \cdot \log(2) - (m + 1) \cdot \log(m + 1) + m \cdot \log m + I^\alpha(\mathbb{P}, \mathbf{Q}) + m \cdot I^{1-\alpha}(\mathbf{Q}, \mathbb{P})$$

where $\alpha := 1/(m + 1)$ and $I^\alpha(\mathbb{P}, \mathbf{Q})$ is the α -skewed Kullback-Leibler divergence between \mathbb{P} and \mathbf{Q} defined through $I^\alpha(\mathbb{P}, \mathbf{Q}) := I(\mathbb{P}, \alpha \mathbb{P} + (1 - \alpha)\mathbf{Q})$. By inequality (27) in [394] one gets $I^\alpha(\mathbb{P}, \mathbf{Q}) \geq -\log\left(1 - \frac{\alpha^2}{4} \cdot V(\mathbf{Q}, \mathbb{P})^2\right)$. Since $(m + 1) \cdot \log(2) - (m + 1) \cdot \log(m + 1) + m \cdot \log m$ is non-negative for all $m > 0$ and takes its minimal value 0 for $m = 1$, we obtain $\inf_{m > 0} J(m \cdot \mathbf{Q}, \mathbb{P}) \geq \inf_{m > 0} K(m)$ where $K(m) := -\log\left(1 - \frac{1}{4(m+1)^2} \cdot V(\mathbf{Q}, \mathbb{P})^2\right) - m \cdot \log\left(1 - \frac{m^2}{4(m+1)^2} \cdot V(\mathbf{Q}, \mathbb{P})^2\right)$. Since $-\log(1 - x) \geq x$ for all $x < 1$ and both $\frac{1}{4(m+1)^2} \cdot V(\mathbf{Q}, \mathbb{P})^2$ and $\frac{m^2}{4(m+1)^2} \cdot V(\mathbf{Q}, \mathbb{P})^2$ are less than 1, it follows that $K(m) \geq \frac{V(\mathbf{Q}, \mathbb{P})^2}{4} \cdot \frac{m^3 + 1}{(m+1)^2}$ where the right-hand side attains its minimal value on $]0, \infty[$ at $m^* = \sqrt{3} - 1 \approx 0.73$. Hence, we obtain $\inf_{m > 0} J(m \cdot \mathbf{Q}, \mathbb{P}) \geq \frac{V(\mathbf{Q}, \mathbb{P})^2}{4} \cdot (2\sqrt{3} - 3) > 0.116 v^2$. Now by (19) in [394], for any \mathbf{Q} one gets $J(\mathbf{Q}, \mathbb{P}) \leq \frac{1}{4} \underline{J}(\mathbf{Q}, \mathbb{P})$ where $\underline{J}(\mathbf{Q}, \mathbb{P}) := I(\mathbf{Q}, \mathbb{P}) + I(\mathbb{P}, \mathbf{Q})$ is the Jensen divergence (also called symmetrized Kullback-Leibler divergence) between \mathbf{Q} and \mathbb{P} . Since (see [395]) $I(\mathbb{P}, \mathbf{Q}) \leq \sum_{k=1}^K \sqrt{\frac{p_k}{q_k}} \cdot |q_k - p_k|$, it follows that $J(\mathbf{Q}^*, \mathbb{P}) \leq \frac{1}{2} \sqrt{L} \cdot V(\mathbf{Q}^*, \mathbb{P})$ which provides $0.116 \cdot v^2 \leq \inf_{m > 0} J(m \cdot \mathbf{Q}^*, \mathbb{P}) = J(m(\mathbf{Q}^*), \mathbb{P}) \leq J(\mathfrak{Q}, \mathbb{P}) \leq J(\mathbf{Q}^*, \mathbb{P}) \leq \frac{1}{2} \sqrt{L} \cdot (v + h)$. For small v , the difference between the RHS and the LHS in the above display is $cst \cdot v + o(v) + \frac{1}{2} \sqrt{L} \cdot h$ which proves that the bounds are sharp locally, with non-trivial lower bound. Other upper bounds can be adapted to sets \mathfrak{Q} defined through tighter conditions on $\sup_{\mathbf{Q} \in \mathfrak{Q}} \sup_{k=1, \dots, K} \frac{p_k}{q_k}$ and $\sup_{\mathbf{Q} \in \mathfrak{Q}} \sup_{k=1, \dots, K} \frac{q_k}{p_k}$ (of e.g. [395]).

³²notice that $V(\mathbf{Q}, \mathbb{P})$ always takes values in the interval $[0, 2[$

H. Case 8

As a continuation of the discussion at the beginning of Case 6b (cf. Subsection XII-F2), we give another φ -divergence which also allows for vectors with negative components (i.e. φ is finite especially for some range of negative values $t < 0$). Indeed, for $\beta \in]0, 1]$, $\tilde{c} \in]0, \infty[$ and $]a_{F_{bw,\beta,\tilde{c}}}, b_{F_{bw,\beta,\tilde{c}}}[=]1 - \frac{1}{\beta}, \infty[$ let us define

$$F_{bw,\beta,\tilde{c}}(t) := \begin{cases} \frac{\tilde{c}}{2\beta} \cdot \left(1 - \frac{1}{(\beta \cdot t + 1 - \beta)^2}\right), & \text{if } t \in]1 - \frac{1}{\beta}, \infty[, \\ -\infty, & \text{if } t \in]-\infty, 1 - \frac{1}{\beta}]. \end{cases}$$

Clearly, $\mathcal{R}(F_{bw,\beta,\tilde{c}}) =]-\infty, \frac{\tilde{c}}{2\beta}[$ and $F_{bw,\beta,\tilde{c}} \in \mathcal{F}$. Since $F_{bw,\beta,\tilde{c}}(1) = 0$, let us choose the natural anchor point $c := 0$, which leads to $] \lambda_-, \lambda_+[= \text{int}(\mathcal{R}(F_{bw,\beta,\tilde{c}})) =]-\infty, \frac{\tilde{c}}{2\beta}[$, $]t_-^{sc}, t_+^{sc}[=]a_{F_{bw,\beta,\tilde{c}}}, b_{F_{bw,\beta,\tilde{c}}}[=]1 - \frac{1}{\beta}, \infty[$ and $]a, b[=]1 - \frac{1}{\beta}, \infty[$.

By using $F_{bw,\beta,\tilde{c}}^{-1}(x) = \frac{1}{\beta} \cdot \left\{ \frac{1}{\sqrt{1 - 2\beta \cdot x/\tilde{c}}} + \beta - 1 \right\}$ for $x \in \text{int}(\mathcal{R}(F_{bw,\beta,\tilde{c}}))$, from formula (131),(133) we can deduce

$$\varphi_{bw,\beta,\tilde{c}}(t) := \varphi_{bw,\beta,\tilde{c}}^{(0)}(t) := \begin{cases} \tilde{c} \cdot \frac{(t-1)^2}{2(\beta \cdot t + 1 - \beta)} \in [0, \infty[, & \text{if } t \in]1 - \frac{1}{\beta}, \infty[, \\ \infty, & \text{if } t \in]-\infty, 1 - \frac{1}{\beta}]. \end{cases} \quad (160)$$

Note that $1 - \frac{1}{\beta} < 0$ so that negative t are allowed here. For $t \geq 0$, $\varphi_{bw,\beta,\tilde{c}}(t)$ is known as *Rukhin's generator* (cf. [396], see e.g. also [397], [12]). From the generator $\varphi_{bw,\beta,\tilde{c}}$ given in (160), we build the corresponding divergence (cf. (4))

$$D_{\varphi_{bw,\beta,\tilde{c}}}(\mathbf{Q}, \mathbf{P}) = \tilde{c} \cdot \sum_{k=1}^K p_k \cdot \frac{\left(\frac{q_k}{p_k} - 1\right)^2}{2\left(\beta \cdot \frac{q_k}{p_k} + 1 - \beta\right)} = \frac{\tilde{c}}{2} \cdot \sum_{k=1}^K \frac{(q_k - p_k)^2}{\beta \cdot q_k + (1 - \beta) \cdot p_k},$$

if $\mathbf{P} \in \mathbb{R}_{\geq 0}^K$ and $\mathbf{Q} \in \mathbb{R}^K$ with $\mathbf{Q} \in]\mathbf{P} \cdot (1 - \frac{1}{\beta}), \infty[$ component-wise;

for the special subcase $\tilde{c} = 1$ and $\mathbf{Q} \in \mathbb{R}_{> 0}^K$, $D_{\varphi_{bw,\beta,1}}(\mathbf{Q}, \mathbf{P})$ can be interpreted as — “non-probability version” of — the well-known *blended weight chi-square divergence of Lindsay* [398] (see e.g. also [399]–[401]). The special case $\tilde{c} = 1$ and $\beta = \frac{1}{2}$ for probability vectors, i.e. $D_{\varphi_{bw,1/2,1}}(\mathbf{Q}, \mathbf{P})$, is equal to (a multiple of the matrix-vector-converted (cf. Remark 17)) *Sanghvi's genetic difference measure* [402] and equal to the double of the so-called (*squared*) *Vincze-Le Cam distance* (cf. [403],[404], see also e.g. [347] — who used the alternative naming *triangular discrimination* — and [349]); $D_{\varphi_{bw,1/2,1}}(\mathbf{Q}, \mathbf{P})$ has been used e.g. in [405] for a machine learning context where \mathbf{Q} and \mathbf{P} are appropriate histograms of RGB color.

Remark 37: (a) By straightforward calculations, one can show that $\varphi_{bw,1,\tilde{c}}$ is equal to the power-divergence generator $\varphi_{\gamma,\tilde{c}} = \tilde{c} \cdot \varphi_{\gamma}$ (cf. (40)) with $\gamma = -1$; the corresponding divergence $D_{\varphi_{bw,1,\tilde{c}}}(\mathbf{Q}, \mathbf{P})$ is thus equal to the power divergence $D_{\tilde{c},\varphi_{-1}}(\mathbf{Q}, \mathbf{P})$ (cf. (41)) which is nothing but the $\tilde{c}/2$ -fold — “non-probability version” — of *Neyman's chi-square divergence*. (b) For the case $\beta = 0$ — which has been excluded above for technical brevity — the divergence generator $\varphi_{bw,0,\tilde{c}}$ corresponds to the power-divergence generator $\varphi_{\gamma,\tilde{c}}$ with $\gamma = 2$; the corresponding divergence $D_{\varphi_{bw,0,\tilde{c}}}(\mathbf{Q}, \mathbf{P})$ is thus equal to the power divergence $D_{\tilde{c},\varphi_2}(\mathbf{Q}, \mathbf{P})$ (cf. (41)) which is nothing but the $\tilde{c}/2$ -fold — “non-probability version” — of Pearson's (i.e. the *classical*) chi-square divergence.

To continue with our general considerations, we can derive from formula (130) (see also (132)) for all $\beta \in]0, 1]$

$$\Lambda_{bw,\beta,\tilde{c}}(z) := \Lambda_{bw,\beta,\tilde{c}}^{(0)}(z) = \begin{cases} -\left(\frac{1}{\beta} - 1\right) \cdot z + \frac{\tilde{c}}{\beta^2} \cdot \left\{1 - \sqrt{1 - \frac{2\beta}{\tilde{c}} \cdot z}\right\}, & \text{if } z \in]-\infty, \frac{\tilde{c}}{2\beta}], \\ \infty, & \text{if } z \in \left]\frac{\tilde{c}}{2\beta}, \infty\right[. \end{cases}$$

which is the cumulant generating function of a probability distribution $\mathcal{Q}[\cdot] = \mathbb{P}[\check{W} \in \cdot]$ of a random variable \check{W} , which can be constructed as follows: $\check{W} := \frac{W}{\beta} - \left(\frac{1}{\beta} - 1\right)$, where W is the random variable constructed in Case 1 (cf. Subsection XII-A) with $\gamma = -1$ and with \tilde{c} replaced by $\frac{\tilde{c}}{\beta^2}$ (recall that W has a tilted stable distribution). In other words, \mathcal{Q} is a special kind of *modified tilted stable distribution*. Notice that \mathcal{Q} is an infinitely divisible (cf. Proposition 27) continuous distribution with density $f_{\check{W}}(u) := \beta \cdot f_W(\beta \cdot u + 1 - \beta) \cdot \mathbf{1}_{]-(\frac{1}{\beta}-1), \infty[}(u)$ ($u \in \mathbb{R}$), where $f_W(\cdot)$ is given in (137) with $\gamma = -1$ and with \tilde{c} replaced by $\frac{\tilde{c}}{\beta^2}$. Moreover, $\mathcal{Q}([0, \infty[) = \mathbb{P}[\check{W} > 0] > 0$. Concerning Remark 9(i), for i.i.d. copies $(\check{W}_i)_{i \in \mathbb{N}}$ of \check{W} , the probability distribution $\mathcal{Q}^{*n_k}[\cdot] := \mathbb{P}[\check{W} \in \cdot]$ of $\check{W} := \sum_{i \in I_k^{(n)}} \check{W}_i = \frac{1}{\beta} \cdot \sum_{i \in I_k^{(n)}} W_i - n_k \cdot \left(\frac{1}{\beta} - 1\right)$ has density $f_{\check{W}}^{\check{W}}(u) := \beta \cdot f_{\check{W}}(\beta \cdot u + (1 - \beta) \cdot n_k) \cdot \mathbf{1}_{]-n_k \cdot (\frac{1}{\beta}-1), \infty[}(u)$ ($u \in \mathbb{R}$), where $f_{\check{W}}(\cdot)$ is given in (138) with $\gamma = -1$ and with \tilde{c} replaced by $\frac{\tilde{c}}{\beta^2}$; accordingly, the determination of the corresponding $\tilde{U}_k^{*n_k}$ and \tilde{ISF}_k works analogously.

I. Case 9

Let us fix any $z_1, z_2 \in \mathbb{R}$, $p \in]0, 1[$ which satisfy $z_1 < 1 < z_2$ and $z_1 \cdot p + z_2 \cdot (1 - p) = 1$ (and thus $p = \frac{z_2 - 1}{z_2 - z_1}$). On $]a_{F_{twop}}, b_{F_{twop}}[:=]z_1, z_2[$ we define

$$F_{twop}(t) := \frac{1}{z_2 - z_1} \cdot \log \left(\frac{(t - z_1) \cdot p}{(z_2 - t) \cdot (1 - p)} \right) = \frac{1}{z_2 - z_1} \cdot \log \left(\frac{(t - z_1) \cdot (z_2 - 1)}{(z_2 - t) \cdot (1 - z_1)} \right), \quad t \in]z_1, z_2[,$$

where for the last equality we have used the above constraint (in order to obtain a two-parameter representation). Straightforwardly, we have $\mathcal{R}(F_{twop}) =]-\infty, \infty[$ and $F_{twop} \in \mathcal{F}$. Since $F_{twop}(1) = 0$, let us choose the natural anchor point $c := 0$, which leads to $] \lambda_-, \lambda_+ [= \text{int}(\mathcal{R}(F_{twop})) =]-\infty, \infty[,]t_-^{sc}, t_+^{sc}[=]z_1, z_2[=]a, b[$. By using

$$F_{twop}^{-1}(x) = \frac{p \cdot z_1 + (1 - p) \cdot z_2 \cdot e^{(z_2 - z_1) \cdot x}}{p + (1 - p) \cdot e^{(z_2 - z_1) \cdot x}}, \quad x \in]-\infty, \infty[,$$

from formula (131) (see also (133)) we deduce

$$\varphi_{twop}(t) := \varphi_{twop}^{(0)}(t) := \begin{cases} \frac{t - z_1}{z_2 - z_1} \cdot \log \left(\frac{(t - z_1) \cdot (z_2 - 1)}{(z_2 - t) \cdot (1 - z_1)} \right) - \log \left(\frac{z_2 - 1}{z_2 - t} \right) \in]0, \infty[, & \text{if } t \in]z_1, z_2[, \\ \log \left(\frac{z_2 - z_1}{z_2 - 1} \right) \in]0, \infty[, & \text{if } t = z_1, \\ \log \left(\frac{z_2 - z_1}{1 - z_1} \right) \in]0, \infty[, & \text{if } t = z_2, \\ \infty, & \text{if } t \in]-\infty, z_1[\cup]z_2, \infty[; \end{cases}$$

notice that this generator is finite only on a *finite* interval. From it, we build the new divergence (cf. (4))

$$D_{\varphi_{twop}}(\mathbf{Q}, \mathbf{P}) = \sum_{k=1}^K \frac{q_k - z_1 \cdot p_k}{z_2 - z_1} \cdot \log \left(\frac{(z_2 - 1) \cdot (q_k - z_1 \cdot p_k)}{(1 - z_1) \cdot (z_2 \cdot p_k - q_k)} \right) - \sum_{k=1}^K p_k \cdot \log \left(\frac{(z_2 - 1) \cdot p_k}{z_2 \cdot p_k - q_k} \right)$$

for $\mathbf{P} \in \mathbb{R}_{>0}^K$, $\mathbf{Q} \in [z_1 \mathbf{P}, z_2 \mathbf{P}]$ component-wise.

Furthermore, we derive from formula (130) (see also (132))

$$\Lambda_{twop}(z) := \Lambda_{twop}^{(0)}(z) := \int_0^z F_{twop}^{-1}(u) du = \log \left(p \cdot e^{z_1 \cdot z} + (1 - p) \cdot e^{z_2 \cdot z} \right), \quad z \in]-\infty, \infty[,$$

which is the well-known cumulant generating function of the two-point probability distribution $\zeta = p \cdot \delta_{z_1} + (1 - p) \cdot \delta_{z_2}$, where $z_1 < 1 < z_2$ and $p = \frac{z_2 - 1}{z_2 - z_1}$. Of course, ζ is a discrete distribution with frequencies $\mathbb{P}[W = z_1] = p$, $\mathbb{P}[W = z_2] = 1 - p$ (and zero elsewhere). Moreover, $\mathbb{P}[W > 0] = 1$ iff $z_1 > 0$, $\mathbb{P}[W = 0] > 0$ iff $z_1 = 0$. Concerning the important Remark 9(i), for i.i.d. copies $(W_i)_{i \in \mathbb{N}}$ of W the probability distribution $\zeta^{*n_k}[\cdot] := \mathbb{P}[\check{W} \in \cdot]$ of $\check{W} := \sum_{i \in I_k^{(n)}} W_i$ is the distribution of the $\text{card}(I_k^{(n)})$ -th step of a generalized random walk starting at zero; this has a nice explicit (“binomial-type”) expression in the special case $z_1 = -z_2$, namely $\sum_{\ell=0}^{\text{card}(I_k^{(n)})} \binom{\text{card}(I_k^{(n)})}{\ell} \cdot p^{\text{card}(I_k^{(n)}) - \ell} \cdot (1 - p)^\ell \cdot \delta_{z_2 \cdot (2\ell - \text{card}(I_k^{(n)}))}$. Within the context of Subsection X-A, for the concrete simulative estimation $D_{\frac{1}{M_{\mathbf{P}}}} \widehat{\varphi_{twop}}(\mathbf{\Omega}, \mathbf{P})$ via (116) and (118), we obtain — in terms of $M_{\mathbf{P}} := \sum_{i=1}^K p_i > 0$, $n_k = n \cdot \tilde{p}_k \in \mathbb{N}$ and \tilde{q}_k^* from proxy method 1 or 2 — that $\tilde{U}_k^{*n_k}$ is the distribution of the $\text{card}(I_k^{(n)})$ -th step of a generalized random walk starting at zero, but with p replaced by $\check{p} := \frac{p \cdot \exp(z_1 \cdot \tau_k)}{p \cdot \exp(z_1 \cdot \tau_k) + (1 - p) \cdot \exp(z_2 \cdot \tau_k)}$ where $\tau_k = \frac{1}{M_{\mathbf{P}}} \cdot F_{twop} \left(\frac{\tilde{q}_k^*}{\tilde{p}_k} \right)$ for $\tilde{q}_k^* \in]z_1 \cdot \tilde{p}_k, z_2 \cdot \tilde{p}_k[$; moreover, \widehat{ISF}_k can be straightforwardly computed by (117).

J. Case 10

It is known that some types of robustness properties of minimum-divergence estimators are connected with the *boundedness* of the derivative φ' of the divergence generator φ ; this property is satisfied for the following cases, which lead to the new classes of divergences (163) and (165).

1) Case 10a:

For any parameter-quadrupel $\alpha, \beta_1, \beta_2, \tilde{c} \in]0, \infty[$ we choose

$$]a_F, b_F[:=]a_{F_{\alpha, \beta_1, \beta_2, \tilde{c}}}, b_{F_{\alpha, \beta_1, \beta_2, \tilde{c}}}[:=]-\infty, \infty[$$

and define with $\check{\theta} := 1 + \alpha \cdot \left(\frac{1}{\beta_2} - \frac{1}{\beta_1} \right)$ (which is in $] -\infty, 1[$ for $\beta_1 < \beta_2$, respectively, in $]1, 1 + \frac{\alpha}{\beta_2}[$ for $\beta_1 > \beta_2$)

$$F_{\alpha, \beta_1, \beta_2, \tilde{c}}(t) := \begin{cases} \tilde{c} \cdot \frac{\beta_1 - \beta_2}{2} + \frac{\tilde{c}}{\frac{1-t}{\alpha} + \frac{1}{\beta_2} - \frac{1}{\beta_1}} \cdot \left(1 - \frac{1}{2} \cdot \sqrt{4 + \left(\frac{1-t}{\alpha} + \frac{1}{\beta_2} - \frac{1}{\beta_1} \right)^2 \cdot (\beta_1 + \beta_2)^2} \right), & \text{if } t \in]a_F, b_F[\setminus \{\check{\theta}\}, \\ \tilde{c} \cdot \frac{\beta_1 - \beta_2}{2}, & \text{if } t = \check{\theta}. \end{cases}$$

One has the continuity $\lim_{t \rightarrow \check{\theta}} F_{\alpha, \beta_1, \beta_2, \tilde{c}}(t) = \tilde{c} \cdot \frac{\beta_1 - \beta_2}{2}$. Moreover, one can see in a straightforward way that $F_{\alpha, \beta_1, \beta_2, \tilde{c}}(\cdot)$ is strictly increasing and that $\mathcal{R}(F_{\alpha, \beta_1, \beta_2, \tilde{c}}) =]-\tilde{c} \cdot \beta_2, \tilde{c} \cdot \beta_1[$. Furthermore, $F_{\alpha, \beta_1, \beta_2, \tilde{c}}(\cdot)$ is smooth on $]a_F, b_F[$, and thus $F_{\alpha, \beta_1, \beta_2, \tilde{c}} \in$

\mathfrak{f} . Since $F_{\alpha,\beta_1,\beta_2,\tilde{c}}(1) = 0$, let us choose the natural anchor point $c := 0$, which leads to $] \lambda_-, \lambda_+[= \text{int}(\mathcal{R}(F_{\alpha,\beta_1,\beta_2,\tilde{c}})) =] - \tilde{c} \cdot \beta_2, \tilde{c} \cdot \beta_1[$ and $] t_-^{sc}, t_+^{sc}[=] a_F, b_F[=] -\infty, \infty[=] a, b[$. Also, the corresponding inverse is

$$F_{\alpha,\beta_1,\beta_2,\tilde{c}}^{-1}(x) = 1 + \alpha \cdot \left(\frac{1}{\beta_2} - \frac{1}{\beta_1} \right) - \alpha \cdot \frac{\frac{1}{\beta_2} - \frac{1}{\beta_1} - \frac{2x}{\tilde{c} \cdot \beta_1 \cdot \beta_2}}{1 + \frac{x}{\tilde{c}} \cdot \left(\frac{1}{\beta_2} - \frac{1}{\beta_1} \right) - \frac{x^2}{\tilde{c}^2 \cdot \beta_1 \cdot \beta_2}}, \quad x \in \text{int}(\mathcal{R}(F_{\alpha,\beta_1,\beta_2,\tilde{c}})); \quad (161)$$

from this and (131) (see also (133)) we can deduce

$$\begin{aligned} \varphi_{\alpha,\beta_1,\beta_2,\tilde{c}}(t) := \varphi_{\alpha,\beta_1,\beta_2,\tilde{c}}^{(0)}(t) &= \tilde{c} \cdot \alpha \cdot \left\{ \frac{\sqrt{4 + (\beta_1 + \beta_2)^2 \cdot \left(\frac{1-t}{\alpha} + \frac{1}{\beta_2} - \frac{1}{\beta_1} \right)^2} - \left(\frac{1-t}{\alpha} + \frac{1}{\beta_2} - \frac{1}{\beta_1} \right) \cdot (\beta_1 - \beta_2) - 2}{2} \right. \\ &\quad \left. + \log \frac{\sqrt{4 + (\beta_1 + \beta_2)^2 \cdot \left(\frac{1-t}{\alpha} + \frac{1}{\beta_2} - \frac{1}{\beta_1} \right)^2} - 2}{\beta_1 \beta_2 \cdot \left(\frac{1-t}{\alpha} + \frac{1}{\beta_2} - \frac{1}{\beta_1} \right)^2} \right\} \in [0, \infty[, \quad t \in] -\infty, \infty[. \quad (162) \end{aligned}$$

Notice that $\varphi_{\alpha,\beta_1,\beta_2,\tilde{c}}(1) = 0$, $\varphi'_{\alpha,\beta_1,\beta_2,\tilde{c}}(1) = 0$, $\varphi_{\alpha,\beta_1,\beta_2,\tilde{c}}(-\infty) = \infty$ and $\varphi_{\alpha,\beta_1,\beta_2,\tilde{c}}(\infty) = \infty$. Moreover, $\varphi'_{\alpha,\beta_1,\beta_2,\tilde{c}}(-\infty) = \varphi'_{\alpha,\beta_1,\beta_2,\tilde{c}}(a_F) = -\tilde{c} \cdot \beta_2$ and $\varphi'_{\alpha,\beta_1,\beta_2,\tilde{c}}(\infty) = \tilde{c} \cdot \beta_1$. From (162), we construct the corresponding divergence (cf. (4))

$$\begin{aligned} D_{\varphi_{\alpha,\beta_1,\beta_2,\tilde{c}}}(\mathbf{Q}, \mathbf{P}) &= \sum_{k=1}^K p_k \cdot \varphi_{\alpha,\beta_1,\beta_2,\tilde{c}}\left(\frac{q_k}{p_k}\right) \\ &= \sum_{k=1}^K p_k \cdot \left[\tilde{c} \cdot \alpha \cdot \left\{ \frac{\sqrt{4 + (\beta_1 + \beta_2)^2 \cdot \left(\frac{1-\frac{q_k}{p_k}}{\alpha} + \frac{1}{\beta_2} - \frac{1}{\beta_1} \right)^2} - \left(\frac{1-\frac{q_k}{p_k}}{\alpha} + \frac{1}{\beta_2} - \frac{1}{\beta_1} \right) \cdot (\beta_1 - \beta_2) - 2}{2} \right. \right. \\ &\quad \left. \left. + \log \frac{\sqrt{4 + (\beta_1 + \beta_2)^2 \cdot \left(\frac{1-\frac{q_k}{p_k}}{\alpha} + \frac{1}{\beta_2} - \frac{1}{\beta_1} \right)^2} - 2}{\beta_1 \beta_2 \cdot \left(\frac{1-\frac{q_k}{p_k}}{\alpha} + \frac{1}{\beta_2} - \frac{1}{\beta_1} \right)^2} \right\} \right], \quad \text{if } \mathbf{P} \in \mathbb{R}_{\geq 0}^K, \mathbf{Q} \in \mathbb{R}^K. \quad (163) \end{aligned}$$

Notice that we can particularly include the case where $p_k = 0$, since for $q_k = 0$ we have $0 \cdot \varphi_{\alpha,\beta_1,\beta_2,\tilde{c}}\left(\frac{0}{0}\right) = 0$ by the convention right after (4), and for $q_k \neq 0$ we have $\lim_{t \rightarrow 0_+} t \cdot \varphi_{\alpha,\beta_1,\beta_2,\tilde{c}}\left(\frac{1}{t}\right) = \tilde{c} \cdot \beta_1$ and $\lim_{t \rightarrow 0_-} t \cdot \varphi_{\alpha,\beta_1,\beta_2,\tilde{c}}\left(\frac{1}{t}\right) = -\tilde{c} \cdot \beta_2$ which are both finite, and hence $p_k \cdot \varphi_{\alpha,\beta_1,\beta_2,\tilde{c}}\left(\frac{q_k}{p_k}\right) = q_k \cdot \frac{p_k}{q_k} \cdot \varphi_{\alpha,\beta_1,\beta_2,\tilde{c}}\left(\frac{q_k}{p_k}\right)$ stays finite as p_k tends to zero. To proceed with our general investigations, with the help of (161) we can derive from (130),(132)

$$\Lambda_{\alpha,\beta_1,\beta_2,\tilde{c}}(z) := \Lambda_{\alpha,\beta_1,\beta_2,\tilde{c}}^{(0)}(z) = \begin{cases} \check{\theta} \cdot z - \tilde{c} \cdot \alpha \cdot \log \left(1 + \frac{z}{\tilde{c}} \cdot \left(\frac{1}{\beta_2} - \frac{1}{\beta_1} \right) - \frac{z^2}{\tilde{c}^2 \cdot \beta_1 \cdot \beta_2} \right), & \text{if } z \in] - \tilde{c} \cdot \beta_2, \tilde{c} \cdot \beta_1[, \\ \infty, & \text{if } z \in] -\infty, -\tilde{c} \cdot \beta_2] \cup [\tilde{c} \cdot \beta_1, \infty[. \end{cases}$$

Notice that $\Lambda_{\alpha,\beta_1,\beta_2,\tilde{c}}(0) = 0$, $\lim_{z \rightarrow -\tilde{c} \cdot \beta_2} \Lambda_{\alpha,\beta_1,\beta_2,\tilde{c}}(z) = \infty$ and $\lim_{z \rightarrow \tilde{c} \cdot \beta_1} \Lambda_{\alpha,\beta_1,\beta_2,\tilde{c}}(z) = \infty$. Moreover, $\Lambda'_{\alpha,\beta_1,\beta_2,\tilde{c}}(-\tilde{c} \cdot \beta_2) = -\infty = a_F$ and $\Lambda'_{\alpha,\beta_1,\beta_2,\tilde{c}}(\tilde{c} \cdot \beta_1) = \infty = b_F$ (which have to be interpreted as limits, as usual).

2) Case 10b:

The analysis for the case $\beta_1 = \beta_2 =: \beta$ can be obtained by taking $\lim_{\beta_1 \rightarrow \beta_2}$ in Case 10a. Alternatively, one can start afresh. Due to its importance and its particularities, we nevertheless state the corresponding results explicitly. To begin with, for any parameter-triple $\alpha, \beta, \tilde{c} \in]0, \infty[$ we choose $] a_F, b_F[:=] a_{F_{\alpha,\beta,\tilde{c}}}, b_{F_{\alpha,\beta,\tilde{c}}}[:=] -\infty, \infty[$ and define with $\check{\theta} := 1$

$$F_{\alpha,\beta,\tilde{c}}(t) := \begin{cases} \frac{\tilde{c} \cdot \alpha}{1-t} \cdot \left(1 - \sqrt{1 + \left(\frac{1-t}{\alpha} \right)^2 \cdot \beta^2} \right), & \text{if } t \in] a_F, b_F[\setminus \{\check{\theta}\}, \\ 0, & \text{if } t = \check{\theta}. \end{cases}$$

Clearly, one has the continuity $\lim_{t \rightarrow \check{\theta}} F_{\alpha,\beta,\tilde{c}}(t) = 0$. Moreover, one can see in a straightforward way that $F_{\alpha,\beta,\tilde{c}}(\cdot)$ is strictly increasing and that $\mathcal{R}(F_{\alpha,\beta,\tilde{c}}) =] -\tilde{c} \cdot \beta, \tilde{c} \cdot \beta[$. Furthermore, $F_{\alpha,\beta,\tilde{c}}(\cdot)$ is smooth on $] a_F, b_F[$, and thus $F_{\alpha,\beta,\tilde{c}} \in \mathfrak{f}$. Since $F_{\alpha,\beta,\tilde{c}}(1) = 0$, let us choose the natural anchor point $c := 0$, which leads to the choice $] \lambda_-, \lambda_+[= \text{int}(\mathcal{R}(F_{\alpha,\beta,\tilde{c}})) =] -\tilde{c} \cdot \beta, \tilde{c} \cdot \beta[$ and $] t_-^{sc}, t_+^{sc}[=] a_F, b_F[=] -\infty, \infty[=] a, b[$. The inverse in (161) collapses to

$$F_{\alpha,\beta,\tilde{c}}^{-1}(x) = 1 + \alpha \cdot \frac{\frac{2x}{\tilde{c} \cdot \beta^2}}{1 - \frac{x^2}{\tilde{c}^2 \cdot \beta^2}}, \quad x \in \text{int}(\mathcal{R}(F_{\alpha,\beta,\tilde{c}}));$$

from this, we can derive from formula (130) (see also (132))

$$\Lambda_{\alpha,\beta,\tilde{c}}(z) := \Lambda_{\alpha,\beta,\tilde{c}}^{(0)}(z) = \begin{cases} \check{\theta} \cdot z - \tilde{c} \cdot \alpha \cdot \log \left(1 - \frac{z^2}{\tilde{c}^2 \cdot \beta^2} \right), & \text{if } z \in] -\tilde{c} \cdot \beta, \tilde{c} \cdot \beta[, \\ \infty, & \text{if } z \in] -\infty, -\tilde{c} \cdot \beta] \cup [\tilde{c} \cdot \beta, \infty[. \end{cases}$$

Notice that $\Lambda_{\alpha,\beta,\tilde{c}}(0) = 0$, $\lim_{z \rightarrow -\tilde{c}\cdot\beta} \Lambda_{\alpha,\beta,\tilde{c}}(z) = \infty$ and $\lim_{z \rightarrow \tilde{c}\cdot\beta} \Lambda_{\alpha,\beta,\tilde{c}}(z) = \infty$. Furthermore, $\lim_{z \rightarrow -\tilde{c}\cdot\beta} \Lambda'_{\alpha,\beta,\tilde{c}}(z) = -\infty$ and $\lim_{z \rightarrow \tilde{c}\cdot\beta} \Lambda'_{\alpha,\beta,\tilde{c}}(z) = \infty$. To proceed, the formula (162) collapses to

$$\varphi_{\alpha,\beta,\tilde{c}}(t) := \varphi_{\alpha,\beta,\tilde{c}}^{(0)}(t) = \tilde{c} \cdot \alpha \cdot \left\{ \sqrt{1 + \beta^2 \cdot \left(\frac{1-t}{\alpha}\right)^2} - 1 + \log \frac{2 \cdot \left(\sqrt{1 + \beta^2 \cdot \left(\frac{1-t}{\alpha}\right)^2} - 1\right)}{\beta^2 \cdot \left(\frac{1-t}{\alpha}\right)^2} \right\} \in [0, \infty[, \quad t \in]-\infty, \infty[. \quad (164)$$

Notice that $\varphi_{\alpha,\beta,\tilde{c}}(1) = 0$, $\varphi'_{\alpha,\beta,\tilde{c}}(1) = 0$, $\varphi_{\alpha,\beta,\tilde{c}}(-\infty) = \infty$ and $\varphi_{\alpha,\beta,\tilde{c}}(\infty) = \infty$. Moreover, $\varphi'_{\alpha,\beta,\tilde{c}}(-\infty) = -\tilde{c} \cdot \beta$ and $\varphi'_{\alpha,\beta,\tilde{c}}(\infty) = \tilde{c} \cdot \beta$. From (164), we construct the corresponding divergence (cf. (4))

$$\begin{aligned} D_{\varphi_{\alpha,\beta,\tilde{c}}}(\mathbf{Q}, \mathbf{P}) &= \sum_{k=1}^K p_k \cdot \varphi_{\alpha,\beta,\tilde{c}}\left(\frac{q_k}{p_k}\right) \\ &= \tilde{c} \cdot \alpha \cdot \sum_{k=1}^K p_k \cdot \left\{ \sqrt{1 + \beta^2 \cdot \left(\frac{1 - \frac{q_k}{p_k}}{\alpha}\right)^2} - 1 + \log \frac{2 \cdot \left(\sqrt{1 + \beta^2 \cdot \left(\frac{1 - \frac{q_k}{p_k}}{\alpha}\right)^2} - 1\right)}{\beta^2 \cdot \left(\frac{1 - \frac{q_k}{p_k}}{\alpha}\right)^2} \right\}, \quad \text{if } \mathbf{P} \in \mathbb{R}_{\geq 0}^K, \mathbf{Q} \in \mathbb{R}^K. \end{aligned} \quad (165)$$

3) Simulation distributions for the Cases 10a,b:

As far as the identification of the corresponding simulation laws ζ (cf. (6)) is concerned, let us first notice that for $\alpha, \beta_1, \beta_2, \tilde{c} \in]0, \infty[$, and anchor point $c = 0$ one can see that — in terms of $\check{\theta} := 1 + \alpha \cdot \left(\frac{1}{\beta_2} - \frac{1}{\beta_1}\right)$ — the derived quantity

$$\Lambda_{\alpha,\beta_1,\beta_2,\tilde{c}}(z) = \check{\theta} \cdot z - \tilde{c} \cdot \alpha \cdot \log \left(1 + \frac{z}{\tilde{c}} \cdot \left(\frac{1}{\beta_2} - \frac{1}{\beta_1}\right) - \frac{z^2}{\tilde{c}^2 \cdot \beta_1 \cdot \beta_2} \right), \quad z \in]-\tilde{c} \cdot \beta_2, \tilde{c} \cdot \beta_1[$$

is the cumulant generating function of a *generalized asymmetric Laplace distribution* $\zeta[\cdot] = \mathbb{P}[W \in \cdot]$ of a random variable $W := \check{\theta} + Z_1 - Z_2$, where Z_1 and Z_2 are auxiliary random variables which are independent and $GAM(\tilde{c} \cdot \beta_1, \tilde{c} \cdot \alpha)$ -distributed respectively $GAM(\tilde{c} \cdot \beta_2, \tilde{c} \cdot \alpha)$ -distributed; for the special case $\tilde{c} = 1$, $\alpha = 1$, $\beta_1 = \beta_2 =: \beta$ (and hence, $\check{\theta} = 1$) one gets that ζ is a *classical Laplace distribution* (two-tailed exponential distribution, bilateral exponential law) with location parameter 1 and scale parameter $\frac{1}{\beta}$. Returning to the general constellation, notice that ζ is an infinitely divisible (cf. Proposition 27) continuous distribution with density

$$f(u) := \frac{\sqrt{2} \cdot \exp\left\{\frac{1}{\sigma \cdot \sqrt{2}} \cdot \left(\frac{1}{\kappa} - \kappa\right) \cdot (u - \theta)\right\}}{\sqrt{\pi} \cdot \sigma^{\tau+1/2} \cdot \Gamma(\tau)} \cdot \left(\frac{\sqrt{2} \cdot |u - \theta|}{\kappa + \frac{1}{\kappa}}\right)^{\tau-1/2} \cdot K_{\tau-1/2}\left(\frac{1}{\sigma \cdot \sqrt{2}} \cdot \left(\kappa + \frac{1}{\kappa}\right) \cdot |u - \theta|\right), \quad u \in \mathbb{R} \setminus \{\theta\}, \quad (166)$$

where $(\theta, \kappa, \sigma, \tau)$ is given in Remark 38 below and K_λ is the modified Bessel function of the third kind with index λ ; for the above-mentioned special case of the classical Laplace distribution, this considerably simplifies to $f(u) := \frac{\beta}{2} \exp\{-\beta \cdot |u - 1|\}$. Moreover, note that in the general case one has $\zeta[]0, \infty[= \mathbb{P}[W > 0] = \int_0^\infty f(u) du \in]0, 1[$, $\zeta[\{0\}] = \mathbb{P}[W = 0] = 0$. Concerning the important Remark 9(i), for i.i.d. copies $(W_i)_{i \in \mathbb{N}}$ of W , the probability distribution $\zeta^{*n_k}[\cdot] := \mathbb{P}[\check{W} \in \cdot]$ of $\check{W} := \sum_{i \in I_k^{(n)}} W_i$ is the same as that of a random variable $\check{W} := \check{\theta} \cdot \text{card}(I_k^{(n)}) + \check{Z}_1 - \check{Z}_2$, where \check{Z}_1 and \check{Z}_2 are auxiliary random variables which are independent and $GAM(\tilde{c} \cdot \beta_1, \tilde{c} \cdot \alpha \cdot \text{card}(I_k^{(n)}))$ -distributed respectively $GAM(\tilde{c} \cdot \beta_2, \tilde{c} \cdot \alpha \cdot \text{card}(I_k^{(n)}))$ -distributed. Within the context of Subsection X-A, for the concrete simulative estimation $D_{\alpha,\beta_1,\beta_2,\tilde{c}}(\Omega, \mathbf{P})$ via (116) and (118), we obtain — in terms of $M_{\mathbf{P}} := \sum_{i=1}^K p_i > 0$, $n_k = n \cdot \tilde{p}_k \in \mathbb{N}$ and \tilde{q}_k^* from proxy method 1 or 2 — that the distribution $\check{U}_k^{*n_k}$ is the same as that of a random variable $\check{W} := \check{\theta} \cdot \text{card}(I_k^{(n)}) + \check{Z}_1 - \check{Z}_2$, where \check{Z}_1 and \check{Z}_2 are auxiliary random variables which are independent and $GAM(\tilde{c} \cdot M_{\mathbf{P}} \cdot \beta_1 - \tau_k, \tilde{c} \cdot M_{\mathbf{P}} \cdot \alpha \cdot \text{card}(I_k^{(n)}))$ -distributed respectively $GAM(\tilde{c} \cdot M_{\mathbf{P}} \cdot \beta_2 + \tau_k, \tilde{c} \cdot M_{\mathbf{P}} \cdot \alpha \cdot \text{card}(I_k^{(n)}))$ -distributed; here, $\tau_k = F_{\alpha,\beta_1,\beta_2,\tilde{c}}\left(\frac{\tilde{q}_k^*}{p_k}\right)$ for $\tilde{q}_k^* \in \mathbb{R}$. Moreover, \widetilde{ISF}_k can be straightforwardly computed by (117).

Remark 38: In the book [406] one can find a very comprehensive study on generalized asymmetric Laplace distributions (also known as Bessel function distributions, McKay distributions), their close relatives (such as e.g. the financial-econometric *variance gamma model* of [407]) as well as their applications; see also e.g. [408] for connections with some other Gamma difference distributions. [406] uses a different parametrization $(\theta, \kappa, \sigma, \tau)$ which is one-to-one with our parametrization $(\check{\theta}, \alpha, \beta_1, \beta_2, \tilde{c} = 1)$, as follows: $\theta = \check{\theta}$, $\tau = \tilde{c} \cdot \alpha$, $\sigma = \frac{1}{\tilde{c}} \cdot \sqrt{\frac{2}{\beta_1 \cdot \beta_2}}$, $\kappa = \sqrt{\frac{\beta_1}{\beta_2}}$. In particular, this implies that we cover *all* generalized asymmetric Laplace distributions with mean 1. For better comparability, we have used the parametrization $(\theta, \kappa, \sigma, \tau)$ in the above-mentioned representation (166) of the density (due to [406]).

XIII. CONCLUDING REMARKS

This paper presents a new approach for the optimization of various different non-linear — deterministic respectively statistical — functionals on \mathbb{R}^K under fairly general (e.g. non-convex, highly disconnected, large-dimensional) constraints, in terms of an appropriately constructed dimension-free bare simulation method which is straightforward to implement and which converges. Algorithms and numerous detailed cases — e.g. for φ -divergences, Renyi divergences, generalized entropies, integer-programming relaxations etc. — are presented with explicit solutions pertaining to the simulations to be performed. As argued, our newly developed optimization procedure does *not* rely on the search for minimizers as a first stage, in contrast with existing methods for similar problems which thus require some involved regularity on the constraint set and also suffer from the so-called “curse of dimensionality”.

Extensions of our new method can be pursued e.g. in two directions. Firstly, the search for optimizers can be handled, either by making use of the simulated values which have been used in order to get the approximation of the optimal value of the objective function, or through dichotomous search. Secondly, the case where the objective function is defined on an infinite-dimensional space (rather than \mathbb{R}^K) is of interest; for instance, in the statistical context this amounts to consider adequacy between a probability distribution P and a model Ω which consists of *continuous* distributions. The basic asymptotics which are used in this paper can be extended to these situations, both in the deterministic case and in the statistical context. This will be part of a follow-up paper.

APPENDIX A PROOFS — PART 1

Proof of Theorem 8. This is a straightforward application of the classical Cramer-type Large Deviation Theorem in the vector case (see Theorem 2.2.30 and Corollary 6.1.6 in [409]). Recall that above we have transformed the original problem into a context where the second argument in $D_\varphi(\cdot, \cdot)$ is a probability vector, as follows: in terms of $M_{\mathbf{P}} := \sum_{i=1}^K p_i > 0$ we normalized $\tilde{\mathbb{P}} := \mathbf{P}/M_{\mathbf{P}}$, and $\tilde{\mathbf{Q}} := \mathbf{Q}/M_{\mathbf{P}}$ for \mathbf{Q} in Ω . With $\tilde{\varphi} \in \Upsilon([a, b])$ defined through $\tilde{\varphi} := M_{\mathbf{P}} \cdot \varphi$, we have obtained

$$D_\varphi(\mathbf{Q}, \mathbf{P}) = \sum_{k=1}^K p_k \cdot \varphi\left(\frac{q_k}{p_k}\right) = \sum_{k=1}^K M_{\mathbf{P}} \cdot \tilde{p}_k \cdot \varphi\left(\frac{M_{\mathbf{P}} \cdot \tilde{q}_k}{M_{\mathbf{P}} \cdot \tilde{p}_k}\right) = D_{\tilde{\varphi}}(\tilde{\mathbf{Q}}, \tilde{\mathbb{P}}) \quad (\text{cf. (12)}).$$

It has followed that the solution of (8) coincides with the one of the problem of finding

$$\tilde{\Phi}_{\tilde{\mathbb{P}}}(\tilde{\Omega}) := \inf_{\tilde{\mathbf{Q}} \in \tilde{\Omega}} D_{\tilde{\varphi}}(\tilde{\mathbf{Q}}, \tilde{\mathbb{P}}), \quad \text{with } \tilde{\Omega} := \Omega/M_{\mathbf{P}} \quad (\text{cf. (13)}).$$

So let us continue by tackling (13). From the assumptions on $\tilde{\varphi}$ and the requirement (15) one can see that

$$\tilde{W}_1 \text{ has moment generating function } MGF_{\tilde{\zeta}}(z) = E_{\mathbb{P}}[e^{z \cdot \tilde{W}_1}] \text{ which is finite on a non-void neighborhood of 0,} \quad (167)$$

$$E_{\mathbb{P}}[\tilde{W}_1] = 1,$$

since $\tilde{\varphi}(1) = 0 = \tilde{\varphi}'(1)$. With the help of these, we obtain the following

Proposition 39: Under the assumptions of Theorem 8, for any set $\tilde{\Omega} \subset \mathcal{M} := \mathbb{R}^K$ with (7) one has

$$\begin{aligned} - \inf_{\tilde{\mathbf{Q}} \in \text{int}(\tilde{\Omega})} D_\varphi(\tilde{\mathbf{Q}}, \tilde{\mathbb{P}}) &\leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P} \left[\boldsymbol{\xi}_n^{\tilde{\mathbf{W}}} \in \tilde{\Omega} \right] \\ &\leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P} \left[\boldsymbol{\xi}_n^{\tilde{\mathbf{W}}} \in \tilde{\Omega} \right] \leq - \inf_{\tilde{\mathbf{Q}} \in \text{cl}(\tilde{\Omega})} D_\varphi(\tilde{\mathbf{Q}}, \tilde{\mathbb{P}}). \end{aligned} \quad (168)$$

Proof of Proposition 39. Recall that $n_k := \text{card}(I_k^{(n)})$ denotes the number of the elements of the block $I_k^{(n)}$ defined right after (16) ($k = 1, \dots, K$). We follow the line of proof of Theorem 2.2.30 in [409], which states the large deviation principle (LDP) for the vector of partial sums of random vectors in \mathbb{R}^K , where we also use Corollary 6.1.6 in [409] in relation with condition (167). Indeed, since by definition the k -th component of the vector $\boldsymbol{\xi}_n^{\tilde{\mathbf{W}}}$ is equal to $\frac{1}{n} \sum_{i \in I_k^{(n)}} \tilde{W}_i$, the current proof will follow from a similar treatment as for the standard Cramer LDP in \mathbb{R}^K . The only difference lies in two facts: firstly, the number of the summands for the k -th coordinate is n_k instead of n in the standard case; secondly, we will need to substitute n_k by its equivalent $n \cdot \tilde{p}_k$, which adds an approximation step. For the upper bound, the proof is based on the corresponding

result for $\mathbf{B} := B_1 \times \cdots \times B_K$ where the B_k 's are open bounded intervals on \mathbb{R}^+ . Since $\lim_{n \rightarrow \infty} \frac{n_k}{n} = \tilde{p}_k$ (cf. (16)), there holds

$$\frac{1}{n} \cdot \log \mathbb{P} \left[\boldsymbol{\xi}_n^{\widetilde{\mathbf{W}}} \in \mathbf{B} \right] = \frac{1}{n} \log \mathbb{P} \left[\bigcap_{k=1}^K \left(\frac{1}{n} \sum_{i \in I_k^{(n)}} \widetilde{W}_i \in B_k \right) \right] = \frac{1}{n} \sum_{k=1}^K \log \mathbb{P} \left[\frac{1+o(1)}{n_k} \sum_{i \in I_k^{(n)}} \widetilde{W}_i \in \frac{1}{\tilde{p}_k} B_k \right], \quad (169)$$

$$\begin{aligned} \text{and hence } \limsup_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P} \left[\boldsymbol{\xi}_n^{\widetilde{\mathbf{W}}} \in \mathbf{B} \right] &\leq \sum_{k=1}^K \tilde{p}_k \cdot \limsup_{n_k \rightarrow \infty} \frac{1}{n_k} \log \mathbb{P} \left[\frac{1}{n_k} \sum_{i \in I_k^{(n)}} \widetilde{W}_i \in \frac{1}{\tilde{p}_k} B_k \right] \\ &\leq - \sum_{k=1}^K \inf_{x_k \in cl(B_k)} \tilde{p}_k \cdot \varphi \left(\frac{x_k}{\tilde{p}_k} \right). \end{aligned} \quad (170)$$

To deduce (170) from (169), we have used (i) the fact that for all k the random variables $\frac{1}{n_k} (1+o(1)) \cdot \sum_{i \in I_k^{(n)}} \widetilde{W}_i$ and $\frac{1}{n_k} \sum_{i \in I_k^{(n)}} \widetilde{W}_i$ are exponentially equivalent in the sense that their difference Δ_{n_k} satisfies $\limsup_{n_k \rightarrow \infty} \frac{1}{n_k} \log \mathbb{P} [|\Delta_{n_k}| > \eta] = -\infty$, making use of the Chernoff inequality for all positive η , as well as (ii) Theorem 4.2.13 in [409]. Now the summation and the inf-operations can be permuted in (170) which proves the claim for the hyper-rectangle \mathbf{B} . As in [409], for a compact set $\widetilde{\Omega}$ we consider its finite covering by such open hyper-rectangles \mathbf{B} and conclude; for Ω being a closed set, a tightness argument holds, following [409] Theorem 2.2.30 verbatim. For the lower bound consider the same hyper-rectangle \mathbf{B} . The argument which locates the tilted distribution at the center of \mathbf{B} , together with the use of the LLN for the corresponding r.v.'s as in [409], in combination with the same approximations as above to handle the approximation of n_k by $n \cdot \tilde{p}_k$, complete the proof. We omit the details. ■

Let us continue with the proof of Theorem 8, by giving the following two helpful lemmas for

$$\Phi_{\mathbb{P}}(\mathbf{A}) := \inf_{\mathbf{Q} \in \mathbf{A}} D_{\varphi}(\mathbf{Q}, \mathbb{P}), \quad \mathbf{A} \subset \mathcal{M} := \mathbb{R}^K. \quad (171)$$

Lemma 40: For any open set $\mathbf{A} \subset \mathcal{M} := \mathbb{R}^K$ one has $\Phi_{\mathbb{P}}(\mathbf{A}) = \Phi_{\mathbb{P}}(cl(\mathbf{A}))$.

This is clear from the continuity of $\Phi_{\mathbb{P}}$.

Lemma 41: For any $\mathbf{A} \subset \mathcal{M} := \mathbb{R}^K$ satisfying (7) one has $\Phi_{\mathbb{P}}(cl(\mathbf{A})) = \Phi_{\mathbb{P}}(\mathbf{A}) = \Phi_{\mathbb{P}}(int(\mathbf{A}))$.

Proof of Lemma 41. Assume first that $\Phi_{\mathbb{P}}(\mathbf{A})$ is finite. Then suppose that \mathbf{A} satisfies (7) and $\Phi_{\mathbb{P}}(cl(\mathbf{A})) < \Phi_{\mathbb{P}}(int(\mathbf{A}))$. The latter implies the existence of a point $\mathbf{a} \in cl(\mathbf{A})$ such that $\mathbf{a} \notin int(\mathbf{A})$ and $D_{\varphi}(\mathbf{a}, \mathbb{P}) = \Phi_{\mathbb{P}}(cl(\mathbf{A}))$. But then, by Lemma 40 and (7) one gets $\Phi_{\mathbb{P}}(int(\mathbf{A})) = \Phi_{\mathbb{P}}(cl(int(\mathbf{A}))) = \Phi_{\mathbb{P}}(cl(\mathbf{A})) = D_{\varphi}(\mathbf{a}, \mathbb{P})$ which leads to a contradiction. When $\Phi_{\mathbb{P}}(\mathbf{A}) = \infty$ then $\Phi_{\mathbb{P}}(cl(\mathbf{A})) = \Phi_{\mathbb{P}}(int(\mathbf{A})) = \Phi_{\mathbb{P}}(\mathbf{A}) = \infty$. ■

Finally, the asymptotic assertion (18) follows from (168), (7) and Lemma 41. This completes the proof of Theorem 8. ■

APPENDIX B PROOFS — PART 2

Before we tackle the proof of Theorem 12, let us introduce the following

Lemma 42: If $\Omega \subset \mathbb{S}^K$ satisfies condition (7), then $\widetilde{\Omega} := \bigcup_{m \neq 0} cl(m \cdot \Omega)$ has the property (7).

This can be deduced in a straightforward way: the assumption implies that $cl(\Omega)$ satisfies (7), and thus also $m \cdot cl(\Omega)$ satisfies (7). But this implies the validity of (7) for the ‘‘cone’’ $\bigcup_{m \neq 0} m \cdot cl(\Omega)$ which is nothing but $\bigcup_{m \neq 0} cl(m \cdot \Omega)$.

Proof of Theorem 12. Recall the interpretations of the two vectors $\boldsymbol{\xi}_{n, \mathbf{X}}^{\mathbf{W}}$ respectively $\boldsymbol{\xi}_{n, \mathbf{X}}^{w\mathbf{W}}$ given in (31) respectively (33), and that the sum of their k components are $\sum_{k=1}^K \frac{1}{n} \sum_{i \in I_k^{(n)}} W_i = \frac{1}{n} \sum_{i=1}^n W_i$ respectively $\sum_{k=1}^K \frac{\sum_{i \in I_k^{(n)}} W_i}{\sum_{k=1}^K \sum_{i \in I_k^{(n)}} W_i} = 1$ (in case of $\sum_{i=1}^n W_i \neq 0$). In the light of these, for $\Omega \subset \mathbb{S}^K$ one gets the set identification

$$\left\{ \boldsymbol{\xi}_{n, \mathbf{X}}^{w\mathbf{W}} \in \Omega \right\} = \bigcup_{m \neq 0} \left\{ \boldsymbol{\xi}_{n, \mathbf{X}}^{\mathbf{W}} \in m \cdot \Omega, \frac{1}{n} \sum_{i=1}^n W_i = m \right\}$$

since $\{\sum_{i=1}^n W_i = 0\}$ amounts to $m = 0$, which cannot hold when $\{\xi_{n,\mathbf{X}}^{w\mathbf{W}} \in \Omega\}$. Now

$$\begin{aligned} \mathbb{P}_{X_1^n}[\xi_{n,\mathbf{X}}^{w\mathbf{W}} \in \Omega] &= \mathbb{P}_{X_1^n}\left[\bigcup_{m \neq 0} \left\{ \xi_{n,\mathbf{X}}^{\mathbf{W}} \in m \cdot \Omega, \frac{1}{n} \sum_{i=1}^n W_i = m \right\}\right] \\ &= \mathbb{P}_{X_1^n}\left[\bigcup_{m \neq 0} \left\{ \xi_{n,\mathbf{X}}^{\mathbf{W}} \in m \cdot \Omega \right\}\right] = \mathbb{P}_{X_1^n}\left[\xi_{n,\mathbf{X}}^{\mathbf{W}} \in \bigcup_{m \neq 0} m \cdot \Omega\right] \end{aligned}$$

since $\{\xi_{n,\mathbf{X}}^{\mathbf{W}} \in m \cdot \Omega\} \subset \{\frac{1}{n} \sum_{i=1}^n W_i = m\}$. Therefore

$$\frac{1}{n} \log \mathbb{P}_{X_1^n}[\xi_{n,\mathbf{X}}^{w\mathbf{W}} \in \Omega] = \frac{1}{n} \log \mathbb{P}_{X_1^n}\left[\xi_{n,\mathbf{X}}^{\mathbf{W}} \in \bigcup_{m \neq 0} m \cdot \Omega\right]. \quad (172)$$

Analogously to the proof of Proposition 39 — applied to $\widetilde{\Omega} := \bigcup_{m \neq 0} m \cdot \Omega$ — one gets in terms of (171)

$$\begin{aligned} -\Phi_{\mathbb{P}}\left(\text{int}\left(\bigcup_{m \neq 0} m \cdot \Omega\right)\right) &\leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_{X_1^n}\left[\xi_{n,\mathbf{X}}^{\mathbf{W}} \in \bigcup_{m \neq 0} m \cdot \Omega\right] \\ &\leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_{X_1^n}\left[\xi_{n,\mathbf{X}}^{\mathbf{W}} \in \bigcup_{m \neq 0} m \cdot \Omega\right] \leq -\Phi_{\mathbb{P}}\left(\text{cl}\left(\bigcup_{m \neq 0} m \cdot \Omega\right)\right). \end{aligned} \quad (173)$$

$$\text{But } \Phi_{\mathbb{P}}\left(\text{int}\left(\bigcup_{m \neq 0} m \cdot \Omega\right)\right) \leq \Phi_{\mathbb{P}}\left(\bigcup_{m \neq 0} \text{int}(m \cdot \Omega)\right) = \inf_{m \neq 0} \Phi_{\mathbb{P}}(\text{int}(m \cdot \Omega)) \quad (174)$$

$$\text{and } \Phi_{\mathbb{P}}\left(\text{cl}\left(\bigcup_{m \neq 0} m \cdot \Omega\right)\right) \geq \Phi_{\mathbb{P}}\left(\bigcup_{m \neq 0} \text{cl}(m \cdot \Omega)\right) = \inf_{m \neq 0} \Phi_{\mathbb{P}}(\text{cl}(m \cdot \Omega)). \quad (175)$$

In fact, the inequality in (174) is straightforward because of $\bigcup_{m \neq 0} \text{int}(m \cdot \Omega) \subset \text{int}\left(\bigcup_{m \neq 0} m \cdot \Omega\right)$ (since the latter is the largest open set contained in $\bigcup_{m \neq 0} m \cdot \Omega$); the inequality in (175) follows from

$$\Phi_{\mathbb{P}}\left(\text{cl}\left(\bigcup_{m \neq 0} m \cdot \Omega\right)\right) \geq \Phi_{\mathbb{P}}\left(\text{cl}\left(\bigcup_{m \neq 0} \text{cl}(m \cdot \Omega)\right)\right) = \Phi_{\mathbb{P}}\left(\bigcup_{m \neq 0} \text{cl}(m \cdot \Omega)\right).$$

An application of Lemma 41 yields $\Phi_{\mathbb{P}}(\text{int}(m \cdot \Omega)) = \Phi_{\mathbb{P}}(m \cdot \Omega) = \Phi_{\mathbb{P}}(\text{cl}(m \cdot \Omega))$ for all $m \neq 0$, and hence

$$\inf_{m \neq 0} \Phi_{\mathbb{P}}(\text{int}(m \cdot \Omega)) = \inf_{m \neq 0} \Phi_{\mathbb{P}}(m \cdot \Omega) = \inf_{m \neq 0} \Phi_{\mathbb{P}}(\text{cl}(m \cdot \Omega)). \quad (176)$$

By combining (172) to (176), one arrives at

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_{X_1^n}[\xi_{n,\mathbf{X}}^{w\mathbf{W}} \in \Omega] &= \lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_{X_1^n}\left[\xi_{n,\mathbf{X}}^{\mathbf{W}} \in \bigcup_{m \neq 0} m \cdot \Omega\right] \\ &= - \inf_{m \neq 0} \Phi_{\mathbb{P}}(m \cdot \Omega) = - \inf_{m \neq 0} \inf_{\mathbf{Q} \in m \cdot \Omega} D_{\varphi}(\mathbf{Q}, \mathbb{P}) = - \inf_{m \neq 0} \inf_{\mathbf{Q} \in \Omega} D_{\varphi}(m \cdot \mathbf{Q}, \mathbb{P}), \end{aligned}$$

where in the second last equality we have “reverted” the notation (171). Note that we did not assume (7) for $\bigcup_{m \neq 0} m \cdot \Omega$. ■

The remaining APPENDICES C ff. are placed in the Supplementary Material of this paper.

ACKNOWLEDGMENT

The authors would like to thank the reviewers and the communicating Associate Editor for their great patience to carefully read through the significantly longer first-submission version (which can be found on arXiv:2107.01693 under a minorly different title), and for their very helpful suggestions on reordering respectively outsourcing some parts thereof; this lead to a more comfortable readability, indeed. W. Stummer is grateful to the Sorbonne Université Paris for its multiple partial financial support and especially to the LPSM for its multiple great hospitality. M. Broniatowski thanks very much the FAU Erlangen-Nürnberg for its partial financial support and hospitality. Moreover, W. Stummer would like to thank Rene Schilling for an interesting discussion on complex-valued foundations of the Bernstein-Widder theorem.

REFERENCES

- [1] I. Csiszár, "Eine informationstheoretische Ungleichung und ihre Anwendung auf den Beweis der Ergodizität von Markoffschen Ketten," *Publ. Math. Inst. Hungar. Acad. Sci.*, Vol. 8, pp. 85–108, 1963.
- [2] M.S. Ali and D. Silvey, "A general class of coefficients of divergence of one distribution from another," *J. Roy. Statist. Soc. Series B (Methodological)*, Vol. 28, no. 1, pp. 131–140, 1966.
- [3] T. Morimoto, "Markov processes and the H-theorem," *J. Phys. Soc. Jpn.*, Vol. 18, no. 3, pp. 328–331, 1963.
- [4] F. Liese and I. Vajda, "On divergences and informations in statistics and information theory," *IEEE Trans. Inf. Theory*, Vol. 52, no. 10, pp. 4394–4412, 2006.
- [5] M. Broniatowski and A. Keziou, "Parametric estimation and tests through divergences and the duality technique," *J. Multiv. Anal.*, Vol. 100, no. 1, pp. 16–36, 2009.
- [6] S. Kullback and R.A. Leibler, *On information and sufficiency*. "Asymptotic distribution of (h, ϕ) -entropies," *Ann. Math. Statistics*, Vol. 22, pp. 79–86, 1951.
- [7] F. Liese and I. Vajda, *Convex Statistical Distances*. Leipzig, Germany: Teubner, 1987.
- [8] T.R.C. Read and N.A.C. Cressie, *Goodness-of-Fit Statistics for Discrete Multivariate Data*. New York, USA: Springer, 1988.
- [9] I. Vajda, *Theory of Statistical Inference and Information*. Dordrecht, NL: Kluwer, 1989.
- [10] I. Csiszár and P.C. Shields, *Information Theory and Statistics: a Tutorial*. Hanover, MA, USA: now Publishers, 2004.
- [11] W. Stummer, *Exponentials, Diffusions, Finance, Entropy and Information*. Aachen, Germany: Shaker, 2004.
- [12] L. Pardo, *Statistical Inference Based on Divergence Measures*. Boca Raton, USA: Chapman & Hall/CRC, 2006.
- [13] F. Liese and K.J. Miescke, *Statistical Decision Theory: Estimation, Testing, and Selection*. New York, USA: Springer, 2008.
- [14] I. Vajda and E.C. van der Meulen, "Goodness-of-fit criteria based on observations quantized by hypothetical and empirical percentiles," in: Z.A. Karian ZA and E.J. Dudewicz (eds.), *Handbook of Fitting Statistical Distributions with R*, pp. 917 – 994. Heidelberg, Germany: CRC, 2010.
- [15] M.D. Reid and R.C. Williamson, "Information, divergence and risk for binary experiments," *J. Machine Learn. Res.*, Vol. 12, pp. 731–817, 2011.
- [16] M. Basseville, "Divergence measures for statistical data processing - an annotated bibliography," *Signal Process.*, Vol. 93, pp. 621–633, 2013.
- [17] W. Stummer and I. Vajda, "On Bregman distances and divergences of probability measures," *IEEE Trans. Inf. Theory*, Vol. 58, no. 3, pp. 1277–1288, 2012.
- [18] M. Broniatowski and I. Vajda, "Several applications of divergence criteria in continuous families," *Kybernetika* Vol. 48, no. 4, pp. 600–636, 2012.
- [19] W. Stummer and A.-L. Kiblinger, "Some new flexibilizations of Bregman divergences and their asymptotics," In: F. Nielsen and F. Barbaresco (eds.), *Geometric Science of Information GSI 2017*, Lecture Notes in Computer Science, vol. 10589, pp. 514–522. Cham, Switzerland: Springer International Publishing, 2017.
- [20] M. Broniatowski and W. Stummer, "Some universal insights on divergences for statistics, machine learning and artificial intelligence," in: F. Nielsen (ed.), *Geometric Structures of Information*, Ser. Signals and Communications Technology, pp. 149–211. Cham, Switzerland: Springer Nature Switzerland, 2019.
- [21] M. Broniatowski and W. Stummer, "A unifying framework for some directed distances in statistics," in: F. Nielsen, A.S.R.S. Rao, C.R. Rao (eds.), *Geometry and Statistics*, Handbook of Statistics, Vol. 46, pp. 145–223. Cambridge MA, USA: Academic Press, 2022.
- [22] A. Bhattacharyya, "On a measure of divergence between two statistical populations defined by their probability distributions," *Bull. Calcutta Math. Soc.*, Vol. 35, pp. 99–109, 1943.
- [23] A. Rényi, "On measures of entropy and information," in: J. Neyman (ed.), *Proc. 4th Berkeley Symp. Math. Stat. Probab. Vol.1*, pp. 547–561. Berkeley, CA, USA: Univ. of California Press, 1961.
- [24] T. van Erven and P. Harremoës, "Rényi divergence and Kullback-Leibler divergence," *IEEE Trans. Inf. Theory*, Vol. 60, no. 7, pp. 3797–3820, 2014.
- [25] A. Bhattacharyya, "On a measure of divergence between two multinomial populations," *Sankhya*, Vol. 7, no. 4, pp. 401–406, 1946.
- [26] A. Bhattacharyya, "On some analogues of the amount of information and their use in statistical estimation (contd.)," *Sankhya*, Vol. 8, no. 3, pp. 201–218, 1947.
- [27] K. Matusita, "On the theory of statistical functions," *Ann. Inst. Statist. Math. Tokyo*, Vol. 3, pp. 17–35, 1951; corrections in *Ann. Inst. Statist. Math. Tokyo*, Vol. 4, pp. 51–53, 1952.
- [28] M.M. Deza and E. Deza, *Encyclopedia of Distances*, 4th ed. Berlin, Germany: Springer, 2016.
- [29] R. Sundaresan, "A measure of discrimination and its geometric properties," in: *Proc. 2002 IEEE Int. Symp. Inf. Theory (ISIT 2002)*, Lausanne, Switzerland, p. 264.
- [30] R. Sundaresan, "Guessing Under Source Uncertainty," *IEEE Trans. Inf. Theory*, Vol. 53, no. 1, pp. 269–287, 2007.
- [31] C. Burbea and C.R. Rao, "On the convexity of some divergence measures based on entropy functions," *IEEE Trans. Inf. Theory*, Vol. 28, no. 3, pp. 489–495, 1982.
- [32] I. Csiszár, "A class of measures of informativity of observation channels," *Periodica Mathem. Hungar.*, Vol. 2, no. 1–4, pp. 191–213, 1972.
- [33] M. Ben-Bassat, "f-entropies, probability of error, and feature selection," *Information and Control*, Vol. 39, pp. 227–242, 1978.
- [34] A. Ben-Tal and M. Teboulle, "Rate-distortion theory with generalized information measures via convex programming duality," *IEEE Trans. Inf. Theory*, Vol. 32, no. 5, pp. 630–641, 1986.
- [35] H.K. Kesavan and J.N. Kapur, "The generalized maximum entropy principle," *IEEE Trans. Syst. Man Cyb.*, Vol. 19, no. 5, pp. 1042–1052, 1989.
- [36] D. Dacunha-Castelle and F. Gamboa, "Maximum d'entropie et probleme des moments," *Ann. Inst. Henri Poincare*, Vol. 26, no. 4, pp. 567–596, 1990.
- [37] M. Teboulle and I. Vajda, "Convergence of best ϕ -entropy estimates," *IEEE Trans. Inf. Theory*, Vol. 39, no. 1, pp. 297–301, 1993.
- [38] F. Gamboa and E. Gassiat, "Asymptotic distribution of (h, ϕ) -entropies," *Ann. Stat.*, Vol. 25, no. 1, pp. 328–350, 1997.
- [39] I. Vajda and J. Zvarova, "On generalized entropies, Bayesian decisions and statistical diversity," *Kybernetika*, Vol. 43, no. 5, pp. 675–696, 2007.
- [40] M. Salicru, M.L. Menendez, D. Morales and L. Pardo, "Asymptotic distribution of (h, ϕ) -entropies," *Commun. Statist. - Theory Meth.*, Vol. 22, no. 7, pp. 2015–2031, 1993.
- [41] C.E. Shannon, "A mathematical theory of communication," *Bell System Technical Journal*, Vol. 27, no. 3, pp. 379–423, 1948.
- [42] J. Havrda and F. Charvat, "Quantification method of classification process," *Kybernetika*, Vol. 3, pp. 30–34, 1967.
- [43] C. Tsallis, "Possible generalization of Boltzmann-Gibbs Statistics," *Journal of Statistical Physics*, Vol. 52, no. 1/2, pp. 479–487, 1988.
- [44] S. Arimoto, "Information-theoretical considerations on estimation problems," *Information and Control*, Vol. 19, pp. 181–194, 1971.
- [45] B.D. Sharma and D.P. Mittal, "New nonadditive measures of entropy for discrete probability distributions," *J. Math. Sci. (Delhi)*, Vol. 10, pp. 28–40, 1975.
- [46] M.O. Hill, "Diversity and evenness: a unifying notation and its consequences," *Ecology*, Vol. 54, no. 2, pp. 427–431, 1973.
- [47] P. Harremoës and F. Topsøe, "Inequalities between entropy and index of coincidence derived from information diagrams," *IEEE Trans. Inf. Theory*, Vol. 47, no. 7, pp. 2944–2960, 2001.
- [48] P. Harremoës and I. Vajda, "On the Bahadur-efficient testing of uniformity by means of the entropy," *IEEE Trans. Inf. Theory*, Vol. 54, no. 1, pp. 321–331, 2008.
- [49] G.P. Patil and C. Taillie, "Diversity as a concept and its measurement," *J. Amer. Statist. Assoc.*, Vol. 77, no. 379, pp. 548–561, 1982.
- [50] J.C.A. van der Lubbe, "An axiomatic theory of heterogeneity and homogeneity," *Metrika*, Vol. 33, pp. 223–245, 1986.
- [51] T.K. Nayak, "On diversity measures based on entropy functions," *Commun. Statist. - Theory Meth.*, Vol. 14, no. 1, pp. 203–215, 1985.
- [52] L. Jost, "Entropy and diversity," *OIKOS*, Vol. 113, no. 2, pp. 363–375, 2006.

- [53] E.T. Jaynes, "Information theory and statistical mechanics I," *Phys. Rev.*, Vol. 106, pp. 620–630, 1957.
- [54] E.T. Jaynes, "Information theory and statistical mechanics II," *Phys. Rev.*, Vol. 108, pp. 171–190, 1957.
- [55] J.N. Kapur, *Maximum-entropy models in science and engineering*. New York, NY, USA: John Wiley & Sons, 1989.
- [56] J.N. Kapur and H.K. Kesavan, *Entropy Optimization Principles With Applications*. San Diego, CA, USA: Academic Press, 1992.
- [57] C. Arndt, *Information Measures*. Berlin, Germany: Springer, 2001.
- [58] H. Gzyl, S. Mayoral and E. Gomes-Goncalves, *Loss Data Analysis. The Maximum Entropy Approach*. Berlin, Germany: de Gruyter, 2018.
- [59] B.G. Lindsay, "Statistical distances as loss functions in assessing model adequacy," in: M.P. Taper and S.R. Lele (eds.), *The Nature of Scientific Evidence*, pp. 439–487. Chicago, IL, USA: The University of Chicago Press, 2004. This includes comments by D.R. Cox and S.P. Ellner as well as a rejoinder by the author.
- [60] B.G. Lindsay, M. Markatou, S. Ray, K. Yang, and S.-C. Chen, "Quadratic distances on probabilities: a unified foundation," *Ann. Statist.*, Vol. 36, no. 2, pp. 983–1006, 2008.
- [61] M. Markatou and E. Sofikitou, "Non-quadratic distances in model assessment," *Entropy*, Vol. 20, No. 464, 2018; doi:10.3390/e20060464.
- [62] M. Markatou and Y. Chen, "Statistical distances and the construction of evidence functions for model adequacy," *Front. Ecol. Evol.*, Vol. 7, No. 447, 2019; doi: 10.3389/fevo.2019.00447.
- [63] I. Bilik and P. Khomchuk, "Minimum divergence approaches for robust classification of ground moving targets," *IEEE Trans. Aero. Elec. Sys.*, Vol. 48, no. 1, pp. 581–603, 2012.
- [64] I. Csiszár, "I-divergence geometry of probability distributions and minimization problems," *Ann. Probab.*, Vol. 3, no. 1, pp. 146–158, 1975.
- [65] I. Csiszár and F. Matúš, "Information projections revisited," *IEEE Trans. Inf. Theory*, Vol. 49, no. 6, pp. 1474–1490, 2003.
- [66] M. Broniatowski and A. Keziou, "Minimization of ϕ -divergences on sets of signed measures," *Stud. Scient. Math. Hungar.*, Vol. 43, pp. 403–442, 2006.
- [67] C. Léonard, "Minimization of entropy functionals," *J. Math. Anal. Appl.*, Vol. 346, pp. 183–204, 2008.
- [68] B. Pelletier, "Inference in φ -families of distributions," *Statistics*, Vol. 45, no. 3, pp. 223–236, 2011.
- [69] M.A. Kumar and I. Sason, "Projection theorems for the Renyi divergence on α -convex sets," *IEEE Trans. Inf. Theory*, Vol. 62, no. 9, pp. 4924–4935, 2016.
- [70] M.A. Kumar and R. Sundaresan, "Minimization problems based on relative α -entropy I: forward projection," *IEEE Trans. Inf. Theory*, Vol. 61, no. 9, pp. 5063–5080, 2015.
- [71] M.A. Kumar and R. Sundaresan, "Minimization problems based on relative α -entropy II: reverse projection," *IEEE Trans. Inf. Theory*, Vol. 61, no. 9, pp. 5081–5095, 2015.
- [72] M. Broniatowski and A. Decurninge, "Estimation for models defined by conditions on their L-moments," *IEEE Trans. Inf. Theory*, Vol. 62, no. 9, pp. 5181–5198, 2016.
- [73] J. Liu and B.G. Lindsay, "Building and using semiparametric tolerance regions for parametric multinomial models," *Ann. Statist.*, Vol. 37, no. 6A, pp. 3644–3659, 2009.
- [74] A. Ghosh and A. Basu, "A new family of divergences originating from model adequacy tests and applications to robust statistical inference," *IEEE Trans. Inf. Theory*, Vol. 64, no. 8, pp. 5581–5591, 2018.
- [75] M. Broniatowski, "A weighted bootstrap procedure for divergence minimization problems," in: J. Antoch et al. (eds.), *Analytical Methods in Statistics*, Springer Proc. Math. Stat. 193, pp. 1–22. Cham, Switzerland: Springer International, 2017.
- [76] Y. Qiao and N. Minematsu, "A study on invariance of f-divergence and its application to speech recognition," *IEEE Trans. Signal Process.*, Vol. 58, no. 7, pp. 3884–3890, 2010.
- [77] X. Nguyen, M.J. Wainwright, and M.I. Jordan, "Estimating divergence functionals and the likelihood ratio by convex risk minimization," *IEEE Trans. Inf. Theory*, Vol. 56, no. 11, pp. 5847–5861, 2010.
- [78] M. Feixas, A. Bardera, J. Rigau, Q. Xu and M. Sbert, *Information Theory Tolls for Image Processing*. San Rafael, CA, USA: Morgan & Claypool, 2014.
- [79] X. Luo, Q. Xu, M. Sbert and Klaus Schoeffmann, "F-divergences driven video key frame extraction," in: *Proc. 2014 IEEE Int. Conf. on Multimedia and Expo (ICME)*, 2014, 6 pages.
- [80] A.-L. Kißlinger and W. Stummer, "New model search for nonlinear recursive models, regressions and autoregressions," in: F. Nielsen and F. Barbaresco (eds.), *Geometric Science of Information GSI 2015*, Lecture Notes in Computer Science, vol. 9389, pp. 693–701. Berlin, Germany: Springer, 2015.
- [81] Z. Mahboubi and M.J. Kochenderfer, "Learning traffic patterns at small airports from flight tracks," *IEEE Trans. Intell. Transp.*, Vol. 18, no. 4, pp. 917–926, 2017.
- [82] S. Guo, W. Huang and Y. Qiao, "Improving scale invariant feature transform with local color contrastive descriptor for image classification," *J. Electron. Imaging*, Vol. 26, No. 1, 013015, 2017; doi:10.1117/1.JEI.26.1.013015.
- [83] I. Csiszár and T. Breuer, "Expected value minimization in information theoretic multiple priors models," *IEEE Trans. Inf. Theory*, Vol. 64, no. 6, pp. 3957–3974, 2018.
- [84] A.-L. Kißlinger and W. Stummer, "A new toolkit for robust distributional change detection," *Appl. Stochastic Models Bus. Ind.*, Vol. 34, pp. 682–699, 2018.
- [85] Y.-K. Noh, B.-T. Zhang and D.D. Lee, "Generative local metric learning for nearest neighbor classification," *IEEE Trans. Pattern Anal. Mach. Intell.*, Vol. 40, no. 1, pp. 106–118, 2018.
- [86] X. Yu, H. Zhang, C. Luo, H. Qi and P. Ren, "Oil spill segmentation via adversarial f-divergence learning," *IEEE Trans. Geosci. Remote Sens.*, Vol. 56, no. 9, pp. 4973–4988, 2018.
- [87] Ö. Arslan, "Statistical coverage control of mobile sensor networks," *IEEE Trans. Robot.*, Vol. 35, no. 4, pp. 889–908, 2019.
- [88] I. Sason, "On data-processing and majorization inequalities for f-Divergences with applications," *Entropy*, Vol. 21, No. 1022, 2019; doi:10.3390/e21101022.
- [89] O. Ciftci, M. Mehrdash, and A. Kargarian, "Data-driven nonparametric chance-constrained optimization for microgrid energy management," *IEEE Trans. Ind. Inform.*, Vol. 16, no. 4, pp. 2447–2457, 2020.
- [90] W. Stummer, "Optimal transport with some directed distances," in: F. Nielsen and F. Barbaresco (eds.), *Geometric Science of Information GSI 2021*, Lecture Notes in Computer Science, vol. 12829, pp. 829–840. Cham, Switzerland: Springer Nature Switzerland, 2021.
- [91] W. Stummer and I. Vajda, "On divergences of finite measures and their applicability in statistics and information theory," *Statistics*, Vol. 44, no. 2, pp. 169–187, 2010.
- [92] C. Gietl and F.P. Ruffel, "Continuity of f -projections and applications to the iterative proportional fitting procedure," *Statistics*, Vol. 51, no. 3, pp. 668–684, 2017.
- [93] G. Avlogiaris, A. Micheas and K. Zografos, "On local divergences between two probability measures," *Metrika*, Vol. 79, pp. 303–333, 2016.
- [94] G. Avlogiaris, A. Micheas and K. Zografos, "On testing local hypotheses via local divergence," *Statist. Methodol.*, Vol. 31, pp. 20–42, 2016.
- [95] I. Csiszár, "Sanov property, generalized I-projection and a conditional limit theorem," *Ann. Probab.*, Vol. 12, no. 3, pp. 768–793, 1984.
- [96] A. György and T. Linder, "Optimal entropy-constrained scalar quantization of a uniform source," *IEEE Trans. Inf. Theory*, Vol. 46, no. 7, pp. 2704–2711, 2000.
- [97] R.R. Tucci, "Method for sampling probability distributions using a quantum computer," *United States Patent*, Patent No. US 8543627 B1, 24th Sep. 2013.
- [98] J.S. Teh, A. Samsudin, M. Al-Mazrooe and A. Akhavan, "GPU and chaos: a new true random number generator," *Nonlinear Dyn.*, Vol. 82, pp. 1913–1922, 2015.
- [99] C. Aghamohammadi and J.P. Crutchfield, "Thermodynamics of random number generation," *Phys. Rev. E*, Vol. 95, pp. 062139-1–062139-11, 2017.

- [100] M. Herrero-Collantes and J.C. Garcia-Escartin, “Quantum random number generators,” *Rev. Mod. Phys.*, Vol. 89, no. 1, pp. 015004-1 – 015004-48, 2017.
- [101] K.A. Balygin, V.I. Zaitsev, A.N. Klimov, S.P. Kulik and S.N. Molotov, “A quantum random number generator based on the 100-Mbit/s Poisson photocount statistics,” *J. Exp. Theor. Phys.*, Vol. 126, no. 6, pp. 728–740, 2018.
- [102] B. Dang, J. Sun, T. Zhang, S. Wang, M. Zhao, K. Liu, L. Xu, J. Zhu, C. Cheng, L. Bao, Y. Yang, H. Wang, Y. Hao and R. Huang, “Physically transient true random number generators based on paired threshold switches enabling Monte Carlo method applications,” *IEEE Electron. Device Lett.*, Vol. 40, no. 7, pp. 1096–1099, 2019.
- [103] L. Gong, J. Zhang, H. Liu, L. Sang and Y. Wang, “True random number generators using electrical noise,” *IEEE Access*, Vol. 7, pp. 125796–125805, 2019.
- [104] S.T. Chandrasekaran, V.E.G. Karnam and A. Sanya, “0.36-mW, 52-Mbps true random number generator based on a stochastic delta-sigma modulator,” *IEEE Solid-State Lett.*, Vol. 3, pp. 190–193, 2020.
- [105] D. Drahi, N. Walk, M.J. Hoban, A.K. Fedorov, R. Shakhovoy, A. Feimov, Y. Kurochkin, W.S. Kolthammer, J. Nunn, J. Barrett and I.A. Walmsley, “Certified quantum random numbers from untrusted light,” *Phys. Rev. X*, Vol. 10, pp. 041048-1 – 041048-32, 2020.
- [106] C. Kollmitzer, S. Schauer, S. Rass and B. Rainer (eds.), *Quantum Random Number Generation*. Cham, Switzerland: Springer Nature, 2020.
- [107] Y. Liu, C. Chen, D.D. Yang, Q. Li and X. Li, “Fast true random number generator based on chaotic oscillation in self-feedback weakly coupled superlattices,” *IEEE Access*, Vol. 8, pp. 182693–182703, 2020.
- [108] I. Fischer and D.J. Gauthier, “High-speed harvesting of random numbers,” *Science*, Vol. 371, 26 February 2021, pp. 889–890, 2021.
- [109] K. Kim, S. Bittner, Y. Zeng, S. Guazzotti, O. Hess, Q.J. Wang and Hui Cao, “Massively parallel ultrafast random bit generation with a chip-scale laser,” *Science*, Vol. 371, 26 February 2021, pp. 948–952, 2021.
- [110] S. Stoller and K.A. Campbell, “Demonstration of three true random number generator circuits using memristor created entropy and commercial off-the-shelf components,” *Entropy*, Vol. 23, No. 371, 2021; doi:10.3390/e23030371.
- [111] A. Schrijver, *Combinatorial Optimization*, Vol.A,B,C. Heidelberg, Germany: Springer, 2003.
- [112] D. Bertsimas and R. Weismantel, *Optimization over Integers*. Belmont, MA, USA: Dynamic Ideas, 2005.
- [113] D.-S. Chen, R.G. Batson and Y. Dang, *Applied Integer Programming*. Hoboken, NJ, USA: Wiley, 2010.
- [114] S. Onn, *Nonlinear Discrete Optimization*. Zürich, Switzerland: European Math. Society Publ. House, 2010.
- [115] B. Korte and J. Vygen, *Combinatorial Optimization*, 6th ed. Berlin, Germany: Springer, 2018.
- [116] L.A. Wolsey, *Integer Programming*, 2nd ed. Hoboken, NJ, USA: Wiley, 2021.
- [117] I. Vajda and K. Vasek, “Majorizations, concave entropies, and comparison of experiments,” *Problems of Control and Information Theory*, Vol. 14, no. 2, pp. 105–115, 1985.
- [118] B. Chen, J. Hu, L. Pu, and Z. Sun, “Stochastic gradient algorithm under (h, ϕ) -entropy,” *Circuits Syst. Signal Process.*, Vol. 26, pp. 941–960, 2007.
- [119] M. Ren, J. Zhang, M. Jiang, M. Yu, and J. Xu, “Minimum (h, ϕ) -entropy control for non-Gaussian stochastic networked control systems and its application to a networked DC motor control system,” *IEEE Trans. Control Syst. Technol.*, Vol. 23, no. 1, pp. 406–411, 2015.
- [120] M. Broniatowski and A. Keziou, “Divergences and duality for estimation and test under moment condition models,” *J. Statistical Planning and Inference*, Vol. 142, pp. 2554–2573, 2012.
- [121] I. Vajda, “About perceptron realizations of Bayesian decisions,” in: *IEEE International Conference on Neural Networks, Washington 1996*, pp. 253–257, 1996.
- [122] P.N. Rathie and P. Kannappan, “A directed-divergence function of type β ,” *Information and Control*, Vol. 20, pp. 38–45, 1972.
- [123] N. Cressie and T.R.C. Read, “Multinomial goodness-of-fit tests,” *J. R. Statist. Soc. B*, Vol. 46, no. 3, pp. 440–464, 1984.
- [124] C. Tsallis, “Generalized entropy-based criterion for consistent testing,” *Phys. Rev. E*, Vol. 58, No. 2, pp. 1442–1445, 1998;
- [125] M. Shiino, “H-Theorem with generalized relative entropies and the Tsallis statistics,” *J. Phys. Soc. Japan*, Vol. 67, no. 11, pp. 3658–3660, 1998.
- [126] S.-I. Amari, *Differential-Geometrical Methods in Statistics*. Berlin, Germany: Springer, 1985.
- [127] Y. Matsuyama, “The α -EM algorithm: surrogate likelihood maximization using α -logarithmic information measures,” *IEEE Trans. Inf. Theory*, Vol. 49, no. 3, pp. 692–706, 2003.
- [128] Y. Matsuyama, “The Alpha-HMM estimation algorithm: prior cycle guides fast paths,” *IEEE Trans. Signal Process.*, Vol. 65, no. 13, pp. 3446–3461, 2017.
- [129] Y. Matsuyama, “Divergence family attains blockchain applications via α -EM algorithm,” in: *Proc. 2019 IEEE Int. Symp. Inf. Theory (ISIT 2019), Paris, France*, p. 727–731.
- [130] C.-J. Ku and T.L. Fine, “Testing for stochastic independence: application to blind source separation,” *IEEE Trans. Signal Process.*, Vol. 53, no. 5, pp. 1815–1826, 2005.
- [131] W. Stummer and I. Vajda, “Optimal statistical decisions about some alternative financial models,” *J. Econometrics*, Vol. 137, no.2, pp. 441–471, 2007.
- [132] W. Stummer and W. Lao, “Limits of Bayesian decision related quantities of binomial asset price models,” *Kybernetika*, Vol. 48, no.4, pp. 750–767, 2012.
- [133] D. Berend, P. Harremoës and A. Kontorovich, “Minimum KL-divergence on complements of L_1 balls,” *IEEE Trans. Inf. Theory*, Vol. 60, no. 6, pp. 3172–3177, 2014.
- [134] J. Verrelst, J.P. Rivera, G. Leonenko, L. Alonso, and J. Moreno, “Optimizing LUT-based RTM inversion for semiautomatic mapping of crop biophysical parameters from Sentinel-2 and -3 Data: role of cost functions,” *IEEE Trans. Geosci. Remote Sens.*, Vol. 52, no. 1, pp. 257–269, 2014.
- [135] O. Salem, A. Serhrouchni, A. Mehaoua, and R. Boutaba, “Event detection in wireless body area networks using Kalman filter and power divergence,” *IEEE Trans. Netw. Serv. Manag.*, Vol. 15, no. 3, pp. 1018–1034, 2018.
- [136] Q. Fu, S.B. Villas-Boas and G. Judge, “Entropy-based China income distributions and inequality measures,” *China Econ. J.*, Vol. 12, no. 3, pp. 352–368, 2019.
- [137] A. Iqbal and A.-K. Seghouane, “An α -divergence-based approach for robust dictionary learning,” *IEEE Trans. Image Process.*, Vol. 28, no. 11, pp. 5729–5739, 2019.
- [138] S. Krömer and W. Stummer, “A new toolkit for mortality data analytics,” in: A. Steland, E. Rafajlowicz, and O. Okhrin (eds.), *Stochastic Models, Statistics and Their Applications*, pp. 393–407. Cham, Switzerland: Springer Nature Switzerland, 2019.
- [139] L. Cai, Y. Chen, N. Cai, W. Cheng and H. Wang, “Utilizing Amari-alpha divergence to stabilize the training of generative adversarial networks,” *Entropy*, Vol. 22, No. 410, 2020; doi:10.3390/e22040410.
- [140] T. Fu, C. Xiao, L.M. Glass and J. Sun, “ α -MOP: molecule optimization with α -divergence,” in: *Proc. 2020 IEEE Int. Conf. on Bioinformatics and Biomedicine (BIBM)*, pp. 240–244, 2020; doi:10.1109/BIBM49941.2020.9313409.
- [141] N.B. Kammerer and W. Stummer, “Some dissimilarity measures of branching processes and optimal decision making in the presence of potential pandemics,” *Entropy*, Vol. 22(8), No. 874, 2020 (123 pages); doi:10.3390/e22080874.
- [142] D.T. Kanapram, F. Patrone, P. Marin-Plaza, M. Marchese, E.L. Bodanese, L. Marcenaro, D.M. Gomez and C. Regazzoni, “Collective awareness for abnormality detection in connected autonomous vehicles,” *IEEE Internet Things J.*, Vol. 7, no. 5, pp. 3774–3789, 2020.
- [143] M. Kumbhakar, “Streamwise velocity profile in open-channel flow based on Tsallis relative entropy,” *Chaos*, Vol. 30, 073136, 2020; doi:10.1063/1.5144867.
- [144] A.B. Dharmawan, S. Mariana, G. Scholz, P. Hörmann, T. Schulze, K. Triyana, Mayra Garces-Schröder, I. Rustenbeck, K. Hiller, H.S. Wasisto and A. Waag, “Nonmechanical parafocal and autofocus features based on wave propagation distribution in lensfree holographic microscopy,” *Sci. Rep.*, Vol. 11, 3213, 2021; doi:10.1038/s41598-021-81098-7.

- [145] X. Liu and S. Sun, "Alpha-divergence minimization with mixed variational posterior for Bayesian neural networks and its robustness against adversarial examples," *Neurocomputing*, Vol. 423, pp. 427–434, 2021.
- [146] A.M. Rekevandi, A.-K. Seghouane and R.J. Evans, "Robust subspace detectors based on α -divergence with application to detection in imaging," *IEEE Trans. Image Process.*, Vol. 30, pp. 5017-5031, 2021.
- [147] A.-K. Seghouane and N. Shokouhi, "Adaptive learning for robust radial basis function networks," *IEEE Trans. Cybernetics*, Vol. 51, No. 5, pp. 2847–2856, 2021.
- [148] Y. Wang, P. Wang, Z. Liu and L.Y. Zhang, "A new item similarity based on α -divergence for collaborative filtering in sparse data," *Expert Syst. Appl.*, Vol. 166, 114074, 2021; doi:10.1016/j.eswa.2020.114074.
- [149] J. S. Sigmon et al., "Content and Performance of the MiniMUGA Genotyping Array: A New Tool To Improve Rigor and Reproducibility in Mouse Research," *Genetics*, Vol. 216, pp. 905–930, 2020.
- [150] W. Ha, E.Y. Sidky, R.F. Barber, T.G. Schmidt and X. Pan, "Estimating the spectrum in computed tomography via Kullback-Leibler divergence constrained optimization," *Med. Phys.*, Vol. 46, no. 1, pp. 81–92, 2019.
- [151] S. Bekhet and A. Ahmed, "Evaluation of similarity measures for video retrieval," *Multimed. Tools Appl.*, Vol. 79, pp. 6265–6278, 2020.
- [152] L.T. Luppino, F.M. Bianchi, G. Moser and S.N. Anfinsen, "Unsupervised image regression for heterogeneous change detection," *IEEE Trans. Geosci. Remote Sens.*, Vol. 57, no. 12, pp. 9960–9975, 2019.
- [153] J. Görtler, T. Spinner, D. Streeb, D. Weiskopf and O. Deussen, "Uncertainty-aware principal component analysis," *IEEE Trans. Visual. Comput. Graph.*, Vol. 26, no. 1, pp. 822–833, 2020.
- [154] T. Zhang, J. Cheng, H. Fu, Z. Gu, Y. Xiao, K. Zhou, S. Gao, R. Zheng and J. Liu, "Noise adaptation generative adversarial network for medical image analysis," *IEEE Trans. Med. Imag.*, Vol. 39, no. 4, pp. 1149–1159, 2020.
- [155] B. Chen, X. He, B. Pan, X. Zou and N. You, "Comparison of beta diversity measures in clustering the high-dimensional microbial data," *PLoS ONE*, Vol. 16, no. 2, e0246893, 2021. doi:10.1371/journal.pone.0246893.
- [156] M. Broniatowski, E. Miranda and W. Stummer, "Testing the number and the nature of the components in a mixture distribution," in: F. Nielsen and F. Barbaresco (eds.), *Geometric Science of Information GSI 2019*, Lecture Notes in Computer Science, vol. 11712, pp. 309–318. Cham, Switzerland: Springer Nature Switzerland, 2019.
- [157] A.-L. Kießlinger and W. Stummer, "Robust statistical engineering by means of scaled Bregman distances," in: C. Agostinelli, A. Basu, P. Filzmoser and D. Mukherjee (eds.), *Recent Advances in Robust Statistics – Theory and Applications*, pp. 81–113. New Delhi, India: Springer, 2016.
- [158] M. Broniatowski and W. Stummer, "A precise bare simulation approach to the minimization of some distances. Foundations," *arXiv:2107.01693v1*, 94 pages, July 2021; typo-correction in *arXiv:2107.01693v2*, 93 pages, October 2022.
- [159] S. Guisau and C. Reischer, "The relative information generating function," *Inf. Sci.*, Vol. 35, no. 3, pp. 235–241, 1985.
- [160] D.E. Clark, "Local entropy statistics for point processes," *IEEE Trans. Inf. Theory*, Vol. 66, no. 2, pp. 1155–1163, 2020.
- [161] H. Chernoff, "A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations," *Ann. Math. Statistics*, Vol. 23, pp. 493–507, 1952.
- [162] E. Torgersen, *Comparison of Statistical Experiments*. Cambridge, UK: Cambridge University Press, 1991.
- [163] G.T. Toussaint, "Probability of error, expected divergence and the affinity of several distributions," *IEEE Trans. Syst. Man Cyb.*, Vol. 8, no. 6, pp. 482–485, 1978.
- [164] X.-B. Peng and Z.-J. Li, "Target scale adaptive control based on comparing Bhattacharyya coefficient," *Adv. Mat. Res.*, Vol. 971–973, pp. 1772–1777, 2014.
- [165] I.B. Ayed, K. Punithakumar and S. Li, "Distribution matching with the Bhattacharyya similarity: a bound optimization framework," *IEEE Trans. Pattern Anal. Mach. Intell.*, Vol. 37, no. 9, pp. 1777–1791, 2015.
- [166] B.K. Patra, R. Launonen, V. Ollikainen and S. Nandi, "A new similarity measure using Bhattacharyya coefficient for collaborative filtering in sparse data," *Knowl.-Based Syst.*, Vol. 82, pp. 163–177, 2015.
- [167] Y. El Merabet, Y. Ruichek, S. Ghaffarian, Z. Samir, T. Boujija, R. Messoussi, R. Touahni and A. Sbihi, "Maximal similarity based region classification method through local image region descriptors and Bhattacharyya coefficient-based distance: application to horizon line detection using wide-angle camera," *Neurocomputing*, Vol. 265, pp. 28–41, 2017.
- [168] P.-H. Chiu, P.-H. Tseng and K.-T. Feng, "Interactive mobile augmented reality system for image and hand motion tracking," *IEEE Trans. Veh. Technol.*, Vol. 67, no. 10, pp. 9995–10009, 2018.
- [169] Y.-K. Noh, J. Hamm, F.C. Park, B.-T. Zhang and D.D. Lee, "Fluid dynamic models for Bhattacharyya-based discriminant analysis," *IEEE Trans. Pattern Anal. Mach. Intell.*, Vol. 40, no. 1, pp. 92–105, 2018.
- [170] L. Bai, A. Velichko and B.W. Drinkwater, "Grain scattering noise modeling and its use in the detection and characterization of defects using ultrasonic arrays," *IEEE Trans. Ultrason. Ferr. Freq. Control*, Vol. 66, no. 11, pp. 1798–1813, 2019.
- [171] V.S. Dixit and P. Jain, "Proposed similarity measure using Bhattacharyya coefficient for context aware recommender system," *J. Intell. Fuzzy Syst.*, Vol. 36, pp. 3105–3117, 2019.
- [172] W. Guan, Z. Liu, S. Wen, H. Xie and X. Zhang, "Visible light dynamic positioning method using improved camshift-Kalman algorithm," *IEEE Photonics J.*, Vol. 11, no. 6, 7906922, 2019; doi:10.1109/jphot.2019.29444080.
- [173] Z. Lin, H. Qin and S. C. Chan, "A new probabilistic representation of color image pixels and its applications," *IEEE Trans. Image Process.*, Vol. 28, no. 4, pp. 2037–2050, 2019.
- [174] C. Chen, C. Zhou, P. Liu and D. Zhang, "Iterative reweighted Tikhonov-regularized multihypothesis prediction scheme for distributed compressive video sensing," *IEEE Trans. Circ. Syst. Video Techn.*, Vol. 30, no. 1, pp. 1–10, 2020.
- [175] A. Jain, S. Nagar, P.K. Singh and J. Dhar, "EMUCF: enhanced multistage user-based collaborative filtering through non-linear similarity for recommendation systems," *Expert Syst. Appl.*, Vol. 161, 113724, 2020; doi:10.1016/j.eswa.2020.113724.
- [176] R. Pascuzzo, N.P. Oxtoby, A.L. Young, J. Blevins, G. Castelli, S. Garbarino, M.L. Cohen, L.B. Schonberger, P. Gambetti, B.S. Appleby, D.C. Alexander and A. Bizzi, "Prion propagation estimated from brain diffusion MRI is subtype dependent in sporadic Creutzfeldt-Jakob disease," *Acta Neuropathologica*, Vol. 140, pp. 169–181, 2020.
- [177] W. Sun, J. Wang and F. Jin, "An automatic coordinate unification method of multitemporal point clouds based on virtual reference datum detection," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, Vol. 13, pp. 3942–3950, 2020.
- [178] S. Valpione, E. Galvani, J. Tweedy, P.A. Mundry, A. Banyard, P. Middlehurst, J. Barry, S. Mills, Z. Salih, J. Weightman, A. Gupta, G. Gremel, F. Baenke, N. Dhomen, P.C. Lorigan and R. Marais, "Immune-awakening revealed by peripheral T cell dynamics after one cycle of immunotherapy," *Nature Cancer*, Vol. 1, No. 2, pp. 210–221, 2021.
- [179] X. Wang, J. Wang, C. Miao and K. Zeng, "Forewarning method of downburst based on feature recognition and extrapolation," *Natural Hazards*, Vol. 103, pp. 903–921, 2020.
- [180] Y. Xu, Y. Xue, G. Hua and J. Cheng, "An adaptive distributed compressed video sensing algorithm based on normalized Bhattacharyya coefficient for coal mine monitoring video," *IEEE Access*, Vol. 8, pp. 158369–158379, 2020.
- [181] S. Zhao, G. Ding, Y. Gao, X. Zhao, Y. Tang, J. Han, H. Yao and Q. Huang, "Discrete probability distribution prediction of image emotions with shared sparse learning," *IEEE Trans. Affect. Comput.*, Vol. 11, no. 4, pp. 574–587, 2020.
- [182] D. Chen, J. Spencer, J.-M. Mirebeau, K. Chen, M. Shu and L.D. Cohen, "A generalized asymmetric dual-front model for active contours and image segmentation," *IEEE Trans. Image Process.*, Vol. 30, pp. 5056–5071, 2021.

- [183] E.C.L. De Oliveira, K. Santana, L. Josino, A.H. Lima e Lima and C. de Souza de Sales Junior, "Predicting cell-penetrating peptides using machine learning algorithms and navigating in their chemical space," *Sci. Rep.*, Vol. 11, 7628, 2021; doi:10.1038/s41598-021-87134-w.
- [184] A. Eshaghi, A.L. Young, P.A. Wijeratne, F. Prados, D.L. Arnold, S. Narayanan, C.R.G. Guttman, F. Barkhof, D.C. Alexander, A.J. Thompson, D. Chard and O. Ciccarelli, "Identifying multiple sclerosis subtypes using unsupervised machine learning and MRI data," *Nature Commun.*, Vol. 12, No. 2078, 2021; doi:10.1038/s41467-021-22265-2.
- [185] C. Feng, H. Zhao, Y. Li, Z. Cheng and J. Wena, "Improved detection of focal cortical dysplasia in normal-appearing FLAIR images using a Bayesian classifier," *Med. Phys.*, Vol. 48, no. 2, pp. 912–925, 2021.
- [186] W. Hanli, C. Hao, W. Lei, H. Jingguang and L. Zhenxing, "A novel pilot protection scheme for transmission lines based on current distribution histograms and their Bhattacharyya coefficient," *Electr. Pow. Syst. Res.*, Vol. 194, 107056, 2021; doi:10.1016/j.epsr.2021.107056.
- [187] R. Jiang, Q. Wang, S. Shi, X. Mou and S. Chen, "Flow-assisted visual tracking using event camera," *CAAI Trans. Intell. Technol.*, Vol. 6, pp. 192–202, 2021.
- [188] A. Lysiak and M. Szmajda, "Empirical comparison of the feature evaluation methods based on statistical measures," *IEEE Access*, Vol. 9, pp. 27868–27883, 2021.
- [189] T. Joel and R. Sivakumar, "Nonsampled contourlet transform with cross-guided bilateral filter for despeckling of medical ultrasound images," *Int. J. Imaging Syst. Technol.*, Vol. 31, pp. 763–777, 2021.
- [190] D. Reising, J. Cancellieri, T.D. Loveless, F. Kandah and A. Skjellum, "Radio identity verification-based IoT security using RF-DNA fingerprints and SVM," *IEEE Internet Things J.*, Vol. 8, no. 10, pp. 8356–8371, 2021.
- [191] T. Skrbic, A. Maritan, A. Giacometti and J.R. Banavar, "Local sequence-structure relationships in proteins," *Protein Sci.*, Vol. 30, pp. 818–829, 2021.
- [192] S. Tsiapoki, O. Bahrami, M.W. Häckell, J.P. Lynch and R. Rolfes, "Combination of damage feature decisions with adaptive boosting for improving the detection performance of a structural health monitoring framework: validation on an operating wind turbine," *Struct. Health. Monit.*, Vol. 20, no. 2, pp. 637–660, 2021.
- [193] P. van Molle, T. Verbelen, B. Vankeirsbilck, J. De Vylder, B. Diricx, T. Kimpe, P. Simoens and B. Dhoedt, "Leveraging the Bhattacharyya coefficient for uncertainty quantification in deep neural networks," *Neural Comput. Appl.*, Vol. 33, no. 16, pp. 10259–10275, 2021; doi:10.1007/s00521-021-05789-y.
- [194] Z. Yang, W. Yan, F. Li, Z. Yu and L. Guo, "Evaluating onset times of acoustic emission signals using histogram distances," *IEEE Trans. Ind. Electron.*, Vol. 68, no. 6, pp. 5237–5247, 2021.
- [195] Y. Zhou and Y. Yu, "Human visual search follows a suboptimal Bayesian strategy revealed by a spatiotemporal computational model and experiment," *Commun. Biol.*, Vol. 4, No. 34, 2021; doi:10.1038/s42003-020-01485-0.
- [196] T. Kailath, "The divergence and Bhattacharyya distance measures in signal selection," *IEEE Trans. Commun. Techn.*, Vol. 15, no. 1, pp. 52–60, 1967.
- [197] D. Zhao, S. Liu, C. Rong, A. Zhong, and S. Liu, "Toward understanding the isomeric stability of fullerenes with density functional theory and the information-theoretic approach," *ACS Omega*, Vol. 3, pp. 17986–17990, 2018.
- [198] C. Bunte and A. Lapidoto, "Encoding tasks and Renyi entropy," *IEEE Trans. Inf. Theory*, Vol. 60, no. 9, pp. 5065–5076, 2014.
- [199] I. Sason, "On the Renyi divergence, joint range of relative entropies, and a channel coding theorem," *IEEE Trans. Inf. Theory*, Vol. 62, no. 1, pp. 23–34, 2016.
- [200] M. A. Kumar, A. Sunny, A. Thakre, A. Kumar and G.D. Manohar, "A unified framework for problems on guessing, source coding and tasks partitioning," in: *Proc. 2022 IEEE Int. Symp. Inf. Theory (ISIT 2022), Espoo, Finland*, pp. 3339–3344.
- [201] T. Prest, "Sharper bounds in lattice-based cryptography using the Renyi divergence," in: T. Takagi and T. Peyrin (eds.), *ASIACRYPT 2017*, Part I, Lecture Notes in Computer Science, vol. 10624, pp. 347–374. Cham, Switzerland: Springer Nature, 2017.
- [202] S. Bai, T. Lepoint, A. Roux-Langlois, A. Sakzad, D. Stehle and R. Steinfeld, "Improved security proofs in lattice-based cryptography: using the Renyi divergence rather than the statistical distance," *J. Cryptol.*, Vol. 31, pp. 610–640, 2018.
- [203] Y. He, C. Song, C. Jing and X. Lei, "Robot active olfaction search in turbulent flow and infotaxis search based on Renyi divergence," in: *Proc. OCEANS 2017 - Aberdeen*, 2017, 9 pages. doi: 10.1109/OCEANSE.2017.8084882.
- [204] A. Momeni, K. Rouhi, H. Rajabalipanah and A. Abdolali, "An information theory-inspired strategy for design of reprogrammable encrypted graphene-based coding metasurfaces at terahertz frequencies," *Sci. Rep.*, Vol. 8, 6200, 2018; doi:10.1038/s41598-018-24553-2.
- [205] A.D. Staszowska, P. Fox-Roberts, L.M. Hirvonen, C.J. Peddie, L.M. Collinson, G.E. Jones and S. Cox, "The Renyi divergence enables accurate and precise cluster analysis for localization microscopy," *Bioinformatics*, Vol. 34, no. 23, pp. 4102–4111, 2018.
- [206] L. Yu and V.Y.F. Tan, "Wyner's common information under Renyi divergence measures," *IEEE Trans. Inf. Theory*, Vol. 64, no. 5, pp. 3616–3632, 2018; Comments and Corrections: Vol. 66, no. 4, pp. 2599–2508, 2020
- [207] G. Zhang, F. Lian, C. Han, H. Chen and N. Fu, "Two novel sensor control schemes for multitarget tracking via delta generalised labelled multi-Bernoulli filtering," *IET Signal Process.*, Vol. 12, no. 9, pp. 1131–1139, 2018.
- [208] L. Yu and V.Y.F. Tan, "Simulation of random variables under Renyi divergence measures of all orders," *IEEE Trans. Inf. Theory*, Vol. 65, no. 6, pp. 3349–3383, 2019.
- [209] J. Blanchet, F. He and K. Murthy, "On distributionally robust extreme value analysis," *Extremes*, Vol. 23, pp. 317–347, 2020.
- [210] H. Cai, Y. Yang, S. Gehly, C. He, and M. Jah, "Sensor tasking for search and catalog maintenance of geosynchronous space objects," *Acta Astronaut.*, Vol. 175, pp. 234–248, 2020.
- [211] R. Gholami and G.A. Hodtani, "A more general information theoretic study of wireless location verification system," *IEEE Trans. Veh. Technol.*, Vol. 69, no. 9, pp. 9938–9950, 2020.
- [212] M.T. Seweryn, M. Pietrzak and Q. Mab, "Application of information theoretical approaches to assess diversity and similarity in single-cell transcriptomics," *Comput. Struct. Biotechn. J.*, Vol. 18, pp. 1830–1837, 2020.
- [213] L. Zhou, "Multiple private key generation for continuous memoryless sources with a helper," *IEEE Trans. Foren. Sec.*, Vol. 15, pp. 2629–2640, 2020.
- [214] K. Makkawi, N. Ait-Tmazirte M. El Badaoui El Najjar and N. Moubayed, "Adaptive diagnosis for fault tolerant data fusion based on α -Renyi divergence strategy for vehicle localization," *Entropy*, Vol. 23, No. 463, 2021; doi:10.3390/e23040463
- [215] Y. Mao, W. Hong, H. Wang, Q. Li, and S. Zhong, "Privacy-preserving computation offloading for parallel deep neural networks training," *IEEE Trans. Paralle. Distr. Syst.*, Vol. 32, no. 7, pp. 1777–1788, 2021.
- [216] A. Tarighati and J. Jalden, "Optimality of rate balancing in wireless sensor networks," *IEEE Trans. Signal Process.*, Vol. 64, no. 14, pp. 3735–3749, 2016.
- [217] S. Bi, M. Broggi and M. Beer, "The role of the Bhattacharyya distance in stochastic model updating," *Mech. Syst. Signal Proc.*, Vol. 129, pp. 437–452, 2019.
- [218] S. Bi, M. Broggi, P. Wei and M. Beer, "The Bhattacharyya distance: enriching the P-box in stochastic sensitivity analysis," *Mech. Syst. Signal Proc.*, Vol. 129, pp. 265–281, 2019.
- [219] Y. Fu and Z. He, "Bhattacharyya distance criterion based multibit quantizer designs for cooperative spectrum sensing in cognitive radio networks," *Wirel. Netw.*, Vol. 25, pp. 2665–2674, 2019.
- [220] A. Cohen, D. Malak, V.B. Bracha and M. Medard, "Adaptive Causal Network Coding With Feedback," *IEEE Trans. Commun.*, Vol. 68, no. 7, pp. 4325–4341, 2020.
- [221] C. Xu, K. Wang, P. Li, R. Xia, S. Guo, and M. Guo, "Renewable energy-aware big data analytics in geo-distributed data centers with reinforcement learning," *IEEE Trans. Netw. Sci. Eng.*, Vol. 7, no. 1, pp. 205–215, 2020.
- [222] M. Xu, V. Jog and P.-L. Loh, "Optimal rates for community estimation in the weighted stochastic block model," *Ann. Stat.*, Vol. 48, no. 1, pp. 183–204, 2020.

- [223] M. Arrigoni and G.K.H. Madsen, "Evolutionary computing and machine learning for discovering of low-energy defect configurations," *npj Comput. Mater.*, Vol. 7, No. 71, 2021; doi:10.1038/s41524-021-00537-1
- [224] S. Fan, Y. Sun and P. Shui, "Region-merging method with texture pattern attention for SAR image segmentation," *IEEE Trans. Geosci. Remote Sens. Lett.*, Vol. 18, no. 1, pp. 112–116, 2021.
- [225] M.A. Mahfouz, A. Shoukry and M.A. Ismail, "EKNN: ensemble classifier incorporating connectivity and density into kNN with application to cancer diagnosis," *Artif. Intell. Med.*, Vol. 111, 101985, 2021; doi:10.1016/j.artmed.2020.101985.
- [226] K.T. Matchev and P. Shyamsundar, "ThickBrick: optimal event selection and categorization in high energy physics. Part I. Signal discovery," *J. High Energy Phys.*, Vol. 2021, No.3, 291, 2021; doi:10.1007/JHEP03(2021)291.
- [227] R. Wang, M. Dang, K. Harada, G. Han, F. Wang, M.P. Pizzi, M. Zhao, G. Tatlonghari, S. Zhang, D. Hao, Y. Lu, S. Zhao, B.D. Badgwell, M. B. Murphy, N. Shanbhag, J.S. Estrella, S. Roy-Chowdhuri, A.A.F. Abdelhakeem, Y. Wang, G. Peng, S. Hanash, G.A. Calin, X. Song, Y. Chu, J. Zhang, M. Li, K. Chen, A.J. Lazar, A. Futreal, S. Song, J.A. Ajani and L. Wang, "Single-cell dissection of intratumoral heterogeneity and lineage diversity in metastatic gastric adenocarcinoma," *Nature Med.*, Vol. 27, 141–151, 2021.
- [228] A.J. Webster, K. Gaitskell, I. Turnbull, B.J. Cairns and R. Clarke, "Characterisation, identification, clustering, and classification of disease," *Sci. Rep.*, Vol. 11, 5405, 2021; doi:10.1038/s41598-021-84860-z.
- [229] T. Xiahou, Z. Zeng and Y. Liu, "Remaining useful life prediction by fusing expert knowledge and condition monitoring information," *IEEE Trans. Ind. Inform.*, Vol. 17, no. 4, pp. 2653–2663, 2021.
- [230] W.K. Wootters, "Statistical distance and Hilbert space," *Phys. Rev. D*, Vol. 23, no. 2, pp. 357–362, 1981.
- [231] N. Ay, J. Jost, H.V. Le, and L. Schwachhöfer, *Information Geometry*. Cham, Switzerland: Springer International, 2017.
- [232] C.R. Rao, "Cluster analysis applied to a study of race mixture in human populations," in: J. van Ryzin (ed.), *Classification and clustering*, pp. 175–197. New York, USA: Academic Press, 1977.
- [233] F. Juhasz, "The Hellinger Distance as Used for the Representation of Serological ABO Distances Among Earlier Human Populations," *Hum. Genet.*, Vol. 63, pp. 228–231, 1983.
- [234] J.A. Martin-Fernandez, C. Barcelo-Vidal, and V. Pawlowsky-Glahn, "Measures of difference for compositional data and hierarchical clustering methods," in: A. Buccianti, G. Nardi and R. Potenza (eds.), *Proceedings in IAMG'98, The Fourth Annual Conference of the International Association for Mathematical Geology*, pp. 526–531. Naples, Italy: De Frede, 1998.
- [235] M. Greenacre, "Weighted metric multidimensional scaling," in: M. Vichi, P. Monari, S. Mignani, and A. Montanari (eds.), *New Developments in Classification and Data Analysis*, pp. 141–149. Berlin, Germany: Springer, 2005.
- [236] S. Jolad, A. Roman, M.C. Shastri, M. Gadgil, and A. Basu, "A new family of bounded divergence measures and application to signal detection," in: M. De Marsico, G. Sanniti di Baja, and A. Fred (eds.), *Proceedings of the 5th International Conference on Pattern Recognition Applications and Methods (ICPRAM 2016)*, pp. 72–83. Setubal, Portugal: Scitepress, 2016.
- [237] A. Ghosh and A. Basu, "A scale-invariant generalization of the Renyi entropy, associated divergences and their optimizations under Tsallis' nonextensive framework," *IEEE Trans. Inf. Theory*, Vol. 67, no. 4, pp. 2141–2161, 2021.
- [238] E. Lutwak, D. Yang, and G. Zhang, "Cramer-Rao and moment-entropy inequalities for Renyi entropy and generalized Fisher information," *IEEE Trans. Inf. Theory*, Vol. 51, no. 2, pp. 473–478, 2005.
- [239] S. Yaghi, Y. Altug and S. Verdu, "Minimax Renyi redundancy," *IEEE Trans. Inf. Theory*, Vol. 64, no. 5, pp. 3715–3733, 2018.
- [240] B.D. Sharma and D.P. Mittal, "New nonadditive measures of entropy for discrete probability distributions," *J. Combin. Inform. System Sci.*, Vol. 2, no. 4, pp. 122–132, 1977.
- [241] D. Morales, L. Pardo, M. Salicru and M.L. Menendez, "Asymptotic properties of divergence statistics in a stratified random sampling and its applications to test statistical hypotheses," *J. Statist. Plan. Inf.*, Vol. 38, pp. 201–222, 1994.
- [242] O. Onicescu, "Energie informationnelle," *C. R. Acad. Sci. Paris Ser. A.*, Vol. 263, no. 3, pp. 841–842, 1966.
- [243] L. Pardo and I.J. Taneja, "Information energy and its applications," *Adv. Electron. El. Phys.*, Vol. 80, pp. 165–241, 1991.
- [244] A. Theodorescu, "Energie informationnelle et notions apparentees," *Trabajos de Estadist. Investigacion Oper.*, Vol. XXVIII, no. 2-3, pp. 183–206, 1977.
- [245] L. Pardo, "Order- α weighted information energy," *Inf. Science*, Vol. 40, pp. 155–164, 1986.
- [246] S.-B. Liu, C.-Y. Rong, Z.-M. Wu and T. Lu, "Renyi entropy, Tsallis entropy and Onicescu information energy in density functional reactivity theory," *Acta Phys.-Chim. Sin.*, Vol. 31, no. 11, pp. 2057–2063, 2015.
- [247] S. Lopez-Rosa, A.L. Martin, J. Antolin and J.C. Angulo, "Electron-pair entropic and complexity measures in atomic systems," *Int. J. Quantum Chem.*, Vol. 119, No. e25861, 2019; doi:10.1002/qua.25861.
- [248] C. Rong, B. Wang, D. Zhao and S. Liu, "Information-theoretic approach in density functional theory and its recent applications to chemical problems," *WIREs Comput Mol Sci.*, Vol. 10, No. e1461, 2020; doi:10.1002/wcms.1461.
- [249] S.W. Golomb, "The information generating function of a probability distribution," *IEEE Trans. Inf. Theory*, Vol. 12, no. 1, pp. 75–77, 1966.
- [250] J.C. Principe, *Information Theoretic Learning*. New York, NY, USA: Springer, 2010.
- [251] A.-M. Acu, G. Bascanbaz-Tunca and I. Rasa, "Information potential for some probability density functions," *Appl. Math. Comput.*, Vol. 389, 125578, 2021; doi:10.1016/j.amc.2020.125578.
- [252] D. Morales, L. Pardo and I. Vajda, "Uncertainty of discrete stochastic systems: general theory and statistical inference," *IEEE Trans. Syst. Man Cyb. - Part A: Syst. Humans*, Vol. 26, no. 6, pp. 681–697, 1996.
- [253] L. Hannah and J.A. Kay, *Concentration in Modern Industry*. London, UK: The Macmillan Press, 1977.
- [254] A. Nunes, T. Trappenberg and M. Alda, "The definition and measurement of heterogeneity," *Transl. Psychiat.*, Vol. 10, No. 299, 2020; doi:10.1038/s41398-020-00986-0.
- [255] L.L. Campbell, "Exponential entropy as a measure of extent of a distribution," *Z. Wahrscheinlichkeit verw. Geb./Probab. Theory Rel. Fields*, Vol. 5, pp. 217–225, 1966.
- [256] V. Greiff, P. Bhat, S.C. Cook, U. Menzel, W. Kang and S.T. Reddy, "A bioinformatic framework for immune repertoire diversity profiling enables detection of immunological status," *Genome Med.*, Vol. 7, 49, 2015; doi:10.1186/s13073-015-0169-8.
- [257] Z. Ma and L. Li, "Measuring metagenome diversity and similarity with Hill numbers," *Mol. Ecol. Resour.*, Vol. 18, pp. 1339–1355, 2018.
- [258] W. Jasinska, M. Manhart, J. Lerner, L. Gauthier, A.W.R. Serohijos and S. Bershtein, "Chromosomal barcoding of E. coli populations reveals lineage diversity dynamics at high resolution," *Nature Ecol. & Evol.*, Vol. 4, pp. 437–452, 2020.
- [259] Z. Ma, L. Li and Y.-P. Zhang, "Defining individual-level genetic diversity and similarity profiles," *Sci. Rep.*, Vol. 10, 5805, 2020; doi:10.1038/s41598-020-62362-8.
- [260] N. Lassance and F. Vrins, "Minimum Renyi entropy portfolios," *Ann. Oper. Res.*, Vol. 299, pp. 23–46, 2021.
- [261] R. Ahlswede, "Identification entropy," in: R. Ahlswede et al. (eds.), *General Theory of Information Transfer and Combinatorics*, Lecture Notes in Computer Science, vol. 4123, pp. 595–613. Berlin, Germany: Springer, 2006.
- [262] R. Ahlswede and N. Cai, "An interpretation of identification entropy," *IEEE Trans. Inf. Theory*, Vol. 52, no. 9, pp. 4198–4207, 2006.
- [263] A.M. Peter and A. Rangarajan, "Information geometry for landmark shape analysis: unifying shape representation and deformation," *IEEE Trans. IEEE Trans. Pattern Anal. Mach. Intell.*, Vol. 31, no. 2, pp. 337–350, 1993.
- [264] G.C. Yalcin and C. Beck, "Generalized statistical mechanics of cosmic rays: application to positron-electron spectral indices," *Sci. Rep.*, Vol. 8, 1764, 2018; doi:10.1038/s41598-018-20036-6.
- [265] T. Wen and W. Jiang, "Measuring the complexity of complex networks by Tsallis entropy," *Physica A*, Vol. 526, 121054, 2019; doi:10.1016/j.physa.2019.121054.

- [266] A. K. Bhandari, "A novel beta differential evolution algorithm-based fast multilevel thresholding for color image segmentation," *Neural Comput. Appl.*, Vol. 32, pp. 4583–4613, 2020.
- [267] T.T. Erguzel, C. Uyulan, B. Unsulver, A. Evrensel, M. Cebi, C.O. Noyan, B. Metin, G. Eryilmaz, G.H. Sayar and N. Tarhan, "Entropy: a promising EEG biomarker dichotomizing subjects with opioid use disorder and healthy controls," *Clin. EEG Neurosci.*, Vol. 51, no. 6, pp. 373–381, 2020.
- [268] M.-S. Kang and K.-T. Kim, "Automatic SAR image registration via Tsallis entropy and iterative search process," *IEEE Sensors J.*, Vol. 20, no. 14, pp. 7711–7720, 2020.
- [269] A. Namdari and Z.S. Li, "An entropy-based approach for modeling Lithium-Ion battery capacity fade," in: *Proc. 2020 Annual Reliability and Maintainability Symposium (RAMS)*, 2020, 7 pages.
- [270] G. Zhang, X. Su and V.P. Singh, "Modelling groundwater-dependent vegetation index using entropy theory," *Ecol. Model.*, Vol. 416, 108916, 2020; doi:10.1016/j.ecolmodel.2019.108916.
- [271] M. Kumbhakar, R.K. Ray, S.K. Chakraborty, K. Ghoshal and V.P. Singh, "Mathematical modelling of streamwise velocity profile in open channels using Tsallis entropy," *Commun. Nonlinear Sc. Numer. Simulat.*, Vol. 94, 105581, 2021; doi:10.1016/j.cnsns.105581.
- [272] Z. Ramezani and A. Pourdarvish, "Transfer learning using Tsallis entropy: an application to Gravity Spy," *Physica A*, Vol. 561, 125273, 2020; doi:10.1016/j.physa.2020.125273.
- [273] T. De Wet, F. Österreicher and M. Thaler, "Tsallis' entropies – axiomatics, associated f -divergences and Fisher's information," *South Afr. Stat. J.*, Vol. 54, no. 2, pp. 163–175, 2020.
- [274] C. Brukner and A. Zeilinger, "Operationally invariant information in quantum measurements," *Phys. Rev. Lett.*, Vol. 83, no. 17, pp. 3354–3357, 1999.
- [275] P. Nath, "On a coding theorem connected with Renyi's entropy," *Information and Control*, Vol. 29, pp. 234–242, 1975.
- [276] E. Arikan, "An inequality on guessing and its application to sequential decoding," *IEEE Trans. Inf. Theory*, Vol. 42, no. 1, pp. 99–105, 1996.
- [277] I. Sason and S. Verdú, "Improving bounds on lossless source coding and guessing moments via Renyi measures," *IEEE Trans. Inf. Theory*, Vol. 64, no. 6, pp. 4323–4346, 2018.
- [278] C.H. Bennett, G. Brassard, C. Crépeau and U.M. Maurer, "Generalized privacy amplification," *IEEE Trans. Inf. Theory*, Vol. 41, no. 6, pp. 1915–1923, 1995.
- [279] M.M. Mayoral, "Renyi's entropy as an index of diversity in sample-stage cluster sampling," *J. Information Sciences*, Vol. 105, pp. 101–114, 1998.
- [280] S. Aviyente, L.A.W. Brakel, R.K. Kushwaha, M. Snodgrass, H. Shevrin, and W.J. Williams, "Characterization of event related potentials using information theoretic distance measures," *IEEE Trans. Bio-Med. Eng.*, Vol. 51, no. 5, pp. 737–743, 2004.
- [281] B. Tao, L. Zhu, H. Ding and Y. Xiong, "An alternative time-domain index for condition monitoring of rolling element bearings – a comparison study," *Reliab. Eng. Syst. Saf.*, Vol. 92, pp. 660–670, 2007.
- [282] J. Jiao, J. Yue, D. Pei and Z. Hu, "Application of feature fusion using coaxial vibration signal for diagnosis of rolling element bearings," *Shock Vibrat.*, Vol. 2020, 8831723, 2020; doi:10.1155/2020/8831723.
- [283] D.-T. Pham, F. Vrins, and M. Verleysen, "On the risk of using Renyi's entropy for blind source separation," *IEEE Trans. Signal Process.*, Vol. 56, no. 10, pp. 4611–4620, 2008.
- [284] I. Sason, "Tight bounds on the Renyi entropy via majorization with applications to guessing and compression," *Entropy*, Vol. 20, No. 896, 2018; doi:10.3390/e20120896.
- [285] P. Carravilla, J. Chojnacki, E. Ruja, S. Insausti, E. Largo, D. Waithe, B. Apellaniz, T. Sicard, J.-P. Julien, C. Eggeling and J.L. Nieva, "Molecular recognition of the native HIV-1 MPER revealed by STED microscopy of single virions," *Nat. Commun.*, Vol. 10, no. 78, 2019; doi:10.1038/s41467-018-07962-9.
- [286] K. Joshi, M.R. de Massy, M. Ismail, J.L. Reading, I. Uddin, A. Woolston, E. Hatipoglu, T. Oakes, R. Rosenthal, T. Peacock, T. Ronel, M. Noursadeghi, V. Turati, A.J.S. Furness, A. Georgiou, Y.N.S. Wong, A.B. Aissa, M.W. Sunderland, M. Jamal-Hanjani, S. Veeriah, N.J. Birkbak, G.A. Wilson, C.T. Hiley, E. Ghorani, J.A. Guerra-Assuncao, J. Herrero, T. Enver, S.R. Hadrup, A. Hackshaw, K.S. Peggs, N. McGranahan, C. Swanton, TRACERx consortium, S.A. Quezada and B. Chain, "Spatial heterogeneity of the T cell receptor repertoire reflects the mutational landscape in lung cancer," *Nature Med.*, Vol. 25, 1549–1559, 2019.
- [287] Z. German-Sallo, "Entropy indices based fault detection," *Proced. Manufact.*, Vol. 46, pp. 549–554, 2020.
- [288] K. Schober, F. Voit, S. Grassmann, T.R. Müller, J. Eggert, S. Jarosch, Bianca Weißbrich, P. Hoffmann, L. Borkner, E. Nio, L. Fanchi, C.R. Clouser, A. Radhakrishnan, L. Mihatsch, P. Lückemeier, J. Leube, G. Dössinger, L. Klein, M. Neuenhahn, J.D. Oduro, L. Cicin-Sain, V.R. Buchholz and D.H. Busch, "Reverse TCR repertoire evolution toward dominant low-affinity clones during chronic CMV infection," *Nature Immunol.*, Vol. 21, pp. 434–441, 2020.
- [289] J.P. Amezcua-Sanchez, "Entropy algorithms for detecting incipient damage in high-rise buildings subjected to dynamic vibrations," *J. Vib. Control*, Vol. 27, no. 3–4, pp. 426–436, 2021.
- [290] P. Barennes, V. Quiniou, M. Shugay, E.S. Egorov, A.N. Davydov, D.M. Chudakov, I. Uddin, M. Ismail, T. Oakes, B. Chain, A. Eugster, K. Kashofer, P.P. Rainer, S. Darko, A. Ransier, D.C. Douek, D. Klatzmann and E. Mariotti-Ferrandiz, "Benchmarking of T cell receptor repertoire profiling methods reveals large systematic biases," *Nature Biotechnol.*, Vol. 39, pp. 236–245, 2021.
- [291] M. Kumar, A.K. Bhandari, N. Singh and A. Ghosh, "A new multilevel histogram thresholding approach using variational mode decomposition," *Multimed. Tools and Appl.*, Vol. 80, pp. 11331–11363, 2021.
- [292] B. Pandey, "Renyi entropy as a measure of cosmic homogeneity," *J. Cosmol. Astropart. Phys.*, Vol. 2021, No. 2, 023, 2021; doi:10.1088/1475-7516/2021/02/023.
- [293] T.D. Popescu and D. Aiordachioaie, "New procedure for change detection operating on Renyi entropy with application in seismic signals Processing," *Circuits Syst. Signal Process.*, Vol. 36, pp. 3778–3798, 2017.
- [294] D. Chen, D.-D. Shi, M. Qin, S.-M. Xu and G.-J. Pan, "Complex network comparison based on communicability sequence entropy," *Phys. Rev. E*, Vol. 98, No. 012319, 2018; doi:10.1103/PhysRevE.98.012319.
- [295] D.-D. Shi, D. Chen and G.-J. Pan, "Characterization of network complexity by communicability sequence entropy and associated Jensen-Shannon divergence," *Phys. Rev. E*, Vol. 101, No. 042305, 2020; doi:10.1103/PhysRevE.101.042305.
- [296] J.P. Bagrow and E.M. Boltt, "An information-theoretic, all-scales approach to comparing networks," *Appl. Network Sci.*, Vol. 4, 45, 2018; doi:10.1007/s41109-019-0156-x
- [297] B. Jena, M.K. Naik, R. Panda and A. Abraham, "Maximum 3D Tsallis entropy based multilevel thresholding of brain MR image using attacking Manta Ray foraging optimization," *Eng. Appl. Artif. Intell.*, Vol. 103, 104293, 2021; doi:10.1016/j.engappai.2021.104293.
- [298] P. Boskoski and D. Juricic, "Fault detection of mechanical drives under variable operating conditions based on wavelet packet Renyi entropy signatures," *Mech. Syst. Signal Process.*, Vol. 31, pp. 369–381, 2012.
- [299] E.C. Pielou, "The measurement of diversity in different types of biological Collections," *J. Theoret. Biol.*, Vol. 13, pp. 131–143, 1966.
- [300] B. Gabriel, C. Medin, J. Alves, R. Nduati, R.K. Bosire, D. Wamalwa, C. Farquhar, G. John-Stewart and B.L. Lohman-Payne, "Analysis of the TCR repertoire in HIV-exposed but uninfected infants," *Sci. Rep.*, Vol. 9, 11954, 2019; doi:10.1038/s41598-019-48434-4.
- [301] R.J.M. Bashford-Rogers, L. Bergamaschi, E.F. McKinney, D.C. Pomball, F. Mescia, J.C. Lee, D.C. Thomas, S.M. Flint, P. Kellam, D.R.W. Jayne, P.A. Lyons and K.G.C. Smith, "Analysis of the B cell receptor repertoire in six immune-mediated diseases," *Nature*, Vol. 574, pp. 122–126, 2019.
- [302] A. Lyubushin, "Seismic noise wavelet-based entropy in Southern California," *J. Seismol.*, Vol. 25, pp. 25–39, 2021.
- [303] A. De Santis, A.G. Gaggia, and U. Vaccaro, "Bounds on entropy in a guessing game," *IEEE Trans. Inf. Theory*, Vol. 47, no. 1, pp. 468–473, 2001.
- [304] M. Johansson and M. Sternad, "Resource allocation under uncertainty using the maximum entropy principle," *IEEE Trans. Inf. Theory*, Vol. 51, no. 12, pp. 4103–4117, 2005.

- [305] S. Marano and M. Franceschetti, "Ray propagation in a random lattice: a maximum entropy, anomalous diffusion process," *IEEE Trans. Antenn. Propag.*, Vol. 53, no. 6, pp. 1888–1896, 2005.
- [306] L. Miao, H. Qi and H. Szu, "A maximum entropy approach to unsupervised mixed-pixel decomposition," *IEEE Trans. Image Process.*, Vol. 16, no. 4, pp. 1008–1021, 2007.
- [307] E.S.C. Rodrigues, F.A. Rodrigues, R.L.A. Rocha and P.L.P. Correa, "Adaptive approach for a maximum entropy algorithm in ecological niche modeling," *IEEE Lat. Amer. Trans.*, Vol. 9, no. 3, pp. 331–338, 2011.
- [308] D. Xiong, M. Zhang and H. Li, "A maximum-entropy segmentation model for statistical machine translation," *IEEE Trans. Audio Speech Lang. Process.*, Vol. 19, no. 8, pp. 2494–2505, 2011.
- [309] R.H. Chan, T.H. Chan, H.M. Yeung and R.W. Wang, "Composition vector method based on maximum entropy principle for sequence comparison," *IEEE/ACM Comput. Biol. Bioinform.*, Vol. 9, no. 1, pp. 79–87, 2012.
- [310] R.P. Mann and R. Garnett R, "The entropic basis of collective behaviour," *J. R. Soc. Interface*, Vol. 12, No. 20150037, 2015; <http://dx.doi.org/10.1098/rsif.2015.0037>.
- [311] A.K. Singh, H.P. Singh and Karmeshu, "Analysis of finite buffer queue: maximum entropy probability distribution with shifted fractional geometric and arithmetic means," *IEEE Commun. Lett.*, Vol. 19, no. 2, pp. 163–166, 2015.
- [312] A. Baddeley, "A statistical commentary on mineral prospectivity analysis," in: B.S.D. Sagar, Q. Cheng, F. Agterberg (eds.), *Handbook of Mathematical Geosciences*, pp. 25–65. Cham, Switzerland: Springer International, 2018.
- [313] G.A. Einicke, H.A. Sabti, D.V. Thiel and M. Fernandez, "Maximum-entropy-rate selection of features for classifying changes in knee and ankle dynamics during running," *IEEE J. Biomed. Health Inform.*, Vol. 22, no. 4, pp. 1097–1103, 2018.
- [314] H. Han, Y. Wang, Y. Zou, J. Liao and Y. Xu, "Three-dimensional substructure imaging of blood cells using maximum entropy tomography based on two non-orthogonal phase images," *Opt. Laser Technol.*, Vol. 136, 106799, 2021; doi:10.1016/j.optlastec.2020.106799.
- [315] R. Burkard, M. Dell'Amico and S. Martello, *Assignment Problems*, rev. repr. Philadelphia, PA, USA: SIAM, 2009.
- [316] D.L. Applegate, R.E. Bixby, V. Chvatal and W.J. Cook, *The Traveling Salesman Problem*. Princeton, NJ, USA: Princeton University Press, 2006.
- [317] G. Gutin and A.P. Punnen (eds.), *The Traveling Salesman Problem and Its Variations*. New York, NY, USA: Springer, 2007.
- [318] W.J. Cook, *In Pursuit of the Traveling Salesman*. Princeton, NJ, USA: Princeton University Press, 2012.
- [319] G. Dantzig, R. Fulkerson, and S. Johnson, "Solution of a large-scale traveling-salesman problem," *Oper. Res.*, Vol. 2, no. 4, pp. 393–410, 1954.
- [320] P. Bertrand, M. Broniatowski, and J.-F. Marcotorchino, "Logical indetermination coupling: a method to minimize drawing matches and its applications," *arXiv:2012.14674v1*, 28 pages, 2020.
- [321] P. Bertrand, M. Broniatowski, and J.-F. Marcotorchino, "Independence versus indetermination: basis of two canonical clustering criteria," *Adv. Data Anal. Classif.*, 2022; doi:10.1007/s11634-021-00484-1.
- [322] S. Yang, S. Tan and J.-X. Xu, "Consensus based approach for economic dispatch problem in a smart grid," *IEEE Trans. Power Syst.*, Vol. 28, no. 4, pp. 4416–4426, 2013.
- [323] V. Loia and A. Vaccaro, "Decentralized economic dispatch in smart grids by self-organizing dynamic agents," *IEEE Trans. Syst. Man Cybern.*, Vol. 44, no. 4, pp. 397–408, 2014.
- [324] A.J. Wood, B.F. Wollenberg and G.B. Sheble, *Power Generation, Operation, and Control*, 3rd ed. Hoboken, NJ, USA: John Wiley & Sons, 2014.
- [325] Y. Xu, W. Zhang, W. Liu and W. Yu, *Distributed Energy Management of Electrical Power Systems*. Hoboken, NJ, USA: John Wiley & Sons, 2021.
- [326] J.S. Sadowsky and J.A. Bucklew, "On Large Deviations Theory and Asymptotically Efficient Monte Carlo Estimation," *IEEE Trans. Inf. Theory*, Vol. 36, no. 3, pp. 579–588, 1990.
- [327] S. Bernstein, "Sur les fonctions absolument monotones," *Acta Math.*, Vol. 52, pp. 1–66, 1929.
- [328] R.L. Schilling, R. Song and Z. Vondracek, *Bernstein Functions*, 2nd ed. Berlin, Germany: de Gruyter, 2012.
- [329] D.V. Widder, "Necessary and sufficient conditions for the representation of a function by a doubly infinite Laplace integral," *Bull. Amer. Math. Soc.*, Vol. 40, no. 4, pp. 321–326, 1934.
- [330] Y.S. Chow and H. Teicher, *Probability Theory*, 3rd ed. New York, NY, USA: Springer, 1997.
- [331] D.V. Widder, *The Laplace Transform*. Princeton, NJ, USA: Princeton University Press, 1941.
- [332] N.I. Akhiezer, *The Classical Moment Problem and Some Related Questions in Analysis*. Edinburgh, UK: Oliver & Boyd, 1965.
- [333] D.S. Shucker, "Extensions and generalizations of a theorem of Widder and of the theory of symmetric local semigroups," *J. Funct. Anal.*, Vol. 58, pp. 291–309, 1984.
- [334] J. Jaksetic and J. Pecaric, "Exponential convexity method," *J. Convex Anal.*, Vol. 20, no. 1, pp. 181–197, 2013.
- [335] N.O. Kotelina and A.B. Pevny, "Exponential convexity and total positivity," *Siberian Electr. Math. Rep.*, Vol. 17, pp. 802–806, 2020; doi:10.33048/semi.2020.17.057.
- [336] D.W. Stroock, *Probability Theory: An Analytic View*, 2nd ed. New York, USA: Cambridge University Press, 2011.
- [337] D. McLeish, "Simulating random variables using moment-generating functions and the saddlepoint approximation," *J. Stat. Comput. Simul.*, Vol. 84, no. 2, pp. 324–334, 2014.
- [338] V.M. Zolotarev, *One-dimensional Stable Distributions*. Providence, USA: American Mathematical Society, 1986.
- [339] L. Devroye, "Random variate generation for exponentially and polynomially tilted stable distributions," *ACM Trans. Model. Comput. Simul.*, Vol. 19, no. 4, Article 18, 20 pages, 2009. doi:10.1145/1596519.1596523.
- [340] L. Devroye and L. James, "On simulation and properties of the stable law," *Stat. Methods Appl.*, Vol. 23, no. 3, 307–343, 2014. doi:10.1007/s10260-014-0260-0.
- [341] O.O. Aalen, "Modeling the heterogeneity in survival analysis by the compound Poisson distribution," *Ann. Appl. Probab.*, Vol. 2, no. 4, pp. 951–972, 1992.
- [342] M.C.K. Tweedie, "Functions of a statistical variate with given means, with special reference to Laplacian distributions," *Proc. Camb. Philos. Soc.*, Vol. 43, pp. 41–49, 1947.
- [343] C.N. Morris, "Natural exponential families with quadratic variance functions," *Ann. Stat.*, Vol. 10, no. 1, pp. 65–80, 1982.
- [344] G. Letac and M. Mora, "Natural real exponential families with cubic variance functions," *Ann. Stat.*, Vol. 18, no. 1, pp. 1–37, 1990.
- [345] J. Lin, "Divergence measures based on the Shannon entropy," *IEEE Trans. Inf. Theory*, Vol. 37, no. 1, pp. 145–151, 1991.
- [346] M.C. Pardo and I. Vajda, "About distances of discrete distributions satisfying the data processing theorem of information theory," *IEEE Trans. Inf. Theory*, Vol. 43, no. 4, pp. 1288–1293, 1997.
- [347] F. Topsøe, "Some inequalities for information divergence and related measures of discrimination," *IEEE Trans. Inf. Theory*, Vol. 46, no. 4, pp. 1602–1609, 2000.
- [348] D.M. Endres and J.E. Schindelin, "A new metric for probability distributions," *IEEE Trans. Inf. Theory*, Vol. 40, no. 7, pp. 1858–1860, 1993.
- [349] I. Vajda, "On metric divergences of probability measures," *Kybernetika*, Vol. 45, no. 6, pp. 885–900, 2009.
- [350] I. Sason, "Tight bounds for symmetric divergence measures and a new inequality relating f -divergences," in: *Proc. 2015 IEEE Information Theory Workshop (ITW)*, 2014, 5 pages.
- [351] R. Sibson, "Information radius," *Probab. Theory Rel. Fields*, Vol. 14, pp. 149–160, 1969.
- [352] D. Kvitsiani, S. Ranade, B. Hangya, H. Taniguchi, J.Z. Huang and A. Kepecs, "Distinct behavioural and network correlates of two interneuron types in prefrontal cortex," *Nature*, Vol. 498, pp. 363–366, 2013.
- [353] Q. Xu, Y. Liu, X. Li, Z. Yang, J. Wang, M. Sbert and R. Scopigno, "Browsing and exploration of video sequences: A new scheme for key frame extraction and 3D visualization using entropy based Jensen divergence," *Inf. Sci.*, Vol. 278, pp. 736–756, 2014.

- [354] G. Jenkinson, E. Pujadas, J. Goutsias and A.P. Feinberg, "Potential energy landscapes identify the information-theoretic nature of the epigenome," *Nature Genetics*, Vol. 49, no. 5, pp. 719–729, 2017.
- [355] F. Martin, J. Carballera, L. Moreno, S. Garrido, and P. Gonzalez, "Using the Jensen-Shannon, density power, and Itakura-Saito divergences to implement an evolutionary-based global localization filter for mobile robots," *IEEE Access*, Vol. 5, pp. 13922–13940, 2017.
- [356] S. Suo, Q. Zhu, A. Saadatpour, L. Fei, G. Guo and G.-C. Yuan, "Revealing the critical regulators of cell identity in the mouse cell atlas," *Cell Rep.*, Vol. 25, pp. 1436–1445, 2018.
- [357] J. Abante, Y. Fang, A.P. Feinberg and J. Goutsias, "Detection of haplotype-dependent allele-specific DNA methylation in WGBS data," *Nature Commun.*, Vol. 11, No. 5238, 2020; doi:10.1038/s41467-020-19077-1.
- [358] A. Afek, H. Shi, A. Rangadurai, H. Sahay, A. Senitzki, S. Khani, M. Fang, R. Salinas, Z. Mielko, M.A. Pufall, G.M.K. Poon, T.E. Haran, M.A. Schumacher, H.M. Al-Hashimi and R. Gordan, "DNA mismatches reveal conformational penalties in protein-DNA recognition," *Nature*, Vol. 587, pp. 291–296, 2020.
- [359] R. Alaiz-Rodriguez and A.C. Parnell, "An information theoretic approach to quantify the stability of feature selection and ranking algorithms," *Knowl.-Based Syst.*, Vol. 195, 105745, 2020; doi:10.106/j.knosys.2020.105745.
- [360] G. Biau, B. Cadre, M. Sangnier and U. Tanielian, "Some theoretical properties of GANs," *Ann. Stat.*, Vol. 48, no. 3, pp. 1539–1566, 2020.
- [361] A. Carre, G. Klausner, M. Edjlali, M. Lerousseau, J. Briand-Diop, R. Sun, S. Ammari, S. Reuze, E. Alvarez Andres, T. Estienne, S. Niyoteka, E. Battistella, M. Vakalopoulou, F. Dhermain, N. Paragios, E. Deutsch, C. Oppenheim, J. Pallud and C. Rober, "Standardization of brain MR images across machines and protocols: bridging the gap for MRI-based radiomics," *Sci. Rep.*, Vol. 10, 12340, 2020; doi:10.1038/s41598-020-69298-z.
- [362] A. Chakraborty, S. Easwaran and S. Sinha, "Uncovering hierarchical structure of international FOREX market by using similarity metric between fluctuation distributions of currencies," *Acta Phys. Pol. A*, Vol. 138, no. 1, pp. 105–115, 2020.
- [363] J. Chong, P. Liu, G. Zhou and J. Xia, "Using MicrobiomeAnalyst for comprehensive statistical, functional, and meta-analysis of microbiome data," *Nature Protocols*, Vol. 15, pp. 799–821, 2020.
- [364] L. Cui, J. Wu, D. Pi, P. Zhang and P. Kennedy, "Dual implicit mining-based latent friend recommendation," *IEEE Trans. Syst. Man Cyber.: Syst.*, Vol. 50, no. 5, pp. 1663–1678, 2020.
- [365] X. Guo and Y. Yuan, "Semi-supervised WCE image classification with adaptive aggregated attention," *Med. Image Anal.*, Vol. 64, 101733, 2020. doi:10.1016/j.media.2020.101733.
- [366] J. Jiang, M. Chen and J.A. Fan, "Deep neural networks for the evaluation and design of photonic devices," *Nat. Rev. Materials*, 2020; doi:10.1038/s41578-020-00260-1.
- [367] Ö. Kartal, M.W. Schmid and U. Grossniklaus, "Cell type-specific genome scans of DNA methylation divergence indicate an important role for transposable elements," *Genome Biol.*, Vol. 21, No. 172, 2020; doi:10.1186/s13059-020-02068-2.
- [368] T. Laszlovszky, D. Schlingloff, P. Hegedüs, T.F. Freund, A. Gulyas, A. Kepecs and B. Hangya, "Distinct synchronization, cortical coupling and behavioral function of two basal forebrain cholinergic neuron types," *Nature Neurosci.*, Vol. 23, pp. 992–1003, 2020.
- [369] K.A. Lawson, C.M. Sousa, X. Zhang, E. Kim, R. Akthar, J.J. Caumanns, Y. Yao, N. Mikolajewicz, C. Ross, K.R. Brown, A.A. Zid, Z.P. Fan, S. Hui, J.A. Krall, D.M. Simons, C.J. Slater, V. De Jesus⁴, L. Tang, R. Singh, J.E. Goldford, S. Martin, Q. Huang, E.A. Francis, A. Habsid, R. Climie, D. Tieu, J. Wei, R. Li, A.H.J. Tong, M. Aregger, K.S. Chan, H. Han, X. Wang, P. Mero, J.H. Brumell, A. Finelli, L. Ailles, G. Bader, G.A. Smolen, G.A. Kingsbury, T. Hart, C. Kung and J. Moffat, "Functional genomic landscape of cancer-intrinsic evasion of killing by T cells," *Nature*, Vol. 586, pp. 120–126, 2020.
- [370] C.H. Li, E.L. Coffey, A. Dall'Agnesse, N.M. Hannett, X. Tang, J.E. Henninger, J.M. Platt, O. Oksuz, A.V. Zamudio, L.K. Afeyan, J. Schuijers, X.S. Liu, S. Markoulaki, T. Lungjangwa, G. LeRoy, D.S. Svoboda, E. Wogram, T.I. Lee, R. Jaenisch and R.A. Young, "MeCP2 links heterochromatin condensates and neurodevelopmental disease," *Nature*, Vol. 586, pp. 440–444, 2020.
- [371] J.A. T. Machado, J.M. Rocha-Neves and J.P. Andrade, "Computational analysis of the SARS-CoV-2 and other viruses based on the Kolmogorov's complexity and Shannon's information theories," *Nonlinear. Dyn.*, Vol. 101, pp. 1731–1750, 2020.
- [372] S. Mohammadi, J. Davila-Velderrain and M. Kellis, "A multiresolution framework to characterize single-cell state landscapes," *Nature Commun.*, Vol. 11, 5399, 2020; doi:10.1038/s41467-020-18416-6.
- [373] A. Mohanty, Q. Li, M.A. Tadayon, S.P. Roberts, G.R. Bhatt, E. Shim, X. Ji, J. Cardenas, S.A. Miller, A. Kepecs and M. Lipson, "Reconfigurable nanophotonic silicon probes for sub-millisecond deep-brain optical stimulation," *Nature Biomed. Eng.*, Vol. 4, no. 2, pp. 223–231, 2020.
- [374] S. Perera, D. Kasthurirathna and M. Bliemer, "Topological rationality of supply chain networks," *Int. J. Prod. Res.*, Vol. 58, no. 10, pp. 3126–3149, 2020.
- [375] F. Pierri, C. Piccardi and S. Ceri, "Topology comparison of Twitter diffusion networks effectively reveals misleading information," *Sci. Rep.*, Vol. 10, 1372, 2020; doi:10.1038/s41598-020-58166-5.
- [376] R. Rabadan, Y. Mohamedi, U. Rubin, T. Chu, A.N. Alghalith, O. Elliott, L. Arnes, S. Cal, A.J. Obaya, A.J. Levine and P.G. Camara, "Identification of relevant genetic alterations in cancer using topological data analysis," *Nature Commun.*, Vol. 11, 3808, 2020; doi:10.1038/s41467-020-17659-7.
- [377] J.G. Reiter, W.-T. Hung, I.-H. Lee, S. Nagpal, P. Giunta, S. Degner, G. Liu, E.C.E. Wassenaar, W.R. Jeck, M.S. Taylor, A.A. Farahani, H.D. Marble, S. Knott, O. Kranenburg, J.K. Lennerz and K. Naxerova, "Lymph node metastases develop through a wider evolutionary bottleneck than distant metastases," *Nature Genetics*, Vol. 52, pp. 692–700, 2020.
- [378] B. Van de Sande, C. Flerin, K. Davie, M. De Waegeneer, G. Hulselmans, S. Aibar, R. Seurinck, W. Saelens, R. Cannoodt, Q. Rouchon, T. Verbeiren, D. De Maeyer, J. Reumers, Y. Saeys and S. Aerts, "A scalable SCENIC workflow for single-cell gene regulatory network analysis," *Nature Protocols*, Vol. 15, pp. 2247–2276, 2020.
- [379] M.A. Skinnider, C.W. Johnston, M. Gunabalasingam, N.J. Merwin, A.M. Kieliszek, R.J. MacLellan, H. Li, M.R.M. Ranieri, A.L.H. Webster, M.P.T. Cao, A. Pfeifle, N. Spencer, Q. H. To, D.P. Wallace, C.A. Dejong 3 and N.A. Magarvey, "Comprehensive prediction of secondary metabolite structure and biological activity from microbial genome sequences," *Nature Commun.*, Vol. 11, 6058, 2020; doi:10.1038/s41467-020-19986-1.
- [380] S. Tuo, H. Liu and H. Chen, "Multipopulation harmony search algorithm for the detection of high-order SNP interactions," *Bioinformatics*, Vol. 36, no. 16, pp. 4389–4398, 2020.
- [381] S. Uttam, A.M. Stern, C.J. Sevinsky, S. Furman, F. Pullara, D. Spagnolo, L. Nguyen, A. Gough, F. Ginty, D.L. Taylor and S.C. Chennubhotla, "Spatial domain analysis predicts risk of colorectal cancer recurrence and infers associated tumor microenvironment networks," *Nature Commun.*, Vol. 11, 3515, 2020; 10.1038/s41467-020-17083-x.
- [382] X. Zhang, C. Delpha, and D. Diallo, "Jensen-Shannon divergence for non-destructive incipient crack detection and estimation," *IEEE Access*, Vol. 8, pp. 116148–116162, 2020.
- [383] T. Zhi, Y. Liu, J. Wang and H. Zhang, "Resist interest flooding attacks via Entropy-SVM and Jensen-Shannon divergence in information-centric networking," *IEEE Syst. J.*, Vol. 14, no. 2, pp. 297–301, 2020.
- [384] P. Acera Mateos, R.F. Balboa, S. Easteal, E. Eyras and H.R. Patel, "PACIFIC: a lightweight deep-learning classifier of SARS-CoV-2 and co-infecting RNA viruses," *Sci. Rep.*, Vol. 11, 3209, 2021; doi:10.1038/s41598-021-82043-4.
- [385] Z. Avsec, M. Weilert, A. Shrikumar, S. Krueger, A. Alexandari, K. Dalal, R. Propf, C. McAnany, J. Gagneur, A. Kundaje and J. Zeitlinger, "Base-resolution models of transcription-factor binding reveal soft motif syntax," *Nature Genetics*, Vol. 53, pp. 354–366, 2021.
- [386] J. Chen, J.E. Markowitz, V. Lilascharoen, S. Taylor, P. Sheurpukdi, J.A. Keller, J.R. Jensen, B.K. Lim, S.R. Datta and L. Stowers, "Flexible scaling and persistence of social vocal communication," *Nature*, Vol. 593, pp. 108–113, 2021.
- [387] M.A. Koldobskiy, G. Jenkinson, J. Abante, V.A. Rodriguez DiBlasi, W. Zhou, E. Pujada, A. Idrizi, R. Tryggvadottir, C. Callahan, C.L. Bonifant, K.R. Rabin, P.A. Brown, H. Ji, J. Goutsias and A.P. Feinberg, "Converging genetic and epigenetic drivers of paediatric acute lymphoblastic leukaemia identified by an information-theoretic analysis," *Nature Biomed. Eng.*, Vol. 5, no. 4, pp. 360–376, 2021.

- [388] C.S. McGinnis, D.A. Siegel, G. Xie, G. Hartoularos, M. Stone, C.J. Ye, Z.J. Gartner, N.R. Roan and S.A. Lee, “No detectable alloreactive transcriptional responses under standard sample preparation conditions during donormultiplexed single-cell RNA sequencing of peripheral blood mononuclear cells,” *BMC Biology*, Vol. 19, 10, 2021; doi:10.1186/s12915-020-00941-x.
- [389] C. Mühlroth and M. Grottko, “Artificial intelligence in innovation: how to spot emerging trends and technologies,” *IEEE Trans. Eng. Manag.*, Vol. 69, no. 2, pp. 493–510, 2022.
- [390] M. Necci, D. Piovesan, CAID Predictors, DisProt Curators and S.C.E. Tosatto, “Critical assessment of protein intrinsic disorder prediction,” *Nature Meth.*, Vol. 18, pp. 472–481, 2021.
- [391] D. Okada, N. Nakamura, K. Setoh, T. Kawaguchi, K. Higasa, Y. Tabara, F. Matsuda and R. Yamada, “Genome-wide association study of individual differences of human lymphocyte profiles using large-scale cytometry data,” *J. Hum. Genet.*, Vol. 66, pp. 557–567, 2021.
- [392] Z. Zhang, J. Ding, J. Xu, J. Tang and F. Guo, “Multi-scale time-series kernel-based learning method for brain disease diagnosis,” *IEEE J. Biomed. Health Inform.*, Vol. 25, no. 1, pp. 209–217, 2021.
- [393] J.N. Kapur, “Four families of measures of entropy,” *Indian J. Pure Appl. Math.*, Vol. 17, pp. 429–449, 1986.
- [394] T. Yamano, “Some bounds for skewed α -Jensen-Shannon divergence,” *Results in Appl. Math.*, Vol. 3, No. 100064, 2019; doi:10.1016/j.rinam.2019.100064.
- [395] S.S. Dragomir, “Upper bounds for the Kullback-Leibler distance and applications,” *Bull. Math. Soc. Sc. Math. Roumanie*, Vol. 43(91), no. 1, pp. 25–37, 2000.
- [396] A.L. Rukhin, “Optimal estimator for the mixture parameter by the method of moments and information affinity,” in: *Trans. 12th Prague Conf. Information Theory, Statistical Decision Functions and Random Processes*, pp. 214–219. Prague, Czech Republic: Czech Acad. Sci., 1994.
- [397] Y. Marhuenda, D. Morales, J.A. Pardo and M.C. Pardo, “Choosing the best Rukhin goodness-of-fit statistics,” *Comp. Statist. & Data Anal.*, Vol. 49, pp. 643–662, 2005.
- [398] B.G. Lindsay, “Efficiency versus robustness: the case for minimum Hellinger distance and related methods,” *Ann. Statist.*, Vol. 22, no. 2, pp. 1081–1114, 1994.
- [399] A. Basu and B.G. Lindsay, “Minimum disparity estimation for continuous models: efficiency, distributions and robustness,” *Ann. Inst. Statist. Math.*, Vol. 46, no. 4, pp. 683–705, 1994.
- [400] L. Györfy and I. Vajda, “A class of modified Pearson and Neyman statistics,” *Statistics & Decisions*, Vol. 19, pp. 239–251, 2001.
- [401] A. Basu, H. Shioya and C. Park, *Statistical Inference: The Minimum Distance Approach*. Boca Raton, USA: CRC Press, 2011.
- [402] L.D. Sanghvi, “Comparison of genetical and morphological methods for a study of biological differences,” *Am. J. Phys. Anthropol.*, Vol. 11, no. 3, pp. 385–404, 1953.
- [403] I. Vincze, “On the concept and measure of information contained in an observation,” in: J. Gani and V.K. Rohatgi (eds.), *Contributions to Probability*, pp. 207–214. New York, NY, USA: Academic Press, 1981.
- [404] L. Le Cam, *Asymptotic Methods in Statistical Decision Theory*. New York, NY, USA: Springer, 1986.
- [405] T. Liu, Z. Yuan, J. Sun, J. Wang, N. Zheng, X. Tang, and H.-Y. Shum, “Learning to Detect a Salient Object,” *IEEE Trans. Pattern Anal. Mach. Intell.*, Vol. 33, no. 2, pp. 353–367, 2011.
- [406] S. Kotz, T.J. Kozubowski and K. Podgorski, *The Laplace Distribution and Generalizations*. Boston, MA, USA: Birkhäuser, 2001.
- [407] D.B. Madan and E. Seneta, “The variance gamma (V.G.) model for share market returns,” *J. Business*, Vol. 63, no. 4, pp. 511–524, 1990.
- [408] B. Klar, “A note on gamma difference distributions,” *Journal of Statistical Computation and Simulation*, vol. 85, no. 18, pp. 3708–3715, 2015.
- [409] A. Dembo and O. Zeitouni, *Large Deviations Techniques and Applications*, 2nd ed (corr. print.). New York, USA: Springer, 2009.
- [410] F.W. Steutel and K. van Harn, *Infinite Divisibility of Probability Distributions On The Real Line*. New York, USA: Marcel Dekker Inc., 2004.

A precise bare simulation approach to the minimization of some distances. I. Foundations — Supplementary Material —

Michel Broniatowski³³ and Wolfgang Stummer³⁴

APPENDIX C
PROOFS — PART 3

Proof of Lemma 14.

From (41) one gets straightforwardly for arbitrary $\tilde{c} > 0$

$$D_{\tilde{c}, \varphi_\gamma}(m \cdot \mathbf{Q}, \mathbb{P}) := \begin{cases} \frac{\tilde{c} \cdot (m^\gamma \cdot H_\gamma - m \cdot A \cdot \gamma + \gamma - 1)}{\gamma \cdot (\gamma - 1)}, & \text{if } \gamma \in] - \infty, 0[, \mathbb{P} \in \mathbb{S}_{\geq 0}^K, \mathbf{Q} \in A \cdot \mathbb{S}_{> 0}^K \text{ and } m > 0, \\ \tilde{c} \cdot (-\log m + \tilde{I} - 1 + m \cdot A), & \text{if } \gamma = 0, \mathbb{P} \in \mathbb{S}_{\geq 0}^K, A \cdot \mathbf{Q} \in \mathbb{S}_{> 0}^K \text{ and } m > 0, \\ \frac{\tilde{c} \cdot (m^\gamma \cdot H_\gamma - m \cdot A \cdot \gamma + \gamma - 1)}{\gamma \cdot (\gamma - 1)}, & \text{if } \gamma \in]0, 1[, \mathbb{P} \in \mathbb{S}_{\geq 0}^K, \mathbf{Q} \in A \cdot \mathbb{S}_{\geq 0}^K \text{ and } m \geq 0, \\ \tilde{c} \cdot (A \cdot m \cdot \log m + m \cdot (I - A) + 1), & \text{if } \gamma = 1, \mathbb{P} \in \mathbb{S}_{\geq 0}^K, \mathbf{Q} \in A \cdot \mathbb{S}_{\geq 0}^K \text{ and } m \geq 0, \\ \frac{\tilde{c} \cdot (m^\gamma \cdot H_\gamma \cdot 1_{[0, \infty[}(m) - m \cdot A \cdot \gamma + \gamma - 1)}{\gamma \cdot (\gamma - 1)}, & \text{if } \gamma \in]1, 2[, \mathbb{P} \in \mathbb{S}_{> 0}^K, \mathbf{Q} \in A \cdot \mathbb{S}^K \text{ and } m \in] - \infty, \infty[, \\ \frac{\tilde{c} \cdot (m^2 \cdot H_2 - m \cdot A \cdot 2 + 2 - 1)}{2 \cdot (2 - 1)}, & \text{if } \gamma = 2, \mathbb{P} \in \mathbb{S}_{> 0}^K, \mathbf{Q} \in A \cdot \mathbb{S}^K \text{ and } m \in] - \infty, \infty[, \\ \frac{\tilde{c} \cdot (m^\gamma \cdot H_\gamma \cdot 1_{[0, \infty[}(m) - m \cdot A \cdot \gamma + \gamma - 1)}{\gamma \cdot (\gamma - 1)}, & \text{if } \gamma \in]2, \infty[, \mathbb{P} \in \mathbb{S}_{> 0}^K, \mathbf{Q} \in A \cdot \mathbb{S}^K \text{ and } m \in] - \infty, \infty[, \\ \infty, & \text{else,} \end{cases} \quad (177)$$

where we have used the three m -independent abbreviations

$$H_\gamma := \sum_{k=1}^K (q_k)^\gamma \cdot (p_k)^{1-\gamma} = 1 + \gamma \cdot (A - 1) + \frac{\gamma \cdot (\gamma - 1)}{\tilde{c}} \cdot D_{\tilde{c}, \varphi_\gamma}(\mathbf{Q}, \mathbb{P}), \quad (\text{cf. (42)})$$

$$I := \sum_{k=1}^K q_k \cdot \log \left(\frac{q_k}{p_k} \right) = \frac{1}{\tilde{c}} \cdot D_{\tilde{c}, \varphi_1}(\mathbf{Q}, \mathbb{P}) + A - 1, \quad (\text{cf. (43)})$$

$$\tilde{I} := \sum_{k=1}^K p_k \cdot \log \left(\frac{p_k}{q_k} \right) = \frac{1}{\tilde{c}} \cdot D_{\tilde{c}, \varphi_0}(\mathbf{Q}, \mathbb{P}) + 1 - A. \quad (\text{cf. (44)})$$

To proceed, let us fix an arbitrary constant $\tilde{c} > 0$.

(i) Case $\gamma \cdot (1 - \gamma) \neq 0$.

(ia) Let us start with the subcase $\gamma \in] - \infty, 0[$. From the first and the last line of (177), it is clear that the corresponding m -infimum can not be achieved for $m \leq 0$; since $H_\gamma > 0$ one gets the unique minimizer $m_{\min} = \left(\frac{H_\gamma}{A} \right)^{1/(1-\gamma)} > 0$ and the minimum $D_{\tilde{c}, \varphi_\gamma}(m_{\min} \cdot \mathbf{Q}, \mathbb{P}) = \frac{\tilde{c}}{\gamma} \cdot \left(1 - \frac{H_\gamma^{1/(1-\gamma)}}{A^{\gamma/(1-\gamma)}} \right)$. Hence, (45) is established. The assertions (46) and (47) follow immediately by monotonicity inspection of $x \rightarrow \frac{\tilde{c}}{\gamma} \cdot \left[1 - \frac{1}{A^{\gamma/(1-\gamma)}} \cdot \left[1 + \gamma \cdot (A - 1) + \frac{\gamma \cdot (\gamma - 1)}{\tilde{c}} \cdot x \right]^{-1/(\gamma-1)} \right]$ for $x \geq 0$ such that $1 + \gamma \cdot (A - 1) + \frac{\gamma \cdot (\gamma - 1)}{\tilde{c}} \cdot x \geq 0$.

(ib) The subcase $\gamma \in]0, 1[$ (cf. the third line of (177)) works analogously if $H_\gamma > 0$; furthermore, if $H_\gamma = 0$ — which can only appear when \mathbb{P}, \mathbf{Q} have disjoint supports (singularity) — then $\inf_{m > 0} D_{\tilde{c}, \varphi_\gamma}(m \cdot \mathbf{Q}, \mathbb{P}) = \frac{\tilde{c}}{\gamma}$ which is (the corresponding special case of) (45).

(ic) In the subcase $\gamma \in]1, \infty[$ (cf. the fifth, sixth and seventh line of (177)) it is straightforward to see that the desired infimum can not be achieved for $m < 0$. Hence, one can proceed analogously to subcase (ia).

(id) The assertions (48) to (51) are straightforward.

(ii) Case $\gamma = 1$. From the fourth line of (177), one obtains the unique minimizer $m_{\min} = \exp\{-I/A\}$ and the minimum $D_{\tilde{c}, \varphi_1}(m_{\min} \cdot \mathbf{Q}, \mathbb{P}) = \tilde{c} \cdot (1 - A \cdot m_{\min})$, which leads to (52). The monotonicity of $x \rightarrow \tilde{c} \cdot (1 - \exp\{-x/\tilde{c}\})$ for $x \geq 0$

³³LPSM, Sorbonne Université, 4 place Jussieu, 75252 Paris, France. ORCID 0000-0001-6301-5531.

³⁴Department of Mathematics, University of Erlangen-Nürnberg (FAU), Cauerstrasse 11, 91058 Erlangen, Germany; e-mail: stummer@math.fau.de. ORCID 0000-0002-7831-4558. Corresponding author.

implies immediately (53) and (54); moreover, (55) and (56) are immediate.

(iii) Case $\gamma = 0$. The second line of (177) implies the unique minimizer $m_{min} = 1/A$, the minimum $D_{\tilde{c} \cdot \varphi_0}(m_{min} \cdot \mathbf{Q}, \mathbb{P}) = \tilde{c} \cdot (\tilde{I} + \log A)$, and hence (57). The assertions (58) to (61) are obvious. ■

APPENDIX D PROOFS — PART 4

Proof of Proposition 27. The assertion follows straightforwardly from the following two facts:

(i) a moment generating function MGF is infinitely divisible if and only if MGF^c is a moment generating function for all $c > 0$ (cf. e.g. (the MGF-version of) Prop. IV.2.5 of [410]).

(ii) $z \mapsto MGF(z)$ is a moment generating function if and only if $z \mapsto MGF(\check{c} \cdot z) =: MGF_{\check{c}}(z)$ is a moment generating function for all $\check{c} > 0$.

Notice that for each $c > 0, \check{c} > 0$ one has $\text{int}(\text{dom}(MGF)) = \text{int}(\text{dom}(MGF^c))$ and $\text{int}(\text{dom}(MGF_{\check{c}})) = \frac{1}{\check{c}} \cdot \text{int}(\text{dom}(MGF))$, and hence the light-tailedness remains unchanged: $0 \in \text{int}(\text{dom}(MGF))$ if and only if $0 \in \text{int}(\text{dom}(MGF^c))$ if and only if $0 \in \text{int}(\text{dom}(MGF_{\check{c}}))$. Since $\varphi \in \Upsilon(]a, b[)$, we have

$$\varphi(t) = \sup_{z \in]\lambda_-, \lambda_+[} \left(z \cdot t - \log \left(\int_{\mathbb{R}} e^{z \cdot y} d\zeta(y) \right) \right), \quad t \in]a, b[, \quad (178)$$

and thus for the exponential of its Fenchel-Legendre transform

$$e^{\varphi_*(z)} = \int_{\mathbb{R}} e^{z \cdot y} d\zeta(y), \quad z \in]\lambda_-, \lambda_+[. \quad (179)$$

Now, let $\tilde{\varphi} := \tilde{c} \cdot \varphi \in \Upsilon(]a, b[)$ for arbitrarily fixed $\tilde{c} > 0$. From the application of (6) to $\tilde{\varphi}$ we obtain

$$\tilde{\varphi}(t) = \sup_{\tilde{z} \in]\tilde{\lambda}_-, \tilde{\lambda}_+[} \left(\tilde{z} \cdot t - \log \left(\int_{\mathbb{R}} e^{\tilde{z} \cdot \tilde{y}} d\tilde{\zeta}_{\tilde{c}}(\tilde{y}) \right) \right), \quad t \in]a, b[, \quad (180)$$

for some unique probability distribution $\tilde{\zeta}_{\tilde{c}}$ on \mathbb{R} . Here, we have used $\tilde{\lambda}_- := \inf_{t \in]a, b[} \tilde{\varphi}'(t) = \tilde{c} \cdot \lambda_-$ and $\tilde{\lambda}_+ := \sup_{t \in]a, b[} \tilde{\varphi}'(t) = \tilde{c} \cdot \lambda_+$. Dividing (180) by \tilde{c} , we arrive at

$$\begin{aligned} \varphi(t) = \frac{\tilde{\varphi}(t)}{\tilde{c}} &= \sup_{\tilde{z} \in]\tilde{c} \cdot \lambda_-, \tilde{c} \cdot \lambda_+[} \left(\frac{\tilde{z}}{\tilde{c}} \cdot t - \log \left(\left(\int_{\mathbb{R}} e^{\frac{\tilde{z}}{\tilde{c}} \cdot \tilde{y} \cdot \tilde{c}} d\tilde{\zeta}_{\tilde{c}}(\tilde{y}) \right)^{1/\tilde{c}} \right) \right), \\ &= \sup_{z \in]\lambda_-, \lambda_+[} \left(z \cdot t - \log \left(\left(\int_{\mathbb{R}} e^{z \cdot \tilde{y} \cdot \tilde{c}} d\tilde{\zeta}_{\tilde{c}}(\tilde{y}) \right)^{1/\tilde{c}} \right) \right), \quad t \in]a, b[, \end{aligned}$$

and hence for the exponential of its Fenchel-Legendre transform

$$e^{\varphi_*(z)} = \left(\int_{\mathbb{R}} e^{z \cdot \tilde{y} \cdot \tilde{c}} d\tilde{\zeta}_{\tilde{c}}(\tilde{y}) \right)^{1/\tilde{c}}, \quad z \in]\lambda_-, \lambda_+[. \quad (181)$$

From (179) and (181) we deduce the relation $(MGF_{\zeta}(z))^{\tilde{c}} = MGF_{\tilde{\zeta}_{\tilde{c}}}(\tilde{c} \cdot z)$ which (with the help of (i) and (ii)) implies the infinite divisibility of ζ .

For the reverse direction, let us assume that $\varphi \in \Upsilon(]a, b[)$ and that the corresponding ζ is infinitely divisible. Recall that $]a, b[= \text{int}(\text{dom}(\varphi))$. Moreover, we fix an arbitrary constant $\tilde{c} > 0$. Of course, there holds $\tilde{c} \cdot \varphi \in \tilde{\Upsilon}(]a, b[)$ and $\text{dom}(\tilde{c} \cdot \varphi) = \text{dom}(\varphi)$. Furthermore, by multiplying (178) with $\tilde{c} > 0$ and by employing (i), (ii) we get

$$\begin{aligned} \tilde{c} \cdot \varphi(t) &= \sup_{z \in]\lambda_-, \lambda_+[} \left(\tilde{c} \cdot z \cdot t - \log \left(\left(\int_{\mathbb{R}} e^{\tilde{c} \cdot z \cdot \frac{y}{\tilde{c}}} d\zeta(y) \right)^{\tilde{c}} \right) \right) = \sup_{\tilde{z} \in]\tilde{c} \cdot \lambda_-, \tilde{c} \cdot \lambda_+[} \left(\tilde{z} \cdot t - \log \left(\left(\int_{\mathbb{R}} e^{\frac{\tilde{z}}{\tilde{c}} \cdot y} d\zeta(y) \right)^{\tilde{c}} \right) \right) \\ &= \sup_{\tilde{z} \in]\tilde{c} \cdot \lambda_-, \tilde{c} \cdot \lambda_+[} \left(\tilde{z} \cdot t - \log \left(\int_{\mathbb{R}} e^{\tilde{z} \cdot y} d\tilde{\zeta}_{\tilde{c}}(y) \right) \right), \quad t \in]a, b[, \end{aligned}$$

for some probability distribution $\tilde{\zeta}_{\tilde{c}}$ on \mathbb{R} . ■

Proof of Proposition 28. It is well known that a candidate function $M :]-\infty, 0[\mapsto]0, \infty[$ is the moment-generating function of an infinitely divisible probability distribution if and only if $(\log M)'$ is absolutely monotone (see e.g. Theorem 5.11 of [328]). By applying this to $M(z) := e^{-a \cdot z + \varphi^*(z)}$ respectively $M(z) := e^{b \cdot z + \varphi^*(-z)}$, one gets straightforwardly the assertion (a) respectively (b); notice that $b = \infty$ respectively $a = -\infty$ can be deduced from the fact that the support of an infinitely divisible distribution is always (one-sided or two-sided) unbounded. For the third case $a = -\infty, b = \infty$ one can use the assertion (cf. e.g. [343], p.73) that a candidate function $M :]\lambda_-, \lambda_+[\mapsto]0, \infty[$ is the moment-generating function of an infinitely divisible probability distribution if the connected function $z \mapsto (\log M)''(z)/(\log M)''(0)$ is the moment-generating function of some auxiliary probability distribution; but the latter is equivalent to exponential convexity (cf. Widder's theorem). By applying this to $M(z) := e^{\varphi^*(z)}$, one ends up with (c). ■

APPENDIX E
PROOFS — PART 5

Proof of Theorem 22. (i) Clearly, on $]\lambda_-, \lambda_+[$ the function Λ is differentiable with strictly increasing derivative

$$\Lambda'(z) = F^{-1}(z + c) + 1 - F^{-1}(c), \quad z \in]\lambda_-, \lambda_+[. \quad (182)$$

Hence, Λ is strictly convex and smooth (because of the smoothness of F^{-1}), and satisfies $\Lambda(0) = 0$ as well as $\Lambda'(0) = 1$. Also, the corresponding extensions of Λ to $z = \lambda_-$ and $z = \lambda_+$ are continuous.

(ii) It is straightforward to see that on $]t_-^{sc}, t_+^{sc}[$ the function φ is differentiable with strictly increasing derivative

$$\varphi'(t) = F(t + F^{-1}(c) - 1) - c, \quad t \in]t_-^{sc}, t_+^{sc}[. \quad (183)$$

Hence, φ is strictly convex and smooth (because of the smoothness of F), and satisfies $\varphi(1) = 0$ as well as $\varphi'(1) = 0$. Also, the corresponding extensions of φ to $t = t_-^{sc}$ and $t = t_+^{sc}$ are continuous.

To prove that $\text{int}(\text{dom}(\varphi)) =]a, b[$, let us first notice that obviously there holds $a \leq t_-^{sc}$ and $t_+^{sc} \leq b$. Moreover, the validity of $\varphi(t) < \infty$ for all $t \in]t_-^{sc}, t_+^{sc}[$ is clear from (131) since $t + F^{-1}(c) - 1 \in]a_F, b_F[= \text{int}(\text{dom}(F))$ and the involved integral over the continuous function F^{-1} is taken over a compact interval.

For the subcase $t_-^{sc} = -\infty = a$ we have thus shown $\text{dom}(\varphi) \cap]-\infty, 1] =]-\infty, 1] =]a, 1]$, whereas for the subcase $t_+^{sc} = \infty = b$ we have verified $\text{dom}(\varphi) \cap [1, \infty[= [1, \infty[= [1, b[$.

Let us next examine the subcase “ $t_-^{sc} > -\infty$ and $\varphi(t_-^{sc}) < \infty$ ”: if $\lambda_- > -\infty$ then $a = -\infty$ and (131) implies $\varphi(t) = \varphi(t_-^{sc}) + \lambda_- \cdot (t - t_-^{sc}) < \infty$ for all $t \in]-\infty, t_-^{sc}] =]a, t_-^{sc}]$, which leads to $\text{dom}(\varphi) \cap]-\infty, 1] =]-\infty, 1] =]a, 1]$; in contrast, if $\lambda_- = -\infty$ then $a = t_-^{sc}$ and (131) implies $\varphi(t) = \varphi(t_-^{sc}) + \lambda_- \cdot (t - t_-^{sc}) = \infty$ for all $t \in]-\infty, t_-^{sc}[=]-\infty, a[$, which leads to $\text{dom}(\varphi) \cap]-\infty, 1] = [a, 1]$.

In the subcase “ $t_-^{sc} > -\infty$ and $\varphi(t_-^{sc}) = \infty$ ”, due to the strict convexity of φ one always has $\lim_{t \downarrow t_-^{sc}} \varphi'(t) = -\infty$; this implies, by $\lim_{t \downarrow t_-^{sc}} \varphi'(t) = \lambda_-$, that $\lambda_- = -\infty$ and thus $a = t_-^{sc}$; from (131) we derive $\varphi(t) = \varphi(t_-^{sc}) + \lambda_- \cdot (t - t_-^{sc}) = \infty$ for all $t \in]-\infty, t_-^{sc}[=]-\infty, a[$, which leads to $\text{dom}(\varphi) \cap]-\infty, 1] =]a, 1]$.

As a further step, we deal with the subcase “ $t_+^{sc} < \infty$ and $\varphi(t_+^{sc}) < \infty$ ”: if $\lambda_+ < \infty$ then $b = \infty$ and (131) implies $\varphi(t) = \varphi(t_+^{sc}) + \lambda_+ \cdot (t - t_+^{sc}) < \infty$ for all $t \in [t_+^{sc}, \infty[= [t_+^{sc}, b[$, which leads to $\text{dom}(\varphi) \cap [1, \infty[= [1, \infty[= [1, b[$; in contrast, if $\lambda_+ = \infty$ then $b = t_+^{sc}$ and (131) implies $\varphi(t) = \varphi(t_+^{sc}) + \lambda_+ \cdot (t - t_+^{sc}) = \infty$ for all $t \in]t_+^{sc}, \infty[=]b, \infty[$, which leads to $\text{dom}(\varphi) \cap [1, \infty[= [1, b]$.

In the subcase “ $t_+^{sc} < +\infty$ and $\varphi(t_+^{sc}) = \infty$ ”, due to the strict convexity of φ one always gets $\lim_{t \uparrow t_+^{sc}} \varphi'(t) = \infty$; this implies, by $\lim_{t \uparrow t_+^{sc}} \varphi'(t) = \lambda_+$, that $\lambda_+ = \infty$ and thus $b = t_+^{sc}$; from (131) we deduce $\varphi(t) = \varphi(t_+^{sc}) + \lambda_+ \cdot (t - t_+^{sc}) = \infty$ for all $t \in]t_+^{sc}, \infty[=]b, \infty[$, which leads to $\text{dom}(\varphi) \cap [1, \infty[= [1, b]$.

Putting things together, we have proved that $\text{int}(\text{dom}(\varphi)) =]a, b[$.

(iii) From (182) and (183) one gets easily

$$\Lambda'^{-1}(t) = F(t + F^{-1}(c) - 1) - c = \varphi'(t), \quad t \in]t_-^{sc}, t_+^{sc}[, \quad (184)$$

as well as $\Lambda'^{-1}(1) = 0$. From this, we derive

$$\begin{aligned} & t \cdot \Lambda'^{-1}(t) - \Lambda(\Lambda'^{-1}(t)) \\ &= t \cdot [F(t + F^{-1}(c) - 1) - c] + [F^{-1}(c) - 1] \cdot [F(t + F^{-1}(c) - 1) - c] \\ &\quad - \int_0^{F(t + F^{-1}(c) - 1) - c} F^{-1}(u + c) du \\ &= \varphi(t), \quad t \in]t_-^{sc}, t_+^{sc}[, \end{aligned}$$

and hence,

$$\varphi(t) = \max_{z \in]\lambda_-, \lambda_+[} (z \cdot t - \Lambda(z)), \quad t \in]t_-^{sc}, t_+^{sc}[,$$

i.e. on $]t_-^{sc}, t_+^{sc}[$ the divergence generator φ is the classical Legendre transform of the restriction of Λ to $] \lambda_-, \lambda_+[$. If “ $\lambda_- > -\infty$, $\Lambda(\lambda_-) \in]-\infty, \infty[$ and $\Lambda'(\lambda_-) \in]-\infty, \infty[$ ” respectively “ $\lambda_+ < +\infty$, $\Lambda(\lambda_+) \in]-\infty, \infty[$ and $\Lambda'(\lambda_+) \in]-\infty, \infty[$ ”, then one can apply classical facts of Fenchel-Legendre transformation to get the corresponding left-hand respectively right-hand linear extensions of φ on the complement of $]t_-^{sc}, t_+^{sc}[$, in order to obtain the desired

$$\varphi(t) = \sup_{z \in]-\infty, \infty[} (z \cdot t - \Lambda(z)), \quad t \in \mathbb{R};$$

notice that $t_-^{sc} = \lim_{z \downarrow \lambda_-} \Lambda'(z)$ and $t_+^{sc} = \lim_{z \uparrow \lambda_+} \Lambda'(z)$.

(iv) This is just the reverse of (iii), by applying standard Fenchel-Legendre-transformation theory. \blacksquare

APPENDIX F FURTHER DETAILS AND PROOFS FOR SUBSECTION X-A

Proof of Lemma 20. By Assumption (OM), one gets for all $\lambda \in cl(\Lambda)$ that $\{\mathbf{x} \in (dom(\tilde{\varphi})^n : T(\mathbf{x}) = \lambda) \cap]t_-^{sc}, t_+^{sc}[^n \neq \emptyset$. Moreover, for any $\mathbf{x} = (x_1, \dots, x_n)$ in \mathbb{R}^n , by the independence of the components of $\tilde{\mathbf{W}}$ we have

$$\begin{aligned} I_{\tilde{\mathbf{W}}}(\mathbf{x}) &= \sup_{\mathbf{z}=(z_1, \dots, z_n) \in \mathbb{R}^n} \left(\langle \mathbf{x}, \mathbf{z} \rangle - \sum_{i=1}^n \Lambda_{\tilde{\varphi}}(z_i) \right) = \sup_{\mathbf{z} \in]\lambda_-, \lambda_+[^n} \left(\sum_{i=1}^n (x_i \cdot z_i - \Lambda_{\tilde{\varphi}}(z_i)) \right) \\ &= \sum_{i=1}^n \left(\sup_{z_i \in]\lambda_-, \lambda_+]} (x_i \cdot z_i - \Lambda_{\tilde{\varphi}}(z_i)) \right) = \sum_{i=1}^n \tilde{\varphi}(x_i) = \sum_{k=1}^K \sum_{i \in I_k^{(n)}} \tilde{\varphi}(x_i) \end{aligned}$$

which is finite if and only if $\mathbf{x} \in (dom(\tilde{\varphi}))^n$ (recall that $\tilde{\varphi}$ is a nonnegative function). Hence, for each $\lambda \in \Lambda$ we obtain

$$I(\lambda) := \inf_{\mathbf{x} \in \mathbb{R}^n : T(\mathbf{x}) = \lambda} I_{\tilde{\mathbf{W}}}(\mathbf{x}) = \inf_{\mathbf{x} \in (dom(\tilde{\varphi}))^n : T(\mathbf{x}) = \lambda} I_{\tilde{\mathbf{W}}}(\mathbf{x}) = \inf_{\mathbf{x} \in (dom(\tilde{\varphi}))^n : T(\mathbf{x}) = \lambda} \sum_{k=1}^K \sum_{i \in I_k^{(n)}} \tilde{\varphi}(x_i) \quad (185)$$

$$= \sum_{k=1}^K n_k \cdot \tilde{\varphi}(\lambda_k) = n \cdot \sum_{k=1}^K \tilde{p}_k \cdot \tilde{\varphi}(\lambda_k) = \inf_{\mathbf{x} \in]t_-^{sc}, t_+^{sc}[^n : T(\mathbf{x}) = \lambda} I_{\tilde{\mathbf{W}}}(\mathbf{x}); \quad (186)$$

here, we have employed the following facts: (i) the right-most infimum in (185) is achieved by minimizing each of the K terms $\sum_{i \in I_k^{(n)}} \tilde{\varphi}(x_i)$ under the linear constraint $\frac{1}{n_k} \cdot \sum_{i \in I_k^{(n)}} x_i = \lambda_k$, and by the strict convexity of $\tilde{\varphi}$ on $]t_-^{sc}, t_+^{sc}[$ the minimum of this generic term is attained when all components x_i are equal to λ_k , and (ii) the outcoming minimum does not depend on the particular (generally non-unique) choice of the x_i 's. Notice that we have used the relation $n_k = n \cdot \tilde{p}_k$ as well. To proceed, let $\underline{\lambda}$ be a minimal rate point (mrp) of Λ , which means that $\underline{\lambda} \in \partial\Lambda$ and $I(\underline{\lambda}) \leq I(\lambda)$ for all $\lambda \in \Lambda$. By Assumption (OM) one can run all the steps in (185) and (186) with $\underline{\lambda}$ instead of λ , and hence

$$I(\underline{\lambda}) = \inf_{\mathbf{x} \in \mathbb{R}^n : T(\mathbf{x}) = \underline{\lambda}} I_{\tilde{\mathbf{W}}}(\mathbf{x}) = \inf_{\mathbf{x} \in]t_-^{sc}, t_+^{sc}[^n : T(\mathbf{x}) = \underline{\lambda}} I_{\tilde{\mathbf{W}}}(\mathbf{x}) = n \cdot \sum_{k=1}^K \tilde{p}_k \cdot \tilde{\varphi}(\underline{\lambda}_k) = n \cdot \sum_{k=1}^K \tilde{p}_k \cdot \tilde{\varphi}(\tilde{q}_k / \tilde{p}_k)$$

where for the last equality we have employed the vector $\tilde{\mathbf{Q}} = \underline{\lambda} \cdot \mathfrak{D}^{-1} \in \partial\tilde{\Omega}$ which we have called the “dominating point of $\tilde{\Omega}$ ”. Also we have proved

$$I(\underline{\lambda}) = n \cdot \inf_{\tilde{\mathbf{Q}} \in \tilde{\Omega}} \sum_{k=1}^K \tilde{p}_k \cdot \tilde{\varphi}(\tilde{q}_k / \tilde{p}_k).$$

The existence of a mrp $\underline{\lambda}$ of Λ (or equivalently, of a mrp $\tilde{\mathbf{Q}}$ of $\tilde{\Omega}$) is a straightforward consequence of the continuity of $D_{\tilde{\varphi}}(\cdot, \tilde{\mathbb{P}})$, the strict increasingness of $[0, 1] \ni t \mapsto D_{\tilde{\varphi}}((1-t) \cdot \tilde{\mathbb{P}} + t \cdot \tilde{\mathbf{Q}}, \tilde{\mathbb{P}})$ for $\tilde{\mathbf{Q}} \in int(\tilde{\Omega})$, and $\partial\tilde{\Omega} \ni \tilde{\mathbf{Q}} \cdot \mathfrak{D} = \underline{\lambda} \in \partial\Lambda$. \blacksquare

On the obtainment of proxies of minimal rate points by proxy method 2:

For the rest of this section, besides (OM) we assume that $dom(\tilde{\varphi}) =]a, b[=]t_-^{sc}, t_+^{sc}[$, and that in case of $a = -\infty$ or $b = +\infty$ the divergence generator $\tilde{\varphi}$ is regularly varying at a or b accordingly, with positive index β , i.e. (with a slight abuse of notation)

- if $a = -\infty$, then for all $\lambda > 0$ there holds

$$\lim_{u \rightarrow -\infty} \frac{\tilde{\varphi}(\lambda \cdot u)}{\tilde{\varphi}(u)} = \lambda^\beta,$$

- if $b = +\infty$, then for all $\lambda > 0$ there holds

$$\lim_{u \rightarrow +\infty} \frac{\tilde{\varphi}(\lambda \cdot u)}{\tilde{\varphi}(u)} = \lambda^\beta;$$

this assumption is denoted by $(H\tilde{\varphi})$.

A proxy of $\tilde{\mathbf{Q}}$ can be obtained by sampling from a distribution on \mathbb{R}^K defined through

$$f(\tilde{\mathbf{Q}}) := C \cdot \exp\left(-\sum_{k=1}^K \tilde{p}_k \cdot \tilde{\varphi}(\tilde{q}_k/\tilde{p}_k)\right) = C \cdot \exp\left(-D_{\tilde{\varphi}}(\tilde{\mathbf{Q}}, \tilde{\mathbb{P}})\right) \quad (\text{cf. (112)})$$

where C is a normalizing constant; strict convexity of $\tilde{\varphi}$ together with $(H\tilde{\varphi})$ prove that f is a well-defined (Lebesgue-) density for a random variable \mathbf{T} on \mathbb{R}^K . We denote by $\mathbb{F}(\cdot) := \mathbb{P}[\mathbf{T} \in \cdot]$ the corresponding distribution on \mathbb{R}^K having density f . The distribution of \mathbf{T} given $(\mathbf{T} \in \tilde{\Omega})$ concentrates on the points in $\tilde{\Omega}$ which minimize $D_{\tilde{\varphi}}(\tilde{\mathbf{Q}}, \tilde{\mathbb{P}})$ as $\tilde{\mathbf{Q}}$ runs in $\tilde{\Omega}$, when $D_{\tilde{\varphi}}(\tilde{\Omega}, \tilde{\mathbb{P}})$ is large. This can be argued as follows. We will consider the case when $\tilde{\Omega}$ is a compact subset in $\mathbb{R}_{>0}^K$ and $\tilde{\varphi}$ satisfies $(H\tilde{\varphi})$ with $b = +\infty$. For the case when $\tilde{\Omega}$ is not compact, or belongs to $\mathbb{R}^K/\{\mathbf{0}\}$, see the Remark 45 hereunder. Consider a compact set Γ in $\tilde{\Omega}$ and let Γ_t be defined as deduced from Γ in a way that makes $D_{\tilde{\varphi}}(\Gamma_t, \tilde{\mathbb{P}})$ increase with t for sufficiently large t . For instance, set

$$\Gamma_t := t \cdot \Gamma. \quad (187)$$

Hence, in case of $b = +\infty$ the divergence

$$D_{\tilde{\varphi}}(\Gamma_t, \tilde{\mathbb{P}}) = \inf_{\mathbf{g}_t \in \Gamma_t} \sum_{k=1}^K \tilde{p}_k \cdot \tilde{\varphi}\left(\frac{(\mathbf{g}_t)_k}{\tilde{p}_k}\right) = \inf_{\mathbf{g} \in \Gamma} \sum_{k=1}^K \tilde{p}_k \cdot \tilde{\varphi}\left(\frac{t \cdot g_k}{\tilde{p}_k}\right)$$

tends to infinity as $t \rightarrow \infty$; the case $a = -\infty$ works analogously with $t \rightarrow -\infty$. In case of $b < \infty$ we may consider

$$\Gamma_t := \{b - \mathbf{g}/t; \mathbf{g} \in \Gamma\} \quad (188)$$

and indeed $D_{\tilde{\varphi}}(\Gamma_t, \tilde{\mathbb{P}}) \rightarrow \infty$ as $t \rightarrow \infty$, with a similar statement when $a > -\infty$.

Assume that Γ has a dominating point $\underline{\mathbf{g}}$. Then Γ_t has dominating point $\underline{\mathbf{g}}_t := t \cdot \underline{\mathbf{g}}$. We prove that \mathbf{T} with distribution (112) cannot be too far away (depending on t) from $\underline{\mathbf{g}}_t$ whenever \mathbf{T} belongs to Γ_t . This argument is valid in the present description of some asymptotics which makes Γ_t as a model for $\tilde{\Omega}$ for large t ; considering the case when $D_{\tilde{\varphi}}(\tilde{\Omega}, \tilde{\mathbb{P}})$ is large is captured through the asymptotic statement

$$\lim_{t \rightarrow \infty} D_{\tilde{\varphi}}(\Gamma_t, \tilde{\mathbb{P}}) = +\infty.$$

There holds the following

Proposition 43: With the above notation and under condition $(H\tilde{\varphi})$, denote by \mathbf{B} a neighborhood of $\underline{\mathbf{g}}$ and $\mathbf{B}_t := t \cdot \mathbf{B}$. Then

$$\mathbb{F}[\Gamma_t \cap \mathbf{B}_t^c | \Gamma_t] = \mathbb{P}[\mathbf{T} \in \Gamma_t \cap \mathbf{B}_t^c | \mathbf{T} \in \Gamma_t] \rightarrow 0$$

as $t \rightarrow \infty$, which proves that simulations under (112) produce proxies of the dominating points $\underline{\mathbf{g}}_t$ in Γ_t .

Before we start with the proof of Proposition 43, we first quote the following

Lemma 44: Let $\tilde{\varphi}$ satisfy $(H\tilde{\varphi})$ with $b = +\infty$. Then for all \mathbf{A} in \mathbb{R}^K such that

$$\check{\alpha} := D_{\tilde{\varphi}}(\mathbf{A}, \tilde{\mathbb{P}}) := \inf_{\mathbf{v} \in \mathbf{A}} \sum_{k=1}^K \tilde{p}_k \cdot \tilde{\varphi}\left(\frac{v_k}{\tilde{p}_k}\right)$$

is finite there holds

$$\lim_{t \rightarrow \infty} \frac{1}{t} \cdot \log \int_{\mathbf{A}} \exp\left(-t \cdot \sum_{k=1}^K \tilde{p}_k \cdot \tilde{\varphi}\left(\frac{v_k}{\tilde{p}_k}\right)\right) dv_1 \dots dv_k = -D_{\tilde{\varphi}}(\mathbf{A}, \tilde{\mathbb{P}}).$$

Proof of Lemma 44. Let us first remark that according to the geometry of the set \mathbf{A} , various combinations for the asymptotics with (187) or (188) may occur; for sake of brevity, we only handle the simplest ones, since all turn to be amenable through the same arguments. Denote for positive r

$$\mathbf{B}(r) := \left\{ \mathbf{v} \in \mathbb{R}^K : \sum_{k=1}^K \tilde{p}_k \cdot \tilde{\varphi}\left(\frac{v_k}{\tilde{p}_k}\right) > r \right\}.$$

It holds, by making the change of variable $r = t \cdot \check{\alpha} + t \cdot s$,

$$\begin{aligned} \int_{\mathbf{A}} \exp \left(-t \cdot \sum_{k=1}^K \tilde{p}_k \cdot \tilde{\varphi} \left(\frac{v_k}{\tilde{p}_k} \right) \right) dv_1 \dots dv_K &= \int \dots \int 1_{\mathbb{R}^+}(r) \cdot 1_{\mathbf{A}}(\mathbf{v}) \cdot 1_{t \cdot \sum_{k=1}^K \tilde{p}_k \cdot \tilde{\varphi} \left(\frac{v_k}{\tilde{p}_k} \right), \infty}[(r)] \cdot e^{-r} dr dv_1 \dots dv_K \\ &= t \cdot e^{-t \cdot \check{\alpha}} \int \dots \int 1_{] -\check{\alpha}, \infty[}(s) \cdot 1_{\mathbf{A}}(\mathbf{v}) \cdot 1_{\mathbf{B}^c(\check{\alpha} + s)}(\mathbf{v}) \cdot e^{-t \cdot s} ds dv_1 \dots dv_K = t \cdot e^{-t \cdot \check{\alpha}} \int_{-\check{\alpha}}^{\infty} \text{Vol}(\mathbf{A} \cap \mathbf{B}^c(\check{\alpha} + s)) \cdot e^{-t \cdot s} ds. \end{aligned}$$

Let $I_t := t \cdot \int_0^{\infty} \text{Vol}(\mathbf{A} \cap \mathbf{B}^c(\check{\alpha} + s)) \cdot e^{-t \cdot s} ds$. We prove that

$$\lim_{t \rightarrow \infty} \frac{1}{t} \cdot \log I_t = 0. \quad (189)$$

When $a = -\infty$ or $b = +\infty$, since $\tilde{\varphi}$ satisfies $(\text{H}\tilde{\varphi})$ there exists a polynomial P such that

$$\text{Vol}(\mathbf{A} \cap \mathbf{B}^c(\check{\alpha} + s)) \leq P(s);$$

whence, assuming without loss of generality that $\text{dom}(\tilde{\varphi}) = \mathbb{R}^+$, we obtain

$$\frac{1}{t} \cdot \log I_t \leq \frac{1}{t} \cdot \log \int_0^{\infty} P \left(\frac{u}{t} \right) \cdot t \cdot e^{-u} du$$

which yields that for large t

$$\frac{1}{t} \cdot \log I_t < 0.$$

When dealing with a context where a or b has finite value and the corresponding sets Γ_t are ‘‘far away’’ from Γ in terms of the distance measure $D_{\tilde{\varphi}}(\cdot, \tilde{\mathbb{P}})$, then $\text{Vol}(\mathbf{A} \cap \mathbf{B}^c(\check{\alpha} + s))$ is bounded. Hence, $\limsup_{t \rightarrow \infty} \frac{1}{t} \cdot \log I_t \leq 0$. Now fix $\varepsilon > 0$. Then, since $\text{Vol}(\mathbf{A} \cap \mathbf{B}^c(a + s))$ is increasing in s , we get

$$\begin{aligned} I_t &\geq t \cdot \int_{\varepsilon}^{\infty} \text{Vol}(\mathbf{A} \cap \mathbf{B}^c(\check{\alpha} + s)) \cdot e^{-t \cdot s} ds \\ &\geq \text{Vol}(\mathbf{A} \cap \mathbf{B}^c(\check{\alpha} + \varepsilon)) \cdot e^{-t \cdot \varepsilon}. \end{aligned}$$

Hence

$$\frac{1}{t} \cdot \log I_t \geq \frac{1}{t} \cdot \log \text{Vol}(\mathbf{A} \cap \mathbf{B}^c(\check{\alpha} + \varepsilon)) - \varepsilon$$

which yields $\liminf_{t \rightarrow \infty} \frac{1}{t} \cdot \log I_t \geq 0$. Therefore (189) holds, which concludes the proof. \blacksquare

We now turn to the

Proof of Proposition 43. Without loss of generality, let $b = +\infty$, Γ_t as in (187) and Condition $(\text{H}\tilde{\varphi})$ hold. Moreover, consider an arbitrary neighborhood \mathbf{B} of \underline{g} and the corresponding neighborhoods $\mathbf{B}_t := t \cdot \mathbf{B}$ of $\underline{g}_t = t \cdot \underline{g}$. There holds

$$\begin{aligned} \frac{1}{\tilde{\varphi}(t)} \cdot \log \mathbb{P}[\mathbf{T} \in \Gamma_t] &= \frac{C}{\tilde{\varphi}(t)} \cdot \log \int_{\Gamma_t} \exp \left(-\sum_{k=1}^K \tilde{p}_k \cdot \tilde{\varphi} \left(\frac{w_k}{\tilde{p}_k} \right) \right) dw_1 \dots dw_K \\ &\stackrel{(1)}{=} \frac{C \cdot K}{\tilde{\varphi}(t)} \cdot \log t + \frac{C}{\tilde{\varphi}(t)} \cdot \log \int_{\Gamma} \exp \left(-t^\beta \cdot \sum_{k=1}^K \tilde{p}_k \cdot \left(\tilde{\varphi} \left(\frac{v_k}{\tilde{p}_k} \right) \cdot (1 + o(1)) \right) \right) dv_1 \dots dv_K \\ &\stackrel{(2)}{=} \frac{C \cdot K}{\tilde{\varphi}(t)} \cdot \log t + \frac{C}{(\tilde{\varphi}(t)/t^\beta)} \cdot \frac{1}{t^\beta} \cdot \log \left((1 + o(1)) \cdot \int_{\Gamma} \exp \left(-t^\beta \cdot \sum_{k=1}^K \tilde{p}_k \cdot \tilde{\varphi} \left(\frac{v_k}{\tilde{p}_k} \right) \right) dv_1 \dots dv_K \right) \\ &\stackrel{(3)}{=} -\frac{C \cdot t^\beta}{\tilde{\varphi}(t)} \cdot D_{\tilde{\varphi}}(\Gamma, \tilde{\mathbb{P}}) \cdot (1 + o(1)) \\ &\stackrel{(4)}{=} -\check{\ell}(t) \cdot D_{\tilde{\varphi}}(\Gamma, \tilde{\mathbb{P}}) \cdot (1 + o(1)) \end{aligned}$$

as t tends to infinity. In the above display, (1) follows from $\tilde{\varphi}(t \cdot x) = (t \cdot x)^\beta \cdot \ell(t \cdot x) = t^\beta \cdot x^\beta \cdot \ell(x) \cdot \frac{\ell(t \cdot x)}{\ell(x)} = t^\beta \cdot \tilde{\varphi}(x) \cdot (1 + o(1))$ as t tends to infinity and x lies in a compact subset of $]0, \infty[$, where ℓ is a slowly varying function. The equality (2) follows from compactness of Γ together with the fact that $\tilde{\varphi}$ is a regularly varying function with index β , so that

$$\lim_{t \rightarrow \infty} \frac{\tilde{\varphi}(t \cdot v)}{\tilde{\varphi}(t)} = v^\beta$$

uniformly upon v on compact sets in $]0, \infty[$. The remaining equalities (3) and (4) follow from classical properties of regularly varying functions, where $\ell := 1/\ell$ is a slowly varying function at infinity, together with standard Laplace-Integral approximation.

In the same way we can show

$$\frac{1}{\tilde{\varphi}(t)} \cdot \log \mathbb{P}[\mathbf{T} \in \Gamma_t \cap \mathbf{B}_t^c] = -\check{\ell}(t) \cdot D_{\tilde{\varphi}}(\Gamma \cap \mathbf{B}^c, \tilde{\mathbb{P}}) \cdot (1 + o(1))$$

as t tends to infinity. Since \mathbf{B} is a neighborhood of the unique dominating point \underline{g} of Γ , one gets that $D_{\tilde{\varphi}}(\Gamma \cap \mathbf{B}^c, \tilde{\mathbb{P}}) > D_{\tilde{\varphi}}(\Gamma, \tilde{\mathbb{P}})$. This implies that

$$\mathbb{P}[\mathbf{T} \in \Gamma_t \cap \mathbf{B}_t^c \mid \mathbf{T} \in \Gamma_t] \rightarrow 0 \quad \text{as } t \rightarrow \infty. \quad \blacksquare$$

Remark 45: Firstly, let us quote that the case when $\tilde{\Omega}$ is an unbounded subset in $\mathbb{R}^K \setminus \{0\}$ is somewhat immaterial for applications. Anyhow, if compactness of Γ is lost, then in order to use the same line of arguments as above, it is necessary to strengthen the assumptions (H $\tilde{\varphi}$) e.g. as follows: when $b = +\infty$ then $\tilde{\varphi}$ has to be asymptotically homogeneous with degree $\beta > 0$, in the sense that $\tilde{\varphi}(t \cdot x) = t^\beta \cdot \tilde{\varphi}(x) \cdot (1 + o(1))$ as $t \rightarrow \infty$; for the subcase $a = -\infty$ one employs an analogous assumption as $t \rightarrow -\infty$. The case when $\tilde{\Omega}$ is a compact set in $\mathbb{R}^K \setminus \{0\}$ can be treated as above, by combining the asymptotics in t in the neighborhood of a and b accordingly.

APPENDIX G PROOF FOR SUBSECTION X-B

Proof of Proposition 21. Recall the weighted empirical measure

$$\xi_{n, \mathbf{X}}^{\mathbf{V}} := \left(\frac{1}{n} \sum_{i \in I_1^{(n)}} V_i, \dots, \frac{1}{n} \sum_{i \in I_K^{(n)}} V_i \right)$$

which satisfies the K linear constraints defined in (122) through

$$E_S[\xi_{n, \mathbf{X}}^{\mathbf{V}}] = \xi_{M, \mathbf{X}}^{\mathbf{W}^*} = \overline{W}^* \cdot \xi_{M, \mathbf{X}}^{w \mathbf{W}^*}$$

where $\mathbf{Q}^* := (q_1^*, \dots, q_K^*) = \xi_{M, \mathbf{X}}^{w \mathbf{W}^*} \in \text{int}(\Omega)$ and $\overline{W}^* = \frac{1}{M} \sum_{j=1}^M W_j^*$. The probability distribution S defined on \mathbb{R}^n is the Kullback-Leibler-divergence projection of $\zeta^{\otimes n}$ on the class of all probability distributions on \mathbb{R}^n which satisfy (122). We prove that $\liminf_{n \rightarrow \infty} S[\xi_{n, \mathbf{X}}^{w \mathbf{V}} \in \Omega] > 0$. To start with, we define for strictly positive δ the set

$$A_{n, \delta} := \left\{ \left| \frac{1}{n} \sum_{i=1}^n V_i - \overline{W}^* \right| \leq \delta \right\}$$

and write

$$S[\xi_{n, \mathbf{X}}^{w \mathbf{V}} \in \Omega] = S[\{\xi_{n, \mathbf{X}}^{w \mathbf{V}} \in \Omega\} \cap A_{n, \delta}] + S[\{\xi_{n, \mathbf{X}}^{w \mathbf{V}} \in \Omega\} \cap A_{n, \delta}^c] =: I + II.$$

By the law of large numbers, the second term II tends to 0 as n tends to infinity. Moreover, one can rewrite

$$I = S \left[\bigcup_{m \in [\overline{W}^* - \delta, \overline{W}^* + \delta]} \{\xi_{n, \mathbf{X}}^{\mathbf{V}} \in m \cdot \Omega\} \right]$$

which entails

$$I \geq S \left[\frac{1}{n_k} \sum_{i \in I_k^{(n)}} V_i \in \mathcal{V}_\eta \left(\overline{W}^* \cdot \frac{q_k^*}{p_k} \right) \text{ for all } k \in \{1, \dots, K\} \right],$$

where $\mathcal{V}_\eta \left(\overline{W}^* \cdot \frac{q_k^*}{p_k} \right)$ denotes a neighborhood of $\overline{W}^* \cdot \frac{q_k^*}{p_k}$ with radius η being small when δ is small, for large enough n , making use of the a.s. convergence of n_k/n to p_k . Now, for any $k \in \{1, \dots, K\}$ one has

$$S \left[\frac{1}{n_k} \sum_{i \in I_k^{(n)}} V_i \notin \mathcal{V}_\eta \left(\overline{W}^* \cdot \frac{q_k^*}{p_k} \right) \right] \leq \exp \left(-n_k \cdot \inf_{x \in \mathcal{V}_\eta \left(\overline{W}^* \cdot \frac{q_k^*}{p_k} \right)^c} \varphi(x) \right) \quad (190)$$

since any margin of S with index in $I_k^{(n)}$ is a corresponding Kullback-Leibler-divergence projection of ζ on the set of all distributions on \mathbb{R} with expectation $\overline{W}^* \cdot \frac{q_k^*}{p_{M,k}^{emp}}$ — where $p_{M,k}^{emp}$ denotes the fraction of the X_i 's (within X_1, \dots, X_M) which are equal to d_k (cf. (29)) — and therefore has a moment generating function which is finite in a non-void neighborhood of 0, which yields (190) by the Markov Inequality. Note that the event $\{\xi_{M, \mathbf{X}}^{w \mathbf{W}^*} \in \text{int}(\Omega)\}$ is regenerative, so that M can be chosen large enough to make $p_{M,k}^{emp}$ close to p_k for all $k \in \{1, \dots, K\}$. This proves the claim. \blacksquare