

GLIVENKO–CANTELLI FOR f -DIVERGENCE

HAOMING WANG AND LEK-HENG LIM

ABSTRACT. We extend the celebrated Glivenko–Cantelli theorem, sometimes called the fundamental theorem of statistics, from its standard setting of total variation distance to all f -divergences. A key obstacle in this endeavor is to define f -divergence on a subcollection of a σ -algebra that forms a π -system but not a σ -subalgebra. This is a side contribution of our work. We will show that this notion of f -divergence on the π -system of rays preserves nearly all known properties of standard f -divergence, yields a novel integral representation of the Kolmogorov–Smirnov distance, and has a Glivenko–Cantelli theorem. We will also discuss the prospects of a Vapnik–Chervonenkis theory for f -divergence.

1. INTRODUCTION

The Glivenko–Cantelli theorem [22, 8] is a cornerstone of empirical process theory. It is likely the best known statement regarding the asymptotic behavior of stochastic processes formed by empirical measures [16]. It is also crucial in nonparametric statistics and forms the basis for statistical consistency in many estimation problems [16]. In the Kolmogorov–Smirnov test, the theorem ensures that the test statistic has desirable asymptotic properties [46]. In statistical resampling methods like the bootstrap, it guarantees that the empirical distribution derived from resampled data will approximate the true distribution as the number of samples grows [18]. Because of its many central roles, the Glivenko–Cantelli theorem is often called the fundamental theorem of statistics [49, 27, 14, 44].

In machine learning, the Glivenko–Cantelli theorem can be extended to the Vapnik–Chervonenkis theorem [55], also known as the fundamental theorem of learning theory [45]. It is used to show the consistency of the principle of empirical risk minimization (ERM) [54] and to derive bounds on generalization error by ensuring the difference between empirical and true distributions diminishes with larger samples, aiding in over-fitting control [6].

Given its pivotal role in both statistics and machine learning, it is surprising that the Glivenko–Cantelli theorem is somewhat limited in scope—it only works with the *total variation norm* restricted to a Glivenko–Cantelli class. Modern AI applications, on the other hand, have ushered in a plethora of alternative distances between two probability distributions μ and ν , most of them easier to compute than the total variation norm. The most prominent of these are the f -divergences

$$(1) \quad D_f(\mu||\nu) = \int_{\Omega} f\left(\frac{d\mu}{d\nu}\right) d\nu = \int_{\Omega} f(g(x)) d\nu(x),$$

where g is the Radon–Nikodym derivative with $d\mu(x) = g(x) d\nu(x)$. The total variation norm is itself an f -divergence with $f(t) = |t - 1|/2$ but many other f -divergence with yet other f 's have played a prominent role in important AI applications recently. In [31], the Kullback–Leibler (KL) divergence is used to regularize the posterior distribution in the Variational Autoencoder (VAE), aligning it with a chosen prior. In [23], Generative Adversarial Networks (GANs) are trained by minimizing the Jensen–Shannon (JS) divergence between the model and real data distributions. In [33], Rényi α -divergences provide a smooth interpolation from the evidence lower bound to the log

2020 *Mathematics Subject Classification.* 28A33, 28A50, 46E27, 60E05, 60E15, 60F17, 94A17.

Key words and phrases. probability measures, Kolmogorov–Smirnov distance, total variation distance, f -divergence, Glivenko–Cantelli theorem.

respectively. Their union is the class of *rays*, denoted

$$\mathcal{R} = \mathcal{R}_c \cup \mathcal{R}_o.$$

These are all π -systems [17, p. 202] but not σ -algebras. In this article, we restrict ourselves to “left” rays, i.e., where the left-end point of the interval is always $-\infty$. In topology, \mathcal{R}_o is called the left-ray topology or left-order topology.

Glivenko [22] showed that \mathcal{R}_c is a Glivenko–Cantelli class whereas Cantelli [8] showed that \mathcal{R}_o is a Glivenko–Cantelli class but we will soon see that it makes no difference whether we use \mathcal{R}_c , \mathcal{R}_o , or \mathcal{R} ; all three statements are equivalent. Obviously, the name “Glivenko–Cantelli class” came about much later [46], earlier statements of the result were more of the following form:

Theorem 1.4 (Glivenko–Cantelli). *Let ν be a Borel probability measure and ν_n be the corresponding empirical measure, $n \in \mathbb{N}$. Then, almost surely,*

$$(2) \quad \sup_{A \in \mathcal{R}} (\nu_n(A) - \nu(A)) \rightarrow 0.$$

There is no substantive difference if we instead use \mathcal{R}_o or \mathcal{R}_c in place of \mathcal{R} in Theorem 1.4 as

$$(3) \quad \sup_{A \in \mathcal{R}_c} (\mu(A) - \nu(A)) = \sup_{A \in \mathcal{R}_o} (\mu(A) - \nu(A)) = \sup_{A \in \mathcal{R}} (\mu(A) - \nu(A))$$

for any probability measures μ and ν on $(\mathbb{R}, \mathcal{B})$, a fact that follows easily from the continuity of probability measures. The common value in (3) defines a distance between probability measures called the Kolmogorov–Smirnov distance, although it is usually defined in terms of \mathcal{R}_c [30].

The *total variation distance* is defined by

$$(4) \quad D_{\text{TV}}(\mu \parallel \nu) = \sup_{A \in \mathcal{B}} (\mu(A) - \nu(A)).$$

So the Kolmogorov–Smirnov distance is just the total variation distance “restricted” to $\mathcal{R} \subseteq \mathcal{B}$, i.e., a partial variation over a smaller class of sets than the full Borel σ -algebra. As we have alluded to earlier, the total variation is an example of f -divergence, an extensive class of distance between probability measures. Note that whether we take absolute values or not in (4) makes no difference to the value because $A \in \mathcal{B}$ iff $A^c \in \mathcal{B}$. Throughout this article, an integral sign \int unadorned with upper and lower limits is taken to mean $\int_{\mathbb{R}}$.

Definition 1.5 (f -divergence). Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a continuous convex function with $f(1) = 0$. Let μ and ν be probability measures on $(\mathbb{R}, \mathcal{B})$ be such that $\mu \ll \nu$, i.e., μ is absolutely continuous with respect to ν , so that we may speak of Radon–Nikodym derivative $d\mu/d\nu$. Then the f -divergence [41, 13, 37, 1] between μ and ν is given by

$$(5) \quad D_f(\mu \parallel \nu) = \int f\left(\frac{d\mu}{d\nu}\right) d\nu.$$

Different choices of f yield various well-known divergences, including Kullback–Leibler [3], Le Cam [32], Jensen–Shannon [34], Jeffreys [29], Chernoff [10], Pearson χ^2 [39], Hellinger squared [26], exponential [36], and alpha–beta [19] divergences, and so on. For $f(t) = |t - 1|/2$, we get [52]:

$$\int \frac{1}{2} \left| \frac{d\mu}{d\nu} - 1 \right| d\nu = \sup_{A \in \mathcal{B}} (\mu(A) - \nu(A)),$$

the total variation distance. To obtain a Glivenko–Cantelli theorem for an arbitrary f -divergence, we first need an analog of Kolmogorov–Smirnov distance corresponding to an f -divergence.

Question 1. Is there a notion of f -divergence over \mathcal{R} for general f that (a) reduces to the Kolmogorov–Smirnov distance in (3) when $f(t) = |t - 1|/2$; and (b) reduces to the standard f -divergence in (5) when \mathcal{R} is replaced by \mathcal{B} ?

The answer to Question 1 is yes, provided by what we will call an f -divergence over \mathcal{R} and denoted $D_f^{\mathcal{R}}(\mu||\nu)$. We will establish the existence of f -divergence over \mathcal{R} in Section 3 and develop some of its properties in Section 4, showing in particular that this new notion preserves some of the best-known relations between different f -divergences. In Section 5, we will prove the result claimed in the title of our article, namely, $D_f^{\mathcal{R}}(\nu_n||\nu)$ converges to zero almost surely. Note that when $f(t) = |t - 1|/2$, this reduces to the Glivenko–Cantelli theorem in Theorem 1.4.

As we will see, it takes some effort to establish the notion of f -divergence over \mathcal{R} . For the total variation distance, getting from (4) to (3) is a matter of simply replacing \mathcal{B} by \mathcal{R} . But for general f -divergence, \mathcal{B} does not even appear directly in (5) and it is no longer a matter of simply replacing \mathcal{B} by \mathcal{R} . Indeed we know of no straightforward answer to Question 1.

In particular, we emphasize that the groundbreaking work of Vapnik and Chervonenkis [55] sheds no light on whether the Glivenko–Cantelli theorem remains true when total variation distance is replaced by an f -divergence. They showed that if the VC-dimension of a class $\mathcal{C} \subseteq \Sigma$, i.e., the maximum number of points that \mathcal{C} can shatter, is finite, then \mathcal{C} is a Glivenko–Cantelli class, i.e., $\sup_{A \in \mathcal{C}} (\nu_n(A) - \nu(A)) \rightarrow 0$ almost surely. More generally, existing Vapnik–Chervonenkis theory is still cast entirely in terms of the total variation distance [54]. Indeed, our work in this article suggests that a “Vapnik–Chervonenkis theory for f -divergence” is likely within reach. We will establish some initial results in Section 6.

Lastly, readers familiar with Choquet integrals [11, 15] may think that it provides a plausible answer to Question 1, as we did initially; we will explain why it does not in Section 7.

2. NOTATIONS AND TERMINOLOGIES

2.1. Functions. We reserve the letter f for f -divergence in this article. Other functions will be denoted g, h, φ, ψ . Depending on context, the notation $X_k \uparrow X$ could mean an increasing sequence $X_1 \leq X_2 \leq \dots \leq X_k \in \mathbb{R}$ converging to a real value $X \in \mathbb{R}$ or a nested sequence of sets $X_1 \subseteq X_2 \subseteq \dots \subseteq X_k \subseteq \mathbb{R}$ converging to a set $X \subseteq \mathbb{R}$. Similarly for $X_k \downarrow X$.

We write $\|\cdot\|_p$ for the L^p -norm for all $p \in [1, \infty]$ except $p = 2$. The L^2 -norm will just be denoted as $\|\cdot\|$. We write $L^2(\nu) := L^2(\mathbb{R}, \mathcal{B}, \nu)$ for the space of measurable functions with finite L^2 -norm and

$$L_+^2(\nu) := \{g \in L^2(\nu) : g \geq 0\}$$

for the cone of functions nonnegative almost everywhere. As usual, by “function” we mean an equivalence class of functions that differ at most on a measure-zero set and conditions like “ $h \geq g$ ” will always be in the almost everywhere sense without specification. On the other hand, we will specify the measure every time since we often have to deal with multiple measures.

Composition of functions will always be denoted by \circ and pointwise product of real-valued functions will be denoted by \cdot if necessary for emphasis or otherwise left unmarked.

Let $\emptyset \neq M \subseteq L^2(\nu)$ be a closed convex subset. The *metric projection* onto M is the operator $\text{proj}_M : L^2(\nu) \rightarrow L^2(\nu)$ that takes $g \in L^2(\nu)$ to the closest point $g_* \in M$, i.e., $\|g - g_*\| = \min_{h \in M} \|h - g\|$. The existence and uniqueness of g_* is guaranteed by the conditions on M and so proj_M is well-defined. It is also well-known that proj_M is continuous [2, p. 52].

2.2. Sequences. We will remind the readers of the three notions of convergence used in our article: For a sequence of measurable functions g_n converging to g in $L^2(\nu)$. The convergence is said to be (i) ν -almost surely if g_n converges to g pointwise except on a ν -null set; (ii) in ν -measure if $\lim_{n \rightarrow \infty} \nu(|g_n - g| > \varepsilon) = 0$ for all $\varepsilon > 0$; (iii) in L^2 -norm if $\|g_n - g\| \rightarrow 0$. The following lemma, adapted from [12, pp. 80–83], summarizes the relations between them.

Lemma 2.1. *Let ν be a Borel probability measure.*

- (a) *If $g_n \rightarrow g$ in L^2 -norm, then $g_n \rightarrow g$ in ν -measure.*
- (b) *If $g_n \rightarrow g$ in ν -measure, then there is a subsequence (g_{n_k}) such that $g_{n_k} \rightarrow g$ almost surely.*

(c) If $g_n \rightarrow g$ almost surely and there exists $g \in L^2(\nu)$ that dominates g_n , then $g_n \rightarrow g$ in L^2 -norm.

We will also need an observation that follows from $|\int g_n - g \, d\nu| \leq \int |g_n - g| \, d\nu \leq (\int (g_n - g)^p \, d\nu)^{1/p}$.

Lemma 2.2. *Let $g_n \in L^p(\nu)$. If $g_n \rightarrow g$ in L^p -norm for some $1 \leq p < \infty$, then*

$$\int g \, d\nu = \lim_{n \rightarrow \infty} \int g_n \, d\nu.$$

2.3. Measures. We write 2^Ω for the power set of Ω . A collection of subsets $\mathcal{E} \subseteq 2^\Omega$ is called a π -system if it is closed under finite intersections [17, p. 202]. The following results from [12, p. 38] and [43, p. 58] are reproduced here for easy reference.

Lemma 2.3. *Let μ, ν be finite measures on (Ω, Σ) such that $\mu(\Omega) = \nu(\Omega)$. If \mathcal{A} is a π -system that generates Σ and $\mu(A) = \nu(A)$ for any $A \in \mathcal{A}$, then $\mu = \nu$.*

Lemma 2.4. *Let $g : X \rightarrow Y$ be any function and $\mathcal{E} \subseteq 2^Y$. Then $\sigma(g^{-1}(\mathcal{E})) = g^{-1}(\sigma(\mathcal{E}))$.*

Whenever we write $d\mu/d\nu$, we implicitly assume that $\mu \ll \nu$.

3. f -DIVERGENCE OVER RAYS

The goal of this section is to resolve Question 1 in the affirmative. The answer is an analogue of f -divergence with respect to the class of rays \mathcal{R} that we will call the f -divergence over \mathcal{R} . This appears only at the end of this section in Definition 3.12. It will take the groundwork developed over the course of the whole section before we can show that the notion is indeed well-defined. There are two milestones to watch for: We will see in Theorem 3.15 that the construction of such a divergence would work for any class $\mathcal{C} \subseteq \mathcal{B}$ that satisfies a certain ‘‘Radon–Nikodym property.’’ We will see in Corollary 3.11, which follows from Theorem 3.9, that $\mathcal{C} = \mathcal{R}$ has this property.

We begin by establishing some basic properties related to \mathcal{R} . Throughout this section, ν denotes a Borel probability measure on $(\mathbb{R}, \mathcal{B})$.

Proposition 3.1. *The set*

$$(6) \quad G(\mathcal{R}) := \{g \in L^2(\nu) : \{g > r\} \in \mathcal{R} \text{ for all } r \in \mathbb{R}\}$$

is a closed convex cone in $L^2(\nu)$ comprising all nonincreasing functions in $L^2(\nu)$.

Proof. We first show that $G(\mathcal{R})$ is exactly the set of nonincreasing L^2 -functions. Let g be such a function and let $r \in \mathbb{R}$. Set $a = \sup\{x : g(x) > r\}$. Any $x \in \{g > r\}$ has $x \leq a$ so $\{g > r\} \subseteq (-\infty, a]$. If there is an $x < a$ with $g(x) \leq r$, then there must be some $y > x$ with $g(y) > r \geq g(x)$. This contradicts the assumption that g is nonincreasing. Hence $(-\infty, a) \subseteq \{g > r\} \subseteq (-\infty, a]$. It follows that $\{g > r\} \in \mathcal{R}$ and so $g \in G(\mathcal{R})$. Conversely, let $g \in G(\mathcal{R})$. Suppose we have $x < y$ with $g(x) < g(y)$. Then there is some a with $(-\infty, a) \subseteq \{g > g(x)\} \subseteq (-\infty, a]$. Since $g(y) > g(x)$, it follows that $y \in \{g > g(x)\}$, and so $x < y < a$. Consequently, $x \in \{g > g(x)\}$, i.e., $g(x) > g(x)$, a contradiction. Hence g must be nonincreasing.

$G(\mathcal{R})$ is a closed convex cone because nonnegative linear combinations or L^2 -limits of nonincreasing functions remain nonincreasing. The former is routine. The latter follows from Lemma 2.1: A sequence (g_n) of nonincreasing functions with $\|g_n - g\| \rightarrow 0$ must also converge in ν -measure. So there is a subsequence (g_{n_k}) that converges to g almost surely. So there is a null set A where $\lim_{k \rightarrow \infty} g_{n_k}(x) = g(x)$ for all $x \notin A$. Since each g_{n_k} is nonincreasing, g is almost surely nonincreasing. \square

An immediate consequence of Proposition 3.1 is that $G(\mathcal{R})$ has a well-defined metric projection. For a given $g \in L^2(\nu)$, finding $\text{proj}_{G(\mathcal{R})} g$ is equivalent to finding the nearest nonincreasing function to g in L^2 -norm.

Definition 3.2. A collection $\{A_1, \dots, A_n\}$ is called an *ordered partition* of \mathbb{R} if it satisfies the following conditions:

- (i) disjoint: $A_i \cap A_j = \emptyset$ for any $i \neq j$;
- (ii) cover: $A_1 \cup \dots \cup A_n = \mathbb{R}$;
- (iii) ordered: $\sup A_i \leq \inf A_j$ for each $i \leq j$.

An *ordered simple function* $h : \mathbb{R} \rightarrow \mathbb{R}$ is a simple function defined on an ordered partition $\{A_1, \dots, A_n\}$.

It is well-known that any $g \in L_+^2(\nu)$ can be approximated to arbitrary accuracy by an increasing sequence of nonnegative simple functions (h_n) . We will show in Corollary 3.4 that h_n can be chosen to be ordered simple functions.

Proposition 3.3. *If $g \in L_+^2(\nu)$ is almost surely continuous, then there exists a monotone increasing sequence of nonnegative ordered simple functions (h_n) that converges pointwise to g almost surely.*

Proof. Define the ordered partition

$$\mathcal{P}_n := \left\{ \left[\frac{nk}{2^n}, \frac{n(k+1)}{2^n} \right) : k = -2^n, -2^n + 1, \dots, 2^n - 1 \right\} \cup \{(-\infty, -n), [n, \infty)\}$$

and the ordered simple function $h_n = \sum_{A \in \mathcal{P}_n} \inf_{x \in A} g(x) \mathbb{1}_A$. Since g is almost surely continuous, we can find a null set S such that for any $\varepsilon > 0$ and $x \in \mathbb{R} \setminus S$, there exist $n \in \mathbb{N}$ and $A \in \mathcal{P}_n$ with $x \in A \subseteq B_\varepsilon(x)$. Consequently,

$$\inf_{y \in B_\varepsilon(x)} g(y) \leq \inf_{y \in A} g(y) = h_n(x) \leq g(x),$$

where $\inf_{y \in B_\varepsilon(x)} g(y) \uparrow g(x)$ as $\varepsilon \downarrow 0$ since g is continuous on $\mathbb{R} \setminus S$. Therefore, $h_n(x) \rightarrow g(x)$ almost surely as $n \rightarrow \infty$. For any $n \leq m$, we have $A \subseteq B$ for any $A \in \mathcal{P}_n$ and $B \in \mathcal{P}_m$; so $h_n(x) \leq h_m(x)$ for all $x \in \mathbb{R}$. \square

In Proposition 3.3, h_n is dominated by g ; so we also have $h_n \rightarrow g$ in L^2 -norm by Lemma 2.1. The almost sure continuity in Proposition 3.3 may sometimes be overly restrictive. We derive an alternative version without this requirement.

Corollary 3.4. *Let $g \in L_+^2(\nu)$. Then there exists a sequence of bounded ordered simple functions (h_n) that converges to g in L^2 -norm and almost surely.*

Proof. Given that the space of continuous functions with compact support $C_c(\mathbb{R})$ is dense in $L^2(\nu)$ [28, p. 354], there is a continuous function φ with compact support that approximates g to accuracy $\|g - \varphi\| < \varepsilon/2$ for any $\varepsilon > 0$. On the other hand, the continuous function φ is in $L^2(\nu)$; so by Proposition 3.3 there is an ordered simple function h that approximates it to accuracy $\|\varphi - h\| < \varepsilon/2$. Hence $\|h - g\| < \varepsilon$. This shows that for any $g \in L_+^2(\nu)$, there exists a sequence of ordered simple function h_n such that $h_n \rightarrow g$ in L^2 -norm. By passing through a subsequence if necessary, Lemma 2.1 ensures $h_n \rightarrow g$ almost surely as well. \square

We next see that the metric projection of an ordered simple function onto $G(\mathcal{R})$ remains ordered simple.

Proposition 3.5 (Projection). *Let $h = \sum_{i=1}^n \alpha_i \mathbb{1}_{A_i} \in L_+^2(\nu)$ be an ordered simple function with $\{A_1, \dots, A_n\}$ an ordered partition of \mathbb{R} . Then $\text{proj}_{G(\mathcal{R})} h$ is also an ordered simple function on $\{A_1, \dots, A_n\}$.*

Proof. Suppose that $\text{proj}_{G(\mathcal{R})} h$ is not constant on A_j for some $j = 1, \dots, n$. Define

$$\varphi = \begin{cases} \text{proj}_{G(\mathcal{R})} h & \text{on } A_j^c, \\ c & \text{on } A_j, \end{cases}$$

where $c = \text{proj}_{G(\mathcal{R})} h(\omega_*)$ and $\omega_* = \text{argmin}_{\omega \in A_j} |h(\omega) - \text{proj}_{G(\mathcal{R})} h(\omega)|$. Then $\varphi \in G(\mathcal{R})$ and

$$\int (h - \text{proj}_{G(\mathcal{R})} h)^2 d\nu \geq \int (h - \varphi)^2 d\nu,$$

a contradiction. \square

A consequence of Proposition 3.5 is that metric projection onto $G(\mathcal{R})$ preserves ordering.

Proposition 3.6 (Monotonicity). *Let $g_1, g_2 \in L_+^2(\nu)$. If $g_1 \leq g_2$, then*

$$\text{proj}_{G(\mathcal{R})} g_1 \leq \text{proj}_{G(\mathcal{R})} g_2.$$

Proof. We first show that the statement holds for ordered simple functions. Let $h_1 = \sum_{i=1}^n a_i \mathbb{1}_{A_i}$ and $h_2 = \sum_{j=1}^m b_j \mathbb{1}_{B_j}$ where (A_1, \dots, A_n) and (B_1, \dots, B_m) are ordered partitions. Suppose $h_1 \leq h_2$ but that $\text{proj}_{G(\mathcal{R})} h_1 > \text{proj}_{G(\mathcal{R})} h_2$ on $E = A_i \cap B_j$ for some i and j . Define

$$\begin{aligned} \varphi &= (\text{proj}_{G(\mathcal{R})} h_1) \cdot \mathbb{1}_{E^c} + (\text{proj}_{G(\mathcal{R})} h_2) \cdot \mathbb{1}_E, \\ \psi &= (\text{proj}_{G(\mathcal{R})} h_2) \cdot \mathbb{1}_{E^c} + (\text{proj}_{G(\mathcal{R})} h_1) \cdot \mathbb{1}_E. \end{aligned}$$

Both φ and ψ are clearly nonincreasing and thus \mathcal{R} -measurable. There are six possibilities for the values of these functions on E :

- | | |
|--|--|
| (a) $h_2 \geq h_1 \geq \text{proj}_{G(\mathcal{R})} h_1 \geq \text{proj}_{G(\mathcal{R})} h_2$, | (d) $\text{proj}_{G(\mathcal{R})} h_1 \geq h_2 \geq h_1 \geq \text{proj}_{G(\mathcal{R})} h_2$, |
| (b) $h_2 \geq \text{proj}_{G(\mathcal{R})} h_1 \geq h_1 \geq \text{proj}_{G(\mathcal{R})} h_2$, | (e) $\text{proj}_{G(\mathcal{R})} h_1 \geq h_2 \geq \text{proj}_{G(\mathcal{R})} h_2 \geq h_1$, |
| (c) $h_2 \geq \text{proj}_{G(\mathcal{R})} h_1 \geq \text{proj}_{G(\mathcal{R})} h_2 \geq h_1$, | (f) $\text{proj}_{G(\mathcal{R})} h_1 \geq \text{proj}_{G(\mathcal{R})} h_2 \geq h_2 \geq h_1$. |

For Cases (a), (b), and (c),

$$\|\psi - h_2\| \leq \|\text{proj}_{G(\mathcal{R})} h_2 - h_2\|;$$

for Cases (e) and (f),

$$\|\varphi - h_1\| \leq \|\text{proj}_{G(\mathcal{R})} h_1 - h_1\|;$$

and for Case (d), either of the following

$$\|\psi - h_2\| \leq \|\text{proj}_{G(\mathcal{R})} h_2 - h_2\|, \quad \|\varphi - h_1\| \leq \|\text{proj}_{G(\mathcal{R})} h_1 - h_1\|$$

must hold; all leading to contradictions.

For the general case with $g_1, g_2 \in L_+^2(\nu)$, we set $g_0 := g_2 - g_1 \in L_+^2(\nu)$. By Corollary 3.4, we have two sequences of nonnegative ordered simple functions $(h_{0,n})$ and $(h_{1,n})$ such that $h_{0,n} \rightarrow g_0$ and $h_{1,n} \rightarrow g_1$ in L^2 -norm. Set $h_{2,n} := h_{0,n} + h_{1,n}$, which is nonnegative, ordered simple, and satisfies $h_{2,n} \rightarrow g_2$ in L^2 -norm and $h_{2,n} \geq h_{1,n}$ for each n . By the continuity of the metric projection operator $\text{proj}_{G(\mathcal{R})}$,

$$\text{proj}_{G(\mathcal{R})} g_1 = \lim_{n \rightarrow \infty} \text{proj}_{G(\mathcal{R})} h_{1,n} \leq \lim_{n \rightarrow \infty} \text{proj}_{G(\mathcal{R})} h_{2,n} = \text{proj}_{G(\mathcal{R})} g_2,$$

as required. \square

A slight refinement of the last step of the proof above yields the following.

Corollary 3.7 (Sequential monotonicity). *Let $g \in L_+^2(\nu)$. If (h_n) is a monotone increasing sequence in $L_+^2(\nu)$ converging to g , then $(\text{proj}_{G(\mathcal{R})} h_n)$ is a monotone increasing sequence in $L_+^2(\nu)$ converging to $\text{proj}_{G(\mathcal{R})} g$. Convergence is both in the sense of L^2 -norm and almost surely.*

Proof. Let h_n be a sequence of ordered simple function such that $h_n \rightarrow g$ almost surely and in L^2 -norm. Continuity of $\text{proj}_{G(\mathcal{R})}$ yields

$$\text{proj}_{G(\mathcal{R})} h_n \rightarrow \text{proj}_{G(\mathcal{R})} g$$

in L^2 -norm. By passing to a subsequence if necessary, we may assume that this convergence also holds in the almost sure sense. \square

Proposition 3.8. *Let $A \in \mathcal{R}$. Then*

$$\text{proj}_{G(\mathcal{R})}(g\mathbb{1}_A) = \text{proj}_{G(\mathcal{R})}(g\mathbb{1}_A) \cdot \mathbb{1}_A \leq (\text{proj}_{G(\mathcal{R})} g) \cdot \mathbb{1}_A.$$

Proof. We have

$$\begin{aligned} & \int (g\mathbb{1}_A - \text{proj}_{G(\mathcal{R})}(g\mathbb{1}_A))^2 d\nu \\ &= \int_A (g\mathbb{1}_A - \text{proj}_{G(\mathcal{R})}(g\mathbb{1}_A))^2 d\nu + \int_{A^c} (g\mathbb{1}_A - \text{proj}_{G(\mathcal{R})}(g\mathbb{1}_A))^2 d\nu \\ &= \int_A (g\mathbb{1}_A - \text{proj}_{G(\mathcal{R})}(g\mathbb{1}_A))^2 d\nu + \int_{A^c} (\text{proj}_{G(\mathcal{R})}(g\mathbb{1}_A))^2 d\nu \\ &\geq \int_A (g\mathbb{1}_A - \text{proj}_{G(\mathcal{R})}(g\mathbb{1}_A))^2 d\nu \\ &= \int (g\mathbb{1}_A - \text{proj}_{G(\mathcal{R})}(g\mathbb{1}_A) \cdot \mathbb{1}_A)^2 d\nu. \end{aligned}$$

Since the projection is unique, we have $\text{proj}_{G(\mathcal{R})}(g\mathbb{1}_A) = \text{proj}_{G(\mathcal{R})}(g\mathbb{1}_A) \cdot \mathbb{1}_A$. As $g\mathbb{1}_A \leq g$, the monotonicity in Proposition 3.6 yields

$$\text{proj}_{G(\mathcal{R})}(g\mathbb{1}_A) = \text{proj}_{G(\mathcal{R})}(g\mathbb{1}_A) \cdot \mathbb{1}_A \leq (\text{proj}_{G(\mathcal{R})} g) \cdot \mathbb{1}_A. \quad \square$$

Let $\mathcal{C} \subseteq \mathcal{B}$ be a σ -subalgebra. The Radon–Nikodym theorem states that for σ -finite measures μ and ν on the measurable space $(\mathbb{R}, \mathcal{B})$ with $\mu \ll \nu$, there exists a \mathcal{C} -measurable function $\rho_{\mathcal{C}}$ such that $\mu(A) = \int_A \rho_{\mathcal{C}} d\nu$ for any $A \in \mathcal{C}$.

What happens if \mathcal{C} is not a σ -subalgebra? For us, the most important case to consider is when $\mathcal{C} = \mathcal{R}$, which is not a σ -subalgebra. We will show in Corollary 3.11 that \mathcal{R} satisfies a kind of Radon–Nikodym property: There is a \mathcal{R} -measurable function $\rho_{\mathcal{R}}$ and a subcollection $\mathcal{E} \subseteq \mathcal{R}$ such that $\mu(A) = \int_A \rho_{\mathcal{R}} d\nu$ for any $A \in \mathcal{E}$. The bulk of the work is in establishing the following technical result.

Theorem 3.9. *Let $g \in L_+^2(\nu)$. Define the π -system $\mathcal{E}(g) \subseteq \mathcal{R}$ by*

$$(7) \quad \mathcal{E}(g) := \{\{\text{proj}_{G(\mathcal{R})} g > r\} : r \geq 0\} \cup \{\mathbb{R}\}.$$

Then for any $E \in \mathcal{E}(g)$,

$$(8) \quad \int_E \text{proj}_{G(\mathcal{R})} g d\nu = \int_E g d\nu.$$

Proof. We establish the result in four steps, showing that

- (i) it holds when g is an ordered simple function h ;
- (ii) it holds for $E = \{\text{proj}_{G(\mathcal{R})} g > r\}$ under the assumption that $\nu(\{\text{proj}_{G(\mathcal{R})} g = r\}) = 0$;
- (iii) it extends to all $E = \{\text{proj}_{G(\mathcal{R})} g > r\}$;
- (iv) it holds for $E = \mathbb{R}$.

STEP (i): First we assume that we have an ordered simple function

$$h = \sum_{i=1}^n \alpha_i \mathbb{1}_{A_i}, \quad \text{proj}_{G(\mathcal{R})} h = \sum_{i=1}^n \beta_i \mathbb{1}_{A_i}.$$

So $E = \{\text{proj}_{G(\mathcal{R})} h > r\}$ can be written as $E = A_1 \cup \dots \cup A_k \in \mathcal{E}(h)$ where $\beta_k > r$ and $\beta_{k+1} \leq r$. Write $\alpha = (\alpha_1, \dots, \alpha_n)$, $\beta = (\beta_1, \dots, \beta_n)$, $\omega = (\omega_1, \dots, \omega_n)$ where $\omega_i = \nu(A_i)$, and let $W = \text{diag}(\omega)$. Then β is the solution to the following quadratic programming problem:

$$\begin{aligned} & \text{minimize} && (\alpha - \beta)^\top W (\alpha - \beta) \\ & \text{subject to} && \beta_1 \geq \dots \geq \beta_n \geq 0. \end{aligned}$$

Observe that for

$$\Pi = \begin{bmatrix} 1 & -1 & & \\ & 1 & -1 & \\ & & \ddots & -1 \\ & & & 1 \end{bmatrix}, \quad \Pi^{-1} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ & 1 & \dots & 1 \\ & & \ddots & \vdots \\ & & & 1 \end{bmatrix}.$$

Let $\gamma = \Pi\beta$. Then the quadratic programming problem above is equivalent to

$$\begin{aligned} & \text{minimize} && \gamma^\top \Pi^{-\top} W \Pi^{-1} \gamma - 2\alpha^\top \Pi^{-1} W \gamma \\ & \text{subject to} && -\gamma \preceq 0. \end{aligned}$$

Standard KKT conditions yield

$$\lambda = 2\Pi^{-\top} W (\Pi^{-1} \gamma - \alpha) = 2\Pi^{-\top} W (\beta - \alpha)$$

or

$$(9) \quad \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \vdots \\ \lambda_n \end{bmatrix} = 2 \begin{bmatrix} 1 & & & \\ 1 & 1 & & \\ \vdots & \vdots & \ddots & \\ 1 & 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} \omega_1 & & & \\ & \omega_2 & & \\ & & \ddots & \\ & & & \omega_n \end{bmatrix} \begin{bmatrix} \beta_1 - \alpha_1 \\ \beta_2 - \alpha_2 \\ \vdots \\ \beta_n - \alpha_n \end{bmatrix}.$$

Since $\beta_k = (\Pi^{-1}\gamma)_k = \gamma_k + \gamma_{k+1} + \dots + \gamma_n > r$ and $\beta_{k+1} = \gamma_{k+1} + \gamma_{k+2} + \dots + \gamma_n \leq r$, we must have $\gamma_k > 0$. By the KKT complementary slackness condition,

$$(10) \quad \lambda_k = 2(\Pi^{-\top} W (\beta - \alpha))_k = 2 \sum_{i=1}^k \omega_i (\beta_i - \alpha_i) = 0.$$

Hence we have

$$\int (\text{proj}_{G(\mathbb{R})} h) \cdot \mathbb{1}_E \, d\nu = \sum_{i=1}^k \beta_i \omega_i = \sum_{i=1}^k \alpha_i \omega_i = \int h \mathbb{1}_E \, d\nu.$$

STEP (ii): Next we drop the ordered simple assumption and just require that $g \in L^2_+(\nu)$ be a general function with

$$(11) \quad \nu(\{\text{proj}_{G(\mathbb{R})} g = r\}) = 0.$$

Let (h_n) be a sequence of nonnegative ordered simple functions with $h_n \rightarrow g$ almost surely and in L^2 -norm. So $\text{proj}_{G(\mathbb{R})} h_n \rightarrow \text{proj}_{G(\mathbb{R})} g$ in L^2 -norm and thus in measure. By passing through a subsequence if necessary, we may assume that $\text{proj}_{G(\mathbb{R})} h_n \rightarrow \text{proj}_{G(\mathbb{R})} g$ almost surely and in L^2 -norm. Set $E = \{\text{proj}_{G(\mathbb{R})} g > r\}$ and $E_n = \{\text{proj}_{G(\mathbb{R})} h_n > r\}$, bearing in mind that we assume (11). Then

$$\begin{aligned} \|\mathbb{1}_E - \mathbb{1}_{E_n}\|^2 &= \int (\mathbb{1}_E - \mathbb{1}_{E_n})^2 \, d\nu \\ &= \nu(E) + \nu(E_n) - 2\nu(E \cap E_n) = \nu(E \Delta E_n) = \nu(E \cap E_n^c) + \nu(E_n \cap E^c). \end{aligned}$$

Consider the null set $N = \{\text{proj}_{G(\mathbb{R})} h_n \not\rightarrow \text{proj}_{G(\mathbb{R})} g\}$. If

$$x \in E \cap E_n^c = \{\text{proj}_{G(\mathbb{R})} g > r\} \cap \{\text{proj}_{G(\mathbb{R})} h_n \leq r\}$$

infinitely often, then

$$\liminf_{n \rightarrow \infty} \text{proj}_{G(\mathbb{R})} h_n(x) \leq r < \text{proj}_{G(\mathbb{R})} g(x),$$

and so $x \in N$. Consequently, $\limsup_{n \rightarrow \infty} E \cap E_n^c \subseteq N$. On the other hand, we may decompose

$$\begin{aligned} E_n \cap E^c &= [\{\text{proj}_{G(\mathbb{R})} h_n > r\} \cap \{\text{proj}_{G(\mathbb{R})} g < r\}] \cup [\{\text{proj}_{G(\mathbb{R})} h_n > r\} \cap \{\text{proj}_{G(\mathbb{R})} g = r\}] \\ &=: F_n \cup G_n. \end{aligned}$$

Since G_n is a null set for each n , if $x \in F_n$ infinitely often, then

$$\limsup_{n \rightarrow \infty} \text{proj}_{G(\mathbb{R})} h_n(x) \geq r > \text{proj}_{G(\mathbb{R})} g(x),$$

and so $x \in N$ as well. Consequently, $\limsup_{n \rightarrow \infty} F_n \subseteq N$. We obtain

$$\begin{aligned} \limsup_{n \rightarrow \infty} \|\mathbb{1}_E - \mathbb{1}_{E_n}\|^2 &\leq \limsup_{n \rightarrow \infty} [\nu(E \cap E_n^c) + \nu(F_n) + \nu(G_n)] \\ &\leq \limsup_{n \rightarrow \infty} \nu(E \cap E_n^c) + \limsup_{n \rightarrow \infty} \nu(F_n) \\ &\leq \nu\left(\limsup_{n \rightarrow \infty} E \cap E_n^c\right) + \nu\left(\limsup_{n \rightarrow \infty} F_n\right) = 0, \end{aligned}$$

and by Hölder,

$$\begin{aligned} \|h\mathbb{1}_E - h_n\mathbb{1}_{E_n}\|_1 &= \|h\mathbb{1}_E - h\mathbb{1}_{E_n} + h\mathbb{1}_{E_n} - h_n\mathbb{1}_{E_n}\|_1 \\ &\leq \|h \cdot (\mathbb{1}_E - \mathbb{1}_{E_n})\|_1 + \|(h - h_n) \cdot \mathbb{1}_{E_n}\|_1 \\ &\leq \|h\| \|\mathbb{1}_E - \mathbb{1}_{E_n}\| + \|h - h_n\| \|\mathbb{1}_{E_n}\| \rightarrow 0. \end{aligned}$$

The same argument also yields

$$(\text{proj}_{G(\mathbb{R})} h_n) \cdot \mathbb{1}_{E_n} \rightarrow (\text{proj}_{G(\mathbb{R})} g) \cdot \mathbb{1}_E$$

in L^1 -norm. By Lemma 2.2, we get

$$\int g\mathbb{1}_E d\nu = \lim_{n \rightarrow \infty} \int h_n\mathbb{1}_{E_n} d\nu = \lim_{n \rightarrow \infty} \int (\text{proj}_{G(\mathbb{R})} h_n) \cdot \mathbb{1}_{E_n} d\nu = \int (\text{proj}_{G(\mathbb{R})} g) \cdot \mathbb{1}_E d\nu$$

as required.

STEP (iii): We next omit the condition (11). Let $x = \inf\{\text{proj}_{G(\mathbb{R})} g = r\}$ and assume $x > -\infty$ for nontriviality. For any $y < x$, we have $\text{proj}_{G(\mathbb{R})} g(y) \geq r$ and so

$$(12) \quad \inf_{y \in (-\infty, x)} \text{proj}_{G(\mathbb{R})} g(y) \geq r.$$

If we have strict inequality in (12), then there exists $\varepsilon > 0$ such that $\inf_{y \in (-\infty, x)} \text{proj}_{G(\mathbb{R})} g(y) > r + \varepsilon$. Thus

$$E_r := \{\text{proj}_{G(\mathbb{R})} g > r\} = \{\text{proj}_{G(\mathbb{R})} g > r + \varepsilon\} =: E_{r+\varepsilon}.$$

Since $\text{proj}_{G(\mathbb{R})} g \neq r + \varepsilon$ almost surely, $\nu(\{\text{proj}_{G(\mathbb{R})} g = r + \varepsilon\}) = 0$, and

$$\int_{E_r} \text{proj}_{G(\mathbb{R})} g d\nu = \int_{E_{r+\varepsilon}} \text{proj}_{G(\mathbb{R})} g d\nu = \int_{E_{r+\varepsilon}} g d\nu = \int_{E_r} g d\nu.$$

If we have equality in (12), then for any $k \in \mathbb{N}$, we may pick $x_k \in (-\infty, x)$ such that

$$r < \text{proj}_{G(\mathbb{R})} g(x_k) < r + \frac{1}{k}$$

and $r_k \in (r, \text{proj}_{G(\mathbb{R})} g(x_k))$ such that

$$r < r_k < r + \frac{1}{k} \quad \text{and} \quad \nu(\text{proj}_{G(\mathbb{R})} g = r_k) = 0.$$

Hence $E_k := \{\text{proj}_{G(\mathbb{R})} g > r_k\}$, $k \in \mathbb{N}$, is an increasing sequence of nested sets with $E_k \subseteq E_r$ and so $\bigcup_{k=1}^{\infty} E_k \subseteq E_r$. Conversely, if $y \in E_r$, i.e., $\text{proj}_{G(\mathbb{R})} g(y) > r$, then there exists r_k such that $r_k < \text{proj}_{G(\mathbb{R})} g(y)$; so $y \in E_k$; and so $E_k \uparrow E_r$. Define the finite measures

$$\zeta(A) := \int_A \text{proj}_{G(\mathbb{R})} g d\nu, \quad \xi(A) := \int_A g d\nu$$

for any $A \in \mathcal{B}$. By the continuity of measure from below,

$$\zeta(E_r) = \lim_{k \rightarrow \infty} \zeta(E_k) = \lim_{k \rightarrow \infty} \xi(E_k) = \xi(E_r)$$

as required.

STEP (iv): It remains to treat the case $E = \mathbb{R}$. We will show that $\int \text{proj}_{G(\mathcal{R})} g \, d\nu = \int g \, d\nu$ by first establishing it for an ordered simple function $h = \sum_{i=1}^n \alpha_i \mathbb{1}_{A_i}$. We may assume that $\beta_k > 0$ and $\beta_{k+1} = \beta_{k+2} = \dots = \beta_n = 0$ without loss of generality. By (9) and Step (i),

$$\begin{aligned} \lambda_n &= 2 \sum_{i=1}^n \omega_i(\beta_i - \alpha_i) = 2 \sum_{i=1}^k \omega_i(\beta_i - \alpha_i) + 2 \sum_{j=k+1}^n \omega_j(\beta_j - \alpha_j) \\ &= -\lambda_k - 2(\omega_{k+1}\alpha_{k+1} + \omega_{k+2}\alpha_{k+2} + \dots + \omega_n\alpha_n). \end{aligned}$$

By (10), $\lambda_k = 0$. By KKT dual feasibility, $\lambda_n \geq 0$. Since we also have $\omega_i, \alpha_i \geq 0$, $i = k+1, k+2, \dots, n$,

$$\lambda_n = -2(\omega_{k+1}\alpha_{k+1} + \omega_{k+2}\alpha_{k+2} + \dots + \omega_n\alpha_n) = 0.$$

It follows from (9) that $0 = \lambda_n = 2 \sum_{i=1}^n \omega_i(\beta_i - \alpha_i)$ and so

$$\int \text{proj}_{G(\mathcal{R})} h \, d\nu = \sum_{i=1}^n \beta_i \omega_i = \sum_{i=1}^n \alpha_i \omega_i = \int h \, d\nu.$$

For a general function $g \in L_+^2(\nu)$, pick a sequence of nonnegative ordered simple functions $h_n \rightarrow g$ in L^2 -norm and apply Lemma 2.1 to get

$$\int g \, d\nu = \lim_{n \rightarrow \infty} \int h_n \, d\nu = \lim_{n \rightarrow \infty} \int \text{proj}_{G(\mathcal{R})} h_n \, d\nu = \int \text{proj}_{G(\mathcal{R})} g \, d\nu. \quad \square$$

We state an immediate corollary of Theorem 3.9.

Corollary 3.10. *Let $g \in L_+^2(\nu)$. If $\text{proj}_{G(\mathcal{R})} g(E) \in \mathcal{B}$, then*

$$\int_E \text{proj}_{G(\mathcal{R})} g \, d\nu = \int_E g \, d\nu.$$

Proof. Theorem 3.9 says that (8) holds for any $E \in \mathcal{E}(g)$. The left and right sides of (8) define finite measures ζ and ξ respectively, and $\zeta(\mathbb{R}) = \xi(\mathbb{R})$ as in Step (iii) of the proof of Theorem 3.9. Since $\mathcal{E}(g)$ is a π -system, $\sigma(\mathcal{E}(g)) = (\text{proj}_{G(\mathcal{R})} g)^{-1}(\mathcal{B})$ by Lemmas 2.3 and 2.4. Hence (8) holds for all $E \in (\text{proj}_{G(\mathcal{R})} g)^{-1}(\mathcal{B})$. \square

We now arrive at the Radon–Nikodym property for \mathcal{R} mentioned earlier.

Corollary 3.11 (Radon–Nikodym over rays). *Let μ, ν be Borel probability measures on $(\mathbb{R}, \mathcal{B})$. Then the \mathcal{R} -measurable function*

$$\rho_{\mathcal{R}} = \text{proj}_{G(\mathcal{R})} \frac{d\mu}{d\nu},$$

and the π -system $\mathcal{E}(d\mu/d\nu) = \{\{\rho_{\mathcal{R}} > r\} : r \geq 0\} \cup \{\mathbb{R}\} \subseteq \mathcal{R}$ satisfy

$$\mu(E) = \int_E \rho_{\mathcal{R}} \, d\nu$$

for any $E \in \mathcal{E}(d\mu/d\nu)$.

Proof. For any $E \in \mathcal{E}(d\mu/d\nu)$, it follows from Theorem 3.9 that

$$\mu(E) = \int_E \frac{d\mu}{d\nu} \, d\nu = \int_E \text{proj}_{G(\mathcal{R})} \frac{d\mu}{d\nu} \, d\nu = \int_E \rho_{\mathcal{R}} \, d\nu. \quad \square$$

We will show that the following definition of f -divergence over \mathcal{R} provides an affirmative answer to Question 1: Theorem 3.15 shows that (a) when f is chosen to be $f(t) = |t - 1|/2$, we recover the Kolmogorov–Smirnov distance in (3); and Definition 3.12 shows that by replacing \mathcal{R} by \mathcal{B} , we recover the standard f -divergence in (5). A reminder of our convention in Section 2.3: we assume

that $\mu \ll \nu$ whenever we write $d\mu/d\nu$. This is not an additional imposition; it is also a requirement for standard f -divergence in Definition 1.5.

Definition 3.12 (f -divergence over rays). Let μ and ν be Borel probability measures and \mathcal{R} be the class of rays. Then the f -divergence over \mathcal{R} is defined as

$$(13) \quad D_f^{\mathcal{R}}(\mu \parallel \nu) := \int f\left(\text{proj}_{G(\mathcal{R})} \frac{d\mu}{d\nu}\right) d\nu.$$

Henceforth we will use this terminology in conjunction with standard nomenclatures. For example when we speak of total variation or Hellinger distance over \mathcal{R} we mean the f -divergence over \mathcal{R} with $f(t) = |t-1|/2$ or $f(t) = (\sqrt{t}-1)^2$ respectively. Obviously, if we replace \mathcal{R} in (13) by the Borel σ -algebra \mathcal{B} , then $D_f^{\mathcal{B}}(\mu \parallel \nu) = D_f(\mu \parallel \nu)$ and we recover the standard f -divergence in Definition 1.5. So we may view the standard definition as “ f -divergence over \mathcal{B} .”

We will provide abundant evidence in Section 4 that our f -divergence over \mathcal{R} in Definition 3.12 retains almost all known properties of the standard f -divergence. But our most rudimentary justification is Theorem 3.15, i.e., the total variation distance on \mathcal{R} is exactly the Kolmogorov–Smirnov distance. To get there we need two corollaries.

Corollary 3.13. *Let $g \in L_+^2(\nu)$. Then for any $A \in \mathcal{R}$,*

$$\int_A \text{proj}_{G(\mathcal{R})} g \, d\nu \geq \int_A g \, d\nu.$$

Proof. By Proposition 3.8 and Theorem 3.9, we have

$$\int_A \text{proj}_{G(\mathcal{R})} g \, d\nu = \int (\text{proj}_{G(\mathcal{R})} g) \cdot \mathbb{1}_A \, d\nu \geq \int \text{proj}_{G(\mathcal{R})}(g \mathbb{1}_A) \, d\nu = \int_A g \, d\nu. \quad \square$$

The π -system $\mathcal{E}(g)$ as defined in (7) has the following property.

Corollary 3.14. *Let $g \in L_+^2(\nu)$. Then for any $E \in \mathcal{E}(g)$,*

$$(\text{proj}_{G(\mathcal{R})} g) \cdot \mathbb{1}_E = \text{proj}_{G(\mathcal{R})}(g \mathbb{1}_E).$$

Proof. By Theorem 3.9,

$$\int (\text{proj}_{G(\mathcal{R})} g) \cdot \mathbb{1}_E \, d\nu = \int g \mathbb{1}_E \, d\nu = \int \text{proj}_{G(\mathcal{R})}(g \mathbb{1}_E) \, d\nu.$$

Since $(\text{proj}_{G(\mathcal{R})} g) \cdot \mathbb{1}_E \geq \text{proj}_{G(\mathcal{R})}(g \mathbb{1}_E)$, we must in fact have equality. \square

We now arrive at the main result of this section. The key to Theorem 3.15 is Theorem 3.9. If we have an analogue of Theorem 3.9 for a class of sets \mathcal{C} , then we have an “ f -divergence on \mathcal{C} ” for which an analogue of Theorem 3.15 with \mathcal{C} in place of \mathcal{R} holds.

Theorem 3.15 (Kolmogorov–Smirnov as special case). *Let μ, ν be Borel probability measures on $(\mathbb{R}, \mathcal{B})$. Then*

$$(14) \quad \int \frac{1}{2} \left| \text{proj}_{G(\mathcal{R})} \frac{d\mu}{d\nu} - 1 \right| d\nu = \sup_{A \in \mathcal{R}} (\mu(A) - \nu(A)).$$

Proof. Let $E = \{\text{proj}_{G(\mathcal{R})} d\mu/d\nu > 1\}$. By Theorem 3.9 and Corollary 3.11,

$$\int_A \text{proj}_{G(\mathcal{R})} \frac{d\mu}{d\nu} \, d\nu = \int_A \frac{d\mu}{d\nu} \, d\nu = \mu(A)$$

for any $A \in \mathcal{E}(d\mu/d\nu)$. So

$$\begin{aligned}
\int \frac{1}{2} \left| \text{proj}_{G(\mathcal{R})} \frac{d\mu}{d\nu} - 1 \right| d\nu &= \frac{1}{2} \int_E \text{proj}_{G(\mathcal{R})} \frac{d\mu}{d\nu} - 1 d\nu + \frac{1}{2} \int_{E^c} 1 - \text{proj}_{G(\mathcal{R})} \frac{d\mu}{d\nu} d\nu \\
&= \frac{1}{2} \left[2 \int_E \text{proj}_{G(\mathcal{R})} \frac{d\mu}{d\nu} d\nu - 1 \right] - \frac{1}{2} \left[2 \int_E d\nu - 1 \right] \\
(15) \qquad &= \int_E \text{proj}_{G(\mathcal{R})} \frac{d\mu}{d\nu} d\nu - \nu(E) \\
&= \mu(E) - \nu(E) \leq \sup_{A \in \mathcal{R}} (\mu(A) - \nu(A)),
\end{aligned}$$

noting that $E = \{\text{proj}_{G(\mathcal{R})} d\mu/d\nu > 1\} \in \mathcal{R}$ as $\text{proj}_{G(\mathcal{R})} d\mu/d\nu \in G(\mathcal{R})$. For the reverse inequality, by Corollary 3.13, we have

$$\int_A \left(\text{proj}_{G(\mathcal{R})} \frac{d\mu}{d\nu} - 1 \right) d\nu \geq \mu(A) - \nu(A),$$

for any $A \in \mathcal{R}$. By (15),

$$\begin{aligned}
\int \frac{1}{2} \left| \text{proj}_{G(\mathcal{R})} \frac{d\mu}{d\nu} - 1 \right| d\nu &= \int_E \left(\text{proj}_{G(\mathcal{R})} \frac{d\mu}{d\nu} - 1 \right) d\nu \\
&= \sup_{A \in \mathcal{R}} \left(\int_A \text{proj}_{G(\mathcal{R})} \frac{d\mu}{d\nu} d\nu - \nu(A) \right) \geq \sup_{A \in \mathcal{R}} (\mu(A) - \nu(A))
\end{aligned}$$

as required. \square

4. PROPERTIES OF THE f -DIVERGENCES OVER RAYS

Theorem 3.15 serves as an indication that our notion of f -divergence over rays in Definition 3.12 is the right one because it restricts to the Kolmogorov–Smirnov distance for a special choice of f . In this section, we will see that just about every property that we know holds for f -divergences also holds for f -divergences over \mathcal{R} . Henceforth we will denote the convex cone of functions convex on $[0, \infty)$ and both vanishing and strictly convex at 1 by

$$(16) \qquad \mathcal{F} := \{f : [0, \infty) \rightarrow \mathbb{R} : f \text{ is convex, } f \text{ is strictly convex at 1, and } f(1) = 0\}.$$

The strict convexity at 1 is as in [40, p. 120]: For all $x, y \in [0, \infty)$ and $t \in (0, 1)$ such that $tx + (1-t)y = 1$, we have $tf(x) + (1-t)f(y) > 0$. This condition guarantees that $D_f(\mu\|\nu) = 0$ if and only if $\mu = \nu$. As a precaution, in situations like the Kullback–Leibler divergence where $f(x) = x \log x$, the value $f(0)$ may only be defined in a limiting sense.

Again μ, ν will denote Borel probability measures on $(\mathbb{R}, \mathcal{B})$ throughout this section, with the implicit assumption that $\mu \ll \nu$ whenever we speak of their Radon–Nikodym derivative $d\mu/d\nu$.

Theorem 4.1. *Let $f, g \in \mathcal{F}$. Then the f -divergence over \mathcal{R} has the following properties:*

- (i) *Linearity:* $D_{\alpha f + \beta g}^{\mathcal{R}}(\mu\|\nu) = \alpha D_f^{\mathcal{R}}(\mu\|\nu) + \beta D_g^{\mathcal{R}}(\mu\|\nu)$ for any $\alpha, \beta \geq 0$.
- (ii) *Nonnegativity:* $D_f^{\mathcal{R}}(\mu\|\nu) \geq 0$.
- (iii) *Affine invariance:* If $g(x) = f(x) + c(x-1)$, $c \in \mathbb{R}$, then $D_f^{\mathcal{R}}(\mu\|\nu) = D_g^{\mathcal{R}}(\mu\|\nu)$.
- (iv) *Boundedness:* $D_f^{\mathcal{R}}(\mu\|\nu) \leq D_f(\mu\|\nu)$.
- (v) *Identity:* If $D_f^{\mathcal{R}}(\mu\|\nu) = 0 = D_f^{\mathcal{R}}(\nu\|\mu)$, then $\mu = \nu$.

Proof. Item (i) is routine from definition. For item (ii), since f is convex, applying Jensen’s inequality and Theorem 3.9 gives

$$\int f \left(\text{proj}_{G(\mathcal{R})} \frac{d\mu}{d\nu} \right) d\nu \geq f \left(\int \text{proj}_{G(\mathcal{R})} \frac{d\mu}{d\nu} d\nu \right) = f \left(\int \frac{d\mu}{d\nu} d\nu \right) = f(1) = 0.$$

Item (iii) follows from

$$\begin{aligned} D_g^{\mathcal{R}}(\mu\|\nu) &= \int f\left(\text{proj}_{G(\mathcal{R})} \frac{d\mu}{d\nu}\right) + c\left(\text{proj}_{G(\mathcal{R})} \frac{d\mu}{d\nu} - 1\right) d\nu \\ &= \int f\left(\text{proj}_{G(\mathcal{R})} \frac{d\mu}{d\nu}\right) d\nu + c \int \left(\text{proj}_{G(\mathcal{R})} \frac{d\mu}{d\nu} - 1\right) d\nu = D_f^{\mathcal{R}}(\mu\|\nu), \end{aligned}$$

noting that the last integral vanishes. The proofs for items (iv) and (v) are considerably more involved and will be deferred to Proposition 4.3 and Corollary 4.7 respectively. \square

We caution the reader that special properties of a specific f -divergence may sometimes be lost for its counterpart over \mathcal{R} . For example, it is well-known that while f -divergence is generally not symmetric, i.e., $D_f(\mu\|\nu) \neq D_f(\nu\|\mu)$, the total variation distance is. This symmetry is lost for the total variation distance over \mathcal{R} as

$$D_{\text{TV}}^{\mathcal{R}}(\mu\|\nu) = \int \frac{1}{2} \left| \text{proj}_{G(\mathcal{R})} \frac{d\mu}{d\nu} - 1 \right| d\nu = \sup_{A \in \mathcal{R}} (\mu(A) - \nu(A));$$

although one may recover it by symmetrizing:

$$(17) \quad \widehat{D}_{\text{TV}}^{\mathcal{R}}(\mu\|\nu) := \max\{D_{\text{TV}}^{\mathcal{R}}(\mu\|\nu), D_{\text{TV}}^{\mathcal{R}}(\nu\|\mu)\} = \sup_{A \in \mathcal{R}} |\mu(A) - \nu(A)|.$$

Strictly speaking, it is this symmetrized version $\widehat{D}_{\text{TV}}^{\mathcal{R}}$ that should be called the Kolmogorov–Smirnov *distance* [30] although we have been loosely using the term for the *divergence* $D_{\text{TV}}^{\mathcal{R}}$ as well.

This is also the reason why we need both $D_f^{\mathcal{R}}(\mu\|\nu) = 0$ and $D_f^{\mathcal{R}}(\nu\|\mu) = 0$ in Theorem 4.1(v). It would not be true otherwise. We illustrate this in Figure 1 with the total variation and Hellinger distances over \mathcal{R} . One can see that $D_f^{\mathcal{R}}(\mu\|\nu) = 0$ or $D_f^{\mathcal{R}}(\nu\|\mu) = 0$ does not necessarily imply $\mu = \nu$ but that $D_f^{\mathcal{R}}(\mu\|\nu) = D_f^{\mathcal{R}}(\nu\|\mu) = 0$ does.

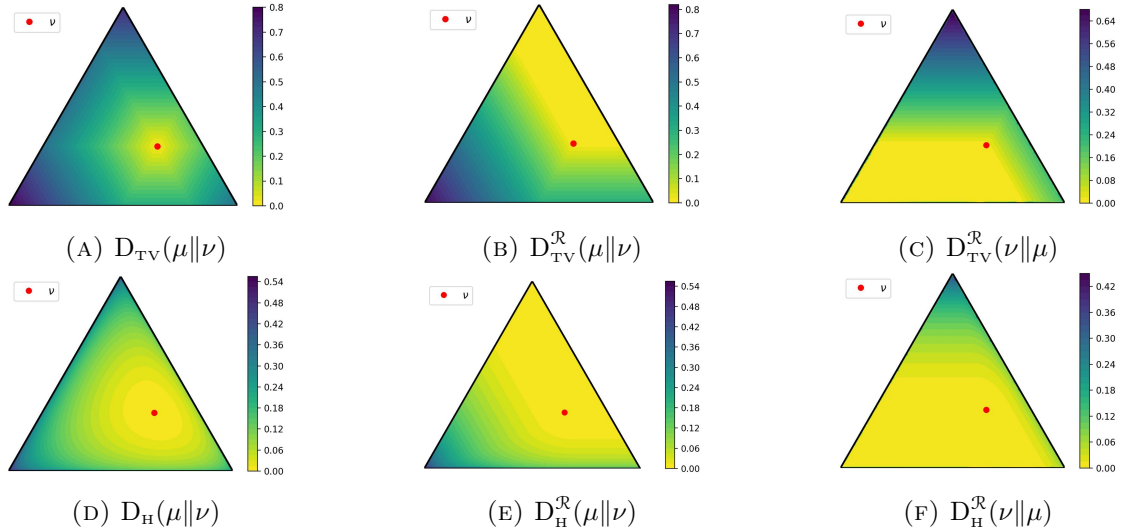


FIGURE 1. 40 Level curves of total variation D_{TV} , total variation over rays $D_{\text{TV}}^{\mathcal{R}}$, Hellinger distance D_{H} and Hellinger distance over rays $D_{\text{H}}^{\mathcal{R}}$ for fixed $\nu = [0.2, 0.5, 0.3]$ as μ ranges over the simplex of distributions on a three-element set.

Evidently, the projection operator $\text{proj}_{G(\mathcal{R})}$ plays a key role in our definition of f -divergence over \mathcal{R} . The next two results show how it interacts with nondecreasing functions and with convex functions.

Lemma 4.2. *Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be nondecreasing. If $g \in L^2_+(\nu)$ is almost surely continuous, then*

$$\int (\text{proj}_{G(\mathbb{R})} g) \cdot (f \circ \text{proj}_{G(\mathbb{R})} g) \, d\nu = \int g \cdot (f \circ \text{proj}_{G(\mathbb{R})} g) \, d\nu.$$

Proof. Let $h = \sum_{i=1}^n \alpha_i \mathbb{1}_{A_i}$ and $\text{proj}_{G(\mathbb{R})} h = \sum_{i=1}^n \beta_i \mathbb{1}_{A_i}$. Suppose $\text{proj}_{G(\mathbb{R})} h$ takes m distinct values $\tau_1 > \dots > \tau_m$ where $m \leq n$. For each $k = 1, \dots, m$, set

$$E_k := \bigcup_{i=1}^k \{\text{proj}_{G(\mathbb{R})} h = \tau_i\} \in \mathcal{E}(h).$$

Let $\zeta_m := f(\tau_m)$ and $\zeta_{m-j} := f(\tau_{m-j}) - \zeta_{m-j+1} - \zeta_{m-j+2} - \dots - \zeta_m$, $j = 1, \dots, m-1$. Then

$$\begin{aligned} \int (\text{proj}_{G(\mathbb{R})} h) \cdot (f \circ \text{proj}_{G(\mathbb{R})} h) \, d\nu &= \int (\text{proj}_{G(\mathbb{R})} h) \cdot \left(\sum_{i=1}^n f(\beta_i) \mathbb{1}_{A_i} \right) \, d\nu \\ &= \int (\text{proj}_{G(\mathbb{R})} h) \cdot \sum_{i=1}^m \zeta_i \mathbb{1}_{E_i} \, d\nu \\ &= \sum_{i=1}^m \zeta_i \int_{E_i} \text{proj}_{G(\mathbb{R})} h \, d\nu = \sum_{i=1}^m \zeta_i \int_{E_i} h \, d\nu \\ &= \int h \cdot \sum_{i=1}^m \zeta_i \mathbb{1}_{E_i} \, d\nu = \int h \cdot (f \circ \text{proj}_{G(\mathbb{R})} h) \, d\nu. \end{aligned}$$

Now let h_n be a nondecreasing sequence of nonnegative ordered simple function that converges to g almost surely and in L^2 -norm. By passing through a subsequence if necessary, $\text{proj}_{G(\mathbb{R})} h_n$ is a nonincreasing sequence that converges to $\text{proj}_{G(\mathbb{R})} g$ almost surely. Since f is nondecreasing,

$$f \circ \text{proj}_{G(\mathbb{R})} h_n \uparrow f \circ \text{proj}_{G(\mathbb{R})} g$$

almost surely. By the monotone convergence theorem,

$$\begin{aligned} \int (\text{proj}_{G(\mathbb{R})} g) \cdot (f \circ \text{proj}_{G(\mathbb{R})} g) \, d\nu &= \lim_{n \rightarrow \infty} \int (\text{proj}_{G(\mathbb{R})} h_n) \cdot (f \circ \text{proj}_{G(\mathbb{R})} h_n) \, d\nu \\ &= \lim_{n \rightarrow \infty} \int h_n \cdot (f \circ \text{proj}_{G(\mathbb{R})} h_n) \, d\nu = \int g \cdot (f \circ \text{proj}_{G(\mathbb{R})} g) \, d\nu. \quad \square \end{aligned}$$

Proposition 4.3. *Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be convex. If $g \in L^2_+(\nu)$ is almost surely continuous, then*

$$\int f \circ \text{proj}_{G(\mathbb{R})} g \, d\nu \leq \int f \circ g \, d\nu.$$

Proof. By convexity,

$$f(\text{proj}_{G(\mathbb{R})} g(x)) - f(g(x)) \leq f'(\text{proj}_{G(\mathbb{R})} g(x)) \cdot (\text{proj}_{G(\mathbb{R})} g(x) - g(x))$$

almost surely. Integrating with respect to ν and applying Lemma 4.2 to f' , we get

$$\int f \circ \text{proj}_{G(\mathbb{R})} g - f \circ g \, d\nu \leq \int f'(\text{proj}_{G(\mathbb{R})} g) \cdot (\text{proj}_{G(\mathbb{R})} g - g) \, d\nu = 0. \quad \square$$

Theorem 4.1(iv) follows from Proposition 4.3 since any $f \in \mathcal{F}$ is convex and so

$$D_f^{\mathbb{R}}(\nu \parallel \mu) = \int f\left(\text{proj}_{G(\mathbb{R})} \frac{d\mu}{d\nu}\right) \, d\nu \leq \int f\left(\frac{d\mu}{d\nu}\right) \, d\nu = D_f(\nu \parallel \mu).$$

It remains to establish Theorem 4.1(v), whose proof will take four intermediate results: Propositions 4.4 and 4.5, Corollaries 4.6 and 4.7, each of independent interest. We begin by showing how we may convert f -divergence over rays to standard f -divergence.

Proposition 4.4. *For any Borel probability measures μ and ν , there exists a Borel probability measure ζ such that for any $f \in \mathcal{F}$,*

$$D_f^{\mathcal{R}}(\mu\|\nu) = D_f(\zeta\|\nu).$$

Proof. We define ζ by

$$\zeta(B) := \int_B \text{proj}_{G(\mathbb{R})} \frac{d\mu}{d\nu} d\nu$$

for any $B \in \mathcal{B}$. By Corollary 3.11, ζ is a probability measure absolutely continuous with respect to ν . Hence

$$\frac{d\zeta}{d\nu} = \text{proj}_{G(\mathbb{R})} \frac{d\mu}{d\nu}$$

and so

$$D_f^{\mathcal{R}}(\mu\|\nu) = \int f\left(\text{proj}_{G(\mathbb{R})} \frac{d\mu}{d\nu}\right) d\nu = \int f\left(\frac{d\zeta}{d\nu}\right) d\nu = D_f(\zeta\|\nu). \quad \square$$

Next we do the reverse: converting standard f -divergence to f -divergence over rays.

Proposition 4.5. *For any Borel probability measures μ and ν , there exist Borel probability measures η, τ such that for any $f \in \mathcal{F}$,*

$$D_f(\mu\|\nu) = D_f^{\mathcal{R}}(\eta\|\tau).$$

Proof. The distribution function¹ $m : [0, \infty) \rightarrow [0, 1]$ of $d\mu/d\nu$ is given by

$$m(\alpha) := \nu\left(\left\{x \in \Omega : \frac{d\mu}{d\nu}(x) > \alpha\right\}\right)$$

for any $\alpha \in [0, \infty)$. Since m may not be injective, we define its generalized inverse $\rho : [0, 1] \rightarrow [0, \infty)$,

$$\rho(u) := \inf\{\alpha \in [0, \infty) : m(\alpha) \leq u\}$$

for any $u \in [0, 1]$. Thus ρ may be viewed a decreasing rearrangement of $d\mu/d\nu$ on $[0, 1]$ satisfying

$$\nu\left(\left\{x \in \Omega : \frac{d\mu}{d\nu}(x) > \alpha\right\}\right) = \lambda(\{u \in [0, 1] : \rho(u) > \alpha\}),$$

with λ the Lebesgue measure on $[0, 1]$ [20, p. 199]. Let $T : [0, 1] \rightarrow \mathbb{R}$ be a strictly increasing measure-preserving transformation. Define the probability measures

$$\tau(A) := \lambda(T^{-1}(A)), \quad \eta(A) := \int_{T^{-1}(A)} \rho d\lambda$$

for any $A \in \mathcal{B}$. Since ρ is nonincreasing and T^{-1} is increasing, the Radon–Nikodym derivative

$$\frac{d\eta}{d\tau} = \rho \circ T^{-1}$$

is also nonincreasing on \mathbb{R} , so $d\eta/d\tau = \text{proj}_{G(\mathbb{R})} d\eta/d\tau$. Hence, for any $f \in \mathcal{F}$,

$$D_f(\mu\|\nu) = \int f\left(\frac{d\mu}{d\nu}\right) d\nu = \int_{[0,1]} f(\rho) d\lambda = \int f(\rho \circ T^{-1}) d\tau = \int f\left(\text{proj}_{G(\mathbb{R})} \frac{d\eta}{d\tau}\right) d\tau = D_f^{\mathcal{R}}(\eta\|\tau)$$

as required. \square

Collectively, Propositions 4.4 and 4.5 lead to the following conclusion: An inequality relation between two divergences hold if and only if that same relation hold for their counterparts over rays.

¹In this article, we use the *decumulative* distribution function instead of the more common cumulative distribution function. They are equivalent — we could have used the latter but it introduces more notational clutter. Another advantage is consistency with our discussion of Choquet integrals in Section 7.

Corollary 4.6 (Preservation of inequality relations). *Let $f, g \in \mathcal{F}$ and $\psi, \varphi : [0, \infty) \rightarrow \mathbb{R}$. Then*

$$\psi(D_f(\mu\|\nu)) \leq \varphi(D_g(\mu\|\nu))$$

holds for all Borel probability measures μ and ν if and only if

$$\psi(D_f^{\mathcal{R}}(\mu\|\nu)) \leq \varphi(D_g^{\mathcal{R}}(\mu\|\nu))$$

holds for all Borel probability measures μ and ν .

Proof. Suppose $\psi(D_f(\mu\|\nu)) \leq \varphi(D_g(\mu\|\nu))$ for all μ and ν . By Proposition 4.4, there exists ζ such that

$$\psi(D_f^{\mathcal{R}}(\mu\|\nu)) = \psi(D_f(\zeta\|\nu)) \leq \varphi(D_g(\zeta\|\nu)) = \varphi(D_g^{\mathcal{R}}(\mu\|\nu)).$$

Conversely, suppose $\psi(D_f^{\mathcal{R}}(\mu\|\nu)) \leq \varphi(D_g^{\mathcal{R}}(\mu\|\nu))$ for all μ and ν . By Proposition 4.5, there exist η and τ such that

$$\psi(D_f(\mu\|\nu)) = \psi(D_f^{\mathcal{R}}(\eta\|\tau)) \leq \varphi(D_g^{\mathcal{R}}(\eta\|\tau)) = \varphi(D_g(\mu\|\nu)). \quad \square$$

As most readers familiar with f -divergences would suspect, the total variation distance plays a somewhat special role among all f -divergences. One striking property [7] is that it serves as a kind of universal lower bound for all f -divergence, in the sense that

$$\bar{f}(D_{\text{TV}}(\mu\|\nu)) \leq D_f(\mu\|\nu),$$

where $\bar{f}(x) := f(1+x) + f(1-x)$ is the symmetrized function of f . Note that $D_{\text{TV}}(\mu\|\nu)$ takes value in $[0, 1]$, so the $f(1-x)$ term is always evaluated within the domain of $f \in \mathcal{F}$. A consequence of our definition of \mathcal{F} in (16) is that \bar{f} is increasing and $\bar{f}(x) > 0$ when $x > 0$. So by Corollary 4.6, the same result holds verbatim for divergences over rays.

Corollary 4.7 (Total variation over rays as universal lower bound). *Let $f \in \mathcal{F}$ and set $\bar{f}(x) := f(1+x) + f(1-x)$ for $x \in [0, 1]$. Then*

$$\bar{f}(D_{\text{TV}}^{\mathcal{R}}(\mu\|\nu)) \leq D_f^{\mathcal{R}}(\mu\|\nu).$$

Theorem 4.1(v) now follows directly from Corollary 4.7: If $D_f^{\mathcal{R}}(\mu\|\nu) = 0$, then $\bar{f}(D_{\text{TV}}^{\mathcal{R}}(\mu\|\nu)) = 0$. So $D_{\text{TV}}^{\mathcal{R}}(\mu\|\nu) = 0$ as \bar{f} is strictly increasing and $\bar{f}(0) = 0$. Since we have also assumed that $D_f^{\mathcal{R}}(\nu\|\mu) = 0$, we get $D_{\text{TV}}^{\mathcal{R}}(\nu\|\mu) = 0$, and thus $\widehat{D}_{\text{TV}}^{\mathcal{R}}(\nu\|\mu) = 0$. By (17), we deduce that $\mu = \nu$ over \mathcal{R} . But \mathcal{R} is a π -system that generates \mathcal{B} . So Lemma 2.3 assures that we have $\mu = \nu$.

We end this section with another demonstration of the utility of Corollary 4.6: It allows us to extend essentially all known relations between divergences to their counterparts over rays.

Theorem 4.8 (Relations between f -divergences over rays). *The following inequalities hold verbatim when each D_f is replaced by $D_f^{\mathcal{R}}$.*

(a) *Total variation and Hellinger:*

$$\frac{1}{2} D_{\text{H}}^2(\mu\|\nu) \leq D_{\text{TV}}(\mu\|\nu) \leq D_{\text{H}}(\mu\|\nu) \sqrt{1 - \frac{D_{\text{H}}^2(\mu\|\nu)}{4}} \leq 1,$$

$$D_{\text{TV}}(\mu\|\nu) \leq \sqrt{-2 \log\left(1 - \frac{D_{\text{H}}^2(\mu\|\nu)}{2}\right)}.$$

(b) *Kullback–Leibler and total variation:*

$$D_{\text{TV}}^2(\mu\|\nu) \leq \frac{1}{2} D_{\text{KL}}(\mu\|\nu),$$

$$\log\left(\frac{1 + D_{\text{TV}}(\mu\|\nu)}{1 - D_{\text{TV}}(\mu\|\nu)}\right) - \frac{2 D_{\text{TV}}(\mu\|\nu)}{1 + D_{\text{TV}}(\mu\|\nu)} \leq D_{\text{KL}}(\mu\|\nu),$$

$$D_{\text{KL}}(\mu\|\nu) \leq \log\left(1 + \frac{1}{2\nu_{\min}} D_{\text{TV}}(\mu\|\nu)^2\right) \leq \frac{1}{2\nu_{\min}} D_{\text{TV}}(\mu\|\nu)^2,$$

where $\nu_{\min} := \min_x \nu(x)$. These are known respectively as Pinsker, strong Pinsker, and reverse Pinsker inequality.

(c) Total variation and χ^2 :

$$\begin{aligned} D_{\text{TV}}^2(\mu\|\nu) &\leq \frac{1}{4} D_{\chi^2}(\mu, \nu), & D_{\text{TV}}(\mu\|\nu) &\leq \max\left\{\frac{1}{2}, \frac{D_{\chi^2}(\mu\|\nu)}{1 + D_{\chi^2}(\mu, \nu)}\right\}, \\ 4 D_{\text{TV}}^2(\mu\|\nu) &\leq \varphi(D_{\text{TV}}(\mu\|\nu)) \leq D_{\chi^2}(\mu\|\nu), \end{aligned}$$

where

$$\varphi(t) = \begin{cases} 4t^2 & t \leq 1/2, \\ t/(1-t) & t \geq 1/2. \end{cases}$$

(d) Kullback–Leibler and Hellinger:

$$\begin{aligned} D_{\text{H}}^2(\mu\|\nu) &\leq 2 \log \frac{2}{2 - D_{\text{H}}^2(\mu\|\nu)} \leq D_{\text{KL}}(\mu\|\nu), \\ D_{\text{KL}}(\mu\|\nu) &\leq \frac{\log(1/\nu_{\min} - 1)}{1 - 2\nu_{\min}} (1 - (1 - D_{\text{H}}^2(\mu\|\nu))^2). \end{aligned}$$

(e) Kullback–Leibler and χ^2 :

$$D_{\text{KL}}(\mu\|\nu) \leq \log(1 + D_{\chi^2}(\mu\|\nu)) \leq D_{\chi^2}(\mu\|\nu).$$

(f) Le Cam and Hellinger:

$$\frac{1}{2} D_{\text{H}}^2(\mu\|\nu) \leq D_{\text{LC}}(\mu\|\nu) \leq D_{\text{H}}^2(\mu\|\nu).$$

(g) Le Cam and Jensen–Shannon:

$$D_{\text{LC}}(\mu\|\nu) \leq D_{\text{JS}}(\mu\|\nu) \leq 2 \log 2 \cdot D_{\text{LC}}(\mu\|\nu).$$

Proof. These are well-known inequalities between the respective f -divergences: the two inequalities in (a) may be found in [40, p. 124] and [21, Theorem 12] respectively; the three inequalities in (b) in [52, p. 88], [53], and [42, Theorem 28] respectively; those in (c), (d), and (e) may be found in [40, p. 133]; (f) in [32, p. 48]; (g) in [50, Theorem 3.2]. It follows from Corollary 4.6 that these relations also hold for the corresponding f -divergences over \mathcal{R} . \square

Ultimately, a definition is justified by the richness of results that can be proved with it. We view the results in this section as evidence that the f -divergence over rays as defined in Definition 3.12 is the “right” one. In this regard, the next section will supply another major piece of evidence.

5. GLIVENKO–CANTELLI THEOREM FOR f -DIVERGENCES OVER RAYS

Put in our context, the original Glivenko–Cantelli theorem in Theorem 1.4 states that empirical distributions converge to the target distribution almost surely in total variation distance over \mathcal{R} , as defined in Definition 3.12. Here we generalize this result to all f -divergences over \mathcal{R} .

As in the case of total variation distance, this convergence does not always hold if we replace \mathcal{R} with \mathcal{B} or some other subsets of \mathcal{B} . Examples 1.1 and 1.2, where $\limsup_{n \rightarrow \infty} D_{\text{TV}}(\nu_n\|\nu) > 0$, can be readily extended to any $f \in \mathcal{F}$ by way of Corollary 4.7:

$$0 < \bar{f}(\limsup_{n \rightarrow \infty} D_{\text{TV}}(\nu_n\|\nu)) = \limsup_{n \rightarrow \infty} \bar{f}(D_{\text{TV}}(\nu_n\|\nu)) \leq D_f(\nu_n\|\nu),$$

i.e., the f -divergence between the empirical distribution and the target does not in general converge to 0 almost surely. Our main result is that it holds for f -divergence over rays: For any $f \in \mathcal{F}$, Theorems 5.2 and 5.4 collectively show that almost surely

$$(18) \quad \lim_{n \rightarrow \infty} D_f^{\mathcal{R}}(\nu_n\|\nu) = 0 = \lim_{n \rightarrow \infty} D_f^{\mathcal{R}}(\nu\|\nu_n).$$

Note that an f -divergence is generally not symmetric in its arguments so we need to establish both equalities.

We will examine more carefully why a sequence of empirical measures can converge on \mathcal{R} but fails to converge on \mathcal{B} . It turns out that a necessary condition is that its corresponding sequence of Radon–Nikodym derivatives does not converge to 1, but the sequence of their projections onto $G(\mathcal{R})$ does.

Proposition 5.1. *Let ν be a Borel probability measure and ν_n be the corresponding empirical measure, $n \in \mathbb{N}$. Suppose the Radon–Nikodym derivative $d\nu_n/d\nu$ exists for all $n \in \mathbb{N}$ and are uniformly bounded. If*

$$\sup_{A \in \mathcal{B}} (\nu_n(A) - \nu(A)) \not\rightarrow 0 \quad \text{and} \quad \sup_{A \in \mathcal{R}} (\nu_n(A) - \nu(A)) \rightarrow 0,$$

then

$$\frac{d\nu_n}{d\nu} \not\rightarrow 1 \quad \text{and} \quad \text{proj}_{G(\mathcal{R})} \frac{d\nu_n}{d\nu} \rightarrow 1,$$

where convergence are all in an almost sure sense.

Proof. Suppose $\lim_{n \rightarrow \infty} d\nu_n/d\nu = 1$. Then

$$0 = \int \frac{1}{2} \left| \lim_{n \rightarrow \infty} \frac{d\nu_n}{d\nu} - 1 \right| d\nu \geq \limsup_{n \rightarrow \infty} \int \frac{1}{2} \left| \frac{d\nu_n}{d\nu} - 1 \right| d\nu = \limsup_{n \rightarrow \infty} D_{\text{TV}}(\nu_n \| \nu) \geq 0.$$

So $\lim_{n \rightarrow \infty} D_{\text{TV}}(\nu_n \| \nu) = 0$, a contradiction, and thus $d\nu_n/d\nu \not\rightarrow 1$. By assumption, $d\nu_n/d\nu \leq c$ for all $n \in \mathbb{N}$ and for some constant $c > 0$. It then follows from Proposition 3.6 that

$$\text{proj}_{G(\mathcal{R})} \frac{d\nu_n}{d\nu} \leq \text{proj}_{G(\mathcal{R})} c = c$$

for all n , i.e., the sequence $(\text{proj}_{G(\mathcal{R})} d\nu_n/d\nu)_n$ is also uniformly bounded. Since $\dim_{\text{VC}}(\mathcal{R}) < \infty$, it follows from the Vapnik–Chervonenkis theorem (see p. 4 for a statement) that

$$\lim_{n \rightarrow \infty} D_{\text{TV}}^{\mathcal{R}}(\nu_n \| \nu) = \lim_{n \rightarrow \infty} \sup_{A \in \mathcal{R}} (\nu_n(A) - \nu(A)) = 0$$

almost surely; taken together with the bounded convergence theorem,

$$\begin{aligned} 0 &= \lim_{n \rightarrow \infty} \sup_{A \in \mathcal{R}} (\nu_n(A) - \nu(A)) = \lim_{n \rightarrow \infty} \int \frac{1}{2} \left| \text{proj}_{G(\mathcal{R})} \frac{d\nu_n}{d\nu} - 1 \right| d\nu \\ &= \int \frac{1}{2} \left| \lim_{n \rightarrow \infty} \text{proj}_{G(\mathcal{R})} \frac{d\nu_n}{d\nu} - 1 \right| d\nu \geq 0. \end{aligned}$$

Hence the last integrand must converge to zero, or, equivalently,

$$(19) \quad \lim_{n \rightarrow \infty} \text{proj}_{G(\mathcal{R})} \frac{d\nu_n}{d\nu} = 1. \quad \square$$

With the observation in Proposition 5.1, we obtain the left half of (18).

Theorem 5.2 (Glivenko–Cantelli theorem for f -divergence I). *Let \mathcal{R} be the class of rays, $f \in \mathcal{F}$ as in (16), ν a Borel probability measure, and ν_n the corresponding empirical measure, $n \in \mathbb{N}$. Then almost surely*

$$\lim_{n \rightarrow \infty} D_f^{\mathcal{R}}(\nu_n \| \nu) = 0.$$

Proof. By (19) and the reverse Fatou lemma,

$$\begin{aligned} 0 &= \int \left| f \left(\lim_{n \rightarrow \infty} \text{proj}_{G(\mathcal{R})} \frac{d\nu_n}{d\nu} \right) \right| d\nu = \int \lim_{n \rightarrow \infty} \left| f \left(\text{proj}_{G(\mathcal{R})} \frac{d\nu_n}{d\nu} \right) \right| d\nu \\ &\geq \limsup_{n \rightarrow \infty} \int \left| f \left(\text{proj}_{G(\mathcal{R})} \frac{d\nu_n}{d\nu} \right) \right| d\nu \geq \limsup_{n \rightarrow \infty} \left| \int f \left(\text{proj}_{G(\mathcal{R})} \frac{d\nu_n}{d\nu} \right) d\nu \right| \geq 0. \end{aligned}$$

Hence

$$\lim_{n \rightarrow \infty} \int f\left(\text{proj}_{G(\mathcal{R})} \frac{d\nu_n}{d\nu}\right) d\nu = \lim_{n \rightarrow \infty} D_f^{\mathcal{R}}(\nu_n \| \nu) = 0. \quad \square$$

The following ‘‘reciprocal’’ of (19) will lead us to the right half of (18).

Proposition 5.3. *Let ν be a Borel probability measure and ν_n be the corresponding empirical measure, $n \in \mathbb{N}$. Suppose the Radon–Nikodym derivative $d\nu_n/d\nu$ exists for all $n \in \mathbb{N}$ and are uniformly bounded. Then almost surely*

$$\lim_{n \rightarrow \infty} \text{proj}_{G(\mathcal{R})} \frac{d\nu}{d\nu_n} = 1.$$

Proof. Let $x_n := |\text{proj}_{G(\mathcal{R})}(d\nu/d\nu_n) - 1|/2$. Since $\dim_{\text{vc}}(\mathcal{R}) < \infty$, it follows from the Vapnik–Chervonenkis theorem that

$$\sup_{A \in \mathcal{R}} (\nu(A) - \nu_n(A)) = \lim_{n \rightarrow \infty} \int x_n d\nu_n = 0.$$

By assumption, $d\nu/d\nu_n < c$ for all $n \in \mathbb{N}$ and for some constant $c > 0$. If $\lim_{n \rightarrow \infty} x_n \neq 0$, then there exists $\varepsilon > 0$ and subsequences (x_{n_k}) and (ν_{n_k}) such that

$$\varepsilon < \int x_{n_k} d\nu = \int x_{n_k} \frac{d\nu}{d\nu_{n_k}} d\nu_{n_k} \leq c \int x_{n_k} d\nu_{n_k} \rightarrow 0,$$

a contradiction, and thus $\lim_{n \rightarrow \infty} x_n = 0$. □

Theorem 5.4 (Glivenko–Cantelli theorem for f -divergence II). *Let \mathcal{R} be the class of rays, $f \in \mathcal{F}$ as in (16), ν a Borel probability measure, and ν_n the corresponding empirical measure, $n \in \mathbb{N}$. Then almost surely*

$$\lim_{n \rightarrow \infty} D_f^{\mathcal{R}}(\nu \| \nu_n) = 0.$$

Proof. Let $y_n := |f(\text{proj}_{G(\mathcal{R})} d\nu/d\nu_n)|$. Since f is continuous at 1,

$$\lim_{n \rightarrow \infty} y_n = \left| f\left(\lim_{n \rightarrow \infty} \text{proj}_{G(\mathcal{R})} \frac{d\nu}{d\nu_n}\right) \right| = 0.$$

So

$$\left| \int f\left(\text{proj}_{G(\mathcal{R})} \frac{d\nu}{d\nu_n}\right) d\nu_n \right| \leq \int \left| f\left(\text{proj}_{G(\mathcal{R})} \frac{d\nu}{d\nu_n}\right) \right| d\nu_n = \nu_n y_n.$$

If $\nu_n y_n \not\rightarrow 0$, then there exists $\varepsilon > 0$ and a subsequence (y_{n_k}) such that

$$\varepsilon < |\nu_{n_k} y_{n_k}| \leq \|y_{n_k}\| \rightarrow 0,$$

a contradiction, and thus $\lim_{n \rightarrow \infty} \nu_n y_n = 0$, giving us the required result. □

The following is an immediate consequence of Theorems 5.2 and 5.4.

Corollary 5.5 (Glivenko–Cantelli theorem for f -divergence III). *Let \mathcal{R} be the class of rays, $f \in \mathcal{F}$ as in (16), ν a Borel probability measure, and ν_n the corresponding empirical measure, $n \in \mathbb{N}$. Then almost surely*

$$\lim_{n \rightarrow \infty} \max\{D_f^{\mathcal{R}}(\nu \| \nu_n), D_f^{\mathcal{R}}(\nu_n \| \nu)\} = 0.$$

One advantage afforded by the version in Corollary 5.5 is that the expression

$$\widehat{D}_f^{\mathcal{R}}(\mu \| \nu) := \max\{D_f^{\mathcal{R}}(\mu \| \nu), D_f^{\mathcal{R}}(\nu \| \mu)\}$$

defines a distance (and not just a divergence) between probability measures. For $f(t) = |t-1|/2$, this is exactly the Kolmogorov–Smirnov distance in (17). We may of course also replace the expression with an L^p variant for any $p \in [1, \infty)$.

6. VAPNIK–CHERVONENKIS THEORY FOR f -DIVERGENCE

A natural follow-up question from the last section is: Given that we have a Glivenko–Cantelli theorem for f -divergence, is a Vapnik–Chervonenkis theory for f -divergence within reach? A first goal of such a theory would be to extend the work in this article from the class of rays \mathcal{R} to more general classes $\mathcal{C} \subseteq \mathcal{B}$. Surprisingly, we found that the original Vapnik–Chervonenkis theorem [55] provides a key.

Replace \mathcal{R} by a class $\mathcal{C} \subseteq \mathcal{B}$ in the definition of $G(\mathcal{R})$ in (6). As Section 3 shows, if $G(\mathcal{C})$ satisfies the conclusions of Proposition 3.1 and Theorem 3.15 with \mathcal{C} in place of \mathcal{R} , then Definition 3.12 with \mathcal{C} in place of \mathcal{R} gives a well-defined f -divergence over \mathcal{C} . We formalize this with a definition.

Definition 6.1 (Pre-Glivenko–Cantelli class). We say that $\mathcal{C} \subseteq \mathcal{B}$ is a pre-Glivenko–Cantelli class if for any Borel probability measure ν ,

$$G(\mathcal{C}) := \{g \in L^2(\nu) : \{g > r\} \in \mathcal{C} \text{ for all } r \in \mathbb{R}\}$$

is a closed convex cone in $L^2(\nu)$; and for any pair of Borel probability measures μ and ν ,

$$\sup_{A \in \mathcal{C}} (\mu(A) - \nu(A)) = \int \frac{1}{2} \left| \text{proj}_{G(\mathcal{C})} \frac{d\mu}{d\nu} - 1 \right| d\nu.$$

We may speak freely of $D_f^{\mathcal{C}}$ whenever we have a pre-Glivenko–Cantelli class \mathcal{C} and an $f \in \mathcal{F}$. We show that the first requirement in Definition 6.1 follows as long as \mathcal{C} is closed under union and intersection, noting that both \mathcal{R} and any σ -subalgebra of \mathcal{B} will have this property.

Lemma 6.2. *If $\mathcal{C} \subseteq \mathcal{B}$ is closed under countable unions and intersections, then $G(\mathcal{C})$ is a closed convex cone in $L^2(\nu)$.*

Proof. The sum of two functions in $G(\mathcal{C})$ remains in it as

$$\{g_1 + g_2 > c\} = \bigcup_{q \in \mathbb{Q}} \{g_1 > c - q\} \cap \{g_2 > q\} \in \mathcal{C}.$$

From this it is routine to verify that $G(\mathcal{C})$ is a convex cone. If $g_n \in G(\mathcal{C})$ is a sequence with $\|g_n - g\| \rightarrow 0$, then $g_n \rightarrow g$ in ν -measure. By passing through a subsequence if necessary, we may assume that $g_n \rightarrow g$ almost surely. It then follows from the definition of pointwise convergence that

$$\{g > c\} = \bigcup_{k \geq 1} \bigcap_{m \geq 1} \bigcup_{n \geq m} \left\{ g_n > c + \frac{1}{k} \right\} \in \mathcal{C}$$

and so $g \in G(\mathcal{C})$. Hence $G(\mathcal{C})$ is closed. \square

The second requirement in Definition 6.1 is usually harder to establish. For \mathcal{B} it is immediate but for \mathcal{R} it took nearly the whole of Section 3. We will extend the former to include all σ -subalgebras of \mathcal{B} .

Lemma 6.3. *Any σ -subalgebra $\mathcal{C} \subseteq \mathcal{B}$ is a pre-Glivenko–Cantelli class.*

Proof. If \mathcal{C} is a σ -subalgebra, then the restrictions of Borel probability measures μ and ν to \mathcal{C} , $\mu|_{\mathcal{C}}$ and $\nu|_{\mathcal{C}}$, remain probability measures [43, p. 30]. Let g be an almost surely nonnegative \mathcal{C} -measurable function, i.e., $\{g > c\} \in \mathcal{C}$ for all $c \in \mathbb{R}$. Then

$$\int_A g d\nu|_{\mathcal{C}} = \int_A g d\nu$$

for any $A \in \mathcal{C}$, i.e., the value of the integral is preserved. If $\mu \ll \nu$, then $\mu|_{\mathcal{C}} \ll \nu|_{\mathcal{C}}$, ensuring that the Radon–Nikodym derivative $d\mu|_{\mathcal{C}}/d\nu|_{\mathcal{C}}$ exists and that

$$\int_A \frac{d\mu|_{\mathcal{C}}}{d\nu|_{\mathcal{C}}} d\nu = \int_A \frac{d\mu|_{\mathcal{C}}}{d\nu|_{\mathcal{C}}} d\nu|_{\mathcal{C}} = \mu|_{\mathcal{C}}(A) = \mu(A) = \int_A \frac{d\mu}{d\nu} d\nu$$

for any $A \in \mathcal{C}$. So $d\mu|_{\mathcal{C}}/d\nu|_{\mathcal{C}}$ is indeed the conditional expectation and

$$(20) \quad \frac{d\mu|_{\mathcal{C}}}{d\nu|_{\mathcal{C}}} = \mathbb{E}_{\mu} \left[\frac{d\mu}{d\nu} \mid \mathcal{C} \right] = \text{proj}_{G(\mathcal{C})} \frac{d\mu}{d\nu}$$

by [4, p. 90]. Hence

$$\begin{aligned} D_{\text{TV}}^{\mathcal{C}}(\mu||\nu) &= \int \frac{1}{2} \left| \text{proj}_{G(\mathcal{C})} \frac{d\mu}{d\nu} - 1 \right| d\nu = \int \frac{1}{2} \left| \frac{d\mu|_{\mathcal{C}}}{d\nu|_{\mathcal{C}}} - 1 \right| d\nu \\ &= \int \frac{1}{2} \left| \frac{d\mu|_{\mathcal{C}}}{d\nu|_{\mathcal{C}}} - 1 \right| d\nu|_{\mathcal{C}} = \sup_{A \in \mathcal{C}} (\mu(A) - \nu(A)). \quad \square \end{aligned}$$

So Definition 6.1 is not vacuous: The class of rays \mathcal{R} and any σ -subalgebra of \mathcal{B} give two examples of pre-Glivenko–Cantelli classes. The reason for introducing this notion is that any finite VC-dimensional pre-Glivenko–Cantelli class gives us a Vapnik–Chervonenkis theorem for f -divergence.

Theorem 6.4 (Vapnik–Chervonenkis theorem for f -divergence). *Let $f \in \mathcal{F}$ and \mathcal{C} be a pre-Glivenko–Cantelli class. If $\dim_{\text{VC}}(\mathcal{C}) < \infty$, then almost surely*

$$(21) \quad \lim_{n \rightarrow \infty} D_f^{\mathcal{C}}(\nu_n||\nu) = 0 = \lim_{n \rightarrow \infty} D_f^{\mathcal{C}}(\nu||\nu_n).$$

Proof. All that we really need to observe is that the whole of Section 5 does not depend on the fact that \mathcal{R} is the class of rays — every result therein holds verbatim with \mathcal{C} in place of \mathcal{R} as long as $D_f^{\mathcal{C}}$ is well-defined, i.e., as long as \mathcal{C} is a pre-Glivenko–Cantelli class. The original Vapnik–Chervonenkis theorem [55], as applied in Propositions 5.1 and 5.3, shows that a pre-Glivenko–Cantelli class \mathcal{C} with finite VC-dimension must satisfy (21) almost surely. \square

With Theorem 6.4 we are led to the following definition in analogy with standard Glivenko–Cantelli class (i.e., with respect to total variation distance).

Definition 6.5 (Glivenko–Cantelli class for f -divergence). *Let $f \in \mathcal{F}$ and $\mathcal{C} \subseteq \mathcal{B}$ be a pre-Glivenko–Cantelli class for f . We say that \mathcal{C} is a Glivenko–Cantelli class for f if (21) holds almost surely.*

Note that we need to restrict ourselves to pre-Glivenko–Cantelli classes so that $D_f^{\mathcal{C}}$ is well-defined. It is nevertheless possible for such a Glivenko–Cantelli class to have infinite VC-dimension, so Theorem 6.4 merely provides a sufficient criterion.

Example 6.6 (Infinite VC-dimensional Glivenko–Cantelli class for f -divergence). *Let $\mathcal{E} = (E_n)_{n \in \mathbb{N}}$ be a countable partition of \mathbb{R} with measurable E_n , $n \in \mathbb{N}$. Let $\mathcal{C} = \sigma(\mathcal{E}) \subseteq \mathcal{B}$ be the σ -algebra generated by \mathcal{E} . So any set in \mathcal{C} can be expressed as a countable union of sets in \mathcal{E} [9, p. 4]. Thus any \mathcal{C} -measurable function must be a simple function $h = \sum_{n=1}^{\infty} \alpha_n \mathbb{1}_{E_n}$. Any pair of Borel probability measures μ and ν restrict to probability measures $\mu|_{\mathcal{C}}$ and $\nu|_{\mathcal{C}}$ with Radon–Nikodym derivative*

$$\frac{d\mu|_{\mathcal{C}}}{d\nu|_{\mathcal{C}}} = \sum_{n=1}^{\infty} \alpha_n \mathbb{1}_{E_n}.$$

For any $m \in \mathbb{N}$,

$$\mu|_{\mathcal{C}}(E_m) = \int_{E_m} \frac{d\mu|_{\mathcal{C}}}{d\nu|_{\mathcal{C}}} d\nu|_{\mathcal{C}} = \int \mathbb{1}_{E_m} \sum_{n=1}^{\infty} \alpha_n \mathbb{1}_{E_n} d\nu|_{\mathcal{C}} = \alpha_m \cdot \nu|_{\mathcal{C}}(E_m).$$

By (20),

$$\text{proj}_{G(\mathcal{C})} \frac{d\mu}{d\nu} = \frac{d\mu|_{\mathcal{C}}}{d\nu|_{\mathcal{C}}} = \sum_{n=1}^{\infty} \mathbb{1}_{E_n} \frac{\mu|_{\mathcal{C}}(E_n)}{\nu|_{\mathcal{C}}(E_n)},$$

and so the f -divergence over \mathcal{C} is given by

$$D_f^{\mathcal{C}}(\nu_n \| \nu) = D_f(\nu_n|_{\mathcal{C}} \| \nu|_{\mathcal{C}}) = \sum_{i=1}^{\infty} f\left(\frac{\nu_n|_{\mathcal{C}}(E_i)}{\nu|_{\mathcal{C}}(E_i)}\right) \nu|_{\mathcal{C}}(E_i).$$

If the sequence $(\nu_n|_{\mathcal{C}}(E_k)/\nu|_{\mathcal{C}}(E_k))_{k \in \mathbb{N}}$ is uniformly bounded, then

$$\begin{aligned} \limsup_{n \rightarrow \infty} D_f^{\mathcal{C}}(\nu_n \| \nu) &= \limsup_{n \rightarrow \infty} \sum_{i=1}^{\infty} f\left(\frac{\nu_n|_{\mathcal{C}}(E_i)}{\nu|_{\mathcal{C}}(E_i)}\right) \nu|_{\mathcal{C}}(E_i) \\ &\leq \sum_{i=1}^{\infty} \limsup_{n \rightarrow \infty} f\left(\frac{\nu_n|_{\mathcal{C}}(E_i)}{\nu|_{\mathcal{C}}(E_i)}\right) \nu|_{\mathcal{C}}(E_i) = \sum_{i=1}^{\infty} f\left(\lim_{n \rightarrow \infty} \frac{\nu_n|_{\mathcal{C}}(E_i)}{\nu|_{\mathcal{C}}(E_i)}\right) \nu|_{\mathcal{C}}(E_i) = 0. \end{aligned}$$

So $\lim_{n \rightarrow \infty} D_f^{\mathcal{C}}(\nu_n \| \nu) = 0$ almost surely. Likewise $\lim_{n \rightarrow \infty} D_f^{\mathcal{C}}(\nu \| \nu_n) = 0$ almost surely. Hence \mathcal{C} satisfies (21) almost surely even though $\dim_{\text{VC}}(\mathcal{C}) = \infty$.

7. CHOQUET INTEGRAL

To address Question 1, our initial thought was to employ the Choquet integral [11, 15], promising because it allows one to work with an arbitrary class $\mathcal{C} \subseteq 2^{\Omega}$ with essentially no condition imposed (unlike π -system or σ -algebra). This is perfect for us as it seems that we could have taken $\mathcal{C} = \mathcal{R}$ but it did not work. We think it might be instructive to record the difficulties encountered as we remain optimistic that it is not a dead end.

The Lebesgue integral defines integration with respect to a measure; Choquet integral extends this to so-called ‘‘capacities.’’ Let Ω be a set and $\mathcal{C} \subseteq 2^{\Omega}$ be a collection of subsets with $\emptyset, \Omega \in \mathcal{C}$. A function $\nu : \mathcal{C} \rightarrow \mathbb{R}$ is called a *capacity* on \mathcal{C} if it satisfies

- (a) normalization: $\nu(\emptyset) = 0$;
- (b) monotonicity: $\nu(A) \leq \nu(B)$ whenever $A \subseteq B$, $A, B \in \mathcal{C}$.

A function $g : \Omega \rightarrow \mathbb{R}$ is \mathcal{C} -measurable if $\{g > t\} \in \mathcal{C}$ for all $t \in \mathbb{R}$. Let $B_+(\mathcal{C})$ be the set of all bounded, nonnegative, \mathcal{C} -measurable functions. The distribution of $g \in B_+(\mathcal{C})$ is the function $m_{\nu, g} : \mathbb{R} \rightarrow [0, \infty)$, $t \mapsto \nu(\{g > t\})$. The Choquet integral of the function g with respect to the capacity ν is then

$$\oint g \, d\nu := \int_0^{\infty} m_{\nu, g}(t) \, dt,$$

where the right-hand side is a Riemann integral. Choquet coincides with Lebesgue when ν is a measure [15, p. 62].

For the class of rays \mathcal{R} , it is easy to see that $\nu_{\mathcal{R}}$, the restriction of ν to \mathcal{R} , is a capacity. So we could simply define f -divergence over \mathcal{R} as

$$D_f^{\mathcal{R}}(\mu \| \nu) := D_f(\mu_{\mathcal{R}} \| \nu_{\mathcal{R}}) = \oint f\left(\frac{d\mu_{\mathcal{R}}}{d\nu_{\mathcal{R}}}\right) d\nu_{\mathcal{R}}.$$

Two immediate problems are that the Radon–Nikodym derivative $d\mu_{\mathcal{R}}/d\nu_{\mathcal{R}}$ cannot be easily defined and the function $f(d\mu_{\mathcal{R}}/d\nu_{\mathcal{R}})$ may not be \mathcal{R} -measurable, i.e., nonincreasing. The first problem has been discussed extensively [24, 25, 38, 48, 51] but all assumed that \mathcal{C} is a σ -algebra and thus inapplicable to \mathcal{R} . One might think that a way around these problems is to extend the domain of $\mu_{\mathcal{R}}$ and $\nu_{\mathcal{R}}$ back to the whole of \mathcal{B} , i.e.,

$$\widehat{\nu}_{\mathcal{R}}(A) := \begin{cases} \nu(A) & \text{if } A \in \mathcal{R}, \\ 0 & \text{if } A \in \mathcal{B} \setminus \mathcal{R}, \end{cases}$$

and likewise for $\widehat{\mu}_{\mathcal{R}}$. But now $\widehat{\mu}_{\mathcal{R}}$ and $\widehat{\nu}_{\mathcal{R}}$ fail monotonicity and are not even capacities.

We are slightly disappointed that the discussions in [24, 25, 38, 48, 51] are all based on the assumption that \mathcal{C} is a σ -algebra (if so, then why not just use standard Lebesgue integral?) but interested readers may have better luck where we and these authors failed:

Open Problem 1. Is it possible to define an f -divergence using Choquet integral for a reasonably general class $\mathcal{C} \subseteq \mathcal{B}$ so that Theorems 4.1 and 4.8 hold true?

By “reasonably general” we expect \mathcal{C} to cover at least the requirements of Lemma 6.2, or perhaps be a π -system or a λ -system, certainly not just a σ -algebra. A positive answer to this open problem should provide an alternative to Definition 3.12 and could potentially be useful for developing the Vapnik–Chervonenkis theory for f -divergence in Section 6.

8. CONCLUSION

As we have noted in Section 1, the total variation distance plays an outsize role in traditional theoretical statistics but modern AI applications often relies on f -divergences other than the total variation distance. The work in this article would hopefully shed light on how one may extend classical results in statistics to other f -divergences.

REFERENCES

- [1] S. M. Ali and S. D. Silvey. A general class of coefficients of divergence of one distribution from another. *Journal of the Royal Statistical Society. Series B (Methodological)*, 28(1):131–142, 1966.
- [2] A. R. Alimov and I. G. Tsar’kov. *Geometric approximation theory*. Springer Monographs in Mathematics. Springer, Cham, 2021.
- [3] C. M. Bishop. *Pattern recognition and machine learning*. Information Science and Statistics. Springer, New York, 2006.
- [4] A. Bobrowski. *Functional analysis for probability and stochastic processes*. Cambridge University Press, Cambridge, 2005.
- [5] F. Bolley, A. Guillin, and C. Villani. Quantitative concentration inequalities for empirical measures on non-compact spaces. *Probab. Theory Related Fields*, 137(3–4):541–593, 2007.
- [6] O. Bousquet and A. Elisseeff. Stability and generalization. *J. Mach. Learn. Res.*, 2(3):499–526, 2002.
- [7] J. Bröcker. A lower bound on arbitrary f -divergences in terms of the total variation, 2009. Preprint, arXiv:0903.1765.
- [8] F. P. Cantelli. Sulla determinazione empirica delle leggi di probabilita. *Giorn. Ist. Ital. Attuari*, 4:421–424, 1933.
- [9] E. Çinlar. *Probability and stochastics*, volume 261 of *Graduate Texts in Mathematics*. Springer, New York, 2011.
- [10] H. Chernoff. A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations. *Ann. Math. Statistics*, 23:493–507, 1952.
- [11] G. Choquet. Theory of capacities. *Ann. Inst. Fourier (Grenoble)*, 5:131–295, 1953/54.
- [12] D. L. Cohn. *Measure theory*. Birkhäuser Advanced Texts: Basel Textbooks. Birkhäuser/Springer, New York, second edition, 2013.
- [13] I. Csiszar. Eine informationstheoretische ungleichung und ihre anwendung auf beweis der ergodizitaet von markoffschen ketten. *Magyer Tud. Akad. Mat. Kutato Int. Koezl.*, 8:85–108, 1964.
- [14] H. Dehling. Almost sure convergence of random variables. In M. Lovric, editor, *International Encyclopedia of Statistical Science*, pages 32–35. Springer, Berlin Heidelberg, 2011.
- [15] D. Denneberg. *Non-additive measure and integral*, volume 27 of *Theory and Decision Library. Series B: Mathematical and Statistical Methods*. Kluwer Academic Publishers Group, Dordrecht, 1994.
- [16] R. M. Dudley. *Uniform central limit theorems*, volume 142 of *Cambridge Studies in Advanced Mathematics*. Cambridge University Press, New York, second edition, 2014.
- [17] E. B. Dynkin. *Markov processes. Vol. II*, volume Band 122 of *Die Grundlehren der mathematischen Wissenschaften*. Springer-Verlag, Berlin-Göttingen-Heidelberg; Academic Press, Inc., Publishers, New York, 1965. Translated with the authorization and assistance of the author by J. Fabius, V. Greenberg, A. Maitra, G. Majone.
- [18] B. Efron and R. J. Tibshirani. *An introduction to the bootstrap*, volume 57 of *Monographs on Statistics and Applied Probability*. Chapman and Hall, New York, 1993.
- [19] S. Eguchi. A differential geometric approach to statistical inference on the basis of contrast functionals. *Hiroshima Math. J.*, 15(2):341–391, 1985.
- [20] G. Folland. *Real Analysis: Modern Techniques and Their Applications*. Pure and Applied Mathematics: A Wiley Series of Texts, Monographs and Tracts. Wiley, 1999.

- [21] G. L. Gilardoni. On Pinsker’s and Vajda’s type inequalities for Csiszár’s f -divergences. *IEEE Trans. Inform. Theory*, 56(11):5377–5386, 2010.
- [22] V. Glivenko. Sulla determinazione empirica delle leggi di probabilita. *Gion. Ist. Ital. Attauri.*, 4:92–99, 1933.
- [23] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014.
- [24] S. Graf. A Radon–Nikodým theorem for capacities. *J. Reine Angew. Math.*, 320:192–214, 1980.
- [25] J. Harding, M. Marinacci, N. T. Nguyen, and T. Wang. Local Radon–Nikodým derivatives of set functions. *Internat. J. Uncertain. Fuzziness Knowledge-Based Systems*, 5(3):379–394, 1997.
- [26] E. Hellinger. Neue Begründung der Theorie quadratischer Formen von unendlichvielen Veränderlichen. *J. Reine Angew. Math.*, 136:210–271, 1909.
- [27] N. Henze. *Asymptotic Stochastics: An Introduction with a View Towards Statistics*. Mathematics Study Resources. Springer, Berlin Heidelberg, 2024.
- [28] J. K. Hunter and B. Nachtergaele. *Applied analysis*. World Scientific Publishing Co., Inc., River Edge, NJ, 2001.
- [29] H. Jeffreys. *Theory of probability*. Oxford Classic Texts in the Physical Sciences. The Clarendon Press, Oxford University Press, New York, 1998.
- [30] M. Kelbert. Survey of distances between the most popular distributions. *Analytics*, 2(1):225–245, 2023.
- [31] D. P. Kingma and M. Welling. Auto-encoding variational bayes. In *2nd International Conference on Learning Representations (ICLR)*, 2014.
- [32] L. Le Cam. *Asymptotic methods in statistical decision theory*. Springer Series in Statistics. Springer, New York, 1986.
- [33] Y. Li and R. E. Turner. Rényi divergence variational inference. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29, pages 1081–1089. Curran Associates, Inc., 2016.
- [34] J. Lin. Divergence measures based on the Shannon entropy. *IEEE Trans. Inform. Theory*, 37(1):145–151, 1991.
- [35] A. Liu, J.-G. Liu, and Y. Lu. On the rate of convergence of empirical measure in ∞ -Wasserstein distance for unbounded density function. *Quart. Appl. Math.*, 77(4):811–829, 2019.
- [36] S. A. Molchanov. *Geometric modeling in probability and statistics* [book review of MR3308142]. *Bull. Amer. Math. Soc. (N.S.)*, 55(1):109–111, 2018.
- [37] T. Morimoto. Markov processes and the h-theorem. *Journal of the Physical Society of Japan*, 18(3):328–331, 1963.
- [38] Y. Narukawa. Distances defined by Choquet integral. In *2007 IEEE International Fuzzy Systems Conference*, 2007.
- [39] K. Pearson. On the criterion that a given system of deviations from the probable in the case of a correlated system of variables is such that it can be reasonably supposed to have arisen from random sampling. In S. Kotz and N. L. Johnson, editors, *Breakthroughs in Statistics: Methodology and Distribution*, pages 11–28. Springer, New York, 1992.
- [40] Y. Polyanskiy and Y. Wu. *Information Theory: From Coding to Learning*. Cambridge University Press, 2024.
- [41] A. Rényi. On measures of entropy and information. In *Proceedings of the fourth Berkeley symposium on mathematical statistics and probability, volume 1: contributions to the theory of statistics*, volume 4, pages 547–562. University of California Press, 1961.
- [42] I. Sason and S. Verdú. f -divergence inequalities. *IEEE Trans. Inform. Theory*, 62(11):5973–6006, 2016.
- [43] R. L. Schilling. *Measures, integrals and martingales*. Cambridge University Press, Cambridge, second edition, 2017.
- [44] C. Schumacher, F. Schwarzenberger, and I. Veselić. A Glivenko–Cantelli theorem for almost additive functions on lattices. *Stochastic Process. Appl.*, 127(1):179–208, 2017.
- [45] S. Shalev-Shwartz and S. Ben-David. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, 2014.
- [46] G. R. Shorack and J. A. Wellner. *Empirical processes with applications to statistics*, volume 59 of *Classics in Applied Mathematics*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2009.
- [47] R. M. Shortt. Empirical measures. *Amer. Math. Monthly*, 91(6):358–360, 1984.
- [48] M. Sugeno. A note on derivatives of functions with respect to fuzzy measures. *Fuzzy Sets and Systems*, 222:1–17, 2013.
- [49] J. C. Taylor. *An introduction to measure and probability*. Springer, New York, 1997.
- [50] F. Topsøe. Some inequalities for information divergence and related measures of discrimination. *IEEE Trans. Inform. Theory*, 46(4):1602–1609, 2000.
- [51] V. Torra, Y. Narukawa, and M. Sugeno. On the f -divergence for non-additive measures. *Fuzzy Sets and Systems*, 292:364–379, 2016.

- [52] A. B. Tsybakov. *Introduction to nonparametric estimation*. Springer Series in Statistics. Springer, New York, 2009.
- [53] I. Vajda. Note on discrimination information and variation. *IEEE Trans. Inform. Theory*, IT-16:771–773, 1970.
- [54] V. N. Vapnik. *Statistical learning theory*. Adaptive and Learning Systems for Signal Processing, Communications, and Control. John Wiley & Sons, Inc., New York, 1998.
- [55] V. N. Vapnik and A. Y. Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theor. Probability Appl.*, 16(2):264–280, 1971.

COMPUTATIONAL AND APPLIED MATHEMATICS INITIATIVE, DEPARTMENT OF STATISTICS, UNIVERSITY OF CHICAGO, CHICAGO, IL 60637

Email address: `haomingwang@uchicago.edu`

COMPUTATIONAL AND APPLIED MATHEMATICS INITIATIVE, DEPARTMENT OF STATISTICS, UNIVERSITY OF CHICAGO, CHICAGO, IL 60637

Email address: `lekheng@uchicago.edu`