

On metrics robust to noise and deformations

William Leeb

School of Mathematics
University of Minnesota, Twin Cities
Minneapolis, MN

Abstract

We study the properties of a family of distances between functions of a single variable. These distances are examples of integral probability metrics, and have been used previously for comparing probability measures on the line; special cases include the Earth Mover’s Distance and the Kolmogorov Metric. We examine their properties for general signals, proving that they are robust to a broad class of deformations. We also establish corresponding robustness results for the induced sliced distances between multivariate functions. Finally, we establish error bounds for approximating the univariate metrics from finite samples, and prove that these approximations are robust to additive Gaussian noise. The results are illustrated in numerical experiments, which include comparisons with Wasserstein distances.

1 Introduction

Many tasks in statistics and machine learning require specification of a metric that measures the similarity between data vectors. For example, the goal of clustering is to group points together that are close and separate points that are far, where “close” and “far” are determined by a certain metric or similarity measure that is designed to capture relevant features [55, 65, 1]. Such a method’s effectiveness depends crucially on the choice of metric. An appropriate metric will be robust to noise and to irrelevant deformations of the input data, so that only the “meaningful” characteristics of each data vector inform the distance.

A widely-used class of metrics are the Wasserstein distances for comparing probability distributions f and g defined over a metric space \mathcal{X} . Informally, the Wasserstein distance between f and g is equal to the minimal cost of transforming f into g by rearranging the mass, where the cost is determined by the metric on \mathcal{X} [66, 67]; we review the precise definition of Wasserstein distances in Section 2.7. The Wasserstein distances are popular metrics in a range of machine learning and statistical applications [49, 50, 57, 9, 51, 56, 37, 6, 54, 45, 12].

The goal of this work is to begin to address the question of whether the favorable properties of Wasserstein distances are shared by other families of metrics. More specifically, we consider two known “robustness” properties exhibited by the Wasserstein distances. The first result states that if there is a smooth bijection $\Phi : \mathcal{X} \rightarrow \mathcal{X}$ and

$$g(x) = f(\Phi(x)) \frac{d\Phi}{dx}(x), \quad (1)$$

where $\frac{d\Phi}{dx}(x)$ is the Radon-Nikodym derivative of Φ , then

$$W_p(f, g) \leq \sup_{x \in \mathcal{X}} d_{\mathcal{X}}(x, \Phi(x)), \quad (2)$$

where W_p denotes the p -Wasserstein distance. This bound is straightforward to prove (indeed, it is almost tautological from the Monge formulation of Wasserstein distances), and a self-contained argument may be found in Section 2.7; see also [36] for another argument when $p = 1$. Informally, (2) tells us that the Wasserstein distances are robust to “small” deformations of a distribution, where the “size” of a deformation is the maximum distance that any point in the domain may be moved.

The second property of Wasserstein distances that we will consider is from the recent paper [52], and suggests that the p -Wasserstein distance W_p is a good choice for clustering tomographic projection images

that arise from cryo-electron microscopy (cryo-EM), a technique for molecular reconstruction that is increasingly used in structural biology [60, 5, 17]. Suppose f and g are functions of two variables that are each projections of a common three-dimensional volume F ; that is,

$$f(x, y) = \int_{\mathbb{R}} F(R_f(x, y, z)) dz, \quad (3)$$

where R_f is an orthogonal transformation; and similarly for g . Then [52] proves that

$$\min_{R \in SO(2)} W_p(f, g \circ R) \leq \theta(R_f, R_g), \quad (4)$$

where $\theta(R_f, R_g)$ is the angle between the projection directions of R_f and R_g . In fact, in Section 2.7 we observe that Wasserstein distances exhibit a somewhat more general robustness property, namely, the Wasserstein distance between projections of two functions is robust to deformations of those functions. Though we have not seen this result stated in the literature, it follows rather trivially from (2) and the analysis in the paper [52]; see Theorem 2.2.

In many problems where Wasserstein distances are used, the transportation problem solved by the distance computation is not of interest in and of itself; rather, one only requires a distance that is robust to distortions of the data. Furthermore, because Wasserstein distances are defined between positive measures, it can be awkward to use or analyze them in settings where the observed signals take on negative values, as is typically the case when observations are corrupted by additive noise. It is natural to ask whether there are other families of metrics with similar robustness properties as the Wasserstein metrics, and which are robust to additive noise. Furthermore, while the motivation in [52] for the bound (4) is comparing two projection images in cryo-electron microscopy, there are, in fact, numerous scientific problems for which the measurement modality only permits observing projections of an object, from which the object itself must be reconstructed [15, 47, 27, 28, 61, 13]. It is therefore of interest to consider the properties of distances used to compare tomographic projections.

In this paper, we approach these questions by studying a family of simple metrics between single-variable functions, including tomographic projections of multivariate functions. These metrics are induced by norms, denoted by $\|f\|_{V^p}$, which are the p -norms of the Volterra operator (the indefinite integral operator) applied to f . We call the norm $\|f\|_{V^p}$ the *Volterra p -norm*, and its induced metric the *Volterra p -distance*, or *Volterra p -metric*. The Volterra distances have been used previously for comparing probability measures on the real line [46, 41]; however, unlike the Wasserstein distances, the Volterra distances are defined between all integrable functions, not just between probability measures; this makes it natural to consider their robustness to additive noise in addition to geometric distortions.

We show that the Volterra distances exhibit robustness properties like (2) and (4), as well as a property generalizing (2) and (4): namely, that when comparing univariate tomographic projections of multivariate functions, the distance is robust to deformations in the higher dimensional space. In fact, the robustness bounds we prove for Volterra metrics are stronger than those for Wasserstein distances; more precisely, when $p > 1$, the Volterra p -metric is bounded by a concave, non-linear function of the deformation's size, suggesting that the Volterra distances are more robust to large deformations. Finally, we analyze discrete approximations to the Volterra distances, showing that they converge to their continuous counterparts for a broad class of “well-behaved” functions, while also being robust to additive Gaussian noise.

Of course, the Volterra metrics are limited in their applicability, as they are only defined between univariate functions; the Wasserstein distances, by contrast, are naturally defined between densities of any number of variables. It is therefore of interest to ask whether there are Volterra-type metrics for multivariate functions that exhibit similar properties. We leverage the theory for Volterra distances to study distances between multivariate functions induced by “slicing”, a technique that has been used for Wasserstein distances [51, 9, 31, 32, 16, 48]. The robustness properties of the Volterra distances immediately induce corresponding robustness properties of the sliced Volterra distances.

Central to our analysis is the variational characterization of Volterra distances, which is well-known (and for which we provide a self-contained proof in Section 3.2). This characterization shows that the Volterra distances are a special case of the class of *maximum mean discrepancies*, also known as *integral probability metrics*, which are typically used for comparing probability densities [26, 11, 4, 46, 41, 2, 63]. If f and g are

two functions on a measure space \mathcal{X} , a maximum mean discrepancy is a metric of the form

$$d_{\mathcal{F}}(f, g) = \sup_{h \in \mathcal{F}} \int_{\mathcal{X}} h(x)(f(x) - g(x))dx, \quad (5)$$

where \mathcal{F} is a specified class of test functions. It is a consequence of the Kantorovich-Rubinstein Theorem [29, 18, 39, 40, 20] that the 1-Wasserstein distance, or Earth Mover’s Distance, between densities f and g is equal to the distance (5) when \mathcal{F} is the set of 1-Lipschitz functions with respect to a metric on \mathcal{X} . By contrast, when \mathcal{X} is an interval, the Volterra p -distance is obtained by taking \mathcal{F} to be the space of functions with derivative in $L^{p/(p-1)}$. The Volterra 1-distance is equal to the 1-Wasserstein metric in one variable, but the Volterra p -distance and the p -Wasserstein metric are not equal for any $p > 1$.

The remainder of the paper is structured as follows:

1. Section 2 reviews basic definitions, notation, and properties, including of the Lebesgue norms, the Volterra operator, push-forwards and deformations, Wasserstein and sliced Wasserstein distances, and tomographic projections. Theorem 2.2 provides a general robustness result on Wasserstein distances, which, though it follows easily from existing work, does not appear to have been stated previously.
2. Section 3 defines the basic objects of study, namely the Volterra distances between univariate functions, the sliced Volterra distances between multivariate functions, and the trapezoidal rule approximations to the Volterra distances. It also contains a self-contained proof of the well-known variational characterization of Volterra distances (Proposition 3.1).
3. Section 4 contains statements of the theorems on the robustness of the Volterra and sliced Volterra distances. Theorem 4.1 is a general robustness result for comparing univariate projections; Theorems 4.2 and 4.3 give stronger bounds for the case of rotations (i.e. changes in projection angle) and monotonically increasing deformations, respectively. Theorems 4.4 and 4.5 show robustness of the sliced Wasserstein distances.
4. Section 5 analyzes the trapezoidal rule approximation to the Volterra distance. Theorems 5.1 and 5.2 show that when samples are taken from sufficiently regular functions, the approximation converges to the true Volterra distance. Theorem 5.3 and Corollary 5.4 show that the approximations are robust to additive Gaussian noise.
5. Section 6 shows the results of numerical experiments illustrating the theoretical results, including comparisons between the Volterra, Wasserstein, and Lebesgue distances.
6. Sections 7 and 8 contain proofs of the results from Sections 4 and 5, respectively. Theorem 7.1 in Section 7 gives a more general statement about the Volterra distances, from which the theorems in Section 4 may all be derived; the variational characterization in Proposition 3.1 is central to the analysis.
7. Section 9 concludes the paper, providing a summary and topics for future research.

2 Preliminaries

This section introduces the basic definitions and notation that will be used in the rest of the paper. Familiarity with basic concepts of measure and integration, e.g. at the level of [21], will be assumed.

2.1 The Lebesgue p -norms

We recall the standard definition of the Lebesgue p -norm. Let $F : \mathbb{R}^d \rightarrow \mathbb{R}$ be a measurable function. If $1 \leq p < \infty$, then the Lebesgue p -norm of F is defined by

$$\|F\|_{L^p} = \left(\int_{\mathbb{R}^d} |F(\mathbf{x})|^p d\mathbf{x} \right)^{1/p}. \quad (6)$$

For $p = \infty$, we define

$$\|F\|_{L^\infty} = \operatorname{ess\,sup}_{\mathbf{x} \in \mathbb{R}^d} |F(\mathbf{x})|. \quad (7)$$

We denote by L^p the set of all functions $F : \mathbb{R}^d \rightarrow \mathbb{R}$ with $\|F\|_{L^p} < \infty$, and by $L^p(A)$ the subset of L^p containing only functions supported on A . As is well-known, if F is supported in a bounded set $A \subset \mathbb{R}^d$, then $\|F\|_{L^p} \leq \|F\|_{L^q} |A|^{1/p-1/q}$ if $p \leq q$, where $|A|$ denotes the Lebesgue measure of A ; in particular, $L^p(A) \subset L^q(A)$. If F is in L^p and G is in L^q , where $1/p + 1/q = 1$, we define their inner product by

$$\langle F, G \rangle = \int_{\mathbb{R}^d} F(\mathbf{x})G(\mathbf{x})d\mathbf{x}. \quad (8)$$

We also define the normalized Lebesgue p -norm $\|\mathbf{x}\|_{\ell_p}$ for vectors \mathbf{x} in \mathbb{R}^n . When $1 \leq p < \infty$,

$$\|\mathbf{x}\|_{\ell_p} = \left(\frac{1}{n} \sum_{j=1}^n |x_j|^p \right)^{1/p}, \quad (9)$$

and when $p = \infty$,

$$\|\mathbf{x}\|_{\ell_\infty} = \max_{1 \leq k \leq n} |x_k|. \quad (10)$$

Note the normalization by n when $p < \infty$. With this convention, $\|\mathbf{x}\|_{\ell_p} \leq \|\mathbf{x}\|_{\ell_q}$ whenever $p \leq q$.

We will denote the unnormalized 2-norm of a vector \mathbf{x} in \mathbb{R}^d by

$$|\mathbf{x}| = \left(\sum_{j=1}^d x_j^2 \right)^{1/2}. \quad (11)$$

Note that we do not normalize by $1/d$ in this case. The normalized norm $\|\mathbf{x}\|_{\ell_p}$ will be used when \mathbf{x} is a vector of samples of a function, and the unnormalized norm $|\mathbf{x}|$ when \mathbf{x} is a variable. We define the inner product between two vectors \mathbf{x} and \mathbf{y} in \mathbb{R}^d by

$$\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{j=1}^d x_j y_j. \quad (12)$$

2.2 Trapezoidal rule approximation

If f is a function on an interval $[a, b]$ and $n \geq 1$ is an integer, the trapezoidal rule approximation to $\int_a^b f(x)dx$ is defined as

$$T_n(f, a, b) = \frac{b-a}{2n} \sum_{k=0}^{n-1} (f(a_k) + f(a_{k+1})), \quad (13)$$

where

$$a_k = a + \frac{k}{n}(b-a), \quad 0 \leq k \leq n. \quad (14)$$

It is well-known that if f is C^3 , then

$$\left| T_n(f, a, b) - \int_a^b f(x)dx \right| \leq \frac{1}{12} \|f''\|_{L^\infty} \frac{(b-a)^3}{n^2}; \quad (15)$$

see, e.g., [14].

2.3 Absolute continuity

Suppose $a < b$. Recall that a function G on $[a, b]$ is said to be absolutely continuous if it can be written as

$$G(x) = G(a) + \int_a^x g(t) dt \quad (16)$$

for a function g in $L^1([a, b])$. If G is absolutely continuous, then it is differentiable almost everywhere, and $G' = g$ where the derivative exists. We denote by \mathcal{A}_0 the set of absolutely continuous functions G satisfying $G(b) = 0$; these functions may be written as

$$G(x) = - \int_x^b g(t) dt \quad (17)$$

where $g = G'$ almost everywhere. For brevity, whenever G is in \mathcal{A}_0 , G' will denote any function such that $G(x) = - \int_x^b G'(t) dt$.

The following result is standard (e.g., see Section 3.5 of [21]):

Theorem 2.1 (Integration by parts). *If F and G are absolutely continuous functions on $[a, b]$, then*

$$\int_a^b (F'(x)G(x) + F(x)G'(x))dx = F(b)G(b) - F(a)G(a). \quad (18)$$

2.4 The Volterra operator

The Volterra operator \mathcal{V} is defined on $L^1([a, b])$ by

$$(\mathcal{V}f)(x) = \int_a^x f(t)dt. \quad (19)$$

We note that this is only the simplest of a large family of related operators that have been widely studied [25]. Importantly, if f is in $L^1([a, b])$, $\mathcal{V}f$ is in L^∞ , with $\|\mathcal{V}f\|_{L^\infty} \leq \|f\|_{L^1}$; furthermore, $\mathcal{V}f$ is, by definition, absolutely continuous when f is in $L^1([a, b])$.

The adjoint transform \mathcal{V}^* is given by

$$(\mathcal{V}^*f)(x) = \int_x^b f(t)dt. \quad (20)$$

This operator satisfies

$$\langle \mathcal{V}f, g \rangle = \langle f, \mathcal{V}^*g \rangle \quad (21)$$

where f and g are two functions in $L^1([a, b])$.

2.5 Push-forwards and ϵ -deformations

Let $\Omega \subset \mathbb{R}^d$ be a non-empty, bounded, open set. Suppose that μ is a finite, signed measure on Ω , and that $\Psi : \Omega \rightarrow \mathbb{R}^d$ is a measurable function. Then Ψ induces a signed measure on $\Psi(\Omega)$, denoted $\Psi_\# \mu$, defined by

$$(\Psi_\# \mu)(E) = \mu(\Psi^{-1}(E)), \quad (22)$$

for measurable sets $E \subset \Psi(\Omega)$. The measure $\Psi_\# \mu$ is referred to as the *push-forward* of μ induced from Ψ [50]. Note that $(\Psi_\# \mu)(\Psi(\Omega)) = \mu(\Psi^{-1}(\Psi(\Omega))) = \mu(\Omega)$; that is, the push-forward preserves the total measure.

Now suppose that μ is induced from a function f supported on Ω ; that is, $\mu(E) = \int_E f(\mathbf{x})d\mathbf{x}$ for all measurable $E \subset \Omega$. Suppose too that Ψ is a diffeomorphism between Ω and $\Psi(\Omega)$; that is, it is C^1 , one-to-one, and $\det(\nabla\Psi(\mathbf{x})) \neq 0$, where $\nabla\Psi(\mathbf{x})$ denotes the Jacobian matrix of Ψ at \mathbf{x} . Let $\Phi = \Psi^{-1}$ denote the

functional inverse of Ψ . Then $\Psi_{\#}\mu$ has density $f(\Phi(\mathbf{x}))|\det(\nabla\Phi(\mathbf{x}))|$. Indeed, using the change of variables $\mathbf{x} = \Phi(\mathbf{u})$ and $d\mathbf{x} = |\det(\nabla\Phi(\mathbf{u}))|d\mathbf{u}$,

$$(\Psi_{\#}\mu)(E) = \mu(\Psi^{-1}(E)) = \int_{\Psi^{-1}(E)} f(\mathbf{x})d\mathbf{x} = \int_E f(\Phi(\mathbf{u}))|\det(\nabla\Phi(\mathbf{u}))|d\mathbf{u}. \quad (23)$$

When convenient, we will write $(\Psi_{\#}f)(\mathbf{x}) = f(\Phi(\mathbf{x}))|\det(\nabla\Phi(\mathbf{x}))|$, or $f_{\Phi}(\mathbf{x}) = f(\Phi(\mathbf{x}))|\det(\nabla\Phi(\mathbf{x}))|$.

When the diffeomorphism Ψ does not move points by more than a value $\epsilon > 0$ (that is, $|\mathbf{x} - \Psi(\mathbf{x})| \leq \epsilon$ for all \mathbf{x} in Ω , or equivalently, $|\mathbf{x} - \Phi(\mathbf{x})| \leq \epsilon$ for all \mathbf{x} in $\Psi(\Omega)$), we will refer to $\Psi_{\#}f$ as an ϵ -deformation of f ; we will also refer to Ψ itself as an ϵ -deformation of Ω . For example, if \mathbf{u} is a fixed unit vector, then the function $\Psi(\mathbf{x}) = \mathbf{x} + \epsilon\mathbf{u}$ is an ϵ -deformation.

2.6 Tomographic projections and the Radon transform

Let $\mathcal{U} = (\mathbf{u}_1, \dots, \mathbf{u}_r) \in \mathbb{S}^{d-1} \times \dots \times \mathbb{S}^{d-1}$ (where $\mathbb{S}^{d-1} \subset \mathbb{R}^d$ is the $(d-1)$ -dimensional unit sphere) denote an ordered collection of r unit vectors in \mathbb{R}^d . Let $\mathbf{u}_{r+1}, \dots, \mathbf{u}_d$ denote $d-r$ orthonormal vectors that are orthogonal to $\mathbf{u}_1, \dots, \mathbf{u}_r$. We define the operator $\mathcal{P}_{\mathcal{U}}$ by

$$(\mathcal{P}_{\mathcal{U}}F)(t_1, \dots, t_r) = \int_{\mathbb{R}^{d-r}} F(t_1\mathbf{u}_1 + \dots + t_r\mathbf{u}_r + s_1\mathbf{u}_{r+1} + \dots + s_{d-r}\mathbf{u}_d)ds_1 \dots ds_{d-r}. \quad (24)$$

We refer to $\mathcal{P}_{\mathcal{U}}F$ as the *tomographic projection* of F onto the subspace spanned by $\mathbf{u}_1, \dots, \mathbf{u}_r$. When $r = 1$, we will denote the tomographic projection of F onto the span of a unit vector \mathbf{u} by $\mathcal{P}_{\mathbf{u}}F$. Note that in this case, the *Radon transform* $\mathcal{R}F : \mathbb{R} \times \mathbb{S}^{d-1}$ of the function F is defined by $(\mathcal{R}F)(t, \mathbf{u}) = (\mathcal{P}_{\mathbf{u}}F)(t)$. For more background on these transforms, see, for example, the references [47, 27].

2.7 Wasserstein distances

If F and G are probability densities on a subset $\Omega \subset \mathbb{R}^d$, their p -Wasserstein distance $W_p(F, G)$ (also known as the Kantorovich distance) is defined as

$$W_p(F, G) = \min_{\Pi \in \mathcal{M}(F, G)} \left(\int_{\Omega} \int_{\Omega} |\mathbf{x} - \mathbf{y}|^p d\Pi(\mathbf{x}, \mathbf{y}) \right)^{1/p}, \quad (25)$$

where $\mathcal{M}(F, G)$ denotes the space of all probability measures on $\Omega \times \Omega$ with marginals equal to F and G , respectively [66, 67]. That is, $\Pi \in \mathcal{M}(F, G)$ if for all measurable $E \subset \Omega$,

$$\Pi(E \times \Omega) = \int_E F(\mathbf{x})d\mathbf{x}, \quad (26)$$

and

$$\Pi(\Omega \times E) = \int_E G(\mathbf{y})d\mathbf{y}. \quad (27)$$

Informally, $W_p(F, G)$ is the minimal cost of rearranging a unit of mass with distribution F into one with distribution G , where the cost of moving mass between locations \mathbf{x} and \mathbf{y} is $|\mathbf{x} - \mathbf{y}|^p$. The distance $W_1(F, G)$ is also known as the *Earth Mover's Distance (EMD)* between the probability measures F and G [66, 67], which we will also denote by $\text{EMD}(F, G)$. The Wasserstein distances and their variants have been widely used in statistics, machine learning, image processing, and related areas [49, 50, 57, 9, 51, 56, 37, 6, 54, 45, 12].

The Wasserstein distance is a relaxation of the Monge distance, defined by

$$M_p(F, G) = \min_{\Phi \in \mathcal{T}(F, G)} \left(\int_{\Omega} |\mathbf{x} - \Phi(\mathbf{x})|^p F(\mathbf{x})d\mathbf{x} \right)^{1/p} \quad (28)$$

where $\mathcal{T}(F, G)$ is the set of all functions $\Phi : \Omega \rightarrow \Omega$ such that $\int_E G(\mathbf{x})d\mathbf{x} = \int_{\Phi^{-1}(E)} F(\mathbf{x})d\mathbf{x}$; that is, all functions Φ which push F onto G , in the sense described in Section 2.5. Indeed, any function Φ in $\mathcal{T}(F, G)$ induces a measure Π_Φ in $\mathcal{M}(F, G)$, with

$$\int_{\Omega} \int_{\Omega} |\mathbf{x} - \mathbf{y}|^p d\Pi_\Phi(\mathbf{x}, \mathbf{y}) = \int_{\Omega} |\mathbf{x} - \Phi(\mathbf{x})|^p F(\mathbf{x}) d\mathbf{x}, \quad (29)$$

and hence $W_p(F, G) \leq M_p(F, G)$. (In fact, when $M_p(F, G)$ is finite, equality holds; see [57].) This implies (2); indeed, if Φ is a smooth bijection on Ω satisfying $|\mathbf{x} - \Phi(\mathbf{x})| \leq \epsilon$ for all \mathbf{x} , and $F_\Phi(\mathbf{x}) = F(\Phi(\mathbf{x}))|\det(\nabla\Phi(\mathbf{x}))|$, then Φ is contained in $\mathcal{T}(F, F_\Phi)$, and so

$$W_p(F, F_\Phi) \leq M_p(F, F_\Phi) \leq \left(\int_{\Omega} |\mathbf{x} - \Phi(\mathbf{x})|^p F(\mathbf{x}) d\mathbf{x} \right)^{1/p} \leq \epsilon \left(\int_{\Omega} F(\mathbf{x}) d\mathbf{x} \right)^{1/p} = \epsilon. \quad (30)$$

In fact, a more general robustness result may be shown, which we state now. The proof is nearly identical to that found in [52].

Theorem 2.2. *Suppose F is a probability density supported on a set Ω , and let $\Phi : \Omega \rightarrow \Omega$ be an ϵ -deformation. For $\mathcal{U} = (\mathbf{u}_1, \dots, \mathbf{u}_r)$, where $\mathbf{u}_1, \dots, \mathbf{u}_r$ are orthonormal, let $\mathcal{P} = \mathcal{P}_{\mathcal{U}}$ denote the tomographic projection operator. Then for all $p \geq 1$,*

$$W_p(\mathcal{P}F, \mathcal{P}F_\Phi) \leq \epsilon. \quad (31)$$

Proof. An identical proof to that of Lemma 1 in [52] shows that $W_p(\mathcal{P}F, \mathcal{P}F_\Phi) \leq W_p(F, F_\Phi)$ (note that the left side refers to transportation in \mathbb{R}^{d-1} , and the right side to \mathbb{R}^d). The bound then follows from (30). \square

When $d = 1$, it is known [57] that $W_p(F, G)$ may be written as follows:

$$W_p(F, G) = \|(\mathcal{V}F)^{-1} - (\mathcal{V}G)^{-1}\|_{L^p}. \quad (32)$$

Here, $(\mathcal{V}F)^{-1}$ denotes the functional inverse of $\mathcal{V}F$, defined as

$$(\mathcal{V}F)^{-1}(x) = \inf\{t \in [a, b] : (\mathcal{V}F)(t) \geq x\}. \quad (33)$$

When $p = 1$, it is also well-known that $W_1(F, G) = \|\mathcal{V}F - \mathcal{V}G\|_{L^1}$.

2.8 Sliced Wasserstein distance

The *sliced Wasserstein distance* [51] is defined between two probability densities F and G in \mathbb{R}^d as follows:

$$SW_{p,\eta}(F, G) = \left(\int_{\mathbb{S}^{d-1}} W_p^p(\mathcal{P}_{\mathbf{w}}F, \mathcal{P}_{\mathbf{w}}G) d\eta(\mathbf{w}) \right)^{1/p} \quad (34)$$

where η is a suitable probability measure over \mathbb{S}^{d-1} . That is, $SW_{p,\eta}(F, G)$ is obtained by averaging the distances between the one-dimensional projections of F and G over all directions.

One advantage of these metrics is that each one-dimensional distance $W_p(\mathcal{P}_{\mathbf{w}}F, \mathcal{P}_{\mathbf{w}}G)$ may be computed rapidly by using the formula (32). Sliced Wasserstein distances have been the subject of considerable research activity [51, 9, 31, 32, 16, 48].

3 Volterra distances

In this section, we introduce our basic objects of study, the Volterra distances and their discrete approximations, and the sliced Volterra distances. Section 3.1 defines the Volterra norms and corresponding distances for functions of a single variable; Section 3.2 reviews the variational characterization of these distances; Section 3.3 defines the sliced Volterra distances; and Section 3.4 defines the trapezoidal rule approximations to the Volterra distances.

3.1 The Volterra norms and distances

Let f be in $L^1([a, b])$. For any value $1 \leq p \leq \infty$, we define the following norm, which we will refer to as the *Volterra p -norm*:

$$\|f\|_{V^p} = \|\mathcal{V}f\|_{L^p}. \quad (35)$$

Concretely, when $1 \leq p < \infty$,

$$\|f\|_{V^p} = \left(\int_a^b \left| \int_a^x f(t) dt \right|^p dx \right)^{1/p}, \quad (36)$$

and when $p = \infty$,

$$\|f\|_{V^\infty} = \operatorname{ess\,sup}_{a \leq x \leq b} \left| \int_a^x f(t) dt \right|. \quad (37)$$

Note that, because $\mathcal{V}f$ is in $L^\infty([a, b])$, the Volterra p -norm of f is finite for any function f in $L^1([a, b])$. If f and g are two functions in $L^1([a, b])$, we will refer to $\|f - g\|_{V^p}$ as the *Volterra p -distance*, or *Volterra p -metric*, between f and g .

Remark 1. When $p = \infty$ and f and g are two probability densities, the Volterra ∞ -distance is known as the Kolmogorov Metric between f and g [24]: $\operatorname{KM}(f, g) = \|f - g\|_{V^\infty}$. The KM arises in the context of goodness-of-fit testing in statistics [23].

Remark 2. When $p = 1$ and f and g are two probability densities, the Volterra 1-distance is equal to the Earth Mover's Distance between f and g described in Section 2.7: $\operatorname{EMD}(f, g) = \|f - g\|_{V^1}$. When $p > 1$, however, the p -Wasserstein distance $W_p(f, g)$ is equal to $\|(\mathcal{V}f)^{-1} - (\mathcal{V}g)^{-1}\|_{V^p}$, and is generally *not* equal to the Volterra p -distance $\|\mathcal{V}f - \mathcal{V}g\|_{L^p}$.

3.2 Variational formulation of $\|f\|_{V^p}$

The following result is an alternate formulation of the Volterra norm that will be useful for analysis. It essentially appears as Theorem 1 in [41]; we provide a self-contained proof for the reader's convenience.

Proposition 3.1. *Let $1 \leq p \leq \infty$ and let q be the conjugate exponent:*

$$\frac{1}{p} + \frac{1}{q} = 1. \quad (38)$$

Then for any function f in $L^1([a, b])$,

$$\|f\|_{V^p} = \sup_{G \in \mathcal{A}_0: \|G'\|_{L^q} \leq 1} \langle f, G \rangle. \quad (39)$$

Proof of Proposition 3.1. By duality of L^p and L^q , we have:

$$\|f\|_{V^p} = \|\mathcal{V}f\|_{L^p} = \sup_{g: \|g\|_{L^q} \leq 1} \int_a^b (\mathcal{V}f)(x)g(x)dx = \sup_{g: \|g\|_{L^q} \leq 1} \langle \mathcal{V}f, g \rangle = \sup_{g: \|g\|_{L^q} \leq 1} \langle f, \mathcal{V}^*g \rangle. \quad (40)$$

Any function of the form \mathcal{V}^*g is contained in \mathcal{A}_0 , and any function G in \mathcal{A}_0 is of the form $G = \mathcal{V}^*g$ where $g = G'$ almost everywhere. Consequently:

$$\|f\|_{V^p} = \sup_{g: \|g\|_{L^q} \leq 1} \langle f, \mathcal{V}^*g \rangle = \sup_{G \in \mathcal{A}_0: \|G'\|_{L^q} \leq 1} \langle f, G \rangle, \quad (41)$$

which completes the proof. \square

3.3 Sliced Volterra distances

Analogous to the sliced Wasserstein distances reviewed in Section 2.8, we define the *sliced Volterra distance* between functions f and g on \mathbb{R}^d :

$$SV_{p,\eta}(f, g) = \left(\int_{\mathbb{S}^{d-1}} \|\mathcal{P}_{\mathbf{w}}f - \mathcal{P}_{\mathbf{w}}g\|_{V^p}^p d\eta(\mathbf{w}) \right)^{1/p}, \quad (42)$$

where η is a probability measure over \mathbb{S}^{d-1} . That is, $SV_{p,\eta}(f, g)$ is obtained by averaging the Volterra p -distances between the one-dimensional projections of f and g over all directions.

3.4 Trapezoidal rule approximation to $\|f\|_{V^p}$

Suppose f is a function on $[a, b]$, and we are given samples of f on an equispaced grid of points in $[a, b]$, from which we wish to approximate $\|f\|_{V^p}$. That is, let $a_0 < a_1 < \dots < a_n$ be equispaced points in $[a, b]$ defined by (14), that is,

$$a_k = a + \frac{k}{n}(b - a), \quad 0 \leq k \leq n. \quad (43)$$

Note that $a_0 = a$ and $a_n = b$. We suppose we are given the values of $f(a_k)$, $0 \leq k \leq n$, or possibly noisy estimates of these, and wish to approximate $\|f\|_{V^p}$.

To this end, we introduce some convenient notation. If \mathbf{v} is a vector in \mathbb{R}^{n+1} and $1 \leq p < \infty$, we define the norm

$$\|\mathbf{v}\|_{\tau_p} = \frac{b-a}{2n} \sum_{k=0}^{n-1} (|\mathbf{v}[k]|^p + |\mathbf{v}[k+1]|^p) = \left(\frac{b-a}{n} \sum_{k=1}^{n-1} |\mathbf{v}[k]|^p + \frac{b-a}{2n} |\mathbf{v}[0]|^p + \frac{b-a}{2n} |\mathbf{v}[n]|^p \right)^{1/p}. \quad (44)$$

When $p = \infty$, we define $\|\mathbf{v}\|_{\tau_\infty} = \|\mathbf{v}\|_{\ell_\infty}$, that is,

$$\|\mathbf{v}\|_{\tau_\infty} = \max_{0 \leq k \leq n} |\mathbf{v}[k]|. \quad (45)$$

If f is a function on $[a, b]$ and \mathbf{f} is a vector with entries $\mathbf{f}[k] = f(a_k)$, $0 \leq k \leq n$, then $\|\mathbf{f}\|_{\tau_p}$ is the trapezoidal rule approximation to $\|f\|_{L^p}$ when $1 \leq p < \infty$, and $\|\mathbf{f}\|_{\tau_\infty}$ is an approximation to $\|f\|_{L^\infty}$.

Suppose n is a positive integer. We define the following discrete Volterra operator $\mathbf{V} : \mathbb{R}^{n+1} \rightarrow \mathbb{R}^{n+1}$ on a vector \mathbf{x} by $(\mathbf{V}\mathbf{x})[0] = 0$, and, for $1 \leq k \leq n$,

$$(\mathbf{V}\mathbf{x})[k] = \frac{b-a}{2n} \sum_{j=0}^{k-1} (\mathbf{x}[j] + \mathbf{x}[j+1]). \quad (46)$$

We then define the discrete Volterra p -norm of \mathbf{x} as

$$\|\mathbf{x}\|_{\nu_p} = \|\mathbf{V}\mathbf{x}\|_{\tau_p}. \quad (47)$$

The interpretation of this quantity may be understood as follows. Suppose f is a function on $[a, b]$, and let \mathbf{f} be the vector in \mathbb{R}^{n+1} with entries $\mathbf{f}[k] = f(a_k)$, for $0 \leq k \leq n$, where a_k are defined in (14). Then $(\mathbf{V}\mathbf{f})[k]$ is the trapezoidal rule approximation to $(\mathcal{V}f)(a_k)$, and $\|\mathbf{f}\|_{\nu_p}$ approximates $\|f\|_{V^p}$. The approximation error will be bounded in Section 5. We remark that one can also define an approximate Volterra norm based on n equispaced midpoint samples on $[a, b]$, which will have the same theoretical guarantees as the trapezoidal rule approximation considered here.

4 Properties of Volterra distances

This section analyzes the use of Volterra distances for comparing functions on an interval and univariate projections of functions on \mathbb{R}^d . Theorem 4.1 establishes a property similar to Theorem 2.2 for Wasserstein

distances, namely, that the distance between projections of a function and its ϵ -deformation is controlled by ϵ . Notably, however, the bound for Volterra metrics is a concave, non-linear function of ϵ , indicating that the Volterra metrics may be more robust than the Wasserstein distances to large deformations. The bound in Theorem 4.1 can be strengthened in certain special cases: Theorem 4.2 analyzes the case of rotations in the plane, and Theorem 4.3 analyzes the setting of univariate functions related by a monotonically increasing deformation. Theorems 4.4 and 4.5 state corresponding results for the sliced Volterra distances.

Theorem 4.1. *Let A and B be non-empty, bounded, open sets in \mathbb{R}^d , with $r = \text{diam}(A \cup B)$. Let $F : \mathbb{R}^d \rightarrow \mathbb{R}$ be in L^p and supported on A , $\Phi : B \rightarrow A$ be an ϵ -deformation, and $F_\Phi = \Phi_{\#}^{-1}F$, i.e.*

$$F_\Phi(\mathbf{x}) = F(\Phi(\mathbf{x})) |\det(\nabla \Phi(\mathbf{x}))| \quad (48)$$

on B , and 0 elsewhere. Then for any $\mathbf{u} \in \mathbb{S}^{d-1}$,

$$\|\mathcal{P}_{\mathbf{u}}F - \mathcal{P}_{\mathbf{u}}F_\Phi\|_{V^p} \leq \min \left\{ \epsilon \cdot K_{p,\Phi}(F) \cdot C_{p,d}(r), \epsilon^{1/p} \cdot \|F\|_{L^1} \right\}, \quad (49)$$

where

$$K_{p,\Phi}(F) = \min \{ \|F\|_{L^p}, \|F_\Phi\|_{L^p} \}, \quad (50)$$

and

$$C_{p,d}(r) = 2^{(p-1)/p} r^{(d-1)(p-1)/p}. \quad (51)$$

Remark 3. Theorem 4.1 states that Volterra distances between a projection of a function and a projection of its ϵ -deformation may be bounded by an increasing function of ϵ . However, when $p > 1$ the bound on the Volterra distances becomes either constant or strictly concave for large ϵ , indicating that the Volterra distances are more robust to large deformations.

Next, we will consider special cases of the deformation Φ , for which tighter bounds can be shown.

4.1 Changes in projection angle

In this section, $F : \mathbb{R}^2 \rightarrow \mathbb{R}$ will denote a function in $L^p(\mathbb{B}_R)$, where $\mathbb{B}_R \subset \mathbb{R}^2$ is the disc of radius R and center $(0, 0)$. For a given angle θ , if $\mathbf{u} = (\cos(\theta), \sin(\theta))$, we let $f_\theta = \mathcal{P}_{\mathbf{u}}F$; that is,

$$f_\theta(x) = \int_{\mathbb{R}} F(\cos(\theta)x + \sin(\theta)y, \cos(\theta)y - \sin(\theta)x) dy \quad (52)$$

We then have the following result:

Theorem 4.2. *Let θ and φ be real numbers, with $\delta \equiv |\theta - \varphi| < \pi$. Then for all $1 \leq p \leq \infty$,*

$$\|f_\theta - f_\varphi\|_{V^p} \leq (2 \sin(\delta/2))^{1/p} \cdot \min \left\{ \delta^{1-1/p} \cdot \|F\|_{L^p} \cdot R^{2-1/p}, \|F\|_{L^1} \cdot R^{1/p} \right\}. \quad (53)$$

The proof of Theorem 4.2 may be found in Section 7.3.

Remark 4. When $p = 1$, then scaling the problem so that $R = 1$ and $\|F\|_{L^1} = 1$, Theorem 4.2 states

$$\|f_\theta - f_\varphi\|_{V^p} \leq 2 \sin(\delta/2), \quad (54)$$

which matches the bound for Wasserstein distances from [52]. For all values of p , under the same scaling, the first term of the right side of (53) yields the upper bound

$$\|f_\theta - f_\varphi\|_{V^p} \leq \delta \cdot \|F\|_{L^p}, \quad (55)$$

since $2 \sin(|\theta - \varphi|/2) \leq \delta$.

4.2 Monotonically increasing deformations

The bound in Theorem 4.1 can be sharpened by a constant factor depending on p for univariate functions in the case where the deformation Φ is monotonically increasing.

Theorem 4.3. *Suppose $f : \mathbb{R} \rightarrow \mathbb{R}$ is in $L^p(I)$, where I is a closed interval. Let $\Phi : J \rightarrow I$ be an ϵ -deformation with $\Phi'(x) > 0$ for all x . Let $f_\Phi(x) = (\Phi_\#^{-1}f)(x) = f(\Phi(x))\Phi'(x)$ on J , and 0 elsewhere. Then*

$$\|f - f_\Phi\|_{V^p} \leq \min \left\{ \epsilon \cdot K_{p,\Phi}(f), \epsilon^{1/p} \cdot \|f\|_{L^1} \right\}, \quad (56)$$

where

$$K_{p,\Phi}(f) = \min\{\|f\|_{L^p}, \|f_\Phi\|_{L^p}\}. \quad (57)$$

Remark 5. Note that the factor of $2^{(p-1)/p}$ from Theorem 4.1 is not present in Theorem 4.3, due to the fact that Φ is increasing. To see that the monotonicity of Φ is required for this sharper bound, consider the following example. Fix $\eta > \delta > 0$, and let f be defined by

$$f(x) = \begin{cases} 1, & \text{if } -\eta \leq x < 0, \\ -1, & \text{if } 0 \leq x \leq \eta, \\ 0, & \text{otherwise.} \end{cases} \quad (58)$$

Let $\Phi : [-\delta, \delta] \rightarrow [-\eta, \eta]$ be defined by $\Phi(x) = -(\eta/\delta)x$. Then

$$f_\Phi(x) = \begin{cases} -\eta/\delta, & \text{if } -\delta \leq x < 0, \\ \eta/\delta, & \text{if } 0 \leq x \leq \delta, \\ 0, & \text{otherwise.} \end{cases} \quad (59)$$

Then it is straightforward to verify that $\epsilon = \eta + \delta$, $K_{p,\Phi}(f) = 1$, and $\|f - f_\Phi\|_{V^\infty} = 2\eta$. Hence, the right side of (56), with $p = \infty$, is $\epsilon = \eta + \delta$, which is not bigger than $\|f - f_\Phi\|_{V^\infty} = 2\eta$.

4.3 Applications to sliced Volterra distances

The bounds from Theorems 4.1 and 4.2 immediately yield bounds for the sliced Volterra distances.

Theorem 4.4. *Let A and B be non-empty, bounded, open sets in \mathbb{R}^d , with $r = \text{diam}(A \cup B)$. Let $F : \mathbb{R}^d \rightarrow \mathbb{R}$ be in $L^p(A)$, $\Phi : B \rightarrow A$ be an ϵ -deformation, and $F_\Phi = \Phi_\#^{-1}F$, i.e.*

$$F_\Phi(\mathbf{x}) = F(\Phi(\mathbf{x})) |\det(\nabla \Phi(\mathbf{x}))| \quad (60)$$

on B , and 0 elsewhere. Let $1 \leq p \leq \infty$ and fix any probability distribution η over \mathbb{S}^{d-1} . Then

$$\text{SV}_{p,\eta}(F, F_\Phi) \leq \min \left\{ \epsilon \cdot K_{p,\Phi}(F) \cdot C_{p,d}(r), \epsilon^{1/p} \cdot \|F\|_{L^1} \right\}, \quad (61)$$

where

$$K_{p,\Phi}(F) = \min\{\|F\|_{L^p}, \|F_\Phi\|_{L^p}\}, \quad (62)$$

and

$$C_{p,d}(r) = 2^{(p-1)/p} r^{(d-1)(p-1)/p}. \quad (63)$$

Theorem 4.5. *Let $F : \mathbb{R}^2 \rightarrow \mathbb{R}$ be in $L^p(\mathbb{B}_R)$. Suppose $0 \leq \delta < \pi$, and define F_δ by*

$$F_\delta(x, y) = F(x \cos(\delta) + y \sin(\delta), y \cos(\delta) - x \sin(\delta)). \quad (64)$$

Then for all $1 \leq p \leq \infty$ and any probability distribution η over \mathbb{S}^{d-1} ,

$$\text{SV}_{p,\eta}(F, F_\delta) \leq (2 \sin(\delta/2))^{1/p} \cdot \min \left\{ \delta^{1-1/p} \cdot \|F\|_{L^p} \cdot R^{2-1/p}, \|F\|_{L^1} \cdot R^{1/p} \right\}. \quad (65)$$

5 Asymptotic behavior of the discrete norms

In this section, we consider the behavior of the discrete Volterra norms, defined in Section 3.4, for vectors consisting of samples of a function f on $[a, b]$ from an equispaced grid. It will be convenient to introduce some notation. Let n be a positive integer, and define $a_0 < a_1 < \dots < a_n$ as in (14), namely

$$a_k = a + \frac{k}{n}(b - a), \quad 0 \leq k \leq n. \quad (66)$$

Note that $a_0 = a$ and $a_n = b$. Let \mathbf{f} be the vector in \mathbb{R}^{n+1} with entries $\mathbf{f}[k] = f(a_k)$, for $0 \leq k \leq n$.

Denote the mean of f on $[a, b]$ by

$$\mu(f) = \frac{1}{b-a} \int_a^b f(t) dt, \quad (67)$$

and let $f_{\text{cen}}(x) = f(x) - \mu(f)$.

If \mathbf{w} is a vector in \mathbb{R}^{n+1} , denote its trapezoidal mean by

$$\mathbf{m}(\mathbf{w}) = \frac{1}{2n} \sum_{k=0}^{n-1} (\mathbf{w}[k] + \mathbf{w}[k+1]), \quad (68)$$

and let \mathbf{w}_{cen} in \mathbb{R}^{n+1} have entries $\mathbf{w}_{\text{cen}}[k] = \mathbf{w}[k] - \mathbf{m}(\mathbf{w})$.

5.1 Convergence rates for well-behaved functions

We first prove rates on the convergence of the discrete approximation $\|\mathbf{f}\|_{\nu_p}$ to the Volterra norm $\|f\|_{V^p}$, where f is a reasonably well-behaved function. Theorem 5.1 establishes a convergence rate of $O(1/n)$ for piecewise Lipschitz functions, whereas Theorem 5.2 establishes the faster rate of $O(1/n^2)$ for smoother functions.

Theorem 5.1. *Suppose $a = c_0 < c_1 < \dots < c_r = b$, and f has Lipschitz constant bounded by $L > 0$ on each interval (c_j, c_{j+1}) , $0 \leq j \leq r-1$. Then for all $1 \leq p \leq \infty$,*

$$\left| \|\mathbf{f}\|_{\nu_p} - \|f\|_{V^p} \right| \leq C \frac{(b-a)^{1+1/p}}{n} (L(b-a) + r\|f\|_{L^\infty}), \quad (69)$$

where $C > 0$ is a universal constant. The same bound holds by replacing f with f_{cen} and \mathbf{f} with \mathbf{f}_{cen} .

In other words, for piecewise Lipschitz functions, the discrete Volterra norm based on n subintervals converges to the true Volterra norm at a rate of $O(1/n)$.

The proof of Theorem 5.1 may be found in Section 8.1. If instead of being merely piecewise Lipschitz, the function f is C^2 and not too oscillatory, then the discrete Volterra norms give a higher order approximation to the Volterra norms of f :

Theorem 5.2. *Suppose f is a two times continuously differentiable function on $[a, b]$ with only finitely many zero crossings. Let $1 \leq p \leq \infty$. Then for all n sufficiently large,*

$$\left| \|\mathbf{f}\|_{\nu_p} - \|f\|_{V^p} \right| \leq C \frac{K(f, p, a, b)}{n^2}, \quad (70)$$

where

$$K(f, p, a, b) = (b-a)^{3+1/p} \|f''\|_{L^\infty} + (b-a)^{2+1/p} \|f'\|_{L^\infty} + (b-a)^2 |f(b)| \left(\frac{|\mu(f)|}{\|f\|_{V^p}} \right)^{p-1} \quad (71)$$

when $1 \leq p < \infty$, and

$$K(f, \infty, a, b) = (b-a)^3 \|f''\|_{L^\infty} + (b-a)^2 \|f'\|_{L^\infty}, \quad (72)$$

and where $C > 0$ is a universal constant.

An analogous bound holds by replacing f with f_{cen} and \mathbf{f} with \mathbf{f}_{cen} :

$$\left| \|\mathbf{f}_{\text{cen}}\|_{\nu_p} - \|f_{\text{cen}}\|_{V^p} \right| \leq C \frac{(b-a)^{3+1/p} \|f''\|_{L^\infty} + (b-a)^{2+1/p} \|f'\|_{L^\infty}}{n^2}, \quad (73)$$

for all $1 \leq p \leq \infty$.

In other words, for such functions f , the discrete Volterra norm based on n subintervals converges to the true Volterra norm at a rate of $O(1/n^2)$. The proof of Theorem 5.2 may be found in Section 8.2.

5.2 Robustness to Gaussian noise

Next, we show that the discrete Volterra metrics are robust to additive noise. More precisely, as the number n of subintervals on which samples are taken grows, the effects of additive Gaussian noise on the samples of f vanish at a predictable rate.

Theorem 5.3. *Let $\sigma_0, \sigma_1, \dots, \sigma_n, \dots$ be a sequence of positive numbers, and let $Z = (Z[0], \dots, Z[n])$, where $Z[0], Z[1], \dots, Z[n], \dots$ are independent with $Z[j] \sim N(0, \sigma_j^2)$. Suppose too that $\sigma > 0$ satisfies*

$$\frac{1}{n} \sum_{j=1}^n \sigma_j^2 \leq \sigma^2, \quad (74)$$

for all n . Let $t > 0$. Then for all $1 \leq p \leq \infty$,

$$\mathbb{P} \{ \|Z\|_{\nu_p} \geq t \} \leq A e^{-Bt^2 n / \sigma^2}, \quad (75)$$

where $A > 0$ and $B > 0$ are universal constants;

$$\lim_{n \rightarrow \infty} \|Z\|_{\nu_p} = 0 \quad (76)$$

almost surely; and

$$\mathbb{E} \|Z\|_{\nu_p} \leq C \frac{\sigma}{\sqrt{n}}, \quad (77)$$

where $C > 0$ is a universal constant. Furthermore, (75), (76) and (77) hold with Z replaced by Z_{cen} .

Corollary 5.4. *Suppose f satisfies the conditions of Theorem 5.1, Z satisfies the conditions of Theorem 5.3, and $Y = \mathbf{f} + Z$. Let $t > 0$ and $1 \leq p \leq \infty$. Then for all n sufficiently large,*

$$\mathbb{P} \{ \left| \|Y\|_{\nu_p} - \|f\|_{V^p} \right| \geq t \} \leq A e^{-Bt^2 n / \sigma^2}, \quad (78)$$

where $A > 0$ and $B > 0$ are universal constants;

$$\lim_{n \rightarrow \infty} \|Y\|_{\nu_p} = \|f\|_{V^p} \quad (79)$$

almost surely; and

$$\mathbb{E} \left| \|Y\|_{\nu_p} - \|f\|_{V^p} \right| \leq C \frac{\sigma}{\sqrt{n}}, \quad (80)$$

where $C > 0$ is a universal constant. Furthermore, (78), (79) and (80) hold with f replaced by f_{cen} and Y replaced by Y_{cen} .

The proofs of Theorem 5.3 and Corollary 5.4 are provided in Section 8.3.

Remark 6. In the setting of Corollary 5.4, both the signal vector \mathbf{f} and the noise vector Z have comparable p -norms; consequently, $\|Y\|_{\ell_p}$ does not approach $\|f\|_{L^p}$ as $n \rightarrow \infty$. For example, if $\sigma_j = \sigma$ for all j , then almost surely

$$\lim_{n \rightarrow \infty} \|Y\|_{\ell_2}^2 = \|f\|_{L^2}^2 + \sigma^2. \quad (81)$$

By contrast, (79) states that Z has a negligible effect on the Volterra norm when n is large.

6 Numerical results

In this section, we show the results of numerical experiments comparing the Volterra distances to Wasserstein and Lebesgue distances.

Define the approximate Wasserstein distance as follows. Suppose f is a probability density on $[a, b]$, and let $\mathbf{f} = (f(a_0), \dots, f(a_n))$, where

$$a_k = a + \frac{k}{n}(b - a), \quad 0 \leq k \leq n. \quad (82)$$

We define, for $0 < t < 1$, the approximate inverse CDF by

$$\widehat{(\mathcal{V}f)^{-1}}(t) = \min\{a_k : (\mathbf{V}\mathbf{f})[k] \geq t\}. \quad (83)$$

For integer m , we define the midpoint grid values $t_j = (j - 1/2)/m$, $1 \leq j \leq m$. Let $\widehat{(\mathbf{V}\mathbf{f})^{-1}}[j] = \widehat{(\mathcal{V}f)^{-1}}(t_j)$. Then the approximate p -Wasserstein distance between densities f and g is defined by

$$\widehat{W}_p(f, g) = \|\widehat{(\mathbf{V}\mathbf{f})^{-1}} - \widehat{(\mathbf{V}\mathbf{g})^{-1}}\|_{\ell_p s}. \quad (84)$$

In all the experiments reported below, we set $m = n$.

6.1 Distances under translation

We illustrate Theorem 4.3 on the functions shown in Figure 1; these are translations $f_\epsilon(x) = f(x - \epsilon)$ of the function f on $[0, 1]$ defined by

$$f(x) = C \cos(10x - 1)e^{-16(10x - 3/2)^2}, \quad (85)$$

where C is chosen so that the integral of f is 1. Figure 2 shows the estimated Volterra p -distances (top row), p -Wasserstein (middle row), and Lebesgue p -distances (bottom row) based on samples from $n = 500$ subintervals, plotted as functions of the shift size ϵ .

The Volterra distances exhibit the behavior described by the bound in Theorem 4.3, namely, the distances grow as concave functions of the shift size. When $p = 1$ and $p = 2$, the distances continue to grow with the translation, whereas when $p = \infty$ the distances level off, consistent with the upper bound from Theorem 4.3. The Wasserstein distances all increase linearly with the shift size (which can be seen easily from the Monge formulation of the Wasserstein distance). By contrast with these behaviors, all of the Lebesgue distances quickly saturate to a constant value, independent of ϵ , as soon as the translation is big enough so that the numerical supports of the translated functions do not overlap.

6.2 Distances under dilation

Next, we consider the function f defined on $[0, 1]$ by

$$f(x) = Cx^6(1 - x) \quad (86)$$

where C is chosen so that the integral of f is 1. We consider the family of dilations of f parameterized by $\eta \geq 1$; these are the functions f_η defined by $f_\eta(x) = f(\eta x)\eta$ on $[0, 1/\eta]$, and $f_\eta(x) = 0$ elsewhere. The size ϵ of the dilation is

$$\epsilon = 1 - \frac{1}{\eta}. \quad (87)$$

Figure 3 shows the function f and some of its dilates. Figure 3 shows the estimated Volterra p -distances (top row), p -Wasserstein (middle row), and Lebesgue p -distances (bottom row) based on samples from $n = 500$ subintervals, plotted as functions of ϵ .

The Volterra distances exhibit the behavior described by the bound in Theorem 4.3, namely, the distances grow as concave functions of the deformation size ϵ . When $p = 1$ and $p = 2$, the distances continue to grow

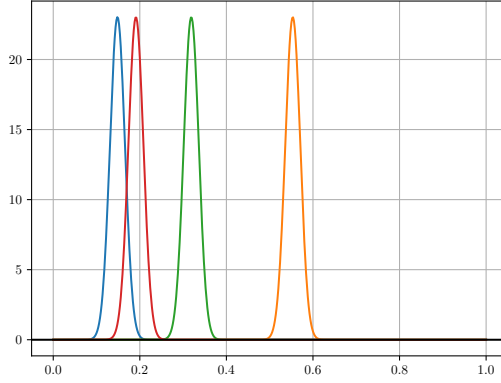


Figure 1: The function (85) (far left, in blue) and its translations, used in the experiment from Section 6.1.

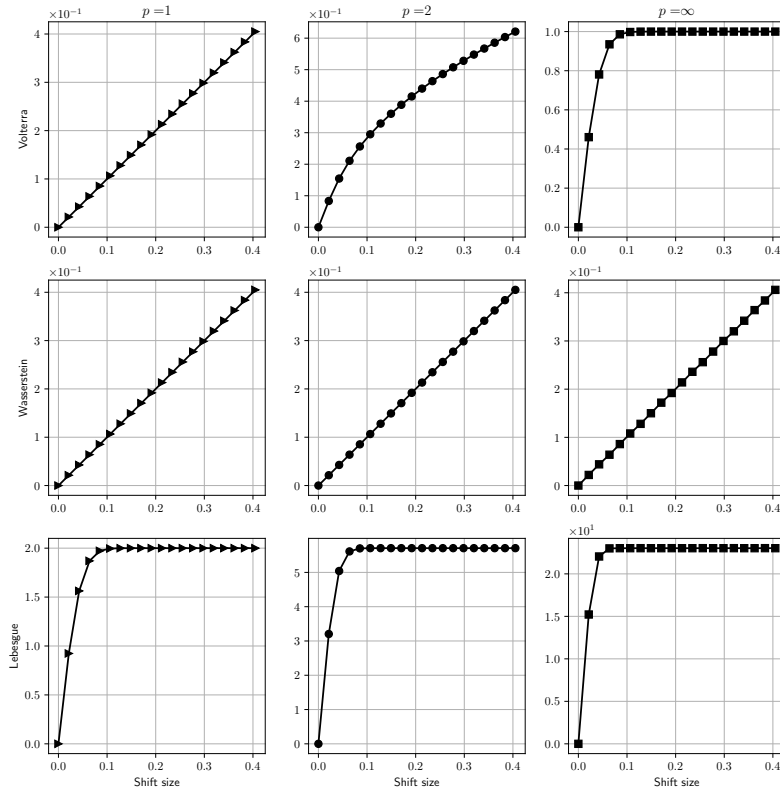


Figure 2: The first row shows the approximated Volterra distances (based on $n = 500$ subintervals) between the function (85) and its shifts, as a function of the shift size. The second row shows the approximated Wasserstein distances, and the third row shows the approximated Lebesgue distances. The values of p (from left to right) are $p = 1, 2, \infty$. See Section 6.1 for details.

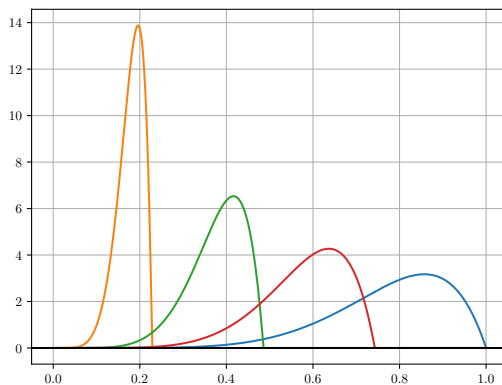


Figure 3: The function (86) (in blue) and its dilations, used in the experiment from Section 6.2.

as ϵ grows, whereas when $p = \infty$ the distances level off, consistent with the upper bound from Theorem 4.3. The Wasserstein distances all increase linearly with the deformation size (which can be seen easily from the Monge formulation of the Wasserstein distance). Because the transformation preserves the integral of f , the L^1 distance levels off when the dilation size is big, since the supports of the function and its dilate are almost disjoint. By contrast, the L^2 and L^∞ distances grow rapidly for large dilation sizes. This is because $\|f_\eta\|_{L^p}$ diverges as η grows, and hence these distances reflect the size of the individual functions and not the relationship between the functions.

6.3 Distances under powers

Next, we consider the function f defined on $[0, 1]$ by

$$f(x) = Cx(1-x)^4, \quad (88)$$

where C is chosen so that the integral of f is 1. We consider the family of deformations $\Phi(x) = x^\alpha$, where $\alpha \geq 1$; the corresponding transformation of f is the function f_α defined by $f_\alpha(x) = f(x^\alpha)\alpha x^{\alpha-1}$ on $[0, 1]$, and $f_\alpha(x) = 0$ elsewhere. The deformation size ϵ is given by

$$\epsilon = \left(\frac{1}{\alpha}\right)^{\frac{1}{\alpha-1}}. \quad (89)$$

Figure 5 shows the function f and its deformations. Figure 5 shows the estimated Volterra p -distances (top row), p -Wasserstein (middle row), and Lebesgue p -distances (bottom row) based on samples from $n = 500$ subintervals, plotted as functions of ϵ .

The Volterra distances exhibit the behavior described by the bound in Theorem 4.3, namely, the distances grow as concave functions of the deformation size ϵ . When $p = 1$ and $p = 2$, the distances continue to grow as ϵ grows, whereas when $p = \infty$ the distances level off, consistent with the upper bound from Theorem 4.3. As in the case of dilations, because the transformation preserves the integral of f , the L^1 distance levels off when the deformation size is big, since the supports of the function and its deformation are almost disjoint. On the other hand, the L^2 and L^∞ distances grow rapidly for large powers.

6.4 Distances between rotated projections I

We illustrate the behavior described by Theorem 4.2 on the function F displayed in Figure 7, given by the formula

$$F(\mathbf{x}) = \frac{1}{2\pi\sigma} \sum_{k=0}^6 h_k e^{-|\mathbf{x}-\mathbf{c}_k|^2/\sigma}, \quad (90)$$

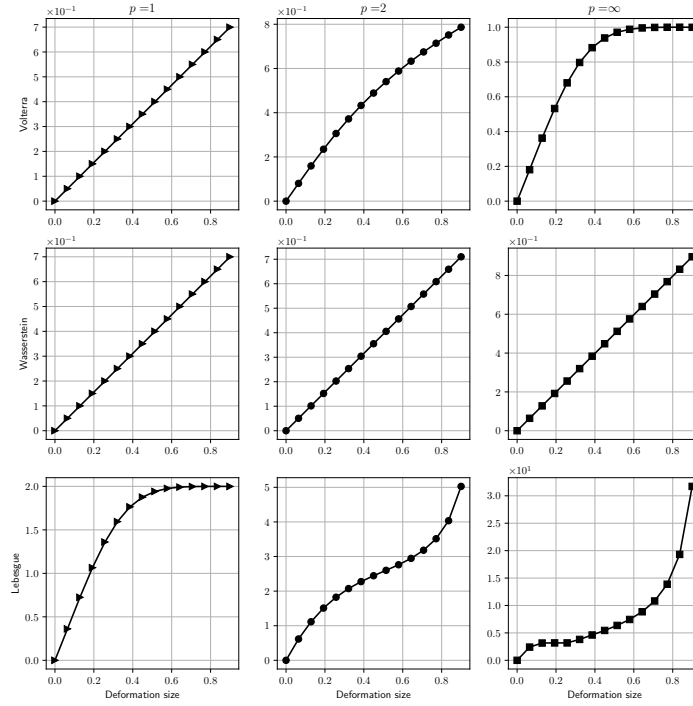


Figure 4: The first row shows the approximated Volterra distances (based on $n = 500$ subintervals) between the function (86) and its dilates, as a function of the deformation size. The second row shows the approximated Wasserstein distances, and the third row shows the approximated Lebesgue distances. The values of p (from left to right) are $p = 1, 2, \infty$. See Section 6.2 for details.

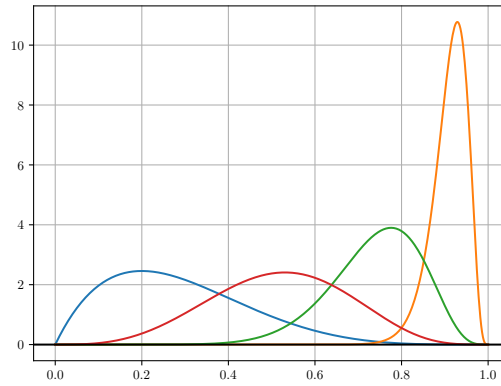


Figure 5: The function (88) (in blue) and its deformations, used in the experiment from Section 6.3.

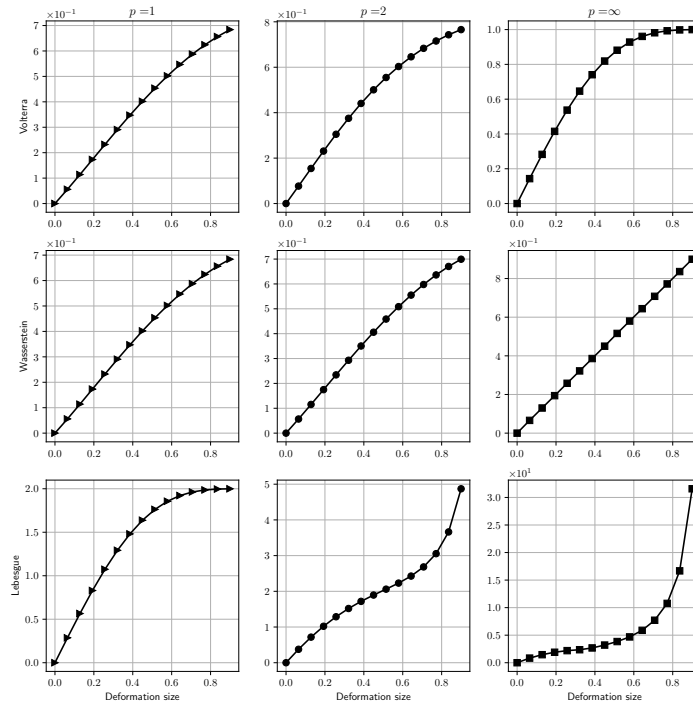


Figure 6: The first row shows the approximated Volterra distances (based on $n = 500$ subintervals) between the function (88) and its deformations, as a function of the deformation size. The second row shows the approximated Wasserstein distances, and the third row shows the approximated Lebesgue distances. The values of p (from left to right) are $p = 1, 2, \infty$. See Section 6.3 for details.

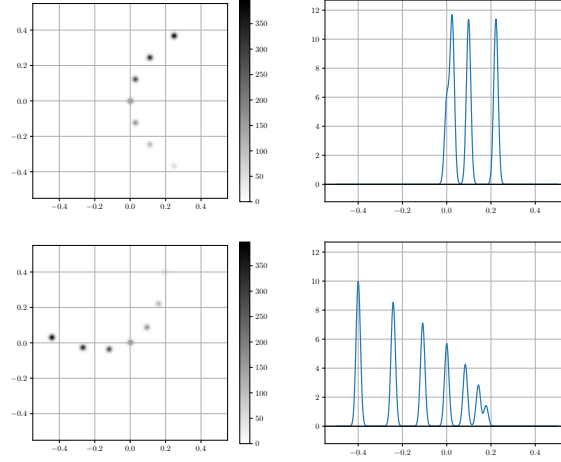


Figure 7: The function F from Section 6.4 (top left) and a rotation (bottom left), with their respective projections onto the x -axis.

where $\sigma = 1/5000$, $\mathbf{c}_k = (x_k, x_k^2)$ and $x_k = -1/3 + k/9$, and $h_k = (k + 1)/24$, $0 \leq k \leq 6$.

We denote by f the projection of F onto the x -axis, and f_θ the projection of F after rotation by θ radians. Figure 7 shows a heatmap of F and a rotation, along with their corresponding projections. Figure 7 shows the estimated Volterra p -distances (top row), p -Wasserstein (middle row), and Lebesgue p -distances (bottom row) based on samples from $n = 500$ subintervals, plotted as functions of the rotation angle. When $p = 2$, the Volterra and Wasserstein distances have very similar behavior, and both vary smoothly with respect to the rotation angle (when $p = 1$, Volterra and Wasserstein are the same, as always). The Lebesgue distances for all p , by contrast, are more irregular functions of the rotation angle.

6.5 Distances between rotated projections II

We next illustrate the behavior described by Theorem 4.2 on the function F displayed in Figure 9, which consists of two nested rings of Gaussian bumps, given by the formula

$$F(\mathbf{x}) = \frac{h}{2\pi\sigma} \sum_{k=0}^4 e^{-|\mathbf{x}-\mathbf{c}_k|^2/\sigma} + \frac{h}{2\pi\sigma} \sum_{k=0}^6 e^{-|\mathbf{x}-\mathbf{d}_k|^2/\sigma} \quad (91)$$

where $\sigma = 1/4000$, $h = 1/12$, and $\mathbf{c}_k = (\cos(\theta_k), \sin(\theta_k))$ with $\theta_k = 2k\pi/5 + (\sqrt{2} + \sqrt{5} + \sqrt{3})\pi$ when $0 \leq k \leq 4$, and $\mathbf{d}_k = (\cos(\varphi_k), \sin(\varphi_k))$ with $\varphi_k = 2k\pi/7 + (\sqrt{2} + \sqrt{5})\pi$ when $0 \leq k \leq 6$.

We denote by f the projection of F onto the x -axis, and f_θ the projection of F after rotation by θ radians. Figure 9 shows a heatmap of F and a rotation, along with their corresponding projections. Figure 9 shows the estimated Volterra p -distances (top row), p -Wasserstein (middle row), and Lebesgue p -distances (bottom row) based on samples from $n = 500$ subintervals, plotted as functions of the rotation angle. As for the example from Section 6.4, when $p = 2$, the Volterra and Wasserstein distances have very similar behavior, and both vary smoothly with respect to the rotation angle, whereas the Lebesgue distances for all p , by contrast, are irregular functions of the rotation angle.

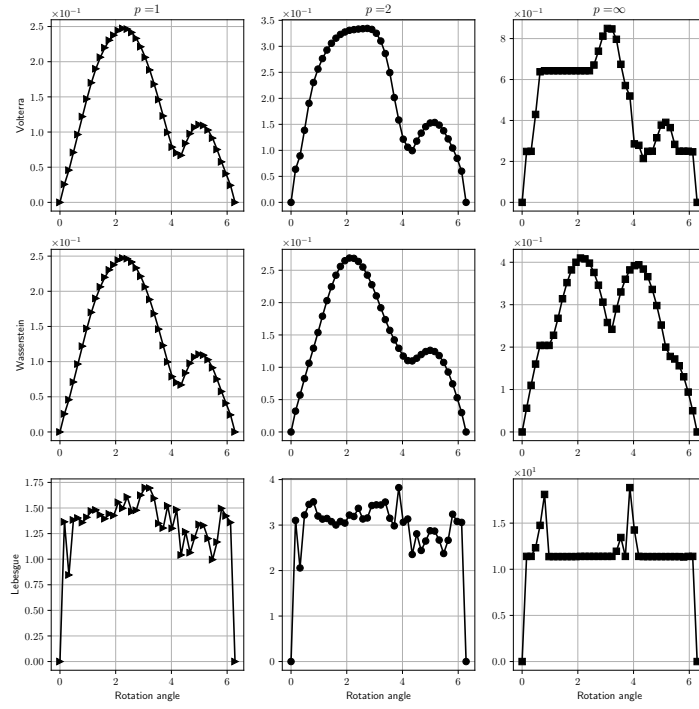


Figure 8: The first row shows the approximated Volterra distances (based on $n = 500$ subintervals) between the projections of F and its rotations, as a function of the rotation angle. The second row shows the approximated Wasserstein distances, and the third row shows the approximated Lebesgue distances. The values of p (from left to right) are $p = 1, 2, \infty$. See Section 6.4 for details.

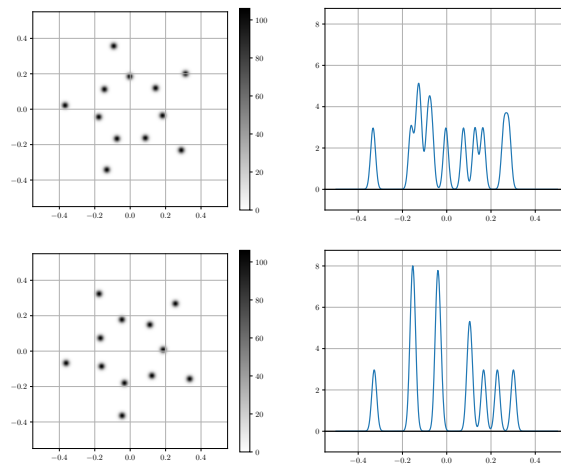


Figure 9: The function F from Section 6.5 (top left) and a rotation (bottom left), with their respective projections onto the x -axis.

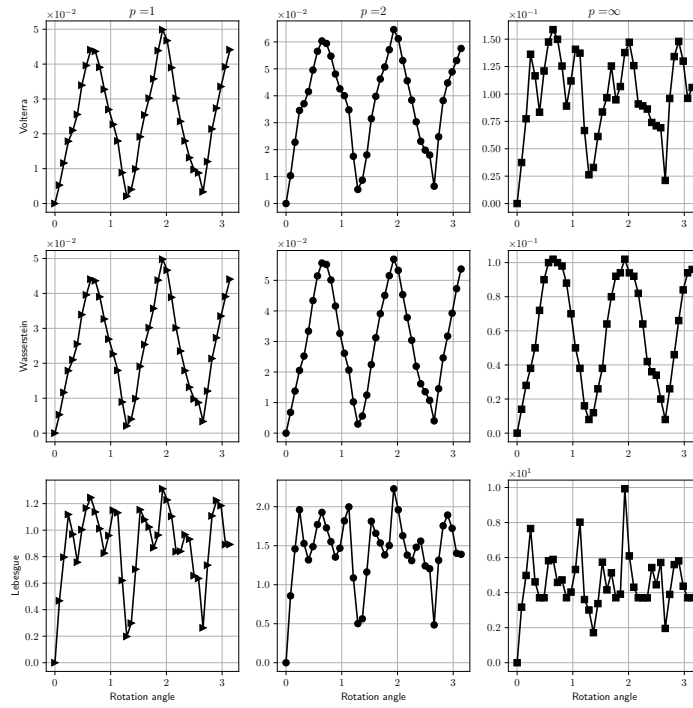


Figure 10: The first row shows the approximated Volterra distances (based on $n = 500$ subintervals) between the projections of F and its rotations, as a function of the rotation angle. The second row shows the approximated Wasserstein distances, and the third row shows the approximated Lebesgue distances. The values of p (from left to right) are $p = 1, 2, \infty$. See Section 6.5 for details.

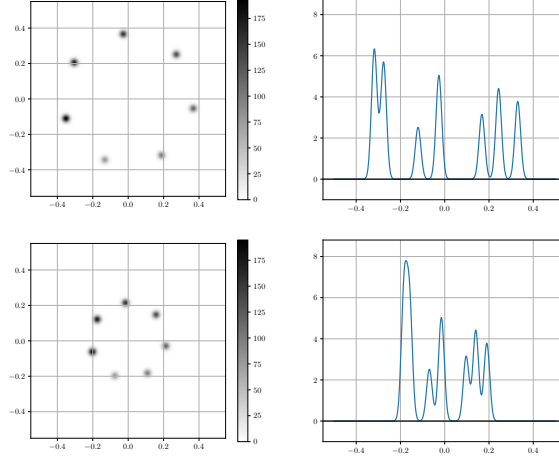


Figure 11: The function F from Section 6.6 (top left) and a shrunken function (bottom left), with their respective projections onto the x -axis.

6.6 Distances under domain shrinking

We next illustrate the behavior of the distances between projections on the function F displayed in Figure 11, given by the formula

$$F(\mathbf{x}) = \frac{1}{2\pi\sigma} \sum_{k=0}^6 h_k e^{-|\mathbf{x}-\mathbf{c}_k|^2/\sigma} \quad (92)$$

where $\sigma = 1/3000$, $h_k = (k+4)/49$, and $\mathbf{c}_k = (\cos(\theta_k), \sin(\theta_k))$ with $\theta_k = 2k\pi/7 + (\sqrt{2} + \sqrt{5} + \sqrt{3})\pi$, $0 \leq k \leq 6$.

The function F is transformed by shrinking the center of each ring towards the origin by an amount ϵ . Figure 11 shows a heatmap of F and a shrunken version, along with their corresponding projections. Figure 11 shows the estimated Volterra p -distances (top row), p -Wasserstein (middle row), and Lebesgue p -distances (bottom row) based on samples from $n = 500$ subintervals, plotted as functions of the shrinkage parameter (proportional to the distance between the Gaussian centers of the original function and the shrunken function). The Wasserstein distances are linear functions of the distance, since each projected Gaussian of the shrunken function is just a translation of the projected Gaussian from the original function. When $p = 2$, the Volterra distance is also a smooth but more concave function. The Lebesgue distances for all p , by contrast, are more irregular functions of the shrinkage parameter.

6.7 Distances under domain squashing

We next illustrate the behavior of the distances between projections on the function F displayed in Figure 13, given by

$$F(\mathbf{x}) = \frac{h}{2\pi\sigma} \sum_{k=0}^{19} e^{-|\mathbf{x}-\mathbf{c}_k|^2/\sigma} \quad (93)$$

where $\sigma = 1/5000$, $h = 1/20$, and $\mathbf{c}_k = (\cos(\theta_k), \sin(\theta_k))$ with $\theta_k = k\pi/10 + (\sqrt{2} + \sqrt{5} + \sqrt{3})\pi$, $0 \leq k \leq 19$.

The function F is transformed by squashing the ring, mapping each center (x, y) to $(\lambda x, y/\lambda)$, and then rotating the result by $\pi/4$. Figure 13 shows a heatmap of F and a squashed version, along with their corresponding projections. Figure 13 shows the estimated Volterra p -distances (top row), p -Wasserstein

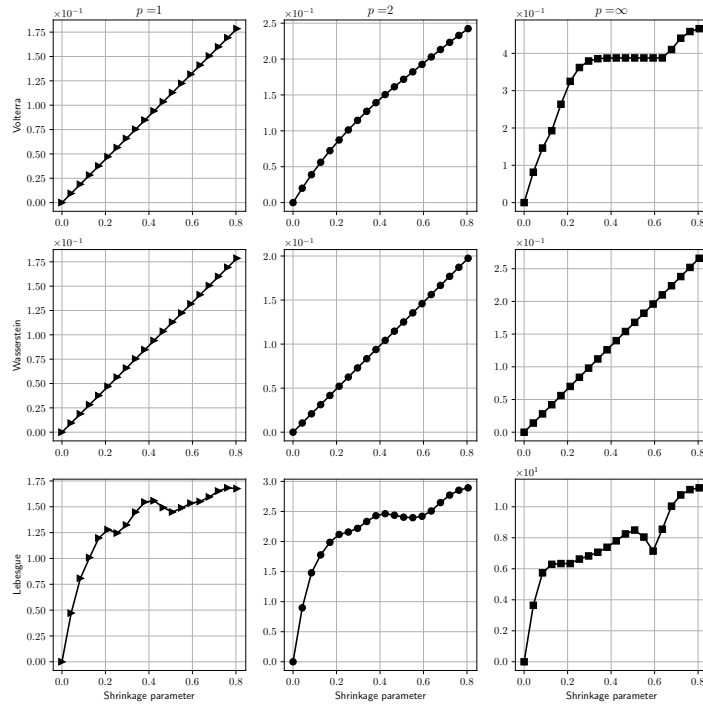


Figure 12: The first row shows the approximated Volterra distances (based on $n = 500$ subintervals) between the projections of F and its shrunken versions, as a function of the shrinkage (proportional to the distance between the Gaussian centers of the original function and the shrunken function). The second row shows the approximated Wasserstein distances, and the third row shows the approximated Lebesgue distances. The values of p (from left to right) are $p = 1, 2, \infty$. See Section 6.6 for details.

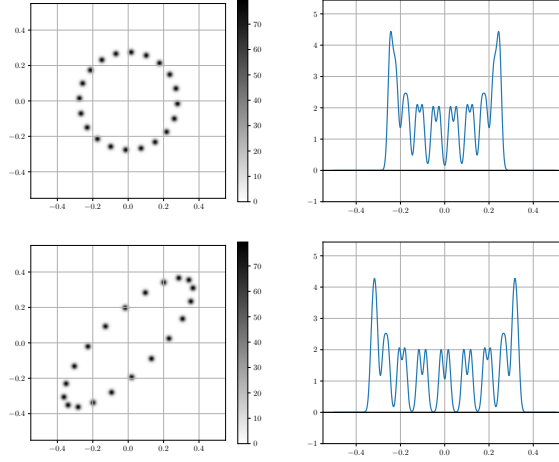


Figure 13: The function F from Section 6.7 (top left) and a squashed function (bottom left), with their respective projections onto the x -axis.

(middle row), and Lebesgue p -distances (bottom row) based on samples from $n = 500$ subintervals, plotted as functions of the squashing parameter λ . When $p = 2$, the Volterra and Wasserstein distances have very similar behavior, and both vary smoothly with respect to the distortion (when $p = 1$, Volterra and Wasserstein are the same, as always). In this example, unlike previous examples, the ∞ -Volterra distance appears to vary more smoothly (for λ close to 1) than the Wasserstein ∞ -distance. As in the other examples, the Lebesgue distances for all p are irregular functions of the distortion.

6.8 Robustness to noise I

To demonstrate the robustness of the Volterra norms under noise described by Corollary 5.4, we run the following experiment. For different values of n , we take a vector \mathbf{f} of $n + 1$ equispaced samples from the function f on $[-1, 1]$ defined by

$$f(x) = xe^{-x^2/4}; \quad (94)$$

A vector Z of iid Gaussian noise with variance .01 is then added to each sample; let $Y = \mathbf{f} + Z$. A plot of a realization of Y , when $n = 512$, is shown in the left panel of Figure 15.

For $p = 1, 2, \infty$, we evaluate the norms $\|Y\|_{\nu_p}$. For each value of n , the experiment is repeated $M = 5000$ times. Denoting the M random signal-plus-noise vectors by Y_1, \dots, Y_M , we record the average absolute error:

$$\text{err}_{n,p} = \frac{1}{M} \sum_{k=1}^M \frac{|\|Y_k\|_{\nu_p} - \|f\|_{V^p}|}{\|f\|_{V^p}}, \quad (95)$$

The right panel of Figure 15 plots $\log_2(\text{err}_{n,p})$ as a function of $\log_2(n)$. The average error decays like $O(1/\sqrt{n})$ as n increases, consistent with Corollary 5.4.

6.9 Robustness to noise II

Next, we examine the simultaneous effects of additive noise and deformation on the distance. We consider the function f defined by

$$f(x) = \sin(2\pi x)\chi_{[0,1/2]}(x), \quad (96)$$

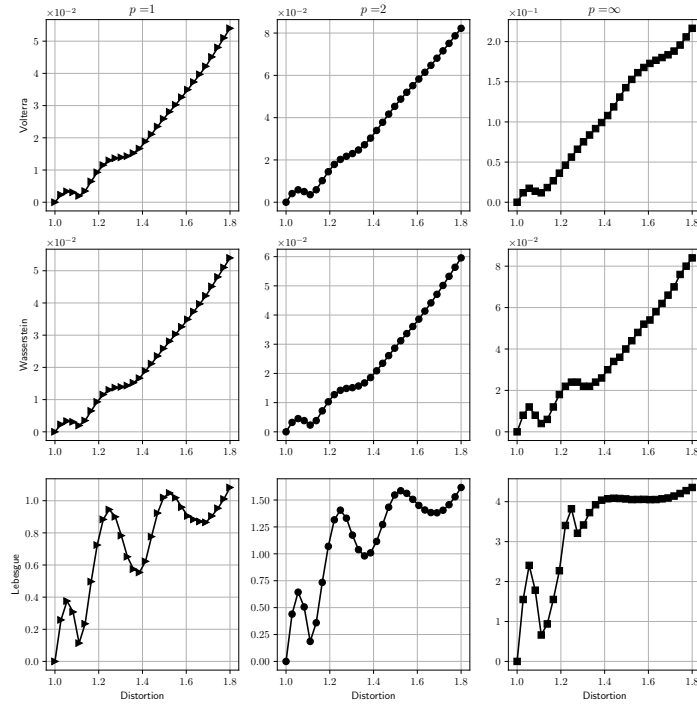


Figure 14: The first row shows the approximated Volterra distances (based on $n = 500$ subintervals) between the projections of F and its squashed versions, as a function of the distortion. The second row shows the approximated Wasserstein distances, and the third row shows the approximated Lebesgue distances. The values of p (from left to right) are $p = 1, 2, \infty$. See Section 6.7 for details.

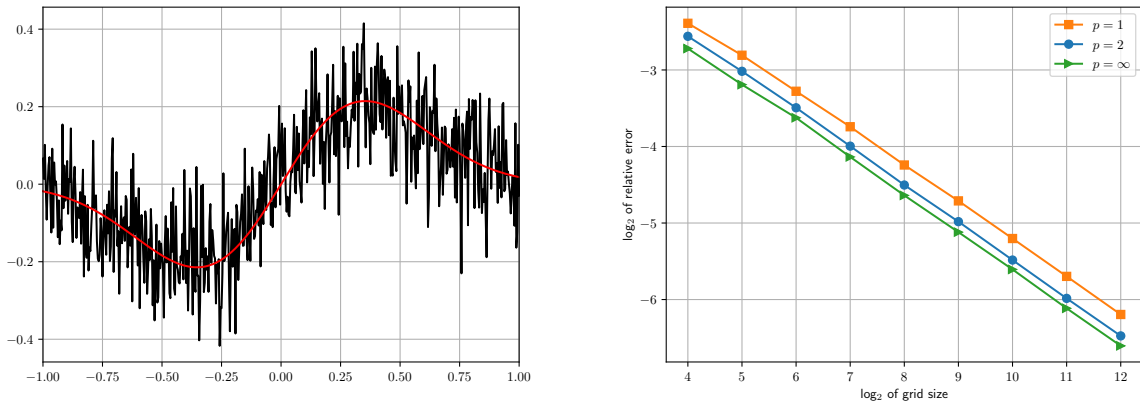


Figure 15: The first panel shows a realization of the noisy draws when $n = 512$, with the noiseless curve graphed in red. The second panel plots $\log_2(\text{err}_{n,p})$ against $\log_2(n)$, for $p = 1, 2, \infty$, where the number of draws is 5000. The slope of each curve is approximately $-1/2$, consistent with the error rate predicted by Corollary 5.4.

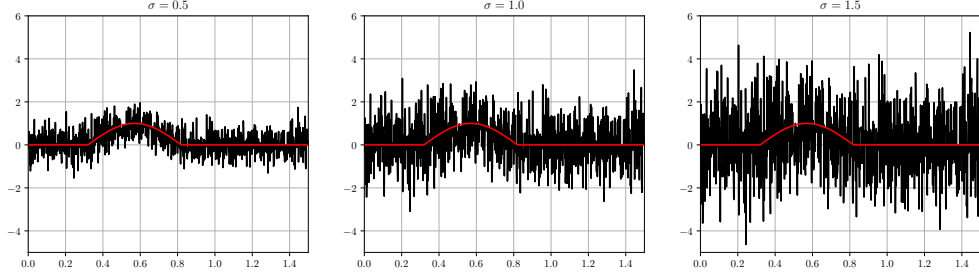


Figure 16: Each panel shows a realization of the noisy curve (in black) from Section 6.9 with the noise levels 0.5, 1.0, and 1.5. The noiseless curve is graphed in red.

and its shifts $f_\epsilon(x) = f(x - \epsilon)$. We sample the functions on a grid with $n = 1000$ subintervals, denoting by \mathbf{x} the vector of samples of f and by \mathbf{x}_ϵ the vector of samples of f_ϵ . For different values of $\sigma \geq 0$, we draw a random vector \mathbf{z} with iid entries $z_j \sim N(0, 1)$ and set $\mathbf{y}_{\epsilon, \sigma} = \mathbf{x}_\epsilon + \sigma \mathbf{z}$. We then compute the distances $\|\mathbf{x} - \mathbf{y}_{\epsilon, \sigma}\|_{\nu_p}$ and $\|\mathbf{x} - \mathbf{y}_{\epsilon, \sigma}\|_{\ell_p}$, for $p = 1, 2, \infty$. The distances are averaged over 5000 independent realizations of the noise vector \mathbf{z} . In Figure 17, we plot these average distances as a function of ϵ , for different noise levels σ . From these plots, we see that the Volterra distances are more robust to noise than the Lebesgue distances, though the effect of the noise is larger for smaller values of p .

7 Proofs from Section 4

We now turn to the proofs of the main theorems from Section 4. Theorems 4.1, 4.2, and 4.3 may all be derived as corollaries of the following general result:

Theorem 7.1. *Let A and B be non-empty, bounded, open sets in \mathbb{R}^d , with $r = \text{diam}(A \cup B)$. Let $F : \mathbb{R}^d \rightarrow \mathbb{R}$ be in L^p and supported on A , $\Phi : B \rightarrow A$ be an ϵ -deformation, and $F_\Phi = \Phi_\#^{-1} F$, i.e.*

$$F_\Phi(\mathbf{x}) = F(\Phi(\mathbf{x})) |\det(\nabla \Phi(\mathbf{x}))| \quad (97)$$

on B , and 0 elsewhere. Then for any $\mathbf{u} \in \mathbb{S}^{d-1}$,

$$\|\mathcal{P}_\mathbf{u} F - \mathcal{P}_\mathbf{u} F_\Phi\|_{V^p} \leq \epsilon^{1/p} \cdot \min \left\{ C(\Psi, r)^{1-1/p} \cdot \|F\|_{L^p}, \|F\|_{L^1} \right\}, \quad (98)$$

where

$$C(\Psi, r) = \max_{|t| \leq r/2} |\{\mathbf{x} \in A : \langle \mathbf{x}, \mathbf{u} \rangle \leq t \leq \langle \Psi(\mathbf{x}), \mathbf{u} \rangle \text{ or } \langle \Psi(\mathbf{x}), \mathbf{u} \rangle \leq t \leq \langle \mathbf{x}, \mathbf{u} \rangle\}|. \quad (99)$$

Theorem 7.1 is proved in Section 7.1. Sections 7.2, 7.3, and 7.4 then contain the proofs for Theorems 4.1, 4.2, and 4.3, respectively.

7.1 Proof of Theorem 7.1

Without loss of generality, suppose $\mathbf{u} = \mathbf{e}_1 = (1, 0, \dots, 0)$, and for brevity, if $G : \mathbb{R}^d \rightarrow \mathbb{R}$ is a function of d variables, let $\mathcal{P}G = \mathcal{P}_{\mathbf{e}_1} G$. That is,

$$(\mathcal{P}G)(t) = \int_{\mathbb{R}^{d-1}} G(t, x_1, \dots, x_{d-1}) dx_1 \cdots dx_{d-1} = \int_{\mathbb{R}^{d-1}} G(t, \mathbf{x}) d\mathbf{x}. \quad (100)$$

Let $R = r/2$. Also without loss of generality, suppose that A and B are contained in $\mathbb{B} = \mathbb{B}_R$, the closed ball of radius R centered at the origin.

First, suppose $p = 1$. We will show that

$$\|\mathcal{P}F - \mathcal{P}F_\Phi\|_{V^1} \leq \epsilon \cdot \|F\|_{L^1}. \quad (101)$$

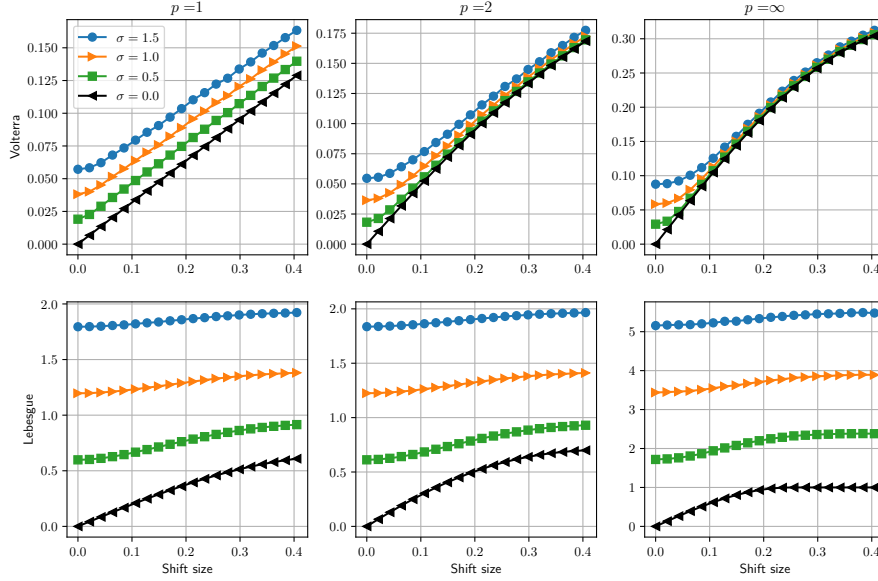


Figure 17: The first row shows the approximated Volterra distances between the function (96) and its noisy shifts, as a function of the shift size, for different noise levels. The second row shows the approximated Lebesgue distances. The values of p (from left to right) are $p = 1, 2, \infty$. See Section 6.9 for details.

By definition,

$$(\mathcal{P}F)(x) = \int_{\mathbb{R}^{d-1}} F(x, \mathbf{y}) d\mathbf{y}, \quad (102)$$

and

$$\begin{aligned} (\mathcal{P}F_{\Phi})(x) &= \int_{\mathbb{R}} F_{\Phi}(x, \mathbf{y}) d\mathbf{y} \\ &= \int_{\mathbf{y}: (x, \mathbf{y}) \in B} F(\Phi(x, \mathbf{y})) |\det(\nabla \Phi(x, \mathbf{y}))| d\mathbf{y}. \end{aligned} \quad (103)$$

Let G on $[-R, R]$ be absolutely continuous, with derivative $g = G'$ satisfying $\|g\|_{L^\infty} \leq 1$.

Performing the change of variables $\mathbf{u} = \Phi(x, \mathbf{y})$ gives

$$\begin{aligned} \int_{-R}^R G(x) (\mathcal{P}F_{\Phi})(x) dx &= \int_{-R}^R G(x) \int_{\mathbf{y}: (x, \mathbf{y}) \in B} F(\Phi(x, \mathbf{y})) |\det(\nabla \Phi(x, \mathbf{y}))| d\mathbf{y} dx \\ &= \int_B G(x) F(\Phi(x, \mathbf{y})) |\det(\nabla \Phi(x, \mathbf{y}))| d\mathbf{y} dx \\ &= \int_A G(\psi_1(\mathbf{u})) F(\mathbf{u}) d\mathbf{u}, \end{aligned} \quad (104)$$

and similarly,

$$\int_{-R}^R G(x) (\mathcal{P}F)(x) dx = \int_A G(x_1) F(\mathbf{x}) d\mathbf{x}. \quad (105)$$

We then have

$$\begin{aligned}
\int_{-R}^R G(x)((\mathcal{P}F)(x) - (\mathcal{P}F_{\Phi})(x)) dx &= \int_A G(x_1)F(\mathbf{x}) d\mathbf{x} - \int_A G(\psi_1(\mathbf{x}))F(\mathbf{x}) d\mathbf{x} \\
&= \int_A (G(x_1) - G(\psi_1(\mathbf{x})))F(\mathbf{x}) d\mathbf{x} \\
&\leq \|F\|_{L^1} \max_{\mathbf{x} \in A} |G(x_1) - G(\psi_1(\mathbf{x}))|.
\end{aligned} \tag{106}$$

Now, because $g = G'$ satisfies $\|g\|_{L^\infty} \leq 1$, we have

$$\begin{aligned}
|G(x_1) - G(\psi_1(\mathbf{x}))| &= \left| \int_{x_1}^{\psi_1(\mathbf{x})} g(t) dt \right| \\
&\leq \|g\|_{L^\infty} |x_1 - \psi_1(\mathbf{x})| \\
&\leq |\mathbf{x} - \Psi(\mathbf{x})| \\
&\leq \epsilon,
\end{aligned} \tag{107}$$

and therefore, taking the supremum over all such G and using Proposition 3.1 shows that

$$\|\mathcal{P}F - \mathcal{P}F_{\Phi}\|_{V^1} \leq \epsilon \cdot \|F\|_{L^1}. \tag{108}$$

This completes the proof when $p = 1$.

We will now prove that

$$\|\mathcal{P}F - \mathcal{P}F_{\Phi}\|_{V^\infty} \leq C(\Psi, r) \|F\|_{L^\infty}. \tag{109}$$

Let $I_{\mathbf{x}}$ be the interval $[x_1, \psi_1(\mathbf{x})]$ when $x_1 \leq \psi_1(\mathbf{x})$, and $[\psi_1(\mathbf{x}), x_1]$ when $x_1 > \psi_1(\mathbf{x})$; and let $\chi(\mathbf{x}, t)$ be 1 if $t \in I_{\mathbf{x}}$, and 0 otherwise.

Now, take an absolutely continuous G on $[-R, R]$ whose derivative $g = G'$ satisfies $\|g\|_{L^1} = 1$. Using (104) and (105) as before, we have

$$\begin{aligned}
\int_{-R}^R G(x)((\mathcal{P}F)(x) - (\mathcal{P}F_{\Phi})(x)) dx &= \int_A G(x_1)F(\mathbf{x}) d\mathbf{x} - \int_A G(\psi_1(\mathbf{x}))F(\mathbf{x}) d\mathbf{x} \\
&= \int_A (G(x_1) - G(\psi_1(\mathbf{x})))F(\mathbf{x}) d\mathbf{x} \\
&\leq \|F\|_{L^\infty} \int_A |G(x_1) - G(\psi_1(\mathbf{x}))| d\mathbf{x}.
\end{aligned} \tag{110}$$

Then, using $g = G'$, we have

$$\begin{aligned}
\int_A |G(x_1) - G(\psi_1(\mathbf{x}))| d\mathbf{x} &= \int_A \left| \int_{I_{\mathbf{x}}} g(t) dt \right| d\mathbf{x} \\
&= \int_A \left| \int_{-R}^{-R} g(t) \chi(\mathbf{x}, t) dt \right| d\mathbf{x} \\
&\leq \int_A \int_{-R}^R |g(t)| \chi(\mathbf{x}, t) dt d\mathbf{x} \\
&= \int_{-R}^R |g(t)| \int_A \chi(\mathbf{x}, t) d\mathbf{x} dt \\
&\leq \|g\|_{L^1} \sup_{|t| \leq R} \int_A \chi(\mathbf{x}, t) d\mathbf{x} \\
&= \sup_{|t| \leq R} \int_A \chi(\mathbf{x}, t) d\mathbf{x} \\
&= C(\Psi, r),
\end{aligned} \tag{111}$$

which completes the proof of (109).

Finally, since the mapping $F \mapsto \int_{-R}^x (\mathcal{P}F - \mathcal{P}F_\Phi)$ is linear, we may now combine the bounds (101) and (109) using the Riesz-Thorin Interpolation Theorem (see, e.g. Theorem 6.27 in [21]) to complete the proof that

$$\|\mathcal{P}_\mathbf{u}F - \mathcal{P}_\mathbf{u}F_\Phi\|_{V^p} \leq \epsilon^{1/p} \cdot C(\Psi, r)^{1-1/p} \cdot \|F\|_{L^p}. \quad (112)$$

Next, we will show that

$$\|\mathcal{P}_\mathbf{u}F - \mathcal{P}_\mathbf{u}F_\Phi\|_{V^p} \leq \epsilon^{1/p} \cdot \|F\|_{L^1}. \quad (113)$$

Let G be absolutely continuous on $[-R, R]$, with derivative $g = G'$ satisfying $\|g\|_{L^q} \leq 1$. From (106),

$$\int_{-R}^R G(x)((\mathcal{P}F)(x) - (\mathcal{P}F_\Phi)(x)) dx \leq \|F\|_{L^1} \max_{\mathbf{x} \in A} |G(x_1) - G(\psi_1(\mathbf{x}))|, \quad (114)$$

and so we must show that for all $\mathbf{x} \in A$,

$$|G(x_1) - G(\psi_1(\mathbf{x}))| \leq \epsilon^{1/p}. \quad (115)$$

As before, let $I_\mathbf{x}$ be the interval $[x_1, \psi_1(\mathbf{x})]$ if $x_1 \leq \psi_1(\mathbf{x})$, and $[\psi_1(\mathbf{x}), x_1]$ if $\psi_1(\mathbf{x}) \leq x_1$, and let $\chi(\mathbf{x}, t)$ be 1 if and only if $t \in I_\mathbf{x}$, and 0 otherwise. Note that

$$\int_{-R}^R \chi(\mathbf{x}, t) dt \leq |x_1 - \psi_1(\mathbf{x})| \leq |\mathbf{x} - \Psi(\mathbf{x})| \leq \epsilon. \quad (116)$$

Using that $G(y) = -\int_y^b g(t)dt$, we may write, for any $\mathbf{x} \in A$,

$$|G(x_1) - G(\psi_1(\mathbf{x}))| = \left| \int_{I_\mathbf{x}} g(t) dt \right|, \quad (117)$$

and Hölder's inequality yields

$$\begin{aligned} \left| \int_{I_\mathbf{x}} g(t) dt \right| &= \left| \int_{-R}^R g(t) \chi(\mathbf{x}, t) dt \right| \\ &\leq \|g\|_{L^q} \left(\int_{-R}^R \chi(\mathbf{x}, t)^p dt \right)^{1/p} \\ &\leq \left(\int_{-R}^R \chi(\mathbf{x}, t) dt \right)^{1/p} \\ &\leq \epsilon^{1/p}, \end{aligned} \quad (118)$$

where the last inequality follows from (116). This completes the proof.

7.2 Proof of Theorem 4.1

Without loss of generality, suppose that A and B are contained in $\mathbb{B} = \mathbb{B}_{r/2}$, and that $\mathbf{u} = \mathbf{e}_1 = (1, 0, \dots, 0)$; and let $\mathcal{P} = \mathcal{P}_{\mathbf{e}_1}$.

We will first show that for all $|t| \leq r/2$,

$$C(\Psi, r) \leq 2r^{d-1}\epsilon. \quad (119)$$

Let $\mathbf{S}_1 = \{\mathbf{x} \in A : x_1 \leq t \leq \psi_1(\mathbf{x})\}$, and let $\mathbf{S}_2 = \{\mathbf{x} \in A : \psi_1(\mathbf{x}) \leq t \leq x_1\}$. Then

$$C(\Psi, r) = |\mathbf{S}_1 \cup \mathbf{S}_2|. \quad (120)$$

To bound the area of \mathbf{S}_1 , observe first that any \mathbf{x} contained in \mathbf{S}_1 must satisfy $t - \epsilon \leq x_1 \leq t$. Indeed, since, by assumption, $\psi_1(\mathbf{x}) - x_1 \leq \epsilon$, we have

$$x_1 \geq \psi_1(\mathbf{x}) - \epsilon \geq t - \epsilon, \quad (121)$$

as claimed. Consequently, since $A \subset \mathbb{B}_{r/2}$,

$$|\mathbf{S}_1| \leq |\{\mathbf{x} \in \mathbb{B}_{r/2} : t - \epsilon \leq x_1 \leq t\}| \leq r^{d-1}\epsilon. \quad (122)$$

Similarly, $|\mathbf{S}_2| \leq r^{d-1}\epsilon$, and hence

$$C(\Psi, r) = |\mathbf{S}_1 \cup \mathbf{S}_2| \leq 2r^{d-1}\epsilon, \quad (123)$$

as claimed.

Consequently, Theorem 7.1 states that

$$\begin{aligned} \|\mathcal{P}F - \mathcal{P}F_\Phi\|_{V^p} &\leq \epsilon^{1/p} \cdot \min \left\{ (2r^{d-1}\epsilon)^{1-1/p} \cdot \|F\|_{L^p}, \|F\|_{L^1} \right\} \\ &= \min \left\{ \epsilon \cdot (2r^{d-1})^{1-1/p} \cdot \|F\|_{L^p}, \epsilon^{1/p} \cdot \|F\|_{L^1} \right\}. \end{aligned} \quad (124)$$

Switching the roles of Ψ and Φ , and using that $(F_\Phi)_\Psi = F$ and $\|F\|_{L^1} = \|F_\Phi\|_{L^1}$, shows the bound

$$\|\mathcal{P}F - \mathcal{P}F_\Phi\|_{V^p} = \|\mathcal{P}(F_\Phi)_\Psi - \mathcal{P}F_\Phi\|_{V^p} \leq \min \left\{ \epsilon \cdot (2r^{d-1})^{1-1/p} \cdot \|F_\Phi\|_{L^p}, \epsilon^{1/p} \cdot \|F\|_{L^1} \right\}. \quad (125)$$

Combining (124) and (125) completes the proof.

7.3 Proof of Theorem 4.2

Without loss of generality, we assume that $\varphi = 0$. Let $c = \cos(\theta)$ and $s = \sin(\theta)$, and let $f = f_\varphi = f_0$. We have $\Phi(x, y) = (cx + sy, cy - sx)$, and $\Psi(x, y) = (cx - sy, cy + sx)$. With $\epsilon = \max_{(x,y) \in \mathbb{B}_R} |\Phi(x, y) - (x, y)|$, we will prove

$$\epsilon \leq 2R \sin(\theta/2), \quad (126)$$

and

$$C(\Psi, 2R) \leq R^2\theta. \quad (127)$$

Assuming these bounds, Theorem 7.1 states that

$$\begin{aligned} \|f - f_\theta\|_{V^p} &\leq \epsilon^{1/p} \cdot \min \left\{ C(\Psi, r)^{1-1/p} \cdot \|F\|_{L^p}, \|F\|_{L^1} \right\} \\ &\leq (2R \sin(\theta/2))^{1/p} \cdot \min \left\{ (R^2\theta)^{1-1/p} \cdot \|F\|_{L^p}, \|F\|_{L^1} \right\} \\ &= (2 \sin(\theta/2))^{1/p} \cdot \min \left\{ \theta^{1-1/p} \cdot \|F\|_{L^p} \cdot R^{2-1/p}, \|F\|_{L^1} \cdot R^{1/p} \right\}, \end{aligned} \quad (128)$$

which is the desired result.

To prove (126), suppose $x^2 + y^2 \leq R^2$. Then from the double angle formula $c = \cos(\theta) = \cos^2(\theta/2) - \sin^2(\theta/2)$, or equivalently $1 - c = 2 \sin^2(\theta/2)$; therefore,

$$\begin{aligned} |\Phi(x, y) - (x, y)|^2 &= (x - cx - sy)^2 + (y - cy + sx)^2 \\ &= (1 - c)^2 x^2 + s^2 y^2 - 2(1 - c)sxy + (1 - c)^2 y^2 + s^2 x^2 + 2(1 - c)sxy \\ &\leq (1 - c)^2 R^2 + s^2 R^2 \\ &\leq (1 - c)^2 R^2 + (1 - c^2) R^2 \\ &= 2(1 - c) R^2 \\ &= 4R^2 \sin^2(\theta/2), \end{aligned} \quad (129)$$

and hence

$$\epsilon \equiv \max_{(x,y) \in \mathbb{B}_R} |\Phi(x,y) - (x,y)| \leq 2R \sin(\theta/2), \quad (130)$$

as desired.

Next, to prove (127), let $I_{x,y}$ be the interval $[x, cx+sy]$ when $x \leq cx+sy$, and the interval $[cx+sy, x]$ when $cx+sy \leq x$; and let $\chi(x,y,t)$ be 1 if $t \in I_{x,y}$ and 0 otherwise. Then $C(\Psi, 2R) = \max_{|t| \leq R} \int_{\mathbb{B}_R} \chi(x,y,t) dx dy$, and so we need to show that for all $t \in [-R, R]$,

$$\int_{\mathbb{B}_R} \chi(x,y,t) dx dy \leq R^2 \theta. \quad (131)$$

It is enough to show this for $R = 1$, since $C(\Psi, 2R) = R^2 C(\Psi, 2)$.

To that end, let $\mathbb{B} = \mathbb{B}_1$, and observe that for all $|t| \leq 1$,

$$\int_{\mathbb{B}} \chi(x,y,t) dx dy = 2|\mathbf{S}_{t,(1,0),(c,s)} \cup \mathbf{S}_{t,(c,s),(1,0)}|, \quad (132)$$

where, for unit vectors \mathbf{v} and \mathbf{w} , $\mathbf{S}_{t,\mathbf{v},\mathbf{w}}$ is the region defined by

$$\mathbf{S}_{t,\mathbf{v},\mathbf{w}} = \{\mathbf{u} \in \mathbb{B} : \langle \mathbf{u}, \mathbf{v} \rangle \leq t \leq \langle \mathbf{u}, \mathbf{w} \rangle\}. \quad (133)$$

By rotational symmetry, the following lemma is immediate:

Lemma 7.2. *If \mathbf{a} and \mathbf{b} are any unit vectors with angle θ , then $|\mathbf{S}_{t,(1,0),(c,y)}| = |\mathbf{S}_{t,\mathbf{a},\mathbf{b}}|$. Furthermore, $|\mathbf{S}_{t,\mathbf{a},\mathbf{b}}| = |\mathbf{S}_{-t,\mathbf{a},\mathbf{b}}|$, and $|\mathbf{S}_{t,\mathbf{a},\mathbf{b}} \cap \mathbf{S}_{-t,\mathbf{a},\mathbf{b}}| = 0$.*

By this lemma, it follows that

$$\int_{\mathbb{B}} \chi(x,y,t) dx dy = 2|\mathbf{S}_{t,\mathbf{v},\mathbf{w}}| = 2|\{\mathbf{u} \in \mathbb{B} : \langle \mathbf{u}, \mathbf{v} \rangle \leq t \leq \langle \mathbf{u}, \mathbf{w} \rangle\}|, \quad (134)$$

where $\mathbf{w} = (\cos(\theta/2), \sin(\theta/2))$ and $\mathbf{v} = (\cos(\theta/2), -\sin(\theta/2))$. Furthermore, we can restrict to $t \geq 0$.

It will be convenient to refer to Figure 18, where \mathbf{w} corresponds to the point labeled B , and \mathbf{v} corresponds to the point labeled E . In the figure, the line AD is perpendicular to OB , and intersects OB at distance t from the origin; consequently, the set of all vectors \mathbf{u} in \mathbb{B} with $\langle \mathbf{u}, \mathbf{w} \rangle \geq t$ is the circular segment through the points A , B and D . Similarly; the line CF is perpendicular to OE , and intersects OE at distance t from the origin; consequently, the set of all vectors \mathbf{u} in \mathbb{B} with $\langle \mathbf{u}, \mathbf{v} \rangle \leq t$ is the circular segment through the points C , A and F . Denote by G the point at the intersection of the lines AD and CF ; then when G lies within the circle, the intersection of these two circular segments is the region bounded by A , C and G . (When G is outside the circle, then the two circular segments are disjoint.)

To evaluate the area of this region, we will first find the area of the full circular segment through A , B and D , and then subtract off the area of the region bounded by C , G and D .

Lemma 7.3. *The area of the circular segment through A , B and D is*

$$\arccos(t) - t\sqrt{1-t^2}, \quad (135)$$

where \arccos takes values in $[0, \pi]$.

Proof. This is immediate from the well-known formula for the area of a circular segment, and the fact that the line segment from O to H has length t . \square

The next lemma is also elementary, and likely known already; however, since we could not find the exact identity in the literature, we provide a self-contained proof.

Lemma 7.4. *When $t \leq \cos(\theta/2)$, the intersection between the circular segment bounded by A , B and D and the circular segment bounded by C , E and F has area*

$$\arcsin\left(\sqrt{1-t^2}\cos(\theta/2) - t\sin(\theta/2)\right) - t\sqrt{1-t^2} + t^2\tan(\theta/2), \quad (136)$$

where \arcsin denotes the inverse of \sin on the interval $[0, \pi/2]$ (and hence takes values in this interval). When $t > \cos(\theta/2)$, the two circular segments are disjoint.

Proof. We begin by showing the second part, namely that when $t > \cos(\theta/2)$, the circular segments are disjoint, or equivalently that the point G lies outside of the circle. Indeed, it is straightforward to show that G is located at the point $(t/\cos(\theta/2), 0)$; hence, G is inside the circle so long as $t/\cos(\theta/2) \leq 1$, or equivalently, $t \leq \cos(\theta/2)$, as desired.

Let us now suppose that $t \leq \cos(\theta/2)$, and evaluate the area of the region bounded by C , G , and D . The line segment from G to D has arc-length parameterization

$$\alpha(s) = t(\cos(\theta/2), \sin(\theta/2)) + s(\sin(\theta/2), -\cos(\theta/2)), \quad (137)$$

and the line segment from C to G has arc-length parameterization

$$\beta(s) = t(\cos(\theta/2), -\sin(\theta/2)) + (\sqrt{1-t^2} + t \cdot \tan(\theta/2) - s)(\sin(\theta/2), \cos(\theta/2)), \quad (138)$$

where

$$t \cdot \tan(\theta/2) \leq s \leq \sqrt{1-t^2}. \quad (139)$$

The counterclockwise arc from D to C has arc-length parameterization

$$\gamma(\varphi) = (\cos(\varphi), \sin(\varphi)), \quad (140)$$

where

$$-\arcsin\left(\sqrt{1-t^2}\cos(\theta/2) - t \cdot \sin(\theta/2)\right) \leq \varphi \leq \arcsin\left(\sqrt{1-t^2}\cos(\theta/2) - t \cdot \sin(\theta/2)\right). \quad (141)$$

When $t \leq \cos(\theta/2)$, we will evaluate the area using Green's Theorem, by computing $\frac{1}{2} \oint (x dy - y dx)$ over each curve. For α , we have

$$\begin{aligned} \frac{1}{2} \oint_{\alpha} x dy &= \frac{1}{2} \int_{t \cdot \tan(\theta/2)}^{\sqrt{1-t^2}} [(t \cdot \cos(\theta/2) + s \cdot \sin(\theta/2))(-\cos(\theta/2))] ds \\ &= \frac{t \cdot \cos^2(\theta/2)}{2} (t \cdot \tan(\theta/2) - \sqrt{1-t^2}) + \frac{\sin(\theta/2)\cos(\theta/2)}{4} (t^2 \cdot \tan^2(\theta/2) - 1 + t^2), \end{aligned} \quad (142)$$

and

$$\begin{aligned} \frac{1}{2} \oint_{\alpha} y dx &= \frac{1}{2} \int_{t \cdot \tan(\theta/2)}^{\sqrt{1-t^2}} [(t \cdot \sin(\theta/2) - s \cdot \cos(\theta/2))(\sin(\theta/2))] ds \\ &= \frac{t \cdot \sin^2(\theta/2)}{2} (\sqrt{1-t^2} - t \cdot \tan(\theta/2)) + \frac{\sin(\theta/2)\cos(\theta/2)}{4} (t^2 \cdot \tan^2(\theta/2) - 1 + t^2), \end{aligned} \quad (143)$$

and hence

$$\frac{1}{2} \oint_{\alpha} (x dy - y dx) = \frac{t}{2} \cdot (t \cdot \tan(\theta/2) - \sqrt{1-t^2}). \quad (144)$$

Similarly,

$$\frac{1}{2} \oint_{\beta} (x dy - y dx) = \frac{t}{2} \cdot (t \cdot \tan(\theta/2) - \sqrt{1-t^2}). \quad (145)$$

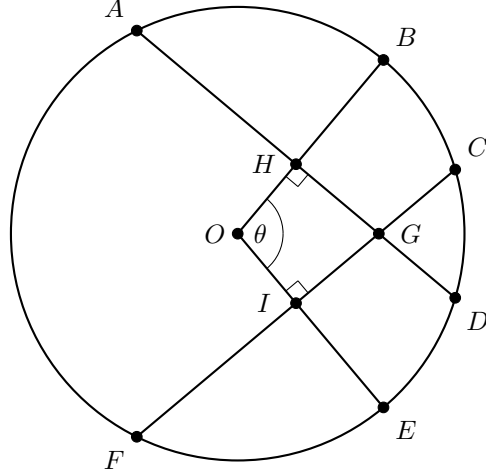


Figure 18: Diagram for the proof of (131). The points labeled B and C are located at $(\cos(\theta/2), \sin(\theta/2))$ and $(\cos(\theta/2), -\sin(\theta/2))$, respectively. The point labeled O is the origin, $(0, 0)$. The line segment OB is orthogonal to the line AD , and the line segment OE is orthogonal to the line FC . The line segments OH and OI each have length t .

Finally, it is straightforward to check that

$$\frac{1}{2} \oint_{\gamma} (x dy - y dx) = \arcsin \left(\sqrt{1-t^2} \cos(\theta/2) - t \cdot \sin(\theta/2) \right). \quad (146)$$

Adding all three integrals together, we find that the area of the region is

$$\frac{1}{2} \oint_{\gamma} (x dy - y dx) = \arcsin \left(\sqrt{1-t^2} \cos(\theta/2) - t \cdot \sin(\theta/2) \right) - t \sqrt{1-t^2} + t^2 \tan(\theta/2), \quad (147)$$

as claimed. \square

From Lemmas 7.3 and 7.4, we find

$$\begin{aligned} & \frac{1}{2} \int_{\mathbb{B}} \chi(x, y, t) dx dy \\ &= \begin{cases} \arccos(t) - t \sqrt{1-t^2} & \text{if } t > \cos(\theta/2); \\ \arccos(t) - \arcsin \left(\sqrt{1-t^2} \cos(\theta/2) - t \cdot \sin(\theta/2) \right) - t^2 \tan(\theta/2), & \text{if } t \leq \cos(\theta/2). \end{cases} \end{aligned} \quad (148)$$

To conclude the proof, we must show that this expression is bounded above by $\theta/2$ for all values of t between 0 and 1. In fact, we will show that (148) is a decreasing function of t , and hence is maximized at $t = 0$. It is immediately apparent that the expression is decreasing in t when $t > \cos(\theta/2)$, since this is the area of the circular segment with chord at distance t from the origin. When $t \leq \cos(\theta)$, we first observe that

$$\begin{aligned} & \frac{d}{dt} \arcsin \left(\sqrt{1-t^2} \cos(\theta/2) - t \cdot \sin(\theta/2) \right) \\ &= \frac{\frac{d}{dt} \left[\sqrt{1-t^2} \cos(\theta/2) - t \cdot \sin(\theta/2) \right]}{\sqrt{1 - \left(\sqrt{1-t^2} \cos(\theta/2) - t \cdot \sin(\theta/2) \right)^2}} \\ &= \frac{-t(1-t^2)^{-1/2} \cos(\theta/2) - \sin(\theta/2)}{\sqrt{1 - \left(\sqrt{1-t^2} \cos(\theta/2) - t \cdot \sin(\theta/2) \right)^2}}, \end{aligned} \quad (149)$$

and the square of the denominator may be written more simply as

$$\begin{aligned}
& 1 - (\sqrt{1-t^2} \cos(\theta/2) - t \cdot \sin(\theta/2))^2 \\
&= 1 - (1-t^2) \cos^2(\theta/2) - t^2 \sin^2(\theta/2) + 2t\sqrt{1-t^2} \cos(\theta/2) \sin(\theta/2) \\
&= 1 - \cos^2(\theta/2) + t^2 \cos^2(\theta/2) - t^2 \sin^2(\theta/2) + 2t\sqrt{1-t^2} \cos(\theta/2) \sin(\theta/2) \\
&= \sin^2(\theta/2) + t^2 \cos^2(\theta/2) - t^2 \sin^2(\theta/2) + 2t\sqrt{1-t^2} \cos(\theta/2) \sin(\theta/2) \\
&= t^2 \cos^2(\theta/2) + (1-t^2) \sin^2(\theta/2) + 2t\sqrt{1-t^2} \cos(\theta/2) \sin(\theta/2) \\
&= \left(t \cdot \cos(\theta/2) + \sqrt{1-t^2} \sin(\theta/2) \right)^2; \tag{150}
\end{aligned}$$

consequently,

$$\begin{aligned}
\frac{d}{dt} \arcsin \left(\sqrt{1-t^2} \cos(\theta/2) - t \cdot \sin(\theta/2) \right) &= \frac{-t(1-t^2)^{-1/2} \cos(\theta/2) - \sin(\theta/2)}{t \cdot \cos(\theta/2) + \sqrt{1-t^2} \sin(\theta/2)} \\
&= \frac{-1}{\sqrt{1-t^2}}. \tag{151}
\end{aligned}$$

Therefore,

$$\begin{aligned}
& \frac{d}{dt} \left[\arccos(t) - \arcsin \left(\sqrt{1-t^2} \cos(\theta/2) - t \cdot \sin(\theta/2) \right) - t^2 \tan(\theta/2) \right] \\
&= \frac{-1}{\sqrt{1-t^2}} + \frac{1}{\sqrt{1-t^2}} - 2t \tan(\theta/2) \\
&= -2t \tan(\theta/2), \tag{152}
\end{aligned}$$

which is negative. Therefore, the maximum value of $\int_{\mathbb{B}} \chi(x, y, t) dx dy$ occurs when $t = 0$, where the value is

$$\begin{aligned}
2 \arccos(0) - 2 \arcsin(\cos(\theta/2)) &= \pi - 2 \arcsin(\sin(\pi/2 + \theta/2)) \\
&= \pi - 2 \arcsin(\sin(\pi/2 - \theta/2)) \\
&= \pi - 2 \left(\frac{\pi}{2} - \frac{\theta}{2} \right) \\
&= \theta; \tag{153}
\end{aligned}$$

where we have used that \sin is even around $\pi/2$. Note that $\pi/2 - \theta/2$ lies between 0 and $\pi/2$ (since θ is between 0 and π), and hence $\arcsin(\sin(\pi/2 - \theta/2)) = \pi/2 - \theta/2$, as claimed. This completes the proof.

7.4 Proof of Theorem 4.3

We will first prove that

$$C(\Psi, r) \leq \epsilon. \tag{154}$$

Let I_x be the interval $[x, \Psi(x)]$ if $x \leq \Psi(x)$, and $[\Psi(x), x]$ if $\Psi(x) \leq x$. Let $\chi(x, t)$ be 1 if $t \in I_x$, and 0 otherwise; then

$$C(\Psi, r) = \max_{|t| \leq r/2} \int_{\mathbb{R}} \chi(x, t) dt. \tag{155}$$

Take $|t| \leq r/2$. Suppose that there is some $x \leq t$ with $t \in I_x$; note that for such x , $I_x = [x, \Psi(x)]$, and so $x \leq \Psi(x)$. Let x^* be the smallest such x . Then $x^* \leq t \leq \Psi(x^*)$. We claim that for all $x > t$, $t \notin I_x$. Indeed, since Ψ is increasing and $x > t \geq x^*$, we have $\Psi(x) > \Psi(x^*) \geq t$. Since both $x > t$ and $\Psi(x) > t$, t does not lie in I_x , as claimed.

Consequently, all x for which t lies in I_x are contained inside the interval $[x^*, t]$. Since $x^* \leq t \leq \Psi(x^*)$ and $|x^* - \Psi(x^*)| \leq \epsilon$, it follows that $|t - x^*| \leq \epsilon$ too. Furthermore, if $x > t$, then $\chi(x, t) = 0$ since $t \notin I_x$; and since x^* is the smallest x for which $t \in I_x$, if $x < x^*$ then $t \notin I_x$, hence $\chi(x, t) = 0$. Therefore,

$$\int_I \chi(x, t) dx \leq \int_{x^*}^t 1 dx = |t - x^*| \leq \epsilon, \quad (156)$$

and so $C(\Psi, r) \leq \epsilon$. Analogous reasoning yields the same bound in the case that there exists $x \geq t$ with $t \in I_x$.

Therefore, Theorem 7.1 states that

$$\begin{aligned} \|F - F_\Phi\|_{V^p} &\leq \epsilon^{1/p} \cdot \min \left\{ \epsilon^{1-1/p} \cdot \|F\|_{L^p}, \|F\|_{L^1} \right\} \\ &= \min \left\{ \epsilon \cdot \|F\|_{L^p}, \epsilon^{1/p} \cdot \|F\|_{L^1} \right\}. \end{aligned} \quad (157)$$

Switching the roles of Ψ and Φ , and using that $(F_\Phi)_\Psi = F$ and $\|F\|_{L^1} = \|F_\Phi\|_{L^1}$, shows the bound

$$\|F - F_\Phi\|_{V^p} = \|(F_\Phi)_\Psi - F_\Phi\|_{V^p} \leq \min \left\{ \epsilon \cdot \|F_\Phi\|_{L^p}, \epsilon^{1/p} \cdot \|F_\Phi\|_{L^1} \right\}. \quad (158)$$

Combining (157) and (158) completes the proof.

8 Proofs for Section 5

First, we introduce some notation that will be useful for the proofs in this section. For a function f on $[a, b]$, define the vector Vf in \mathbb{R}^{n+1} with entries $(Vf)[k] = (\mathcal{V}f)(a_k)$, $0 \leq k \leq n$.

We also state the following simple lemma:

Lemma 8.1. *Suppose $1 \leq p \leq \infty$. Let $\mathbf{1}$ be the all 1's vector in \mathbb{R}^{n+1} . Then $\|\mathbf{1}\|_{\nu_p} \leq (b-a)^{1+1/p}$.*

Proof. For $1 \leq k \leq n$,

$$(\mathbf{V}\mathbf{1})[k] = \frac{b-a}{2n} \sum_{j=0}^{k-1} (1+1) = \frac{k}{n}(b-a), \quad (159)$$

and therefore

$$\|\mathbf{1}\|_{\nu_p} = \left(\frac{b-a}{2n} \sum_{k=0}^{n-1} \left[\left(\frac{k}{n}(b-a) \right)^p + \left(\frac{k+1}{n}(b-a) \right)^p \right] \right)^{1/p} \leq (b-a)^{1+1/p}, \quad (160)$$

as claimed. \square

8.1 Proof of Theorem 5.1

To begin the proof of Theorem 5.1, suppose $1 \leq p < \infty$.

Lemma 8.2. *Suppose f satisfies the assumptions of Theorem 5.1. Let $1 \leq m \leq n$. Then*

$$|T_m(f, a, a_m) - (\mathcal{V}f)(a_m)| \leq \frac{L(b-a)^2}{2n} + \frac{4r\|f\|_{L^\infty}(b-a)}{n}. \quad (161)$$

Proof. Because f is bounded by $\|f\|_{L^\infty}$, for any $0 \leq j \leq m$ we have the bound

$$\left| \frac{b-a}{2n} (f(a_j) + f(a_{j+1})) - \int_{a_j}^{a_{j+1}} f(t) dt \right| \leq \frac{2\|f\|_{L^\infty}(b-a)}{n}. \quad (162)$$

On the other hand, if j is such that there are no c_ℓ in $[a_j, a_{j+1}]$, then f is L -Lipschitz on $[a_j, a_{j+1}]$, and so

$$\begin{aligned}
\left| \frac{b-a}{2n} (f(a_j) + f(a_{j+1})) - \int_{a_j}^{a_{j+1}} f(t) dt \right| &= \frac{1}{2} \left| \int_{a_j}^{a_{j+1}} (f(a_j) - f(t) + f(a_{j+1}) - f(t)) dt \right| \\
&\leq \frac{L}{2} \int_{a_j}^{a_{j+1}} (t - a_j + a_{j+1} - t) dt \\
&= \frac{L}{2} \int_{a_j}^{a_{j+1}} \frac{(b-a)}{n} dt \\
&= \frac{L(b-a)^2}{2n^2}.
\end{aligned} \tag{163}$$

Since there are at most $2r$ intervals $[a_j, a_{j+1}]$ containing a value from among c_0, \dots, c_r (since each c_1, \dots, c_r can be contained in at most 2 such intervals), and there are $m \leq n$ subintervals in total, we have

$$\begin{aligned}
|T_m(f, a, a_m) - (\mathcal{V}f)(a_m)| &= \left| \frac{b-a}{2n} \sum_{j=0}^{m-1} (f(a_j) + f(a_{j+1})) - \int_a^{a_m} f(t) dt \right| \\
&\leq \sum_{j=0}^{m-1} \left| \frac{b-a}{2n} (f(a_j) + f(a_{j+1})) - \int_{a_j}^{a_{j+1}} f(t) dt \right| \\
&\leq \frac{mL(b-a)^2}{2n^2} + \frac{4r\|f\|_{L^\infty}(b-a)}{n}, \\
&\leq \frac{L(b-a)^2}{2n} + \frac{4r\|f\|_{L^\infty}(b-a)}{n},
\end{aligned} \tag{164}$$

as claimed. \square

By Lemma 8.2,

$$|(\mathbf{V}\mathbf{f})[m] - (\mathcal{V}f)[m]| \leq \frac{L(b-a)^2}{2n} + \frac{4r\|f\|_{L^\infty}(b-a)}{n} \tag{165}$$

Consequently,

$$\begin{aligned}
|\|\mathbf{V}\mathbf{f}\|_{\tau_p} - \|\mathcal{V}f\|_{\tau_p}| &\leq \|\mathbf{V}\mathbf{f} - \mathcal{V}f\|_{\tau_p} \\
&\leq (b-a)^{1/p} \|\mathbf{V}\mathbf{f} - \mathcal{V}f\|_{\ell_\infty} \\
&\leq \frac{L(b-a)^{2+1/p}}{2n} + \frac{4r\|f\|_{L^\infty}(b-a)^{1+1/p}}{n}.
\end{aligned} \tag{166}$$

Next, we will show that

$$|\|f\|_{V^p} - \|\mathcal{V}f\|_{\tau_p}| \leq \frac{(b-a)^{1+1/p}}{n} \|f\|_{L^\infty}, \tag{167}$$

Combined with (166), this will conclude the proof. To that end, we have the following lemma:

Lemma 8.3. *Suppose g has Lipschitz constant bounded by A on $[a, b]$, and let $\mathbf{g}[k] = g(a_k)$, where*

$$a_k = a + \frac{k}{n}(b-a), \quad 0 \leq k \leq n. \tag{168}$$

Then for any $1 \leq p \leq \infty$,

$$|\|\mathbf{g}\|_{\tau_p} - \|g\|_{L^p}| \leq \frac{(b-a)^{1+1/p}}{n} A. \tag{169}$$

Proof. First, suppose $1 \leq p < \infty$. For each $0 \leq m \leq n$, let

$$S_m = \left[\left(\frac{b-a}{n} \right) \left(\frac{|g(a_m)|^p + |g(a_{m+1})|^p}{2} \right) \right]^{1/p}, \quad (170)$$

and let

$$R_m = \left(\int_{a_m}^{a_{m+1}} |g(x)|^p dx \right)^{1/p}. \quad (171)$$

Then

$$\|\mathbf{g}\|_{\tau_p} = \left(\sum_{m=0}^{n-1} |S_m|^p \right)^{1/p} \quad (172)$$

and

$$\|g\|_{L^p} = \left(\sum_{m=0}^{n-1} |R_m|^p \right)^{1/p}. \quad (173)$$

The Mean Value Theorem ensures that there is some t_m in the interval $[a_m, a_{m+1}]$ satisfying

$$R_m = \left(\frac{b-a}{n} \right)^{1/p} |g(t_m)|. \quad (174)$$

Then

$$\begin{aligned} |S_m - R_m| &= \left| S_m - \left(\frac{b-a}{n} \right)^{1/p} |g(t_m)| \right| \\ &= \left(\frac{b-a}{n} \right)^{1/p} \left| \left[\left(\frac{|g(a_m)|^p + |g(a_{m+1})|^p}{2} \right) \right]^{1/p} - \left[\left(\frac{|g(t_m)|^p + |g(t_m)|^p}{2} \right) \right]^{1/p} \right| \\ &\leq \left(\frac{b-a}{n} \right)^{1/p} \left(\frac{|g(a_m) - g(t_m)|^p + |g(a_{m+1}) - g(t_m)|^p}{2} \right)^{1/p} \\ &\leq \left(\frac{b-a}{n} \right)^{1/p} \left(\frac{A^p |a_m - t_m|^p + A^p |a_{m+1} - t_m|^p}{2} \right)^{1/p} \\ &\leq \left(\frac{b-a}{n} \right)^{1/p} A |a_{m+1} - a_m| \\ &= A \left(\frac{b-a}{n} \right)^{1+1/p}. \end{aligned} \quad (175)$$

Consequently,

$$\begin{aligned} \left| \|\mathbf{g}\|_{\tau_p} - \|g\|_{L^p} \right| &= \left| \left(\sum_{m=0}^{n-1} |S_m|^p \right)^{1/p} - \left(\sum_{m=0}^{n-1} |R_m|^p \right)^{1/p} \right| \\ &\leq \left(\sum_{m=0}^{n-1} |S_m - R_m|^p \right)^{1/p} \\ &\leq \left(\sum_{m=0}^{n-1} A^p \left(\frac{b-a}{n} \right)^{p+1} \right)^{1/p} \\ &= \frac{(b-a)^{1+1/p}}{n} A. \end{aligned} \quad (176)$$

This completes the proof when p is finite. The proof for $p = \infty$ follows by taking the limit $p \rightarrow \infty$ and using the convergence of the p -norm to the ∞ -norm. \square

Now, $\mathcal{V}f$ has Lipschitz constant bounded by $\|f\|_{L^\infty}$:

$$|(\mathcal{V}f)(x) - (\mathcal{V}f)(y)| = \left| \int_a^x f(t)dt - \int_a^y f(t)dt \right| = \left| \int_x^y f(t)dt \right| \leq L|x - y|. \quad (177)$$

We may therefore apply Lemma 8.3 with $g = \mathcal{V}f$ and $A = \|f\|_{L^\infty}$ to show (167), thereby completing the proof of Theorem 5.1 for f and \mathbf{f} .

To prove the result for f_{cen} and \mathbf{f}_{cen} , from Lemma 8.2 we have

$$|\mu(f) - \mathbf{m}(\mathbf{f})| = \left| \frac{1}{b-a}(\mathcal{V}f)(b) - \frac{1}{b-a}T_n(f, a, b) \right| \leq \frac{L(b-a)}{2n} + \frac{4r\|f\|_{L^\infty}}{n}. \quad (178)$$

Letting $\tilde{\mathbf{f}}$ be the vector in \mathbb{R}^{n+1} with entries $\tilde{\mathbf{f}}[k] = f_{\text{cen}}(a_k) = f(a_k) - \mu(f)$, $0 \leq k \leq n$, applying the result we have already shown to f_{cen} in place of f gives

$$\left| \|\tilde{\mathbf{f}}\|_{\nu_p} - \|f_{\text{cen}}\|_{V^p} \right| \leq C \frac{(b-a)^{1+1/p}}{n} (L(b-a) + r\|f_{\text{cen}}\|_{L^\infty}). \quad (179)$$

Furthermore, for all $0 \leq k \leq n$,

$$\mathbf{f}_{\text{cen}}[k] - \tilde{\mathbf{f}}[k] = \mu(f) - \mathbf{m}(\mathbf{f}), \quad (180)$$

and so, by Lemma 8.1,

$$\begin{aligned} \left| \|\mathbf{f}_{\text{cen}}\|_{\nu_p} - \|\tilde{\mathbf{f}}\|_{\nu_p} \right| &\leq \|\mathbf{f}_{\text{cen}} - \tilde{\mathbf{f}}\|_{\nu_p} \\ &= \|(\mu(f) - \mathbf{m}(\mathbf{f}))\mathbf{1}\|_{\nu_p} \\ &= |\mu(f) - \mathbf{m}(\mathbf{f})| \|\mathbf{1}\|_{\nu_p} \\ &\leq |\mu(f) - \mathbf{m}(\mathbf{f})| (b-a)^{1+1/p} \\ &\leq C \frac{(b-a)^{1+1/p}}{n} (L(b-a) + r\|f_{\text{cen}}\|_{L^\infty}). \end{aligned} \quad (181)$$

The result now follows by combining (179) and (181).

8.2 Proof of Theorem 5.2

We will first prove a general result on approximating the L^p norm:

Proposition 8.4. *Suppose G is a C^3 function on $[a, b]$ with $G(a) = 0$. Suppose $a = c_0 < c_1 < \dots < c_r = b$ are points such that, for $0 \leq j \leq r-1$, $G(c_j) = 0$ and $\text{sign}(G)$ is constant on (c_j, c_{j+1}) . Let $\mathbf{g}[k] = G(a_k)$, $0 \leq k \leq n+1$.*

Let $1 \leq p < \infty$. Then for all n sufficiently large,

$$\left| \|\mathbf{g}\|_{\tau_p} - \|G\|_{L^p} \right| \leq C \frac{(b-a)^2}{n^2} \left((b-a)^{1/p} \|G''\|_{L^\infty} + \frac{|G(b)|^{p-1}}{\|G\|_{L^p}^{p-1}} |G'(b)| \right), \quad (182)$$

where $C > 0$ is a universal constant; and

$$\left| \|\mathbf{g}\|_{\ell_\infty} - \|G\|_{L^\infty} \right| \leq C \frac{(b-a)^2}{n^2} \|G''\|_{L^\infty}, \quad (183)$$

for all n sufficiently large.

Proof. Since the result is trivial if $G \equiv 0$, suppose that G is not constantly zero. First suppose $1 < p < \infty$, and let $H(x) = |G(x)|^p$ for $a \leq x \leq b$. Then H is continuous on $[a, b]$, and, since the sign of G is constant on each interval (c_j, c_{j+1}) , $0 \leq j \leq r-1$, H is C^3 on (c_j, c_{j+1}) .

From (a slightly different form of) the Euler-Maclaurin formula found in [7] and [8], we may write the error of the trapezoidal rule approximation $T_n(H, a, b)$ to $\int_a^b H(x)dx$ as

$$T_n(H, a, b) - \int_a^b H(x)dx = \frac{\delta^2}{2} \sum_{k=0}^{r-1} P_2(t_k) [H'(c_k^-) - H'(c_k^+)] - \frac{\delta^2}{2} \int_a^b H''(x)P_2\left(\frac{x-a}{\delta}\right) dx, \quad (184)$$

where $\delta = (b-a)/n$; P_2 is the second Bernoulli polynomial defined on $[0, 1]$ and extended 1-periodically (that is, $P_2(x) = (x-1/2)^2 - 1/12$ on $[0, 1]$, and $P_2(x+k) = P_2(x)$ for all integers k); $t_k = (c_k - a)/\delta$; and where $H'(c_0^-)$ is understood to denote $H'(b^-)$.

Because $p > 1$, H is differentiable and

$$H'(x) = p|G(x)|^{p-1}\text{sign}(G(x))G'(x). \quad (185)$$

Since $G(c_j) = 0$ when $0 \leq j \leq r-1$, $H'(c_j^+) = \pm p|G(c_j)|^{p-1}G'(c_j^+) = 0$, and for $1 \leq j \leq r-1$, $H'(c_j^-) = 0$ as well. When $j = 0$, $H'(c_0^-) = H'(b^-) = \pm p|G(b)|^{p-1}G'(b)$. Therefore,

$$\left| \frac{\delta^2}{2} \sum_{k=0}^{r-1} P_2(t_k) [H'(c_k^-) - H'(c_k^+)] \right| = \frac{(b-a)^2}{2n^2} p|P_2(1)||G(b)|^{p-1}|G'(b)| = \frac{(b-a)^2 p|G(b)|^{p-1}|G'(b)|}{12n^2}. \quad (186)$$

Next, we will bound the term

$$\frac{\delta^2}{2} \int_a^b H''(x)P_2\left(\frac{x-a}{\delta}\right) dx. \quad (187)$$

For x in each open interval (c_k, c_{k+1}) ,

$$H''(x) = p(p-1)|G(x)|^{p-2}G'(x)^2 + p|G(x)|^{p-1}\text{sign}(G(x))G''(x). \quad (188)$$

Suppose first that $G > 0$ on (c_j, c_{j+1}) , so that $H(x) = G(x)^p$. Let $D(x) = G(x)^{p-1}$; then $D'(x) = (p-1)G(x)^{p-2}G'(x)$, and $D(c_j) = D(c_{j+1}) = 0$. Then, using integration by parts, and $\|P_2\|_{L^\infty} = 1/6$,

$$\begin{aligned} \int_{c_j}^{c_{j+1}} p(p-1)G(x)^{p-2}G'(x)^2 P_2\left(\frac{x-a}{\delta}\right) dx &\leq \frac{p}{6} \int_{c_j}^{c_{j+1}} (p-1)G(x)^{p-2}G'(x)^2 dx \\ &= \frac{p}{6} \int_{c_j}^{c_{j+1}} D'(x)G'(x) dx \\ &= \frac{p}{6} (D(c_{j+1})G'(c_{j+1}) - D(c_j)G'(c_j)) - \frac{p}{6} \int_{c_j}^{c_{j+1}} D(x)G''(x) dx \\ &= -\frac{p}{6} \int_{c_j}^{c_{j+1}} G(x)^{p-1}G''(x) dx. \end{aligned} \quad (189)$$

Consequently, we have the upper bound

$$\left| \int_{c_j}^{c_{j+1}} p(p-1)G(x)^{p-2}G'(x)^2 P_2\left(\frac{x-a}{\delta}\right) dx \right| \leq \frac{p\|G''\|_{L^\infty}}{6} \int_{c_j}^{c_{j+1}} |G(x)|^{p-1} dx; \quad (190)$$

Furthermore, we also have the bound

$$\left| \int_{c_j}^{c_{j+1}} pG(x)^{p-1}G''(x)P_2\left(\frac{x-a}{\delta}\right) dx \right| \leq \frac{p\|G''\|_{L^\infty}}{6} \int_{c_j}^{c_{j+1}} |G(x)|^{p-1} dx. \quad (191)$$

Putting these together shows

$$\left| \int_{c_j}^{c_{j+1}} H''(x)P_2\left(\frac{x-a}{\delta}\right) dx \right| \leq \frac{p\|G''\|_{L^\infty}}{3} \int_{c_j}^{c_{j+1}} |G(x)|^{p-1} dx. \quad (192)$$

The same bound may also be shown if $G < 0$ on (c_j, c_{j+1}) (and it is obvious if $G = 0$ on all of (c_j, c_{j+1})).
Consequently,

$$\begin{aligned}
\left| \frac{\delta^2}{2} \int_a^b H''(x) P_2 \left(\frac{x-a}{\delta} \right) dx \right| &= \left| \frac{\delta^2}{2} \sum_{j=0}^{r-1} \int_{c_j}^{c_{j+1}} H''(x) P_2 \left(\frac{x-a}{\delta} \right) dx \right| \\
&\leq \frac{\delta^2 p \|G''\|_{L^\infty}}{6} \int_a^b |G(x)|^{p-1} dx \\
&\leq \frac{\delta^2 p \|G''\|_{L^\infty}}{6} (b-a)^{1/p} \left(\int_a^b |G(x)|^p dx \right)^{1-1/p} \\
&= \frac{(b-a)^{2+1/p}}{6n^2} p \|G''\|_{L^\infty} \|G\|_{L^p}^{p-1}. \tag{193}
\end{aligned}$$

Applying the Euler-Maclaurin formula then gives

$$\left| \|\mathbf{g}\|_{\tau_p}^p - \int_a^b |G(x)|^p dx \right| \leq \frac{(b-a)^2}{12n^2} p |G(b)|^{p-1} |G'(b)| + \frac{(b-a)^{2+1/p}}{6n^2} p \|G''\|_{L^\infty} \|G\|_{L^p}^{p-1}. \tag{194}$$

It follows that for all n sufficiently large,

$$\|\mathbf{g}\|_{\tau_p}^p \geq \frac{1}{2} \|G\|_{L^p}^p. \tag{195}$$

Now, the function $y \mapsto y^{1/p}$ has derivative $(1/p)y^{1/p-1}$; this is a decreasing function, and hence its maximum value between $\|\mathbf{g}\|_{\tau_p}^p$ and $\|G\|_{L^p}^p$ is bounded above by

$$\frac{1}{p} \left(\frac{1}{2} \|G\|_{L^p}^p \right)^{1/p-1} = \frac{2^{1-1/p}}{p} \|G\|_{L^p}^{1-p}. \tag{196}$$

By the mean value theorem, therefore,

$$\begin{aligned}
\left| \|\mathbf{g}\|_{\tau_p} - \|G\|_{L^p} \right| &\leq \left(\frac{(b-a)^2}{12n^2} p |G(b)|^{p-1} |G'(b)| + \frac{(b-a)^{2+1/p}}{6n^2} p \|G''\|_{L^\infty} \|G\|_{L^p}^{p-1} \right) \cdot \frac{2^{1-1/p}}{p} \|G\|_{L^p}^{1-p} \\
&\leq C \frac{(b-a)^2}{n^2} \left((b-a)^{1/p} \|G''\|_{L^\infty} + |G'(b)| \frac{|G(b)|^{p-1}}{\|G\|_{L^p}^{p-1}} \right), \tag{197}
\end{aligned}$$

where $C > 0$ is universal. This completes the proof when $1 < p < \infty$. The result for $p = 1$ follows by taking the limit $p \rightarrow 1^+$.

When $p = \infty$, let x^* satisfy $\|G\|_{L^\infty} = |G(x^*)|$. If $x^* = b$, then, since $\mathbf{g}[n] = G(b)$, $\|\mathbf{g}\|_{\ell_\infty} = |G(b)| = \|G\|_{L^\infty}$; and since $G(a) = 0$, x^* cannot equal a unless $G \equiv 0$, in which case the result is trivial.

Now suppose $a < x^* < b$. Then $G'(x^*) = 0$, and so a second-order Taylor expansion gives

$$|G(x) - G(x^*)| \leq C \|G''\|_{L^\infty} |x - x^*|^2, \tag{198}$$

where $C > 0$ is universal. Consequently, since there is a grid point a_{k^*} within $(b-a)/n$ of x^* , we have

$$|\mathbf{g}[k^*] - G(x^*)| = |G(a_{k^*}) - G(x^*)| \leq C \frac{(b-a)^2}{n^2} \|G''\|_{L^\infty}, \tag{199}$$

and therefore,

$$\begin{aligned}
\left| \|\mathbf{g}\|_{\ell_\infty} - \|G\|_{L^\infty} \right| &= |G(x^*)| - \|\mathbf{g}\|_{\ell_\infty} \\
&\leq |G(x^*)| - |\mathbf{g}[k^*]| \\
&\leq |G(x^*) - \mathbf{g}[k^*]| \\
&\leq C \frac{(b-a)^2}{n^2} \|G''\|_{L^\infty}, \tag{200}
\end{aligned}$$

which is the desired result. \square

Applying Proposition 8.4 to $G(x) = (\mathcal{V}f)(x)$, and using that $G'(x) = f(x)$ and $G(b) = \mu(f)$, gives the bound

$$\left| \|\mathbf{V}f\|_{\tau_p} - \|f\|_{V^p} \right| \leq C \frac{(b-a)^2}{n^2} \left((b-a)^{1/p} \|f'\|_{L^\infty} + |f(b)| \frac{|\mu(f)|^{p-1}}{\|f\|_{V^p}^{p-1}} \right) \quad (201)$$

when $1 \leq p < \infty$ and for all n sufficiently large,

$$\left| \|\mathbf{V}f\|_{\tau_p} - \|f\|_{V^\infty} \right| \leq C \frac{(b-a)^2}{n^2} \|f'\|_{L^\infty} \quad (202)$$

for all n .

To finish the proof, we will show that

$$\left| \|\mathbf{V}f\|_{\tau_p} - \|\mathbf{f}\|_{\nu_p} \right| \leq C \|f''\|_{L^\infty} \frac{(b-a)^{3+1/p}}{n^2}. \quad (203)$$

Lemma 8.5. *Let $0 \leq m \leq n$. Then*

$$|(\mathcal{V}f)(a_m) - (\mathbf{V}\mathbf{f})[m]| \leq C \|f''\|_{L^\infty} \frac{(b-a)^3}{n^2}, \quad (204)$$

where the constant $C > 0$ is universal.

Proof. When $m = 0$, the left side is 0. When $m \geq 1$, using a standard error estimate for the trapezoidal rule (see Section 2.2) and the fact that $a_m - a = m(b-a)/n$, we get

$$\begin{aligned} |(\mathcal{V}f)(a_m) - (\mathbf{V}\mathbf{f})[m]| &= \left| \int_a^{a_m} f(t) dt - T_m(f, a, a_m) \right| \\ &\leq C \|f''\|_{L^\infty} \frac{(a_m - a)^3}{m^2}, \\ &= C \|f''\|_{L^\infty} \frac{m(b-a)^3}{n^3} \\ &\leq C \|f''\|_{L^\infty} \frac{(b-a)^3}{n^2}, \end{aligned} \quad (205)$$

as claimed. \square

Using Lemma 8.5, we have

$$\begin{aligned} \left| \|\mathbf{V}f\|_{\tau_p} - \|\mathbf{f}\|_{\nu_p} \right| &= \left| \|\mathbf{V}f\|_{\tau_p} - \|\mathbf{V}\mathbf{f}\|_{\tau_p} \right| \\ &\leq \|\mathbf{V}f - \mathbf{V}\mathbf{f}\|_{\tau_p} \\ &\leq (b-a)^{1/p} \|\mathbf{V}f - \mathbf{V}\mathbf{f}\|_{\ell_\infty} \\ &\leq C \|f''\|_{L^\infty} \frac{(b-a)^{3+1/p}}{n^2}. \end{aligned} \quad (206)$$

This completes the proof of Theorem 5.2 for the uncentered function f .

To prove the result for f_{cen} and \mathbf{f}_{cen} , let $\tilde{\mathbf{f}}$ have entries $\tilde{\mathbf{f}}[k] = f_{\text{cen}}(a_k) = f(a_k) - \mu(f)$. Applying the result already shown to f_{cen} in place of f , and noting that $\mu(f_{\text{cen}}) = 0$, gives

$$\left| \|\tilde{\mathbf{f}}\|_{\nu_p} - \|f_{\text{cen}}\|_{V^p} \right| \leq C \frac{(b-a)^{2+1/p}}{n^2} \left((b-a) \|f''_{\text{cen}}\|_{L^\infty} + \|f'_{\text{cen}}\|_{L^\infty} \right). \quad (207)$$

Since f is C^2 ,

$$|\mu(f) - \mathbf{m}(\mathbf{f})| = \left| \frac{1}{b-a} \int_a^b f(x) dx - T_n(f, a, b) \right| \leq C \|f''\|_{L^\infty} \frac{(b-a)^2}{n^2} = C \|f''_{\text{cen}}\|_{L^\infty} \frac{(b-a)^2}{n^2}, \quad (208)$$

where $C > 0$ is universal. Furthermore, for all $0 \leq k \leq n$,

$$\mathbf{f}_{\text{cen}}[k] - \tilde{\mathbf{f}}[k] = \mathbf{m}(\mathbf{f}) - \mu(f), \quad (209)$$

and so, by Lemma 8.1,

$$\begin{aligned} \left| \|\tilde{\mathbf{f}}\|_{\nu_p} - \|\mathbf{f}_{\text{cen}}\|_{\nu_p} \right| &\leq \|\tilde{\mathbf{f}} - \mathbf{f}_{\text{cen}}\|_{\nu_p} \\ &= \|(\mathbf{m}(\mathbf{f}) - \mu(f))\mathbf{1}\|_{\nu_p} \\ &= |\mathbf{m}(\mathbf{f}) - \mu(f)| \|\mathbf{1}\|_{\nu_p} \\ &\leq |\mathbf{m}(\mathbf{f}) - \mu(f)| (b-a)^{1+1/p} \\ &\leq C \|f''\|_{L^\infty} \frac{(b-a)^{3+1/p}}{n^2}. \end{aligned} \quad (210)$$

The result now follows by combining (207) and (210).

8.3 Proofs of Theorem 5.3 and Corollary 5.4

Let $T[0] = 0$, and for $1 \leq k \leq n$, let

$$\begin{aligned} T[k] &= \frac{b-a}{n} \sum_{j=0}^{k-1} \frac{Z[j] + Z[j+1]}{2} \\ &= \frac{b-a}{2n} \sum_{j=0}^{k-1} Z[j] + \frac{b-a}{2n} \sum_{j=1}^k Z[j] \\ &= \frac{1}{2} (S_0[k-1] + S_1[k]), \end{aligned} \quad (211)$$

where

$$S_0[k] = \frac{b-a}{n} \sum_{j=0}^k Z[j], \quad 0 \leq k \leq n-1, \quad (212)$$

and

$$S_1[k] = \frac{b-a}{n} \sum_{j=1}^k Z[j], \quad 1 \leq k \leq n. \quad (213)$$

Lemma 8.6. *For any $t > 0$,*

$$\mathbb{P} \left(\max_{0 \leq k \leq n-1} |S_0[k]| \geq t \right) \leq 2 \exp(-nt^2/2(b-a)^2\sigma^2), \quad (214)$$

and

$$\mathbb{P} \left(\max_{1 \leq k \leq n} |S_1[k]| \geq t \right) \leq 2 \exp(-nt^2/2(b-a)^2\sigma^2). \quad (215)$$

Proof. The method of proof is fairly standard; see, for instance, [53]. Let $\lambda > 0$, and define

$$X[k] = \exp(\lambda S_0[k]), \quad 0 \leq k \leq n-1. \quad (216)$$

Then X is a submartingale, i.e. $\mathbb{E}[X[k] | Z[0], \dots, Z[k-1]] \geq X[k-1]$ for $1 \leq k \leq n-1$. Observe that $S_0[n-1]$ is normally distributed with mean zero and with variance

$$\bar{\sigma}^2 = \frac{(b-a)^2}{n^2} \sum_{j=0}^{n-1} \sigma_j^2 \leq \frac{(b-a)^2}{n} \sigma^2. \quad (217)$$

Consequently, using the standard formula for the Gaussian moment generating function,

$$\mathbb{E}[X[n-1]] = e^{\lambda^2 \sigma^2 / 2} \leq e^{\lambda^2 (b-a)^2 \sigma^2 / 2n}. \quad (218)$$

By Doob's Inequality (e.g. see Theorem 5.4.2 in [19]), for any real number t ,

$$\begin{aligned} \mathbb{P}\left(\max_{0 \leq k \leq n-1} S_0[k] \geq t\right) &= \mathbb{P}\left(\max_{0 \leq k \leq n-1} X[k] \geq \exp(\lambda t)\right) \\ &\leq \mathbb{E}[X[n-1]] \exp(-\lambda t) \\ &\leq e^{\lambda^2 (b-a)^2 \sigma^2 / 2n - \lambda t}. \end{aligned} \quad (219)$$

Taking $\lambda = tn/\sigma^2(b-a)^2$ yields the bound

$$\mathbb{P}\left(\max_{0 \leq k \leq n-1} S_0[k] \geq t\right) \leq \exp(-nt^2/2(b-a)^2\sigma^2). \quad (220)$$

Symmetry and the union bound immediately gives the bound

$$\mathbb{P}\left(\max_{0 \leq k \leq n-1} |S_0[k]| \geq t\right) \leq 2 \exp(-nt^2/2(b-a)^2\sigma^2). \quad (221)$$

An identical argument holds for S_1 , completing the proof. \square

Since $T[0] = 0$ and $T[k] = (S_0[k-1] + S_1[k])/2$ when $k \geq 1$, and since

$$\|Z\|_{\nu_\infty} = \max_{0 \leq k \leq n} |T[k]|, \quad (222)$$

the union bound shows

$$\mathbb{P}(\|Z\|_{\nu_\infty} \geq t) = \mathbb{P}\left(\max_{0 \leq k \leq n} |T[k]| \geq t\right) \leq 4 \exp(-nt^2/2(b-a)^2\sigma^2). \quad (223)$$

This establishes (75) for $p = \infty$; the result then follows for all $p \geq 1$ since $\|Z\|_{\nu_p} \leq \|Z\|_{\nu_\infty}$.

To see that the rest of the theorem follows from (75), observe that since the right side of (75) is summable over n , it follows from the Borel-Cantelli Lemma [10] that $\lim_{n \rightarrow \infty} \|Z\|_{\nu_p} = 0$ almost surely, establishing (76). Furthermore,

$$\begin{aligned} \mathbb{E}[\|Z\|_{\nu_p}] &= \int_0^\infty \mathbb{P}(\|Z\|_{\nu_p} \geq t) dt \\ &\leq 2 \int_0^\infty \exp(-nt^2/2(b-a)^2\sigma^2) dt \\ &= \frac{\sigma(b-a)}{\sqrt{n}} 2 \int_0^\infty \exp(-u^2/2) du, \end{aligned} \quad (224)$$

which establishes (77) and completes the proof of the theorem for Z . The corresponding results for Z_{cen} may be deduced from those of Z and the standard concentration bound for $m(Z) \sim N(0, \sigma^2/n)$:

$$\mathbb{P}\{|m(Z)| > t\} \leq 2e^{-t^2 n / 2\sigma^2}. \quad (225)$$

(See, e.g., Chapter 2 in [68].) This completes the proof of Theorem 5.3.

To prove (78), recall that Theorem 5.1 gives the bound

$$\|f\|_{V^p} - \|f\|_{\nu_p} \leq \frac{C}{n}, \quad (226)$$

where C is a constant not depending on p , t or n . From the triangle inequality we have

$$\begin{aligned} \|Y\|_{\nu_p} - \|f\|_{V^p} &= \|\mathbf{f} + Z\|_{\nu_p} - \|f\|_{V^p} \\ &\leq \|\mathbf{f}\|_{\nu_p} + \|Z\|_{\nu_p} - \|f\|_{V^p} \\ &\leq \frac{C}{n} + \|Z\|_{\nu_p}, \end{aligned} \tag{227}$$

and similarly, since $\|\mathbf{f}\|_{\nu_p} - \|Z\|_{\nu_p} \leq \|Y\|_{\nu_p}$,

$$\begin{aligned} \|f\|_{V^p} - \|Y\|_{\nu_p} &\leq \|f\|_{V^p} - \|\mathbf{f}\|_{\nu_p} + \|Z\|_{\nu_p} \\ &\leq \frac{C}{n} + \|Z\|_{\nu_p}. \end{aligned} \tag{228}$$

Combining (227) and (228) shows

$$\left| \|Y\|_{\nu_p} - \|f\|_{V^p} \right| \leq \frac{C}{n} + \|Z\|_{\nu_p}, \tag{229}$$

If $t - C/n \geq t/2$, which holds for all n sufficiently large, then from Theorem 5.3,

$$\begin{aligned} \mathbb{P} \left\{ \left| \|f\|_{V^p} - \|Y\|_{\nu_p} \right| \geq t \right\} &\leq \mathbb{P} \{ \|Z\|_{\nu_p} \geq t - C/n \} \\ &\leq \mathbb{P} \{ \|Z\|_{\nu_p} \geq t/2 \} \\ &\leq Ae^{-Bn(t/2)^2/\sigma^2}, \end{aligned} \tag{230}$$

which is a bound of the desired form. This completes the proof of (78). The limit (79) follows immediately from (229) and the fact that $\|Z\|_{\nu_p} \rightarrow 0$ almost surely. To prove (80), take expectations of each side of (229) and apply (77). The proofs of the corresponding results for Y_{cen} and f_{cen} are nearly identical.

9 Conclusion

This paper has proven a number of robustness properties of the Volterra distances for functions of a single variable, and the sliced Volterra distances for functions of multiple variables. These results extend previous results known for Wasserstein distances. Our results indicate that the favorable properties of Wasserstein distances may be shared by a wider class of metrics, which may be better suited for certain applications; for instance, the Volterra metrics are defined between functions with negative values and unequal integrals, and are less sensitive to large deformations of the data.

The Volterra metrics are extremely simple: one merely applies a smoothing filter to the input functions, and then computes their Lebesgue distance. It seems likely that one could prove similar results for other families of filters. As such, the present work suggests that constructing metrics with a desired set of robustness properties may only require applying an appropriate collection of filters to the input data, where the filters are designed to smooth the functions with respect to the underlying geometry of their domain. In fact, many metrics that have been considered in recent years are of exactly this type [44, 43, 36, 35, 59, 42].

It is also natural to explore applications of the Volterra distances, sliced Volterra distances, and related metrics to problems where Wasserstein and sliced Wasserstein distances have been used previously. One such area is analysis of data from cryo-electron microscopy (cryo-EM), in which one observes two-variable projections of a three-variable volume (a molecule), at unknown viewing directions, from which the volume is to be determined [60, 5, 17]. Wasserstein metrics have been proposed for clustering images and parameterizing volumes in cryo-EM [62, 52, 69]. Due to their robustness to deformations, Volterra metrics, or other metrics with similar properties, may also be appropriate for heterogeneity analysis in cryo-EM [70, 22, 38, 3, 33, 58, 64, 34], as has been proposed for Wasserstein distances [52]. More generally, it is of interest to explore applications to clustering, nearest neighbor regression, and other metric-based tasks. Questions along these lines will be pursued in future work.

Acknowledgements

I thank Joe Kileel, Amit Moscovich, Rohan Rao, and Amit Singer for stimulating discussions on the papers [52], [30], and [69]. I give additional thanks to Amit Singer for helpful feedback and suggestions on earlier versions of the manuscript. I acknowledge support from NSF BIGDATA award IIS-1837992, BSF award 2018230, and NSF CAREER award DMS-2238821.

References

- [1] Charu C. Aggarwal and Chandan K. Reddy, editors. *Data Clustering: Algorithms and Applications*. CRC Press, 2014.
- [2] Rohit Agrawal and Thibaut Horel. Optimal bounds between f -divergences and integral probability metrics. *The Journal of Machine Learning Research*, 22(1):5662–5720, 2021.
- [3] Yariv Aizenbud and Yoel Shkolnisky. A max-cut approach to heterogeneity in cryo-electron microscopy. *Journal of Mathematical Analysis and Applications*, 479(1):1004–1029, 2019.
- [4] Michael Arbel, Anna Korba, Adil Salim, and Arthur Gretton. Maximum mean discrepancy gradient flow. *Advances in Neural Information Processing Systems*, 32, 2019.
- [5] Tamir Bendory, Alberto Bartesaghi, and Amit Singer. Single-particle cryo-electron microscopy: Mathematical theory, computational challenges, and opportunities. *IEEE Signal Processing Magazine*, 37(2):58–76, 2020.
- [6] Espen Bernton, Pierre E. Jacob, Mathieu Gerber, and Christian P. Robert. On parameter estimation with the Wasserstein distance. *Information and Inference: A Journal of the IMA*, 8(4):657–676, 2019.
- [7] Jean-Paul Berrut. A circular interpretation of the Euler-Maclaurin formula. *Journal of Computational and Applied Mathematics*, 189:375–386, 2006.
- [8] Jean-Paul Berrut and Manfred R. Trummer. Extrapolation quadrature from equispaced samples of functions with jumps. *Numerical Algorithms*, 92:65–88, 2023.
- [9] Nicolas Bonneel, Julien Rabin, Gabriel Peyré, and Hanspeter Pfister. Sliced and Radon Wasserstein barycenters of measures. *Journal of Mathematical Imaging and Vision*, 51:22–45, 2015.
- [10] Tapas Kumar Chandra. *The Borel-Cantelli Lemma*. Springer, 2012.
- [11] Badr-Eddine Chérif-Abdellatif and Pierre Alquier. MMD-Bayes: Robust Bayesian estimation via maximum mean discrepancy. In *Symposium on Advances in Approximate Bayesian Inference*, pages 1–21, 2020.
- [12] Alexander Cloninger, Brita Roy, Carley Riley, and Harlan M. Krumholz. People mover’s distance: Class level geometry using fast pairwise data adaptive transportation costs. *Applied and Computational Harmonic Analysis*, 47(1):248–257, 2019.
- [13] Ronald R. Coifman, Yoel Shkolnisky, Fred J. Sigworth, and Amit Singer. Graph Laplacian tomography from unknown random projections. *IEEE Transactions on Image Processing*, 17(10):1891–1899, 2008.
- [14] Germund Dahlquist and Åke Björck. *Numerical Methods*. Prentice Hall, Inc., 1974.
- [15] Stanley R. Deans. *The Radon Transform and Some Of Its Applications*. 2007.
- [16] Ishan Deshpande, Ziyu Zhang, and Alexander G. Schwing. Generative modeling using the sliced Wasserstein distance. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3483–3491, 2018.
- [17] Allison Doerr. Single-particle cryo-electron microscopy. *Nature Methods*, 13(1):23, 2016.

- [18] Richard M. Dudley. *Real Analysis and Probability*. CRC Press, 2018.
- [19] Rick Durrett. *Probability Theory and Examples*. Cambridge University Press, fourth edition, 2010.
- [20] D.A. Edwards. On the Kantorovich-Rubinstein theorem. *Expositiones Mathematicae*, 29:387–398, 2011.
- [21] Gerald B. Folland. *Real Analysis: Modern Techniques and Their Applications*. John Wiley and Sons, second edition, 1999.
- [22] Joachim Frank and Abbas Ourmazd. Continuous changes in structure mapped by manifold embedding of single-particle data in cryo-EM. *Methods*, 100:61–67, 2016.
- [23] Frank J. Massey, Jr. The Kolmogorov-Smirnov test for goodness of fit. *Journal of the American Statistical Association*, 46(253):68–78, 1951.
- [24] Alison L. Gibbs and Francis Edward Su. On choosing and bounding probability metrics. *International Statistical Review*, 70(3):419–435, 2002.
- [25] Israel Gohberg and Mark Grigorevich Krein. *Theory and Applications of Volterra Operators in Hilbert Space*. AMS, 1970.
- [26] Arthur Gretton, Karsten M. Borgwardt, Malte J. Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *The Journal of Machine Learning Research*, 13(1):723–773, 2012.
- [27] Sigurdur Helgason. *Integral Geometric and Radon Transforms*. 2010.
- [28] Gabor T. Herman. *Fundamentals of Computerized Tomography: Image Reconstruction from Projections*. Springer Science & Business Media, 2009.
- [29] H.G. Kellerer. Duality theorems and probability metrics. In *Proceedings of the Seventh Conference on Probability Theory*, pages 211–220. VNU Press, 1985.
- [30] Joe Kileel, Amit Moscovich, Nathan Zelesko, and Amit Singer. Manifold learning with arbitrary norms. *Journal of Fourier Analysis and Applications*, 27(5):82, 2021.
- [31] Soheil Kolouri, Kimia Nadjahi, Umut Simsekli, Roland Badeau, and Gustavo Rohde. Generalized sliced Wasserstein distances. *Advances in Neural Information Processing Systems*, 32, 2019.
- [32] Soheil Kolouri, Gustavo K. Rohde, and Heiko Hoffmann. Sliced Wasserstein distance for learning Gaussian mixture models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3427–3436, 2018.
- [33] Roy R. Lederman, Joakim Andén, and Amit Singer. Hyper-molecules: on the representation and recovery of dynamical structures for applications in flexible macro-molecules in cryo-EM. *Inverse Problems*, 36(4):044005, 2020.
- [34] Roy R. Lederman and Amit Singer. A representation theory perspective on simultaneous alignment and classification. *Applied and Computational Harmonic Analysis*, 49(3):1001–1024, 2020.
- [35] William Leeb. The mixed lipschitz space and its dual for tree metrics. *Applied and Computational Harmonic Analysis*, 44(3):584–610, 2018.
- [36] William Leeb and Ronald Coifman. Hölder-Lipschitz norms and their duals on spaces with semigroups, with applications to Earth Mover’s Distance. *Journal of Fourier Analysis and Applications*, 22(4):910–953, 2016.
- [37] Elizaveta Levina and Peter Bickel. The Earth Mover’s Distance is the Mallows Distance: Some insights from statistics. In *Proceedings of the Eighth IEEE International Conference on Computer Vision*, volume 2, pages 251–256. IEEE, 2001.

- [38] Hstau Y. Liao, Yaser Hashem, and Joachim Frank. Efficient estimation of three-dimensional covariance and its application in the analysis of heterogeneous samples in cryo-electron microscopy. *Structure*, 23(6):1129–1137, 2015.
- [39] G. Rubinstein L.V. Kantorovich. On a space of completely additive functions. *Vestnik Leningradskogo Universiteta*, 13(7):52–59, 1958.
- [40] G.P. Akilov L.V. Kantorovich. *Functional Analysis*. 2nd edition, 1982.
- [41] Makoto Maejima and Svetlozar T. Rachev. An ideal metric and the rate of convergence to a self-similar process. *The Annals of Probability*, 15(2):708–727, 1987.
- [42] Gal Mishne, Eric Chi, and Ronald Coifman. Co-manifold learning with missing data. In *International Conference on Machine Learning*, pages 4605–4614, 2019.
- [43] Gal Mishne, Ronen Talmon, Israel Cohen, Ronald R. Coifman, and Yuval Kluger. Data-driven tree transforms and metrics. *IEEE Transactions on Signal and Information Processing Over Networks*, 4(3):451–466, 2017.
- [44] Gal Mishne, Ronen Talmon, Ron Meir, Jackie Schiller, Maria Lavzin, Uri Dubin, and Ronald R. Coifman. Hierarchical coupled-geometry analysis for neuronal structure and activity pattern discovery. *IEEE Journal of Selected Topics in Signal Processing*, 10(7):1238–1253, 2016.
- [45] Caroline Moosmüller and Alexander Cloninger. Linear Optimal Transport Embedding: Provable fast Wasserstein distance computation and classification for nonlinear problems. *Information and Inference: A Journal of the IMA*, 12(1):363–389, 2023.
- [46] Alfred Müller. Integral probability metrics and their generating classes of functions. *Advances in Applied Probability*, 29:429–443, 1997.
- [47] Frank Natterer. *The Mathematics of Computerized Tomography*. Wiley, 1986.
- [48] Sloan Nietert, Ziv Goldfeld, Ritwik Sadhu, and Kengo Kato. Statistical, robustness, and computational guarantees for sliced Wasserstein distances. *Advances in Neural Information Processing Systems*, 35:28179–28193, 2022.
- [49] Victor M. Panaretos and Yoav Zemel. Statistical aspects of Wasserstein distances. *Annual Review of Statistics and Its Application*, 6:405–431, 2019.
- [50] Gabriel Peyré and Marco Cuturi. Computational optimal transport: With applications to data science. *Foundations and Trends in Machine Learning*, 11(5-6):355–607, 2019.
- [51] Julien Rabin, Gabriel Peyré, Julie Delon, and Marc Bernot. Wasserstein barycenter and its application to texture mixing. In *International Conference on Scale Space and Variational Methods in Computer Vision*, pages 435–446. Springer, 2011.
- [52] Rohan Rao, Amit Moscovich, and Amit Singer. Wasserstein k-means for clustering tomographic projections. In *NeurIPS 2020, Machine Learning for Structural Biology (MLSB) Workshop*, 2020.
- [53] Daniel Revuz and Marc Yor. *Continuous Martingales and Brownian Motion*. Springer, 2005.
- [54] Philippe Rigollet and Jonathan Weed. Uncoupled isotonic regression via minimum Wasserstein deconvolution. *Information and Inference: A Journal of the IMA*, 8(4):691–717, 2019.
- [55] Lior Rokach. A survey of clustering algorithms. In *Data Mining and Knowledge Discovery Handbook*, pages 269–298. Springer, 2009.
- [56] Yossi Rubner, Carlo Tomasi, and Leonidas J. Guibas. The Earth Mover’s Distance as a metric for image retrieval. *International Journal of Computer Vision*, 40(2):99–121, 2000.

- [57] Filippo Santambrogio. *Optimal Transport for Applied Mathematicians: Calculus of Variations, PDEs, and Modeling*. Birkhäuser, 2015.
- [58] Sjors H. W. Scheres. Processing of structurally heterogeneous cryo-EM data in RELION. *Methods in Enzymology*, 579:125–157, 2016.
- [59] Shirdhonkar Shirdhonkar and David W. Jacobs. Approximate earth mover’s distance in linear time. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8, 2008.
- [60] Amit Singer and Fred J. Sigworth. Computational methods for single-particle electron cryomicroscopy. *Annual Review of Biomedical Data Science*, 3:163–190, 2020.
- [61] Amit Singer and Hau-Tieng Wu. Two-dimensional tomography from noisy projections taken at unknown random directions. *SIAM Journal on Imaging Sciences*, 6(1):136–175, 2013.
- [62] Amit Singer and Ruiyi Yang. Alignment of Density Maps in Wasserstein Distance. *arXiv preprint arXiv:2305.12310v1*, 2023.
- [63] Bharath K. Sriperumbudur, Kenji Fukumizu, Arthur Gretton, Bernhard Schölkopf, and Gert R.G. Lanckriet. On the empirical estimation of integral probability metrics. *Electronic Journal of Statistics*, 6:1550–1599, 2012.
- [64] Bogdan Toader, Fred J. Sigworth, and Roy R. Lederman. Methods for Cryo-EM single particle reconstruction of macromolecules having continuous heterogeneity. *Journal of Molecular Biology*, 435(9):168020, 2023.
- [65] Rosanna Verde and Antonio Irpino. Dynamic clustering of histogram data: using the right metric. In *Selected Contributions in Data Analysis and Classification*, pages 123–134. Springer, 2007.
- [66] Cédric Villani. *Topics in Optimal Transportation*, volume 58 of *Graduate Studies in Mathematics*. AMS, 2001.
- [67] Cédric Villani. *Optimal Transport: Old and New*. Springer, 2008.
- [68] Martin J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge University Press, 2019.
- [69] Nathan Zelesko, Amit Moscovich, Joe Kileel, and Amit Singer. Earthmover-based manifold learning for analyzing molecular conformation spaces. In *IEEE 17th International Symposium on Biomedical Imaging (ISBI)*, pages 1715–1719, 2020.
- [70] Ellen D. Zhong, Tristan Bepler, Bonnie Berger, and Joseph H. Davis. CryoDRGN: reconstruction of heterogeneous cryo-EM structures using neural networks. *Nature Methods*, 18(2):176–185, 2021.