

Mirror Descent on Reproducing Kernel Banach Spaces

Akash Kumar

*Department of Computer Science and Engineering
University of California San Diego*

AKK002@UCSD.EDU

Mikhail Belkin

*Department of Computer Science and Engineering
Halıcıoğlu Data Science Institute
University of California San Diego*

MBELKIN@UCSD.EDU

Parthe Pandit

*Center for Machine Intelligence and Data Science (C-MInDS)
Indian Institute of Technology Bombay, India*

PANDIT@IITB.AC.IN

Abstract

Recent advances in machine learning have led to increased interest in reproducing kernel Banach spaces (RKBS) as a more general framework that extends beyond reproducing kernel Hilbert spaces (RKHS). These works have resulted in the formulation of representer theorems under several regularized learning schemes. However, little is known about an optimization method that encompasses these results in this setting. This paper addresses a learning problem on Banach spaces endowed with a reproducing kernel, focusing on efficient optimization within RKBS. To tackle this challenge, we propose an algorithm based on mirror descent (MDA). Our approach involves an iterative method that employs gradient steps in the dual space of the Banach space using the reproducing kernel.

We analyze the convergence properties of our algorithm under various assumptions and establish two types of results: first, we identify conditions under which a linear convergence rate is achievable, akin to optimization in the Euclidean setting, and provide a proof of the linear rate; second, we demonstrate a standard convergence rate in a constrained setting. Moreover, to instantiate this algorithm in practice, we introduce a novel family of RKBSs with p -norm ($p \neq 2$), characterized by both an explicit dual map and a kernel.

Keywords: reproducing kernel Banach spaces, kernel methods, linear rate, mirror descent, optimization error

1 Introduction

In supervised machine learning, we are given a set of observations $\mathcal{D}_n = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, where the inputs \mathbf{x}_i , $i = 1, 2, \dots, n$ are sampled from a data space \mathcal{X} with corresponding outputs y_i from a label set \mathcal{Y} . The task is to find a function $\hat{f} : \mathcal{X} \rightarrow \mathcal{Y}$, chosen from a *a priori* fixed model class \mathcal{F} , that best predicts $\hat{f}(\mathbf{x}_i) \approx y_i$, $i = 1, 2, \dots$. The choice of the model class \mathcal{F} is usually based on a prior belief in the expressivity of predictive functions that could lead to the optimal classifier f^* (which may not be in \mathcal{F}), as dictated by the nature/environment.

Typically, to find \hat{f} , we consider a minimization problem over the observations D_n with respect to a loss function $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ as follows:

$$\arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i), y_i) \quad (1)$$

Various optimization techniques have been proposed to minimize the optimization error in approximating a solution to Eq. (1), with first-order gradient methods (Nemirovski and Yudin, 1983; Nesterov, 2014; Lee et al., 2016) being the most well-known in the parametric setting, e.g., neural networks.

The choice of \mathcal{F} plays a pivotal role in minimizing the optimization error while approximating the optimal classifier f^* (Bottou and Bousquet, 2007). A fundamental question arises: *Can a function space \mathcal{F} be constructed to minimize approximation error without trading off with optimization error?*

Traditionally, kernel methods, particularly in reproducing kernel Hilbert spaces (RKHSs), have been proposed to bound and lower the approximation error in a regularized empirical setting

$$\min_{f \in \mathcal{F}} \sum_{i=1}^n \ell(y_i, f(x_i)) + \lambda \cdot \Psi(f), \text{ where } \Psi : \mathcal{F} \rightarrow \mathbb{R}_{\geq 0}, \lambda > 0, \quad (2)$$

However, while RKHSs exhibit function space optimality (Scholkopf and Smola, 2001), where the optimal solution to Eq. (2) has a representer that facilitates efficient optimization methods, they may lack expressive power in terms of approximation error.

To overcome this limitation and enrich the diversity of geometric structures and norms, recent research has explored alternative model classes, such as non-Hilbertian *Banach spaces* (characterized by norms that don't adhere to the parallelogram law), aiming for improved approximation capabilities. One noteworthy framework in this context is the *reproducing kernel Banach spaces* (RKBS), as introduced in Zhang et al. (2009); Lin et al. (2022). Within the RKBS framework, several significant problems—such as minimal norm interpolation, regularization networks, support vector machines, sparse learning, and multi-task learning—have been formulated and investigated (Song et al., 2013; Zhang et al., 2009; Ye, 2014; Zhang and Zhang, 2012; Xu and Ye, 2019; Xu, 2023; Wang et al., 2023b).

Additionally, significant efforts have been made to characterize the function spaces learned by neural networks through the perspective of Banach spaces (Bach, 2014; Parhi and Nowak, 2019; Ongie et al., 2020; Parhi and Nowak, 2020; Wright and Gonzalez, 2021). More recently, there has been a shift towards treating this characterization as an optimization problem over RKBS (Spek et al., 2022; Bartolucci et al., 2023; Shilton et al., 2023a; Parhi and Unser, 2023).

Although the aforementioned findings have significantly contributed to our understanding of the approximation capabilities of Banach spaces, the extent to which these formulations address the corresponding *optimization error*, specifically in terms of statistically efficient and/or provable methods for identifying solutions to Eq. (2), remains unexplored in the setting of reproducing kernel Banach spaces.

In this work, we address this gap by proposing an algorithm based on *mirror descent* (Nemirovski and Yudin, 1983). We focus on a general minimization problem defined over the domain of reproducing kernel Banach spaces and demonstrate that the reproducing property allows us to

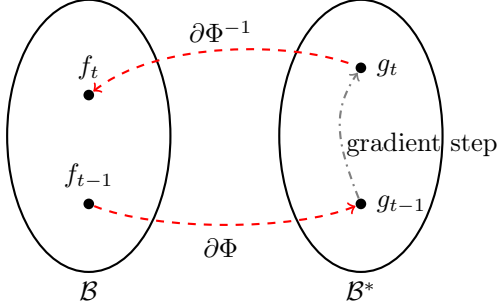


Figure 1: Functional MDA

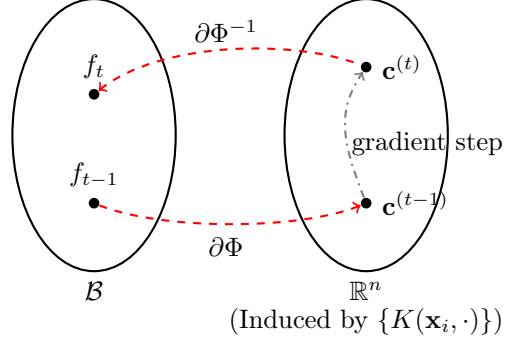


Figure 2: Our MDA for RKBSs

Figure 3: The schematic diagram for the mirror descent algorithm is presented. The first image represents the general functional form of the mirror descent algorithm. In the second image, we illustrate the update rule for Reproducing Kernel Banach Spaces (RKBS), which requires updates in \mathbb{R}^n corresponding to the kernel evaluations on the training data points.

prescribe a gradient update in the dual of the Banach space. Our contributions aim to bridge the gap between approximation and optimization errors in the context of function space design, providing a novel perspective and practical algorithmic results for RKBS.

We summarize the contributions of the work below:

1. **Mirror Descent:** We address a general minimization problem given by Eq. (1), encompassing both regularized and unregularized settings. In this formulation, the function space \mathcal{F} is a non-Hilbertian Banach space, signifying that the norm, denoted as $\|\cdot\|_{\mathcal{F}}$, does not correspond to an inner product. Unlike Hilbert spaces, Banach spaces exhibit a geometric difference in that their dual spaces need not be naturally identifiable with the underlying space. Thus, the gradient of the loss functional ℓ in the first argument does not exist within the function space. Consequently, we turn to the dual of the Banach space, where the gradient steps are executed. These updates in the dual space are then reflected back to the primal Banach space, rendering the algorithm as mirror descent over the Banach space.

Given that derivatives and updates are computed over functions, this approach represents a functional form of mirror descent, distinguishing it from many other works. Furthermore, our focus extends to Banach spaces endowed with a reproducing property, defining them as reproducing kernel Banach spaces. This unique characteristic enables the representation of the mirror descent algorithm as kernel evaluations for updates in the dual space. This, in turn, facilitates simpler and more tractable updates in the dual space expressed in terms of a kernel function.

For example, an instant of the algorithm for optimizing square-error loss functionals with regularization parameter $\lambda > 0$ and using the regularization functional as a mirror map has the form (see Fig. 3):

$$\begin{aligned}
 g_t &\leftarrow (1 - 2\eta\lambda) \cdot g_{t-1} - 2\eta \times (\text{loss differential}) \\
 f_t &\leftarrow \text{InverseMirrormap}(g_t)
 \end{aligned}$$

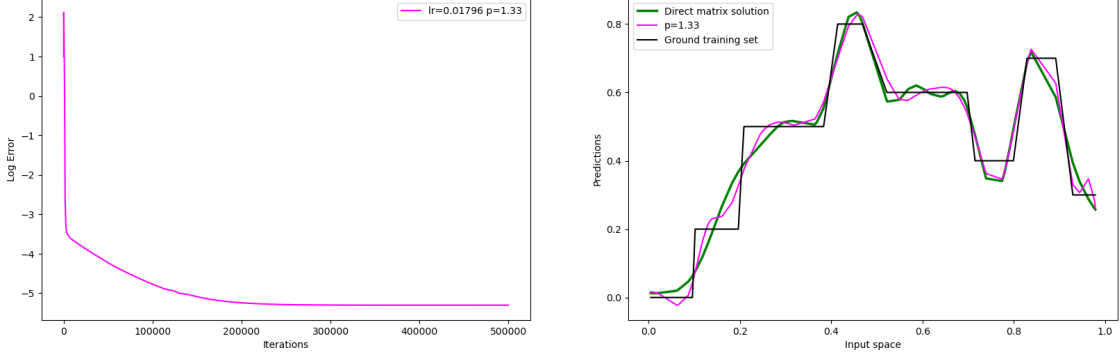


Figure 4: We employ mirror descent on p -norm Reproducing Kernel Banach Spaces (RKBS) for a step function defined on the interval $[0, 1]$. Our training set consists of 80 points, with 15 centers randomly selected from this set. The first image illustrates the logarithm of the training error over iterations. The second image compares the approximation for $p = 1.33$ (depicted by the purple curve) against a solution obtained using NumPy (shown by the green curve). In this comparison, the NumPy solution solves $K_g \cdot \alpha = Y$, where K_g is the similarity matrix computed over the (training, centers) pairs using a Gaussian kernel.

where f_t and g_t are in the primal and dual Banach spaces respectively, and η is the learning rate. Now, if the reproducing kernel K is known, then this could be simplified as

$$\begin{aligned} \mathbf{c}^{(t)} &\leftarrow (1 - 2\eta\lambda)\mathbf{c}^{(t-1)} - 2\eta(f_{t-1}(\mathbf{x}) - Y) \\ f_t &\leftarrow \text{InverseMirrormap}(g_t), \quad g_t = \sum_{i=1}^n \mathbf{c}_i^{(t)} \cdot K(\mathbf{x}_i, \cdot) \end{aligned}$$

where $\mathbf{c}^{(t)} \in \mathbb{R}^n$ (n depends on the training set size), and $f_t(\mathbf{x}) = (f_t(\mathbf{x}_1), f_t(\mathbf{x}_2), \dots, f_t(\mathbf{x}_n))$, $Y = (y_1, y_2, \dots, y_n)$.

We provide the relevant definitions on Banach spaces and the corresponding functional analysis in Section 2, with the MDA discussed in Section 2.2 in details.

Instantiation on an example RKBS: To implement the Mirror Descent Algorithm (MDA), we construct a finite-center based Reproducing Kernel Banach Space (RKBS) (see Section 4). Informally, the Banach space is defined as the set of functions of the following form:

$$\mathcal{B} := \left\{ f : f = \sum_{i=1}^n \alpha_i \cdot H(\cdot, \mathbf{c}_i), \left(\sum_{i=1}^n (|\alpha_i|^p) \right)^{\frac{1}{p}} \lesssim \infty \right\}$$

where $\alpha_i \in \mathbb{R}$, $i = 1, 2, \dots, n$, \mathbf{c}_i are samples/centers in a fixed sample space, and the bivariate function H is defined over pairs of points in that sample space. With minimal assumptions on the linear independence of $H(\cdot, \mathbf{c}_i)$, $i = 1, 2, \dots$, we can demonstrate that it is a reproducing kernel Banach space, and a unique kernel can be explicitly specified.

Additionally, we introduce ℓ_p -type mirror maps, also referred to as dual maps, which, when combined with the kernel, precisely recover the gradient iterations of the Mirror Descent Algorithm (MDA) within this Banach space. In the context of nonregularized learning, the

iterative update steps can be rigorously expressed as follows:

$$\begin{aligned}\beta^{(t)} &\leftarrow \beta^{(t-1)} - \eta \cdot \left(\widehat{H}^\top \left(\widehat{H} \alpha^{(t-1)} - Y \right) \right), \\ \alpha^{(t)} &\leftarrow \frac{\text{sgn}(\beta^{(t)}) |\beta^{(t)}|^{q-1}}{\|\beta^{(t)}\|_q^{q-1}},\end{aligned}$$

where β and α represent the dual and primal parameters, respectively, and \widehat{H} denotes the similarity kernel matrix derived from the training data and centers. This construction of the RKBS is detailed further in Section 4, where we also demonstrate its application to a learning task. An illustration is provided in Fig. 5 in which we show the convergence of the p -norm RKBS for learning a step function where H is an asymmetric Lab-RBF kernel (He et al., 2024).

2. **Convergence of MDA:** More generally, the dual space to a Banach space, comprising real-valued linear transformations defined over the Banach space, may exhibit an unfavorable topology, resulting in the lack of a proper mirroring operation in the mirror descent algorithm. To mitigate such situations, we assume the RKBS is *reflexive*, meaning the dual of the dual space is isometrically isomorphic to the Banach space.

Our main theoretical result establishes the convergence of MDA for reflexive Banach spaces that are Hilbertizable (defined as spaces isomorphic to a Hilbert space¹). This is stated informally as (formally as Theorem 20)

Informal Theorem 1: *If the underlying reflexive Banach space is Hilbertizable, and the loss functional and mirror maps possess properties of (functional notions) μ -strong convexity (see Definition 9) and γ -smoothness (see Definition 10), then in the unconstrained optimization setting (regularized and non-regularized) with a unique minimizer, MDA achieves linear rate depending on the smoothness and strong convexity parameters and a choice of learning rate.*

We show a negative result on the existence of functionals possessing both μ -strong convexity and γ -smoothness for Banach spaces *not* isomorphic to a Hilbert space. Under the same notion, this rate also extends under the (functional form of) *Polyak-Łojasiewicz inequality*, generally studied as a weaker notion than convexity for the loss functional.

Furthermore, under mild assumptions, we extend the analysis to demonstrate the convergence of the algorithm in the constrained setting, achieving a rate of $\mathcal{O}\left(\frac{1}{\sqrt{t}}\right)$, where t is the number of iterations. Formal definitions of (functional) smoothness and convexity of a functional (defined on Banach spaces) are detailed in Section 2.1.

1.1 Related Work

Mirror Descent: The mirror descent method, introduced in Nemirovski and Yudin (1983) for convex optimization problems and later analyzed in Beck and Teboulle (2003), is a well-established first-order method known for its effectiveness in designing algorithms for non-Euclidean geometries. It is particularly recognized for its almost dimension-free rate of convergence, as documented in works like Beck and Teboulle (2003); Ben-Tal et al. (2001); Censor and Zenios (1992); Eckstein (1993); Cesa-Bianchi et al. (2012). More recently, implicit regularization in deep neural networks via mirror descent has been studied in Sun et al. (2022, 2023).

1. *not* isometrically isomorphic

The key principle of the method involves a strongly convex potential function that induces a Bregman divergence (metric) over the problem space, as detailed in [Bauschke et al. \(2017\)](#); [Bubeck \(2015\)](#); [Azizan et al. \(2021\)](#). Subsequently, a gradient descent step is taken in the dual space, which is considered a space of transformations. Previous research has primarily focused on scenarios where the underlying problem space is a Euclidean vector space. In our work, we extend this framework to non-Hilbertian Banach spaces of functions. In this setting, we propose a mirror descent algorithm tailored for a Banach space with a reproducing property, demonstrating favorable convergence properties under standard assumptions.

RKBS: Formally, the concept of reproducing kernel Banach spaces (RKBS) was introduced by [Zhang et al. \(2009\)](#) and later expanded upon in [Lin et al. \(2022\)](#). This framework emerged as a systematic approach to studying learning in Banach spaces, serving as an extension of the well-established reproducing kernel Hilbert spaces (RKHS). In contrast to Hilbert spaces, RKBSs exhibit more intricate geometrical forms and diverse norms. Additionally, RKBS accommodates various kernel functions, including asymmetric kernels ([Zhang and Zhang, 2018](#)) and non-positive definite kernels ([Fukumizu et al., 2011](#)). Recent works have proposed different constructions of RKBSs, such as reflexive RKBS in [Lin et al. \(2022\)](#), semi-inner product (s.i.p) RKBS (refer to [Zhang et al. \(2009\)](#); [Zhang and Zhang \(2011\)](#)), ℓ^1 norm RKBS in [Song et al. \(2013\)](#), and a class of p -norm RKBSs utilizing generalized Mercer kernels ([Xu and Ye, 2019](#)).

While these works have established Representer theorems for various minimization problems in machine learning (detailed in [Lin et al. \(2022\)](#) and [Wang et al. \(2023a\)](#)), a notable gap persists: the absence of a known computational algorithm to find an optimal solution. Our work addresses this gap by presenting an algorithm specifically designed for reflexive RKBSs. We focus on scenarios where an explicit functional form of the point evaluation functional in the dual space of the RKBS can be identified, represented by a unique reproducing kernel.

Learning in Banach spaces: There is a growing interest in understanding and establishing connections between learning problems formulated over Banach spaces ([Bennett and Bredensteiner, 2000](#); [Argyriou et al., 2010](#); [Micchelli and Pontil, 2004, 2007](#); [Unser et al., 2016](#); [Parhi and Nowak, 2019](#); [Srinivasan and Slotine, 2022](#); [Parhi and Unser, 2023](#); [Shilton et al., 2023a](#)).

Numerous significant problems, including p -norm coefficient-based regularization ([Shi et al., 2011](#); [Song et al., 2013](#); [Tong et al., 2010](#)), large-margin classification ([Der and Lee, 2007](#); [Fasshauer et al., 2015](#); [Zhang et al., 2009](#)), lasso in statistics ([Tibshirani, 1996](#)), function spaces of trained neural networks ([Parhi and Nowak, 2021](#); [Shilton et al., 2023a](#)), random Fourier features (RFF) for asymmetric kernels ([qian He et al., 2022](#)), and sparsity in machine learning ([Song et al., 2013](#); [Xu, 2023](#)), have been investigated within the framework of Banach spaces.

Neural Networks and Banach spaces: Recently, there has been a growing interest in understanding deep neural networks through the framework of learning schemes in Banach spaces ([Bartolucci et al., 2023](#); [Shilton et al., 2023b](#); [E et al., 2018](#); [Spek et al., 2022](#); [Parhi and Unser, 2023](#); [Chung and Sun, 2023](#)). A pivotal question in this context is to characterize the function spaces that deep networks learn or represent ([Bach, 2014](#); [Gribonval et al., 2019](#); [Ongie et al., 2020](#); [Parhi and Nowak, 2020](#); [Savarese et al., 2019](#)). In a series of works, [Parhi and Nowak \(2019, 2020, 2021\)](#) established a representer theorem connecting deep ReLU networks with data fitting problems over functions from specific Banach spaces. This line of study was further extended for neural networks with univariate nonlinearity to reproducing kernel Banach spaces ([Bartolucci et al., 2023](#); [Spek et al., 2022](#)), and

later to multivariate nonlinearity in Parhi and Unser (2023). Consequently, any computational algorithm developed for optimization in RKBS would have immediate implications for these results. Additional references on this topic can be found in (Chen et al., 2023; Chung and Sun, 2023; Wright and Gonzalez, 2021).

1.2 Roadmap

In Section 2, we introduce the necessary notation, mathematical formulations, and key definitions related to Reproducing Kernel Banach Spaces (RKBS). Section 2.2 outlines the problem setup and describes the functional mirror descent algorithm. In Section 3, we present our main theoretical results, along with the corresponding proofs of the algorithm’s convergence properties. Finally, Section 4 provides the construction of the p -norm finite-center RKBS and a detailed implementation of the mirror descent algorithm.

2 Preliminaries

Notations: Let a set of letters f, g, h denote elements of a vector space or a dual space. An element of the input space, generally denoted by \mathcal{X} , is represented as \mathbf{x} . Vectors in a Euclidean space are denoted either as α or β . The learning rate is denoted as η . Function spaces, such as Hilbert space, Banach space, among others, are denoted as $\mathcal{H}, \mathcal{B}, \mathcal{C}$. Bivariate functions (over fixed input spaces) or functionals over fixed function spaces are denoted as G, H, K . Bregman divergence over function spaces is denoted by the notation \mathfrak{D} . The dual of an element/space is denoted by $_{-}^*$. There are two sets of notations for non-linear functionals: A, B, C, L (e.g., loss functional) and Φ, Ψ (e.g., mirror maps). All the constants in the theoretical results are denoted by Greek symbols— $\rho, \kappa, \lambda, \gamma$.

Theory of Reproducing kernel Banach spaces: Assume $\mathcal{X} \subseteq \mathbb{R}^d$ is a *locally* compact Hausdorff space (unless stated otherwise). We consider a family of real-valued functions $\mathcal{B} := \{f \mid f : \mathcal{X} \rightarrow \mathbb{R}\}$ over \mathcal{X} . We assume that \mathcal{B} is a complete *vector space* endowed with the norm $\|\cdot\|_{\mathcal{B}}$, *i.e.* $(\mathcal{B}, \|\cdot\|_{\mathcal{B}})$ forms a *Banach space*. Furthermore, we impose the condition that the norms are non-Hilbertian².

For a Banach space \mathcal{B} , we denote its *dual* space as \mathcal{B}^* , *i.e.* a set of real-valued linear transformations on \mathcal{B} , such that for any $g \in \mathcal{B}^*$

$$\|g\|_{\mathcal{B}^*} := \sup_{f \in \mathcal{B}, \|f\|_{\mathcal{B}} \leq 1} g(f)$$

Thus, by definition, the dual space $(\mathcal{B}^*, \|\cdot\|_{\mathcal{B}^*})$ is a Banach space. A transformation $g \in \mathcal{B}^*$ induces a natural action on any $f \in \mathcal{B}$ which is denoted by the bilinear operator, aka duality bracket, $\langle \cdot, \cdot \rangle_{\mathcal{B}} : \mathcal{B} \times \mathcal{B}^* \rightarrow \mathbb{R}$ such that

$$\langle f, g \rangle_{\mathcal{B}} = g(f)$$

Since g is a linear transformation the operation is bilinear in two arguments of the bracket. The **bidual** space of a Banach space \mathcal{B} is the dual of the dual space and is denoted by

$$\mathcal{B}^{**} := (\mathcal{B}^*)^*$$

There is a natural map $\iota = \iota_{\mathcal{B}} : \mathcal{B} \rightarrow \mathcal{B}^{**}$ which assigns to every element $f \in \mathcal{B}$ the linear functional $\iota(f) : \mathcal{B}^* \rightarrow \mathbb{R}$ whose value on any $g \in \mathcal{B}^*$ is obtained by evaluating the linear bounded functional

2. (*i.e.* they don’t satisfy the parallelogram law)

$g : \mathcal{B} \rightarrow \mathbb{R}$ at f , i.e. $g(f)$.

$$\iota(f)(g) = g(f), \quad \forall f \in \mathcal{B}, \forall g \in \mathcal{B}^* \quad (3)$$

As a consequence of Hahn-Banach theorem, the linear map ι is an isometric embedding. Conversely, an interesting class of Banach spaces are ones where \mathcal{B}^{**} can be identified as \mathcal{B} , more formally, they are called *reflexive* spaces as studied in [Buehler and Salamon \(2018\)](#), also defined below

Definition 1 (Reflexive Banach spaces) *A real normed vector space \mathcal{B} is called reflexive if the isometric embedding $\iota : \mathcal{B} \rightarrow \mathcal{B}^{**}$ in Eq. (3) is bijective.*

In other words, a reflexive Banach space \mathcal{B} has the property that $\mathcal{B} \simeq \mathcal{B}^{**}$. In this section, we would use symbol \mathcal{B} for the bidual space as if they *identify* the same.

In order to design a computational algorithm, we are interested in understanding certain properties of a Banach space, in particular, its *reproducing* property, aka a kernel. We adopt the definitions and treatment for a reproducing kernel Banach space (RKBS) as detailed in [Lin et al. \(2022\)](#). We provide formal definitions in the following before stating some useful connections.

Definition 2 (Reproducing kernel Banach spaces (RKBS)) *A reproducing kernel Banach space \mathcal{B} on a prescribed nonempty set \mathcal{X} is a Banach space of functions on \mathcal{X} such that for every $\mathbf{x} \in \mathcal{X}$, the point evaluation functional³ $\delta_{\mathbf{x}} \in \mathcal{B}^*$ on \mathcal{B} is continuous, i.e., there exists a positive constant $C_{\mathbf{x}} \geq 0$ such that*

$$|\delta_{\mathbf{x}}(f)| = |f(\mathbf{x})| \leq C_{\mathbf{x}} \cdot \|f\|_{\mathcal{B}}, \quad \forall f \in \mathcal{B}$$

Definition 3 (Reproducing kernel) *Assume \mathcal{B}_1 is an RKBS on a set Ω_1 . If there exists a Banach space \mathcal{B}_2 of functions on another set Ω_2 , a continuous bilinear form $\langle \cdot, \cdot \rangle_{\mathcal{B}_1 \times \mathcal{B}_2}$, and a function K on $\Omega_1 \times \Omega_2$ such that $K(\mathbf{x}, \cdot) \in \mathcal{B}_2$ for all $\mathbf{x} \in \Omega_1$ and $K(\cdot, \mathbf{y}) \in \mathcal{B}_1$ for all $\mathbf{y} \in \Omega_2$ and*

$$f(\mathbf{x}) = \langle f, K(\mathbf{x}, \cdot) \rangle_{\mathcal{B}_1 \times \mathcal{B}_2} \text{ for all } \mathbf{x} \in \Omega_1 \text{ and for all } f \in \mathcal{B}_1, \quad (4)$$

then we call K a reproducing kernel for \mathcal{B}_1 . If, in addition, \mathcal{B}_2 is also an RKBS on Ω_2 and it holds $K(\cdot, \mathbf{y})$ for all $\mathbf{y} \in \Omega_2$ and

$$g(\mathbf{y}) = \langle K(\cdot, \mathbf{y}), g \rangle_{\mathcal{B}_1 \times \mathcal{B}_2} \text{ for all } \mathbf{y} \in \Omega_2 \text{ and all } g \in \mathcal{B}_2, \quad (5)$$

then we call \mathcal{B}_2 an adjoint RKBS of \mathcal{B}_1 and call \mathcal{B}_1 and \mathcal{B}_2 a pair of RKBSs. In this case, $\tilde{K}(\mathbf{x}, \mathbf{y}) := K(\mathbf{y}, \mathbf{x})$ for $\mathbf{x} \in \Omega_2, \mathbf{y} \in \Omega_1$, is a reproducing kernel for \mathcal{B}_2 .

Eq. (4) and Eq. (5) are called the *reproducing properties* for the kernel K in RKBSs \mathcal{B}_1 and \mathcal{B}_2 , respectively. Note that for different choices of the bilinear form $\langle \cdot, \cdot \rangle_{\mathcal{B}_1 \times \mathcal{B}_2}$ lead to a potentially different reproducing kernel. In this work, we are interested in the case where $\Omega_1 = \mathcal{X}, \Omega_2 = \mathcal{B}$ and $\mathcal{B}_2 = \mathcal{B}^*$. In this regard, if we use the *canonical* bilinear form, i.e. the duality bracket, then we obtain a *unique* reproducing kernel for the RKBS \mathcal{B} . In Section 4, we study this unique kernel corresponding to an RKBS.

3. $\delta_{\mathbf{x}}(f) = f(\mathbf{x})$ for any $f \in \mathcal{B}$ and $\mathbf{x} \in \mathcal{X}$

2.1 Basic Definitions for Functional Optimization

Here, we consider standard notions and definitions for calculus over functionals. It requires a careful treatment to extend the inherent definitions for the Euclidean vectors spaces. We assume that our function space is a Banach space $(\mathcal{B}(\mathcal{X}), \|\cdot\|_{\mathcal{B}})$ of real valued functions on \mathcal{X} . Let $\mathcal{D} \subseteq \mathcal{B}$ be an open set and F a real-valued functional on this domain. First, we provide some necessary definitions for optimization over the function space by defining the notions of differentiability. Then, we talk about the notions of convexity, smoothness, Lipsitzness for F ; extending to the Polyak-Łojasiewicz inequality.

Definition 4 (Gâteaux differential) *Let $F : \mathcal{D} \rightarrow \mathbb{R}$ be a nonlinear transformation. Let $f \in \mathcal{D} \subseteq \mathcal{B}$ and h be arbitrary in \mathcal{B} . If the limit*

$$\partial_f F(h) := \lim_{\gamma \rightarrow 0} \frac{1}{\gamma} [F(f + \gamma \cdot h) - F(f)] \quad (6)$$

exists, it is called the Gâteaux differential of F at f with the increment h . If Eq. (6) exists for each $h \in \mathcal{B}$, then F is called Gâteaux differentiable at f .

For the scope of this work, we have stated the definition of Gâteaux differential in the setting of a Banach space. But it could be easily extended to any vector space that need not have a norm.

Definition 5 (Fréchet derivative) *Consider the transformation F as defined above. If for fixed $f \in \mathcal{D}$ and any $h \in \mathcal{B}$ there exists $\partial_f F(h) \in \mathbb{R}$ which is linear and continuous with respect to h such that*

$$\lim_{\|h\|_{\mathcal{B}} \rightarrow 0} \frac{F(f + h) - F(f) - \partial_f F(h)}{\|h\|_{\mathcal{B}}} = 0, \quad (7)$$

then F is said to be Fréchet differentiable at f and $\partial_f F(h)$ is said to be the Fréchet differential of F at f with increment h .

As is common in the literature of real analysis, we state analogous definitions on the convexity and smoothness of a loss functional.

Definition 6 (Convexity) *We say a functional $F : \mathcal{D} \rightarrow \mathbb{R}$ is convex if for any $f, f' \in \mathcal{B}$ and $\lambda \in [0, 1]$ we have*

$$F(\lambda f + (1 - \lambda)f') \leq \lambda F(f) + (1 - \lambda)F(f')$$

If the inequality is strict we call F a strictly convex functional.

In the functional analysis literature (see [Zălinescu \(1983\)](#)), a slightly stronger notion of convexity is uniform convexity as defined below.

Definition 7 (Uniform convexity) We say a functional $F : \mathcal{D} \rightarrow \mathbb{R}$ is uniformly convex if there exists $\rho : \mathbb{R}_+ \rightarrow \overline{\mathbb{R}}_+$ (with $\rho(t) = 0 \Leftrightarrow t = 0$) such that

$$F(\lambda f + (1 - \lambda)f') \leq \lambda F(f) + (1 - \lambda)F(f') - \lambda(1 - \lambda)\rho(\|f - f'\|)$$

for all $f, f' \in \text{dom } F$ and all $\lambda \in [0, 1]$.

Existence of such functionals strictly implies that the Banach space \mathcal{B} is reflexive (see Theorem 3.5.13 in Zălinescu (1983)). We discuss the requirement of reflexivity of the Banach space in Section 3.

With this, we define the notion of subgradients which is used to define the strong notion of convexity and smoothness of a functional.

Definition 8 (Subgradients) For a functional $F : \mathcal{B} \rightarrow \mathbb{R}$, we define its set of subgradients at any $f_0 \in \mathcal{B}$ as

$$\partial_{f_0} F := \{g \in \mathcal{B}^*; \forall f \in \mathcal{B}, F(f) \geq F(f_0) + \langle f - f_0, g \rangle_{\mathcal{B}}\}$$

Here, we haven't stated any condition on the Fréchet Derivative of F , which if exists renders a singleton set for subgradients.

Definition 9 (μ -strongly convex) We say a functional $F : \mathcal{D} \rightarrow \mathbb{R}$ is μ -strongly convex if for all $f_0 \in \mathcal{B}$, there exists $g \in \partial_{f_0} F$ such that

$$\forall f \in \mathcal{B}, \quad F(f) - F(f_0) \geq \langle f - f_0, g \rangle_{\mathcal{B}} + \frac{\mu}{2} \|f - f_0\|_{\mathcal{B}}^2$$

Note that μ -strong convexity implies that the functional F is uniformly convex for $\rho(t) = \frac{\mu}{2}t^2$.

Definition 10 (γ -smoothness) We say a functional $F : \mathcal{D} \rightarrow \mathbb{R}$ is γ -smooth if for all $f_0 \in \mathcal{B}$, there exists $g \in \partial_{f_0} F$ such that

$$\forall f \in \mathcal{B}, \quad F(f) - F(f_0) \leq \langle f - f_0, g \rangle_{\mathcal{B}} + \frac{\gamma}{2} \|f - f_0\|_{\mathcal{B}}^2 \tag{8}$$

Remark 11 Several notions of smoothness exist in the Hilbert space setting (Nesterov, 2018), and these can be appropriately adapted to the Banach space setting. In this work, we focus on γ -smoothness, which is shown to be equivalent to similar smoothness notions as studied in Theorem 3.1 of Wachsmuth and Wachsmuth (2022). In Section 3, we examine the theoretical properties of the convergence of the specified mirror descent algorithm (see Section 2.2) under various assumptions on the underlying loss functional (see Eq. (1)). It is important to note that any theoretical results, whether assuming γ -smoothness or not, apply to the smoothness concepts discussed in Wachsmuth and Wachsmuth (2022).

In the Euclidean setting, the two constants μ and γ can be directly related as studied in Nesterov (2018). It is straight-forward to note that if a loss functional F is both smooth and strongly convex, then $\mu \leq \gamma$.

Lemma 12 *If a real-valued functional $F : \mathcal{D} \subseteq \mathcal{B} \rightarrow \mathbb{R}$ is μ -strongly convex and γ -smooth, then $\mu \leq \gamma$.*

Now, we state a useful definition on the boundedness of the Fréchet derivative of a functional on \mathcal{B} .

Definition 13 (L-lipshitz) *We say a convex functional $F : \mathcal{D} \rightarrow \mathbb{R}$ is L-lipshitz w.r.t $\|\cdot\|_{\mathcal{B}}$ if for all $f \in \mathcal{B}$, and subdifferential $g \in \partial_f F$,*

$$\|g\|_{\mathcal{B}^*} \leq L$$

Polyak-Łojasiewicz inequality (Polyak, 1963) has been extensively studied notion in the optimization landscape as a relaxation to convexity. We provide an extension of this notion in the functional setting:

Definition 14 (Polyak-Łojasiewicz inequality) *A real-valued functional $F : \mathcal{D} \rightarrow \mathbb{R}$ is called μ -PL if for some $\mu > 0$:*

$$\frac{1}{2} \|\partial_f F\|_{\mathcal{B}^*}^2 \geq \mu(F(f) - F(f^*)), \forall f \in \mathcal{B}$$

where f^* is a global minimizer of F .

2.2 Problem Setup and Algorithmic Insights

Let \mathcal{X} be the data space and \mathcal{Y} a finite label set. We denote a set of n training data points as $D_n := \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\} \subset \mathcal{X} \times \mathcal{Y}$. We consider a modal class \mathcal{B} , a reflexive reproducing kernel Banach space with the goal to find a function $f \in \mathcal{B}$ such that $f(x_i) \approx y_i$ for every example in D_n . For a given function f and data point (\mathbf{x}, y) , we consider a non-negative penalty function $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$, where $\ell(f(\mathbf{x}), y)$ is Fréchet differentiable wrt to f over entire \mathcal{B} unless stated otherwise. Using this penalty function, the total loss on the training set D_n is defined via a loss functional over \mathcal{B} , denoted as $L : \mathcal{B} \rightarrow \mathbb{R}$ as $L(f) = \sum_{i=1}^n \ell(f(x_i), y_i)$. We study the following optimization problem with the loss functional L over \mathcal{B}

$$\min_{f \in \mathcal{B}} L(f) \tag{9}$$

A well-studied choice of loss functional is square loss for which Eq. (9) has the form

$$\min_{f \in \mathcal{B}} \sum_{i=1}^n (f(x_i) - y_i)^2 \tag{10}$$

Equivalently, if the loss functional is regularized with a functional $\Psi_{\geq 0} : \mathcal{B} \rightarrow \mathbb{R}$ with regularization parameter $\lambda > 0$ then we get

$$\min_{f \in \mathcal{B}} \sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \cdot \Psi(f) \tag{11}$$

We solve the optimization problem in Eq. (9) using a *functional form* of mirror descent.

Consider a strongly convex Fréchet differentiable functional $\Phi : \mathcal{B} \rightarrow \mathbb{R}$, also called a *potential* functional in the literature. For the sake of context, we call them mirror maps in this work. The gradient or the Fréchet derivative of Φ can be thought of as an operator from \mathcal{B} to \mathcal{B}^*

$$\partial_{(\cdot)}\Phi : \mathcal{B} \rightarrow \mathcal{B}^* \text{ s.t. } f \mapsto \partial_f\Phi \quad (12)$$

We discuss mirror maps in great details in Section 2.3.

Mirror descent for RKBS \mathcal{B} Assume \mathcal{B} is a reflexive reproducing kernel Banach space of real-valued functions over \mathcal{X} . Assume we a mirror map $\Phi : \mathcal{B} \rightarrow \mathbb{R}$.

We optimize the general minimization problem in Eq. (9) for a given loss functional \mathbf{L} with the following mirror descent algorithm (MDA)

$$g_t := g_{t-1} - \eta \cdot \partial_{f_{t-1}}\mathbf{L} \quad (13a)$$

$$f_t := (\partial\Phi)^{-1}(g_t) \quad (13b)$$

where $g_i \in \mathcal{B}^*$ and $\partial_{f_{t-1}}\mathbf{L}$ is the Fréchet derivative of a loss functional $\mathbf{L} : \mathcal{B} \rightarrow \mathbb{R}$ with respect to f_{t-1} .

We denote the corresponding reproducing kernel to the RKBS \mathcal{B} as $K : \mathcal{X} \times \mathcal{B} \rightarrow \mathbb{R}$. Consider an input \mathbf{x} and the evaluation functional $\delta_{\mathbf{x}} : \mathcal{B} \rightarrow \mathbb{R}$. Now, the Gâteaux differential of $\delta_{\mathbf{x}}(f) = f(\mathbf{x})$ w.r.t f at any function $h \in \mathcal{B}$ is given by

$$\partial_f(\delta_{\mathbf{x}})[h] = \langle h, K(\mathbf{x}, \cdot) \rangle_{\mathcal{B}} \quad (14)$$

Assuming that the loss functional \mathbf{L} is square loss (see Eq. (10)), we can compute the Fréchet derivative of \mathbf{L} wrt f as follows:

$$\begin{aligned} \partial_f\mathbf{L}(\cdot) &= \partial_f \left(\sum_{i=1}^n (f(\mathbf{x}_i) - y_i)^2 \right) = 2 \sum_{i=1}^n (f(\mathbf{x}_i) - y_i) \cdot \partial_f f(\mathbf{x}_i) \\ &= 2 \sum_{i=1}^n (f(\mathbf{x}_i) - y_i) \cdot \partial_f \langle f, K(\mathbf{x}_i, \cdot) \rangle_{\mathcal{B}} \\ &= 2 \sum_{i=1}^n (f(\mathbf{x}_i) - y_i) \cdot K(\mathbf{x}_i, \cdot) \end{aligned} \quad (15)$$

where the last equation follows using Eq. (14). Note that the terms $K(\mathbf{x}_i, \cdot)$ for any i is in the dual space \mathcal{B}^* , and thus the summation is in the dual space.

With this, we provide an *iterative algorithm* using the kernel K :

$$g_t \leftarrow g_{t-1} - 2\eta \sum_{i=1}^n (f_{t-1}(\mathbf{x}_i) - y_i) \cdot K(\mathbf{x}_i, \cdot) \quad (16a)$$

$$f_t \leftarrow (\partial\Phi)^{-1}(g_t) \quad (16b)$$

Notice that if g_0 is initialized as $\sum_{i=1}^n \mathbf{c}_i^{(0)} \cdot K(\mathbf{x}_i, \cdot)$ for some $\mathbf{c}^{(0)} \in \mathbb{R}^n$ then inductively every g_t has the form

$$\begin{aligned}\mathbf{c}^{(t)} &\leftarrow \mathbf{c}^{(t-1)} - 2\eta(f_{t-1}(\mathbf{x}) - Y) \\ g_t &= \sum_{i=1}^n \mathbf{c}_i^{(t)} \cdot K(\mathbf{x}_i, \cdot) \\ f_t &\leftarrow (\partial\Phi)^{-1}(g_t)\end{aligned}$$

where we use the vectorial notation $f_i(\mathbf{x}) = (f_i(\mathbf{x}_1), f_i(\mathbf{x}_2), \dots, f_i(\mathbf{x}_n))$, $Y = (y_1, y_2, \dots, y_n)$. Thus, we can make updates to a vector $\mathbf{c}^{(t)} \in \mathbb{R}^n$ to keep track of gradients in the dual space. Similarly, the MDA can be stipulated in the regularized learning in Eq. (11) where we assume that the mirror map is same as the regularization term Ψ for simplification, i.e. $\Phi = \Psi$

$$g_t \leftarrow (1 - 2\eta\lambda) \cdot g_{t-1} - 2\eta \sum_{i=1}^n (f_{t-1}(\mathbf{x}_i) - y_i) \cdot K(\mathbf{x}_i, \cdot) \quad (18a)$$

$$f_t \leftarrow (\partial\Phi)^{-1}(g_t) \quad (18b)$$

where we have used the fact that $\partial_f\Phi = \partial_f\Psi$. Similarly to the nonregularized setting, this could be further simplified to

$$\mathbf{c}^{(t)} \leftarrow (1 - 2\eta\lambda)\mathbf{c}^{(t-1)} - 2\eta(f_{t-1}(\mathbf{x}) - Y) \quad (19a)$$

$$g_t = \sum_{i=1}^n \mathbf{c}_i^{(t)} \cdot K(\mathbf{x}_i, \cdot) \quad (19b)$$

$$f_t \leftarrow (\partial\Phi)^{-1}(g_t) \quad (19c)$$

Recent works have established representer theorems for different settings in Eq. (9). For example, Zhang et al. (2009) and Lin et al. (2022) studied regularized learning in Eq. (11) for semi-inner product (s.i.p) RKBSs with continuous and convex loss functionals L . They showed that the optimal classifier f in the primal Banach space has a dual representation which can be uniquely written as $g^* = \sum_{i=1}^n c_i^* K(\mathbf{x}_i, \cdot)$ for scalars $c_i^* \in \mathbb{R}$. This emphasizes the novelty of the MDA iterations, as they could potentially achieve the optimal solution to Eq. (9).

Remark 15 Eq. (18) can be easily extended to cases where the Fréchet derivatives of Φ and Ψ map to scaled linear transformations in \mathcal{B}^* . Thus, MDA can be applied to a natural choice studied in the literature where $\Psi = \psi(\|\cdot\|_{\mathcal{B}})$ and $\Phi = \frac{1}{2}\|\cdot\|_{\mathcal{B}}^2$ (or $\phi(\|\cdot\|_{\mathcal{B}}^2)$). However, we can eliminate these assumptions for certain families of RKBSs where the functions are induced by the training points. In such cases, every $g_t \in \mathcal{B}^*$ is induced by $\{K(\mathbf{x}_i, \cdot)\}$, e.g., p -norm RKBSs as constructed in Section 4.

In Section 4, we discuss the construction of an RKBS for which one can stipulate explicit primal and dual space updates for a suitable choice of a mirror map.

Note that MDA for square loss in Eq. (16) and Eq. (18) could be generalized for *any* differentiable loss functional L . The iteration step includes the appropriate Fréchet derivative of the loss functional

$$\partial_f L(f, \{(\mathbf{x}_i, y_i)\}_{i=1}^n) = \sum_{i=1}^n \frac{\partial \ell(z, y_i)}{\partial z} \Big|_{z=f(\mathbf{x}_i)} \cdot \partial_f f(\mathbf{x}_i) = \sum_{i=1}^n \frac{\partial \ell(z, y_i)}{\partial z} \Big|_{z=f(\mathbf{x}_i)} \cdot K(\mathbf{x}_i, \cdot)$$

A natural question is *if the mirror descent algorithm of Eq. (13) is statistically efficient?* We study this in Section 3. This requires appropriate choices of mirror maps to achieve certain convergence guarantees. Below, we provide formal treatment of these maps.

2.3 Mirror Maps

Consider a convex open set of functionals \mathcal{C} such that $\mathcal{B} \subset \bar{\mathcal{C}}$. We are interested in certain real-valued functionals on \mathcal{C} for the mirror descent algorithm of Eq. (16) that could be used to map functions in the primal space \mathcal{B} to transformations in \mathcal{B}^* .

In Section 2.2 we introduced mirror maps to discuss the mirror descent algorithms in Eq. (13). Formally, we define mirror maps as follows:

Definition 16 (Mirror Map) *We say that a functional $\Phi : \mathcal{C} \rightarrow \mathbb{R}$ is a mirror map if it satisfies the following properties:*

1. Φ is strictly convex.
2. Subgradient sets of $\partial_{(\cdot)}\Phi$ don't intersect at non-empty set and $\partial_{(\cdot)}\Phi(\mathcal{B}) = \mathcal{B}^*$.

The Condition 2. above is the key ingredient to the MDA in Eq. (16). The injectivity of subgradient sets makes sure that the algorithm can converge without getting stuck in a loop and the surjectivity makes sure that the algorithm gets back to the primal space reliably. The algorithm would make sense only if the chosen mirror map satisfies this property. We can guarantee this for a wide variety of strictly convex functionals on a reflexive Banach space, including squared p -norms. Note that strict convexity of Φ (in Condition 1.) implies that the subgradient sets at any $f, f' \in \mathcal{B}$ don't intersect, i.e.

$$\partial_f\Phi \cap \partial_{f'}\Phi = \emptyset$$

On the other hand, Condition 2. implies that there exists some f such that for any $g \in \mathcal{B}^*$, $g \in \partial_f\Phi$. In Lemma 17, we establish a strong result on the injectivity and surjectivity of a mirror map.

We assume that the underlying functional is proper, i.e. $F : \mathcal{B} \rightarrow \mathbb{R}$ such that $(\text{dom } F) \neq \emptyset$ and $F(f) \geq -\infty$ for all $f \in \mathcal{B}$. We defer the proof of the lemma to Appendix A.

Lemma 17 *Consider a convex Banach space $(\mathcal{B}, \|\cdot\|_{\mathcal{B}})$. Let $F : \mathcal{B} \rightarrow \mathbb{R}$ be a proper, strictly convex and Gâteaux differentiable functional. Then, for all $f, f' \in \mathcal{B}$,*

$$\partial_f F = \partial_{f'} F \implies f = f' \tag{20}$$

Furthermore, for any linear functional $g^* \in \mathcal{B}^*$, there exists $\hat{f} \in \mathcal{B}$ such that

$$g^* = \partial_{\hat{f}} F$$

In other words, the operator $\partial_{(\cdot)} F : \mathcal{B} \rightarrow \mathcal{B}^*$ is both injective and surjective, where

$$\forall f \in \mathcal{B}, f \longmapsto \partial_{(\cdot)} F(f) := \partial_f F \in \mathcal{B}^*$$

Although, we have stated the lemma for a Gâteaux differentiable mirror app, surjectivity can be achieved in its absence (see the proof in Appendix A).

3 Theoretical results: Convergence of Mirror Descent Algorithm

In this subsection, we consider the unconstrained optimization problem

$$\min_{f \in \mathcal{B}} \mathsf{L}(f) \tag{21}$$

over a Banach space \mathcal{B} and a real loss functional $\mathsf{L} : \mathcal{B} \rightarrow \mathbb{R}$. We would like to understand how well the mirror descent algorithm of Eq. (16) performs if there exists a realizable global minimizer of Eq. (21), *i.e.* there exists $f^* \in \mathcal{B}$ such that $\arg \min_{f \in \mathcal{B}} \mathsf{L}(f) = \{f^*\}$. However, the existence of a global minimizer in the Banach space for Eq. (21) is not guaranteed, a situation that can be resolved by assuming the reflexivity of the space (see Zălinescu (2002, Theorem 2.3.1)).

In this section, we assume that the Banach space \mathcal{B} is reflexive. This assumption serves two purposes: first, it ensures the existence of a global minimizer for Eq. (21) within the space; second, it allows us to study the convergence properties of the mirror descent algorithm under the assumption that the underlying functionals—either the loss functional L or the mirror map Φ —are strongly convex (see Definition 9). The second condition also implies the reflexivity of the space, as stated in the following result.

Theorem 18 (Theorem 3.5.13 Zălinescu (2002)) *Let \mathcal{B} be a Banach space. If there exists a proper, lower semi-continuous, uniformly convex functional $F : \mathcal{B} \rightarrow \mathbb{R}$ whose domain has a nonempty interior, then \mathcal{B} is reflexive.*

3.1 Non-existence of a smooth and strongly convex functional

In the Euclidean setting, the convergence rate for various gradient-based methods applied to the optimization problem

$$\min_{x \in \mathbb{R}^d} f(x), \tag{22}$$

where $f : \mathbb{R}^d \rightarrow \mathbb{R}$, has been extensively studied. To achieve a linear convergence rate with gradient descent (Bubeck, 2015), a common requirement is that the loss function f is both strongly convex and smooth. A natural question arises: *can we achieve a linear rate for the optimization problem in Eq. (21) within reflexive Banach spaces?*

The answer is likely negative for Banach spaces that are not isomorphic to a Hilbert space. We can show that there does not exist a functional $F : \mathcal{B} \rightarrow \mathbb{R}$ that is both strongly convex and γ -smooth on a Banach space that is not isomorphic to a Hilbert space.

We state this negative result regarding the existence of a strongly convex and γ -smooth functional for general Banach spaces as follows. The proof is deferred to Appendix A.2.

Lemma 19 (Existence of strongly convex and γ -smooth functional) *Let \mathcal{B} be a Banach space. If there exists a functional $F : \mathcal{B} \rightarrow \mathbb{R}$ that is both μ -strongly convex and γ -smooth for some $\mu > 0$ and $\gamma < \infty$, then \mathcal{B} is isomorphic to a Hilbert space.*

The proof is based on the characterization of second-order differentiable points of a functional as discussed in Borwein and Noll (1994). Using these points, we can demonstrate the existence of a

Hilbert norm on any separable Banach space. By applying a generalization of the parallelogram law to show the isomorphism of a Banach space onto a Hilbert space (see Kwapień (1972)), the proof extends to general Banach spaces by considering their separable subspaces.

3.2 Conditional linear rate for unconstrained optimization

In the previous section (3.1), we demonstrated that there exist Banach spaces that do not admit functionals which are both strongly convex (see Definition 9) and smooth (see Definition 10). Here, we consider a specific class of RKBSs that satisfy a certain property, which we term ‘Hilbertizable’, and show a linear rate of convergence for unconstrained optimization using the mirror descent algorithm. This rate is achieved for loss functionals that are both strongly convex and γ -smooth.

In Definition 16, we require the mirror map to be strictly convex. To establish convergence guarantees, we further need slightly stronger conditions, namely γ -smoothness and μ -strong convexity.

First, we state an assumption regarding the specific class of Banach spaces that could potentially admit strongly convex and smooth functionals.

Assumption 1 *A Banach space \mathcal{B} is termed Hilbertizable if $(\mathcal{B}, \|\cdot\|_{\mathcal{B}})$ can be isomorphically mapped onto a Hilbert space.*

Although topologically $(\mathcal{B}, \|\cdot\|_{\mathcal{B}})$ is same as a Hilbert space, but note that this is a weaker notion than isometric isomorphism, which, if it exists, implies that the Banach space norm $\|\cdot\|_{\mathcal{B}}$ is a Hilbert norm. Finite sequences with ℓ_p (where $p \in (1, \infty)$) norm are isomorphic to each other. On the other hand, square-summable infinite sequences $\ell^2(\mathbb{N})$ with ℓ_2 and $\ell_2 + \ell_\infty$ norms are isomorphic (but not isometrically) to each other, where one is a Hilbert space and the other Banach space.

Now, we state another assumption that characterizes the existence of a smooth and strongly convex functional.

Assumption 2 *A reflexive Banach space \mathcal{B} is termed smoothly-convex if $(\mathcal{B}, \|\cdot\|_{\mathcal{B}})$ admits a functional $F : \mathcal{B} \rightarrow \mathbb{R}$ which is both μ -strongly convex and γ -smooth (with respect to $\|\cdot\|_{\mathcal{B}}$) for $\mu > 0$ and $\gamma < \infty$.*

Thus, for a Banach space to satisfy Assumption 2, Lemma 19 states that it must satisfy Assumption 1.

Under the aforementioned assumptions on the underlying Banach space, we state and prove the linear rate of MDA as follows:

Theorem 20 (unconstrained optimization) *Consider the optimization problem:*

$$\min_{f \in \mathcal{B}} \mathbf{L}(f)$$

where \mathcal{B} is a reflexive RKBS that satisfies Assumption 1 and Assumption 2. Assume that the loss functional \mathbf{L} is μ -strongly convex and γ -smooth (w.r.t $\|\cdot\|_{\mathcal{B}}$). Furthermore, assume that there exists a mirror map Φ that is ν -strongly convex and ρ -smooth (w.r.t. $\|\cdot\|_{\mathcal{B}}$). Let $f^* \in \mathcal{B}$ be the unique

global minima of the optimization objective, then the mirror descent algorithm of (16) converges to the optimal solution in \mathcal{B} with the learning rate $\eta = \min \left\{ \frac{\nu}{\gamma}, \frac{1}{2\mu\nu\kappa^2} \right\}$. Moreover, the convergence rate is linear, i.e.,

$$\mathsf{L}(f_k) - \mathsf{L}(f^*) \leq (\mathsf{L}(f_0) - \mathsf{L}(f^*)) \cdot e^{-k \cdot \frac{\mu\nu^2\kappa^2}{\gamma}}$$

where κ depends on ρ .

Proof of Theorem 20 In the remainder of the proof, let f_k denote the k -th update in the primal space \mathcal{B} , and let $g_k := \partial_{f_k} \Phi$. Since L or Φ may not be Gâteaux differentiable, we use $\partial_f \mathsf{L}$ or $\partial_f \Phi$ to represent their subgradient sets for any $f \in \mathcal{B}$. In the mirror descent algorithm (MDA), we select one element randomly from these subgradient sets (to resolve ties). So, without loss of generality, we write these sets from the chosen elements. Note that the updates g_k and g_{k-1} in the dual space are related according to the MDA as follows:

$$g_k := g_{k-1} - \eta \cdot \partial_{f_{k-1}} \mathsf{L} \quad (23)$$

Consider the updates f_k and f_{k-1} . Using the strong convexity of Φ , we have

$$\begin{aligned} \Phi(f_k) &\geq \Phi(f_{k-1}) + \langle f_k - f_{k-1}, \partial_{f_{k-1}} \Phi \rangle_{\mathcal{B}} + \frac{\nu}{2} \|f_k - f_{k-1}\|_{\mathcal{B}}^2 \\ \Phi(f_{k-1}) &\geq \Phi(f_k) + \langle f_{k-1} - f_k, \partial_{f_k} \Phi \rangle_{\mathcal{B}} + \frac{\nu}{2} \|f_k - f_{k-1}\|_{\mathcal{B}}^2 \end{aligned}$$

Adding the equations above we get

$$\begin{aligned} 0 &\geq \langle f_k - f_{k-1}, \partial_{f_{k-1}} \Phi - \partial_{f_k} \Phi \rangle_{\mathcal{B}} + \nu \cdot \|f_k - f_{k-1}\|_{\mathcal{B}}^2 \\ &= \langle f_k - f_{k-1}, g_{k-1} - g_k \rangle_{\mathcal{B}} + \nu \cdot \|f_k - f_{k-1}\|_{\mathcal{B}}^2 \\ &= \eta \cdot \langle f_k - f_{k-1}, \partial_{f_{k-1}} \mathsf{L} \rangle_{\mathcal{B}} + \nu \cdot \|f_k - f_{k-1}\|_{\mathcal{B}}^2 \end{aligned} \quad (24)$$

Eq. (24) is useful in the sense that it provides a way to bound $\|f_k - f_{k-1}\|_{\mathcal{B}}^2$ in terms of $\partial_{f_{k-1}} \mathsf{L}$. We would show that the norm of $\|\partial_{f_{k-1}} \mathsf{L}\|_{\mathcal{B}^*}$ could be bounded the other way round in terms of $\|f_k - f_{k-1}\|_{\mathcal{B}}^2$. In order to show that we need the following lemma that establishes the connection of the mirror map Φ to its convex conjugate defined on the dual space.

Proposition 21 ((Zălinescu, 2002, Corollary 3.5.7)) *Let $F : \mathcal{B} \rightarrow \mathbb{R}$ be a continuous convex function and $p, q \in \mathbb{R}$ be such that $1 \leq p \leq 2 \leq q$ and $p^{-1} + q^{-1} = 1$. Then, the following statements are equivalent*

1. $\exists L_2 > 0, \forall f, f' \in \mathcal{B}, \forall g \in \partial_{f'} F:$

$$F(f) \leq F(f') + \langle f - f', g \rangle_{\mathcal{B}} + \frac{L_2}{p} \cdot \|f - f'\|_{\mathcal{B}}^p;$$

2. $\exists L_5 > 0, f, f' \in \mathcal{B},$

$$\langle f - f', \partial_f F - \partial_{f'} F \rangle_{\mathcal{B}} \geq \frac{2}{L_5 q} \cdot \|\partial_f F - \partial_{f'} F\|_{\mathcal{B}^*}^q.$$

We could apply Lemma 21 for the mirror map Φ . Since it is convex and ρ -smooth the condition 1. holds and thus, there exists a scalar $\kappa > 0$ such that $\kappa := \frac{1}{L_3}$ for which condition 2. is satisfied for any functions $f, f' \in \mathcal{B}$.

Using Lemma 21, for the iterates f_k and f_{k-1} we have

$$\begin{aligned} \langle f_k - f_{k-1}, \partial_{f_k} \Phi - \partial_{f_{k-1}} \Phi \rangle_{\mathcal{B}} &\geq \|\partial_{f_k} \Phi - \partial_{f_{k-1}} \Phi\|_{\mathcal{B}^*}^2 \\ \implies \eta \cdot \langle f_k - f_{k-1}, -\partial_{f_{k-1}} \mathbf{L} \rangle_{\mathcal{B}} &\geq \kappa \eta^2 \cdot \|\partial_{f_{k-1}} \mathbf{L}\|_{\mathcal{B}^*}^2 \end{aligned} \quad (25)$$

In the equation above we note that $\partial_{f_k} \Phi = g_k$ for all $k = 0, 1, \dots$. Thus, $\partial_{f_k} \Phi - \partial_{f_{k-1}} \Phi = -\eta \cdot \partial_{f_{k-1}} \mathbf{L}$.

But using Cauchy-Schwartz inequality, we note that

$$\langle f_{k-1} - f_k, \partial_{f_{k-1}} \mathbf{L} \rangle_{\mathcal{B}^*} \leq \|f_{k-1} - f_k\|_{\mathcal{B}} \cdot \|\partial_{f_{k-1}} \mathbf{L}\|_{\mathcal{B}^*}$$

Thus, we can write

$$\|f_k - f_{k-1}\|_{\mathcal{B}} \geq \kappa \eta \cdot \|\partial_{f_{k-1}} \mathbf{L}\|_{\mathcal{B}} \quad (26)$$

Now, using γ -smoothness of \mathbf{L} we note that

$$\begin{aligned} \mathbf{L}(f_k) &\leq \mathbf{L}(f_{k-1}) + \langle f_k - f_{k-1}, \partial_{f_{k-1}} \mathbf{L} \rangle_{\mathcal{B}} + \frac{\gamma}{2} \|f_k - f_{k-1}\|_{\mathcal{B}}^2 \\ &\leq \mathbf{L}(f_{k-1}) + \left(\frac{\gamma}{2} - \frac{\nu}{\eta} \right) \cdot \|f_k - f_{k-1}\|_{\mathcal{B}}^2 \end{aligned} \quad (27)$$

$$\leq \mathbf{L}(f_{k-1}) - \left(\frac{\nu}{\eta} - \frac{\gamma}{2} \right) \cdot \kappa^2 \eta^2 \cdot \|\partial_{f_{k-1}} \mathbf{L}\|_{\mathcal{B}^*}^2 \quad (28)$$

$$= \mathbf{L}(f_{k-1}) - \frac{(2\nu - \gamma\eta)\kappa^2\eta}{2} \cdot \|\partial_{f_{k-1}} \mathbf{L}\|_{\mathcal{B}^*}^2 \quad (29)$$

In Eq. (27), we bound $\langle f_k - f_{k-1}, \partial_{f_{k-1}} \mathbf{L} \rangle_{\mathcal{B}}$ in terms of $\|f_k - f_{k-1}\|_{\mathcal{B}}^2$ using Eq. (24). But then $\|f_k - f_{k-1}\|_{\mathcal{B}}^2$ can be bounded in terms of $\|\partial_{f_{k-1}} \mathbf{L}\|_{\mathcal{B}^*}$ using Eq. (26). Here, note that $\frac{2\nu}{\gamma} > \eta$.

Now, we would try to lower bound $\mathbf{L}(f^*)$ using the definition of the minimizer. Consider the following:

$$\begin{aligned} \mathbf{L}(f^*) &:= \inf_{f \in \mathcal{B}} \mathbf{L}(f) \\ &= \inf_{f \in \mathcal{B}} \mathbf{L}(f_k + f) \\ &\geq \inf_{f \in \mathcal{B}} \mathbf{L}(f_k) + \langle f, \partial_{f_k} \mathbf{L} \rangle_{\mathcal{B}^*} + \frac{\mu}{2} \cdot \|f\|_{\mathcal{B}}^2 \\ &= \mathbf{L}(f_k) + \mu \cdot \left(\inf_{f \in \mathcal{B}} - \left\langle f, -\frac{1}{\mu} \cdot \partial_{f_k} \mathbf{L} \right\rangle_{\mathcal{B}} + \frac{1}{2} \|f\|_{\mathcal{B}}^2 \right) \\ &= \mathbf{L}(f_k) - \frac{1}{2\mu} \cdot \|\partial_{f_{k-1}} \mathbf{L}\|_{\mathcal{B}^*}^2 \end{aligned} \quad (30)$$

In the equation above, we use the fact that the convex conjugate of $\frac{1}{2} \|\cdot\|_{\mathcal{B}}^2$ is $\frac{1}{2} \|\cdot\|_{\mathcal{B}^*}^2$.

Now, we could bound $\mathbf{L}(f_k) - \mathbf{L}(f^*)$ geometrically. Using Eq. (29) and Eq. (30)

$$\begin{aligned} \mathbf{L}(f_k) &\leq \mathbf{L}(f_{k-1}) - \frac{(2\nu - \gamma\eta)\kappa^2\eta}{2} \cdot \|\partial_{f_{k-1}} \mathbf{L}\|_{\mathcal{B}^*}^2 \\ &\leq \mathbf{L}(f_{k-1}) + \mu(2\nu - \gamma\eta)\kappa^2\eta \cdot (\mathbf{L}(f^*) - \mathbf{L}(f_{k-1})) \end{aligned}$$

Now, this could be simplified as follows:

$$\begin{aligned}
\mathbb{L}(f_k) - \mathbb{L}(f^*) &\leq -(\mathbb{L}(f^*) - \mathbb{L}(f_{k-1})) + \mu(2\nu - L\eta)\kappa^2\eta \cdot (\mathbb{L}(f^*) - \mathbb{L}(f_{k-1})) \\
&\leq (1 - \mu(2\nu - \gamma\eta)\kappa^2\eta) (\mathbb{L}(f_{k-1}) - \mathbb{L}(f^*)) \\
&\vdots \\
&= (1 - \mu(2\nu - \gamma\eta)\kappa^2\eta)^k (\mathbb{L}(f_0) - \mathbb{L}(f^*)) \\
&\leq \lambda \cdot \exp(-k \cdot \mu(2\nu - \gamma\eta)\kappa^2\eta) \\
&= \lambda \cdot \exp\left(-k \cdot \frac{\mu\nu^2\kappa^2}{\gamma}\right)
\end{aligned}$$

where we define constant $\lambda := (\mathbb{L}(f_0) - \mathbb{L}(f^*))$. With this, we have completed the proof. \blacksquare

Extension to PŁ condition on the loss functional Now, we consider relaxing the strong convexity condition (see Definition 9) on the loss functional to establish Theorem 22. In Definition 14, we introduced the Polyak-Łojasiewicz (PŁ) condition for the functional setup.

The PŁ inequality (Polyak, 1963; Łojasiewicz, 1963) has been extensively studied as a weaker condition on a loss function compared to strong convexity in minimization problems within the Euclidean setting (see Eq. (22)). This condition is commonly encountered in applications and often leads to fast convergence of numerical algorithms. Under the PŁ condition, a linear rate of convergence for gradient descent in the classical setting above can be demonstrated (Karimi et al., 2016).

Thus, it is natural to explore this in the functional setup. Indeed, it is sufficient for the convergence of the MDA. The rate of convergence remains linear if we replace the μ -strong convexity of the loss functional \mathbb{L} with the μ -PŁ condition. We state the result as follows:

Theorem 22 (Unconstrained Optimization under μ -PŁ Condition) *Consider the optimization problem:*

$$\min_{f \in \mathcal{B}} \mathbb{L}(f), \tag{31}$$

where \mathcal{B} is a reflexive RKBS that satisfies Assumption 1 and Assumption 2. Assume that the loss functional \mathbb{L} is μ -PŁ and γ -smooth (with respect to $\|\cdot\|_{\mathcal{B}}$). Furthermore, assume that there exists a mirror map Φ which is ν -strongly convex and ρ -smooth (with respect to $\|\cdot\|_{\mathcal{B}}$). Let $f^* \in \mathcal{B}$ be the unique global minimizer of the optimization objective. Then, the mirror descent algorithm described by (16) converges to the optimal solution in \mathcal{B} with the learning rate $\eta = \min\left\{\frac{\nu}{\gamma}, \frac{1}{2\mu\nu\kappa^2}\right\}$. Moreover, the convergence rate is linear, i.e.,

$$\mathbb{L}(f_k) - \mathbb{L}(f^*) \leq (\mathbb{L}(f_0) - \mathbb{L}(f^*)) \cdot e^{-k \cdot \frac{\mu\nu^2\kappa^2}{\gamma}},$$

where κ depends on ρ .

The proof of this theorem follows as before, with the addition that Eq. (30) holds due to the PŁ condition of the loss functional \mathbb{L} .

3.3 Constrained Optimization: Convergence of projected MDA

In the previous sections, we study the conditions under which MDA can achieve linear rate of convergence for unconstrained optimization over a reproducing kernel Banach space. Having demonstrated this linear rate in the unconstrained setting, an important question is how MDA, as expressed in Eq. (16), behaves when used in a constrained setting. In this scenario, the algorithm is restricted to choosing functions solely from a constrained subset of the Banach space \mathcal{B} .

We establish that under mild assumptions (not requiring Assumption 1 and Assumption 2) on the mirror map and the loss functional, a rate of $\mathcal{O}\left(\frac{1}{\sqrt{t}}\right)$ can be attained. This paradigm has undergone extensive examination in the classical setting; Bubeck (2015) explored various algorithms employing mirror descent in Banach spaces defined over Euclidean spaces, and with due consideration, these results can be extended to the functional case as well.

Mirror Descent for constrained optimization Let $\mathcal{B}_0 \subseteq \mathcal{B}$ be a compact convex set of functionals. Consider a loss functional $F : \mathcal{B} \rightarrow \mathbb{R}$. We are interested in the following optimization problem:

$$\min_{f \in \mathcal{B}_0} L(f) \quad (32)$$

Now, note that if we try to run the mirror descent algorithm of Eq. (16), then the preimage of g_t , *i.e.* $(\partial\Phi)^{-1}(g_t)$ may not be in \mathcal{B}_0 . A simple solution to this issue could be to project the preimage back to \mathcal{B}_0 . In order to do this, we study the notion of a Bregman divergence.

Consider a real-valued functional $A : \mathcal{B} \rightarrow \mathbb{R}$. We define the *Bregman divergence* \mathfrak{D}_A wrt A as follows:

$$f, f' \in \mathcal{B}, \quad \mathfrak{D}_A(f, f') := A(f) - A(f') - \langle f - f', \partial_{f'} A \rangle_{\mathcal{B}} \quad (33)$$

This gives a straight-forward identity— for any $f, f', h \in \mathcal{B}$

$$\langle f - h, \partial_f A - \partial_{f'} A \rangle_{\mathcal{B}} = \mathfrak{D}_A(f, f') + \mathfrak{D}_A(h, f) - \mathfrak{D}_A(h, f') \quad (34)$$

Consider a convex set of functionals $\mathcal{C} \subseteq \mathcal{B}$ such that closure of \mathcal{C} contains \mathcal{B}_0 , *i.e.* $\mathcal{B}_0 \subseteq \bar{\mathcal{C}}$. We consider mirror maps that satisfy Condition 2. in Definition 16 for the set \mathcal{C} .

Definition 23 (Projection Map) For a mirror map $\Phi : \mathcal{C} \rightarrow \mathbb{R}$, we define a projection map $\Pi_{\mathcal{B}_0}^{\Phi}$ for any $f' \in \mathcal{C}$ as

$$\Pi_{\mathcal{B}_0}^{\Phi}(f') := \arg \min_{f \in \mathcal{B}_0 \cap \mathcal{C}} \mathfrak{D}_{\Phi}(f, f')$$

There are pathological cases where the rhs arg min is an empty set. To avoid such cases we further assume that Φ diverges on the boundary of \mathcal{C} , *i.e.* $\lim_{f \rightarrow \partial\mathcal{C}} \|\partial_f \Phi\|_{\mathcal{B}^*} = +\infty$, thus it is coersive.

With this definition, we propose *projected* mirror descent algorithm to solve the optimization problem in Eq. (32) as follows. Assuming that the loss function L is squared error, *i.e.* $L(f) = \sum_{i=1}^n (f(x_i) - y_i)^2$, we provide a projected iterative algorithm using a reproducing kernel K for \mathcal{B}

and a projection map $\Pi_{\mathcal{B}_0}^\Phi$ (see Definition 23) as follows:

$$g_t \leftarrow g_{t-1} - 2\eta \sum_{i=1}^n (f_{t-1}(x_i) - y_i) \cdot K(\cdot, x_i)$$

$$f_t \leftarrow \Pi_{\mathcal{B}_0}^\Phi \left((\partial\Phi)^{-1}(g_t) \right)$$

In the iteration step, we haven't assumed Gâteaux differentiability of Φ , thus we treat $\partial\Phi$ as a subgradient set. Since a strictly convex functional has non-intersecting subgradient sets, the projection is well-defined. Furthermore, if g_0 is initialized as $\sum_{i=1}^n \mathbf{c}_i^{(0)} \cdot K(\mathbf{x}_i, \cdot)$ for some $\mathbf{c}^{(0)} \in \mathbb{R}^n$ then inductively every g_t has a form similar to Eq. (17). Although the algorithm requires a reproducing kernel from computational point of view, we can show that the projected algorithm would converge even under mild assumptions on a general loss functional L . For the rest of the section, we would consider the following iteration steps to avoid the explicit existence of a reproducing kernel for a reflexive Banach space \mathcal{B} .

$$g_t \leftarrow g_{t-1} - \eta \cdot \partial_{f_{t-1}} \mathsf{L} \tag{36a}$$

$$f_t \leftarrow \Pi_{\mathcal{B}_0}^\Phi \left((\partial\Phi)^{-1}(g_t) \right) \tag{36b}$$

Rest of the section, we will show the convergence of the steps above. First, we show a useful result on the optimality of the projection in the projected mirror descent algorithm with the proof detailed in Appendix A.3. This essentially states that even after the projection the algorithm always points in the direction of descent.

Lemma 24 *Let $f \in \mathcal{B}_0 \cap \mathcal{C}$ and $f' \in \mathcal{C}$, then*

$$\left\langle \Pi_{\mathcal{B}_0}^\Phi(f') - f, \partial_{\Pi_{\mathcal{B}_0}^\Phi(f')} \Phi - \partial_{f'} \Phi \right\rangle_{\mathcal{B}} \leq 0$$

We demonstrate convergence for a class of loss functionals that are L -Lipschitz (Definition 13) and convex. Using the computation of Frechet derivative in Eq. (15), we can show that square and cross-entropy loss functionals satisfy this property if the kernel K and functions $f \in \mathcal{B}$ are bounded. Additionally, the mirror map is required to be μ -strongly convex. A notable example of such mirror maps is the squared p -norms ($1 < p < 2$) on finite-dimensional Banach spaces. In Section 4, we construct a p -norm RKBS and provide a concrete algorithm for mirror maps of the form $\frac{1}{2} \|\cdot\|_p^2$ (see Algorithm 1).

Now, we state the main result of this subsection on the convergence of projected MDA in (36), with the proof detailed below. In the discussion above, we did not assume Φ to be Gâteaux differentiable; however, for ease of discussion and without loss of generality, we assume it in the following statement.

Theorem 25 (constrained optimization) *Assume \mathcal{B} be a reflexive Banach space. Let $\mathcal{B}_0 \subseteq \mathcal{B}$ be a compact convex set of functionals. Consider the optimization problem:*

$$\min_{f \in \mathcal{B}_0} \mathsf{L}(f)$$

Assume that the loss functional L is convex and L -Lipshitz (wrt $\|\cdot\|_{\mathcal{B}}$). Furthermore, assume that there exists a proper mirror map $\Phi : \mathcal{C} \rightarrow \mathbb{R}$ that is Gâteaux differentiable and μ -strongly convex over

$\mathcal{C} \cap \mathcal{B}_0$ (wrt $\|\cdot\|_{\mathcal{B}}$), where \mathcal{C} is the closure of \mathcal{B}_0 . Let $R^2 = \max_{f, f' \in \mathcal{B}_0 \cap \mathcal{C}} |\Phi(f) - \Phi(f')|$. Then, there is a choice of η such that the projected mirror descent algorithm of Eq. (36) converges to the optimal solution in \mathcal{B}_0 at the rate $\tilde{O}(\frac{1}{\sqrt{t}})$.

Proof The proof idea depends on bounding $L(f_t) - L(f^*)$ in terms of Bregman divergences of iterates f_i 's wrt the mirror map Φ . Note that we don't assume anything on the differentiability of the loss functional L , thus $\partial_f L$ needn't have a *unique* subdifferential. Thus, at each iteration step t of projected MDA in Eq. (36), we use a subdifferential $\hat{g}_t \in \partial_{f_t} L$. In the proof below, we would use $\partial_{f_t} L$ for the subdifferential without loss of generality.

Also, since Φ is proper and strongly convex, using Lemma 17 we know that $\partial_{(\cdot)} \Phi$ is invertible. Since \mathcal{B} is reflexive, convex conjugate, denoted as Φ^* , provides the inverse operator from \mathcal{B}^* to \mathcal{B} .

Using convexity of L and Lemma 24 we can write

$$\begin{aligned} L(f_t) - L(f^*) &\leq \langle f_t - f^*, \partial_{f_t} L \rangle_{\mathcal{B}} \\ &\leq \frac{1}{\eta} \langle f_t - f^*, \partial_{f_t} \Phi - g_{t+1} \rangle_{\mathcal{B}} \end{aligned} \quad (37)$$

$$= \frac{1}{\eta} \left\langle f_t - f^*, \partial_{f_t} \Phi - \partial_{\partial_{g_{t+1}} \Phi^*} \Phi \right\rangle_{\mathcal{B}} \quad (38)$$

In Eq. (37), we have used the gradient step in the dual space (see Eq. (36)). Since $\partial_{(\cdot)} \Phi^* : \mathcal{B}^* \rightarrow \mathcal{B}$ is the inverse operator to $\partial_{(\cdot)} \Phi$, we can write $g_{t+1} = \partial_{(\partial_{g_{t+1}} \Phi^*)} \Phi$.

Now, using the identity on the Bregman divergence \mathfrak{D}_{Φ} in Eq. (34), we can write

$$\begin{aligned} &\left\langle f_t - f^*, \partial_{f_t} \Phi - \partial_{\partial_{g_{t+1}} \Phi^*} \Phi \right\rangle_{\mathcal{B}} \\ &= \mathfrak{D}_{\Phi}(f_t, \partial_{g_{t+1}} \Phi^*) + \mathfrak{D}_{\Phi}(f^*, f_t) - \mathfrak{D}_{\Phi}(f^*, \partial_{g_{t+1}} \Phi^*) \\ &\leq \mathfrak{D}_{\Phi}(f_t, \partial_{g_{t+1}} \Phi^*) + \mathfrak{D}_{\Phi}(f^*, f_t) - \mathfrak{D}_{\Phi}(f^*, f_{t+1}) - \mathfrak{D}_{\Phi}(f_{t+1}, \partial_{g_{t+1}} \Phi^*) \\ &= \underbrace{\mathfrak{D}_{\Phi}(f^*, f_t) - \mathfrak{D}_{\Phi}(f^*, f_{t+1})}_I + \underbrace{\mathfrak{D}_{\Phi}(f_t, \partial_{g_{t+1}} \Phi^*) - \mathfrak{D}_{\Phi}(f_{t+1}, \partial_{g_{t+1}} \Phi^*)}_{II} \end{aligned}$$

Now, rest of the proof follows similar steps as in Bubeck (2015) (Theorem 4.2). Summing over $i = 1$ to $i = t$ gives the following bound on (I):

$$\begin{aligned} (I) &= \sum_{i=1}^t \mathfrak{D}_{\Phi}(f^*, f_i) - \mathfrak{D}_{\Phi}(f^*, f_{i+1}) \\ &= \mathfrak{D}_{\Phi}(f^*, f_1) - \mathfrak{D}_{\Phi}(f^*, f_{t+1}) \\ &\leq R^2 \end{aligned}$$

To bound (II), we note the following:

$$\begin{aligned} (II) &= \mathfrak{D}_{\Phi}(f_t, \partial_{g_{t+1}} \Phi^*) - \mathfrak{D}_{\Phi}(f_{t+1}, \partial_{g_{t+1}} \Phi^*) \\ &= \Phi(f_t) - \Phi(f_{t+1}) - \left\langle f_t - f_{t+1}, \partial_{\partial_{g_{t+1}} \Phi^*} \Phi \right\rangle_{\mathcal{B}} \end{aligned} \quad (39)$$

$$\leq \langle f_t - f_{t+1}, \partial_{f_t} \Phi \rangle_{\mathcal{B}} - \frac{\mu}{2} \|f_t - f_{t+1}\|_{\mathcal{B}}^2 - \left\langle f_t - f_{t+1}, \partial_{\partial_{g_{t+1}} \Phi^*} \Phi \right\rangle_{\mathcal{B}} \quad (40)$$

$$\begin{aligned}
&= \left\langle f_t - f_{t+1}, \partial_{f_t} \Phi - \partial_{\partial_{g_{t+1}} \Phi^*} \Phi \right\rangle_{\mathcal{B}} - \frac{\mu}{2} \|f_t - f_{t+1}\|_{\mathcal{B}}^2 \\
&= \langle f_t - f_{t+1}, \eta \cdot \partial_{f_t} \mathsf{L} \rangle_{\mathcal{B}} - \frac{\mu}{2} \|f_t - f_{t+1}\|_{\mathcal{B}}^2 \tag{41}
\end{aligned}$$

$$\leq \eta \cdot \|f_t - f_{t+1}\|_{\mathcal{B}} \cdot \|\partial_{f_t} \mathsf{L}\|_{\mathcal{B}^*} - \frac{\mu}{2} \|f_t - f_{t+1}\|_{\mathcal{B}}^2 \tag{42}$$

$$\leq \eta L \cdot \|f_t - f_{t+1}\|_{\mathcal{B}} - \frac{\mu}{2} \|f_t - f_{t+1}\|_{\mathcal{B}}^2 \tag{43}$$

$$\leq \frac{(\eta L)^2}{2\mu} \tag{44}$$

In Eq. (39), we have used the definition of the Bregman divergence as shown in Eq. (33). We bound the difference $\Phi(f_t) - \Phi(f_{t+1})$ using the μ -strong convexity of the mirror map Φ in Eq. (40). In Eq. (41), we use the mirror descent update of Eq. (36). Now, it remains to bound $\langle f_t - f_{t+1}, \eta \cdot \partial_{f_t} \mathsf{L} \rangle_{\mathcal{B}}$ in terms of $\|f_t - f_{t+1}\|_{\mathcal{B}}$ for which we use the Cauchy-Schwarz inequality, which is what we achieve in Eq. (42). In Eq. (43), we used the L -Lipshitzness of the loss functional L . Finally, we use the following inequality $aq - bq^2 \leq \frac{a^2}{4b}$ for any $q \in \mathbb{R}$ to achieve the final bound in Eq. (44).

Now, adding for t iterations we get

$$\sum_{i=1}^t (\mathsf{L}(f_i) - \mathsf{L}(f^*)) = \frac{1}{\eta} \left(\mathfrak{D}_{\Phi}(f^*, f_1) - \mathfrak{D}_{\Phi}(f^*, f_{t+1}) + \frac{t(\eta L)^2}{2\mu} \right) \leq \frac{R^2}{\eta} + \frac{\eta t L^2}{2\mu}$$

But using convexity

$$\frac{1}{t} \sum_{i=1}^t \mathsf{L}(f_i) - \mathsf{L}(f^*) \leq \mathsf{L} \left(\frac{1}{t} \sum_{i=1}^t f_i \right) - \mathsf{L}(f^*) \leq \frac{R^2}{t\eta} + \frac{\eta L^2}{2\mu}$$

If apriori we set $\eta = \frac{R}{L} \sqrt{\frac{2\mu}{t}}$, then $\frac{R^2}{t\eta} + \frac{\eta L^2}{2\mu} = \frac{\sqrt{2RL}}{\sqrt{\mu t}}$, which gives the stated claim on the rate of convergence. ■

Remark 26 *Theorem 25 is stated without any assumption on the smoothness and differentiability of the loss functional L . But with β -smoothness, one can hope to achieve a rate of $O(\frac{1}{t})$ with a slight variant of MDA by extending the idea of mirror prox as shown in Bubeck (2015) to the functional setting.*

4 Example of MD on an RKBS

In this section, we show a construction of a non-Hilbertian Banach space with a reproducing property. This space is spanned by eigenfunctions centered at chosen datapoints. Using these eigenfunctions, that are written in the form of kernel evaluations, we construct the unique kernel of the Banach space. Furthermore, with a standard ℓ_p type mirror map, we provide the explicit form of the mirror descent in Algorithm 1.

Consider a locally compact Hausdorff set $\mathcal{X} \subseteq \mathbb{R}^d$. Consider a bivariate function $H : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R} \setminus \{-\infty, +\infty\}$. Fix a set of centers $C := \{\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k\} \subset \mathcal{X}$ such that $\{H(\cdot, \mathbf{c}_i)\}_{i=1}^k$ is a set of linearly independent functions over the field \mathbb{R} .

Now, for $p \neq 2$ and $k > 0$, we define the space of functions $\ell_p^k(C)$ as follows:

$$\ell_p^k(C) := \left\{ \sum_{i=1}^k \alpha_i \cdot H(\cdot, \mathbf{c}_i) \mid \alpha \in \mathbb{R}^k, \left(\sum_{i=1}^k |\alpha_i|^p \right)^{\frac{1}{p}} < \infty \right\} \quad (45)$$

The underlying norm is induced by the ℓ_p norm, i.e. for any $f_\alpha \in \ell_p^k$, $\|f_\alpha\|_{\ell_p^k} = \left(\sum_{i=1}^k |\alpha_i|^p \right)^{\frac{1}{p}}$. Note that due to the linear independence assumption, each function $f \in \ell_p^k$ has a unique representation in terms of the basis functions, whereby the norm is uniquely defined. [Nemirovski and Yudin \(1983\)](#) studied a similar Banach space with finite length sequences of size n with p -norm. Here, we differ by using the dot product of the vectors α with the evaluations of the bivariate map H .

In this section, we study the real vector space $\left(\ell_p^k(C), \|\cdot\|_{\ell_p^k} \right)$. For ease of notation we would only write ℓ_p^k unless stated otherwise to signify a different choice of centers or bivariate function H . We will show that $\left(\ell_p^k(C), \|\cdot\|_{\ell_p^k} \right)$ is a reproducing kernel Banach space.

First, we will show that the norm $\|\cdot\|_{\ell_p^k}$ is not induced by an inner product if $p \neq 2$. The key is to observe that any inner product that induces the norm $\|\cdot\|_{\ell_p^k}$ would violate the parallelogram law.

Let $\alpha = (1, 1, 0, \dots, 0)$ and $\alpha' = (1, -1, 0, \dots, 0)$. Note that by parallelogram law,

$$2\|f_\alpha\|_{\ell_p^k}^2 + 2\|f_{\alpha'}\|_{\ell_p^k}^2 = \|f_\alpha + f_{\alpha'}\|_{\ell_p^k}^2 + \|f_\alpha - f_{\alpha'}\|_{\ell_p^k}^2 \implies 4 \cdot 2^{2/p} = 8 \implies p = 2.$$

Note that this space is isometrically isomorphic to the standard ℓ_p space of finite sequences with the $\|\cdot\|_{\ell_p}$ norm. Thus, $\left(\ell_p^k(C), \|\cdot\|_{\ell_p^k} \right)$ is a complete normed vector space, aka a Banach space.

Similarly, for $q = 1 + \frac{q}{p}$, we define the Banach space ℓ_q^k for a set linearly independent functions $\{H(\mathbf{c}_i, \cdot)\}_{i=1}^k$ (over the field \mathbb{R}) with the corresponding norm $\|\cdot\|_{\ell_q^k}$ as follows:

$$\ell_q^k(C) := \left\{ \sum_{i=1}^k \alpha_i \cdot H(\mathbf{c}_i, \cdot) \mid \alpha \in \mathbb{R}^k, \left(\sum_{i=1}^k |\alpha_i|^q \right)^{\frac{1}{q}} < \infty \right\} \quad (46)$$

We would use $(\ell_q^k, \|\cdot\|_{\ell_q^k})$ with the implicit understanding that space in Eq. (46) is assumed.

Remark 27 We assume that the span $\langle H(\mathbf{c}_i, \cdot) : i \in [k] \rangle$ to be dense in ℓ_q^k and span $\langle H(\cdot, \mathbf{c}_i) : i \in [k] \rangle$ to be dense in ℓ_p^k .

In the following, we would show that $\left(\ell_p^k, \|\cdot\|_{\ell_p^k} \right)$ is endowed with the reproducing property (see Definition 3).

Lemma 28 For $1 < p < \infty$, $\left(\ell_p^k, \|\cdot\|_{\ell_p^k} \right)$ is a reproducing kernel Banach space.

Proof We need to show that the evaluation functional on the space ℓ_p^k is continuous, or there exists a positive constant $C_{\mathbf{x}}$ for any given $\mathbf{x} \in X$ such that

$$|\delta_{\mathbf{x}}(f_{\alpha})| = |f_{\alpha}(\mathbf{x})| \leq C_{\mathbf{x}} \cdot \|f_{\alpha}\|_{\ell_p^k}$$

for all $f_{\alpha} \in \ell_p^k$. It is easy to note for $C_{\mathbf{x}} = k \cdot \max_{i \in [k]} |H(\mathbf{c}_i, \mathbf{x})|$ suffices for this as

$$|f_{\alpha}(x)| = \left| \sum_{i=1}^k \alpha_i \cdot H(\mathbf{x}, \mathbf{c}_i) \right| \leq \max_{i \in [k]} |H(\mathbf{x}, \mathbf{c}_i)| \cdot \sum_{i=1}^k |\alpha_i| \leq k \cdot \max_{i \in [k]} |H(\mathbf{x}, \mathbf{c}_i)| \cdot \left(\sum_{i=1}^k |\alpha_i|^p \right)^{\frac{1}{p}}$$

■

Dual space: Now, we consider the dual space of ℓ_p^k , denoted as ℓ_p^{k*} . We would show some interesting properties of this dual space. First, note that by definition, the dual norm is

$$\|F\|_{\ell_p^{k*}} = \sup_{\|f\|_{\ell_p^k}=1} |F(f)| \quad (47)$$

We would show that ℓ_p^{k*} is isometrically isomorphic to ℓ_q^k where $\frac{1}{p} + \frac{1}{q} = 1$.

Lemma 29 Any bounded linear functional $F \in \ell_p^{k*}$ can be uniquely represented as

$$F(f_{\alpha}) := \sum_{i=1}^k \alpha_i \beta_i$$

for all $f_{\alpha} \in \ell_p^k$, where $\sum_{i=1}^k \beta_i \cdot H(\cdot, \mathbf{c}_i) \in \ell_q^k$. Moreover, any function $g_{\beta} \in \ell_q^k$ defines a linear functional F in ℓ_p^{k*} with dual norm

$$\|F\|_{\ell_p^{k*}} = \|g_{\beta}\|_{\ell_q^k} = \left(\sum_{i=1}^k |\beta_i|^q \right)^{\frac{1}{q}}$$

Proof For this result, we assume⁴ that $p > 1$. Consider an arbitrary function $f_{\alpha} \in \ell_p^k$. Note that we can write

$$F(f_{\alpha}) = \sum_{i=1}^k \alpha_i \cdot F(H(\cdot, \mathbf{c}_i)) \quad (\text{by linearity})$$

Define β_i as $F(H(\cdot, \mathbf{c}_i))$ for all $i \in [k]$. Now, we wish to show that $g_{\beta} := \sum_{i=1}^k \beta_i \cdot H(\mathbf{c}_i, \cdot) \in \ell_q^k$, i.e. $\|g_{\beta}\|_{\ell_q^k}$ is bounded.

Consider a choice of α'_i such that

$$\alpha'_i = |\beta_i|^{\frac{q}{p}} \operatorname{sgn}(\beta_i)$$

4. Case where $p = 1$ can be similarly handled.

Thus, the norm of $f_{\alpha'}$ is

$$\|f_{\alpha'}\|_{\ell_p^k} = \left(\sum_{i=1}^k |\beta_i|^q \right)^{\frac{1}{p}}.$$

On the other hand,

$$|\mathbf{F}(f_{\alpha'})| = \sum_{i=1}^k |\beta_i|^{\frac{q}{p}+1}$$

But by definition

$$\begin{aligned} |\mathbf{F}(f_{\alpha'})| &\leq \|\mathbf{F}\|_{\ell_p^{k*}} \cdot \|f_{\alpha'}\|_{\ell_p^k} \\ \implies \sum_{i=1}^k |\beta_i|^{\frac{q}{p}+1} &\leq \|\mathbf{F}\|_{\ell_p^{k*}} \cdot \left(\sum_{i=1}^k |\beta_i|^q \right)^{\frac{1}{p}} \\ \implies \sum_{i=1}^k |\beta_i|^q &\leq \|\mathbf{F}\|_{\ell_p^{k*}} \cdot \left(\sum_{i=1}^k |\beta_i|^q \right)^{\frac{1}{p}} \quad \left(\text{since } \frac{q}{p} + 1 = q \right) \\ \implies \left(\sum_{i=1}^k |\beta_i|^q \right)^{\frac{1}{q}} &\leq \|\mathbf{F}\|_{\ell_p^{k*}} \end{aligned} \tag{48}$$

Since \mathbf{F} is bounded the inequality above implies that $g_{\beta} \in \ell_q^k$. Now, assuming $f_{\alpha''} \in \ell_p^k$ and $g_{\beta} \in \ell_q^k$ such that $\|f_{\alpha''}\|_{\ell_p^k} = 1$, using Hölder's inequality we note that

$$|\mathbf{F}(f_{\alpha''})| = \left| \sum_{i=1}^k \alpha_i'' \beta_i \right| \leq \|f_{\alpha''}\|_{\ell_p^k} \cdot \|g_{\beta}\|_{\ell_q^k} \tag{49}$$

Plugging the bound of Eq. (49) in Eq. (47), we get

$$\|\mathbf{F}\|_{\ell_p^{k*}} \leq \|g_{\beta}\|_{\ell_q^k}$$

But Eq. (48) implies that $\|\mathbf{F}\|_{\ell_p^{k*}} = \|g_{\beta}\|_{\ell_q^k}$.

Showing that any element in ℓ_q^k defines a linear functional is straightforward where we consider a linear functional \mathbf{F} that maps $H(\mathbf{c}_i, \cdot)$ for all $i \in [k]$ to the corresponding parameters. This completes the proof. \blacksquare

Thus, we can treat $(\ell_q^k, \|\cdot\|_{\ell_q^k})$ as the dual of $(\ell_p^k, \|\cdot\|_{\ell_p^k})$ where any element in ℓ_q^k is *identified* as a linear transformation in ℓ_p^{k*} .

Reproducing kernel of ℓ_p^k : In Lemma 28, we showed that $(\ell_p^k, \|\cdot\|_{\ell_p^k})$ has the reproducing property. A question remains if we could find an explicit form for a kernel K for the constructed RKBS. In the following discussion, we show explicit kernels for the primal and dual spaces $(\ell_q^k, \|\cdot\|_{\ell_q^k})$, and $(\ell_p^k, \|\cdot\|_{\ell_p^k})$ respectively, which are both RKBSs. Depending on how we treat ℓ_q^k (dual or adjoint) we can show different kernels for the pair. As shown in Theorem 1 (Zhang et al., 2009), if ℓ_q^k is treated

as a dual then there is a unique kernel, otherwise as an adjoint one can show multiple kernels. In the following, we will show the explicit constructions.

Kernel for adjoint Banach spaces: Here, we study ℓ_p^k and ℓ_q^k as Banach spaces *adjoint* to each other as per Definition 3. We study a choice of bilinear map and show that how H turns out to be the kernel in this setting. Consider the following bilinear form $\langle \cdot, \cdot \rangle_{\mathcal{B}}$ on $\ell_p^k \times \ell_q^k$ defined as

$$\langle f_{\alpha}, g_{\beta} \rangle_{\mathcal{B}} := \sum_{i,j} \alpha_i \beta_j K(\mathbf{c}_j, \mathbf{c}_i) \quad (50)$$

It is easy to check its continuity. Furthermore,

$$\begin{aligned} \langle f_{\alpha}, H(\mathbf{x}, \cdot) \rangle_{\mathcal{B}} &= \sum_{i=1}^k \alpha_i H(\mathbf{x}, \mathbf{c}_i) = f_{\alpha}(\mathbf{x}), \quad \mathbf{x} \in \mathcal{X} \\ \langle H(\cdot, \mathbf{y}), g_{\beta} \rangle_{\mathcal{B}} &= \sum_{i=1}^k \beta_i H(\mathbf{c}_i, \mathbf{y}) = \sum_{i=1}^k \beta_i H(\mathbf{c}_i, \mathbf{y}) = g_{\beta}(\mathbf{y}), \quad \mathbf{y} \in \mathcal{X} \end{aligned}$$

Thus, H forms a reproducing kernel for ℓ_p^k and ℓ_q^k where we consider them as adjoint under a specific choice of bilinear form.

In the following, we study the case when a canonical bilinear form is induced by the dual space.

Canonical bilinear form and its unique kernel: Here, we consider the bilinear form induced by action of the dual space ℓ_p^{k*} on the primal space ℓ_p^k . We would use this canonical action to define a bilinear map on $\ell_p^k \times \ell_q^k$. Any function $g_{\beta} \in \ell_q^k$ is identified as an element in $\mathbb{F}_{\beta} \in \ell_p^{k*}$ as shown in Theorem 29. Thus, we write

$$\langle f_{\alpha}, g_{\beta} \rangle_{\ell_p^k \times \ell_q^k} = \langle f_{\alpha}, \mathbb{F}_{\beta} \rangle_{\ell_p^k \times \ell_p^{k*}} = \mathbb{F}_{\beta}(f_{\alpha}) = \sum_{i=1}^k \alpha_i \beta_i$$

Note, that H can't be the underlying kernel as

$$\langle f_{\alpha}, H(\mathbf{c}_j, \cdot) \rangle_{\ell_p^k \times \ell_q^k} = \alpha_j,$$

whereas the evaluation should be $\sum_{i=1}^k \alpha_i H(\mathbf{c}_j, \mathbf{c}_i)$. Now, consider a bivariate function $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ that is defined for $\mathbf{x} \in \mathcal{X}$ as

$$\begin{aligned} K(\cdot, \mathbf{x}) &:= \sum_{i=1}^k H(\mathbf{c}_i, \mathbf{x}) H(\cdot, \mathbf{c}_i) \\ K(\mathbf{x}, \cdot) &:= \sum_{i=1}^k H(\mathbf{x}, \mathbf{c}_i) H(\mathbf{c}_i, \cdot) \end{aligned}$$

It is easy to check that for all $\mathbf{x} \in \mathcal{X}$ we have $K(\cdot, \mathbf{x}) \in \ell_p^k$ and $K(\mathbf{x}, \cdot) \in \ell_q^k$. Furthermore, we note that

$$\begin{aligned} \langle f_{\alpha}, K(\mathbf{x}, \cdot) \rangle_{\ell_p^k \times \ell_q^k} &= \sum_{i=1}^k \alpha_i H(\mathbf{x}, \mathbf{c}_i) = f_{\alpha}(\mathbf{x}) \\ \langle K(\cdot, \mathbf{x}), g_{\beta} \rangle_{\ell_p^k \times \ell_q^k} &= \sum_{i=1}^k H(\mathbf{c}_i, \mathbf{x}) \beta_i = g_{\beta}(\mathbf{x}) \end{aligned}$$

Remark 30 *The construction of the RKBS as shown in Eq. (45) can be extended to an arbitrary countable set of centers. This has been studied in Xu and Ye (2019). They consider a locally compact Hausdorff space Ω with measure μ in \mathbb{R}^d and eigenfunctions $\phi_n \in L_0(\Omega)$ for all $n \in \mathbb{N}$ and show that the following space*

$$B_K^p(\Omega) := \left\{ f := \sum_{i \in \mathbb{N}} a_i \phi_i \mid (a_i \in \mathbb{N}) \in \ell_p \right\} \quad (51)$$

is an RKBS with the norm $\|\cdot\|_{B_K^p} = \|\cdot\|_{\ell_p}$, where $K \in L_0(\Omega \times \Omega)$ is a generalized Mercer kernel as defined in Xu and Ye (2019). Note that we don't assume any condition on the bivariate function H other than the linear independence conditions over the fixed centers.

Mirror map To instantiate the MDA, we consider a natural choice of mirror map. Consider the functional $\Phi_p : \ell_p^k \rightarrow \mathbb{R}$ defined as follows:

$$\Phi_p(f_\alpha) := \frac{1}{2} \cdot \|f_\alpha\|_{\ell_p^k}^2, \quad \forall f_\alpha \in \ell_p^k$$

Note that the Fréchet differential of Φ_p at f_α with increment $f_{\alpha'}$ is

$$\partial_{f_\alpha} \Phi_p(f_{\alpha'}) = \sum_{i=1}^k \alpha'_i \cdot \frac{\text{sgn}(\alpha_i) |\alpha_i|^{p-1}}{\|\alpha\|_p^{p-2}}$$

Since ℓ_p^{k*} is isometrically isomorphic to ℓ_q^k , using Theorem 29 we have the following correspondence

$$\partial_{f_\alpha} \Phi_p \equiv g_\beta = \sum_{i=1}^k \beta_i \cdot H(\mathbf{c}_i, \cdot), \quad \text{where } \beta = \frac{\text{sgn}(\alpha) |\alpha|^{p-1}}{\|\alpha\|_p^{p-2}} \quad (52)$$

Now, refer to the map

$$\partial_{(\cdot)} \Phi_p : \ell_p^k \rightarrow \ell_p^k, \quad f \mapsto \partial_f \Phi_p \quad (53)$$

As a direct consequence of Lemma 17, it is straight-forward to show that $\partial_{(\cdot)} \Phi_p$ is invertible. Furthermore, the inverse has some desired properties as well.

Corollary 31 *Fix any $p \in [1, \infty)$ and consider the Banach space ℓ_p^k as shown in Eq. (45). Then, the inverse of the Gâteaux derivative of the map $\frac{1}{2} \|\cdot\|_p^2$ on ℓ_p^k is the Gâteaux derivative of $\frac{1}{2} \|\cdot\|_q^2$ on ℓ_q^k , where $\frac{1}{p} + \frac{1}{q} = 1$.*

Proof *Note that $\frac{1}{2} \|\cdot\|_p^2$ is a proper strictly convex, Gâteaux differentiable map on ℓ_p^k . Since the Fenchel conjugate of $\frac{1}{2} \|\cdot\|_p^2$ is $\frac{1}{2} \|\cdot\|_q^2$ (cf Example 3.27 Boyd and Vandenberghe (2004)), the claim of the corollary follows immediately using Lemma 17. \blacksquare*

In Algorithm 1, we provide explicit updates for MDA of Eq. (16) for the optimization problem of Eq. (9) over ℓ_p^k for squared-error loss. The algorithm maps functions in the primal space ℓ_p^k to the dual space ℓ_q^k via a mirror map as shown in Eq. (52). Below, we provide the justification for the update steps in Algorithm 1.

First, we note that the linear combination of the inverse of a function g_β wrt $(\partial_{(\cdot)} \Phi_p)^{-1}$, say α , can be computed in terms of β .

Algorithm 1: Mirror descent on ℓ_p^k

Data: Given a training set $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$, model class ℓ_p^k , similarity kernel \widehat{H} , initial parameters α_0, β_0

Result: Parameters $\widehat{\alpha}$

- 1 Initialize $t = 0$;
 - while** $t \leq T$ **do**
 - 2 $\beta^{(t)} \leftarrow \beta^{(t-1)} - \eta \cdot (\widehat{H}^\top (\widehat{H}\alpha^{(t-1)} - Y))$;
 - 3 $\alpha^{(t)} \leftarrow \frac{\text{sgn}(\beta^{(t)})|\beta^{(t)}|^{q-1}}{\|\beta^{(t)}\|_q^{q-1}}$;
 - 4 $t++$;
 - 5 **return** $\widehat{\alpha} = \alpha_T$;
-

Lemma 32 We have $(\partial_{(\cdot)}\Phi_p)^{-1}g_\beta = f_\alpha$ where $\alpha_i = \text{sgn}(\beta_i) \frac{|\beta_i|^{q-1}}{\|\beta\|_q^{q-1}}$ for all $i \in [k]$ and $\frac{1}{p} + \frac{1}{q} = 1$.

Proof Using Corollary 31, it is clear that $(\partial_{(\cdot)}\Phi_p)^{-1} = \partial_{(\cdot)}\Phi_q$. Now, since $\Phi_q : \ell_q^k \rightarrow \mathbb{R}$, thus we note that

$$\begin{aligned} \partial_{(\cdot)}\Phi_q : \ell_q^k &\rightarrow \ell_p^k \\ g_\beta &\mapsto \sum_{i=1}^k \text{sgn}(\beta_i) \frac{|\beta_i|^{q-1}}{\|\beta\|_q^{q-1}} \cdot H(\cdot, \mathbf{c}_i) \end{aligned}$$

which is immediate using Eq. (52). ■

We consider the following notations to set up the algorithm.

Notations: Coefficients of a mirror descent update at iteration step t in MDA (see Eq. (16)) are denoted as $\beta^{(t)}$ for the function $g_{\beta^{(t)}}$ in ℓ_q^k (similarly $\alpha^{(t)}$ for $f_{\alpha^{(t)}} \in \ell_p^k$). α, β and $Y := (y_1, y_2, \dots, y_n)^T$ are treated as column vectors in \mathbb{R}^k . We define a similarity matrix \widehat{H} of dimension $(n \times k)$ where $\widehat{H}_{ij} := H(\mathbf{x}_i, \mathbf{c}_j)$. For the set of centers $C = \{\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_n\}$, we also use $\bar{\mathbf{c}}$ to denote a vector, e.g. $H(\cdot, \bar{\mathbf{c}})$ denotes a row vector $(H(\cdot, \mathbf{c}_1), H(\cdot, \mathbf{c}_2), \dots, H(\cdot, \mathbf{c}_k))$. Similarly, for a training set $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ we use $\bar{\mathbf{x}}$ to denote row vector $H(\bar{\mathbf{x}}, \cdot) := (H(\mathbf{x}_1, \cdot), H(\mathbf{x}_2, \cdot), \dots, H(\mathbf{x}_n, \cdot))$.

Using the notations above, for clarity we rewrite the update steps of Eq. (16) here for squared error loss:

$$g_{\beta^{(t)}} \leftarrow g_{\beta^{(t-1)}} - 2\eta \sum_{i=1}^n (f_{\alpha^{(t-1)}}(\mathbf{x}_i) - y_i) \cdot K(\mathbf{x}_i, \cdot) \quad (54a)$$

$$f_{\alpha^{(t)}} \leftarrow (\partial\Phi_p)^{-1}(g_{\beta^{(t)}}) \quad (54b)$$

Explicit updates for the first equation Line 2 of Algorithm 1 can be derived as follows:

$$\begin{aligned} \sum_{i=1}^k \beta_i^{(t)} \cdot H(\mathbf{c}_i, \cdot) &= \sum_{i=1}^k \beta_i^{(t-1)} \cdot H(\mathbf{c}_i, \cdot) - 2\eta \sum_{i=1}^n \left(\sum_{j=1}^k (f_{\alpha^{(t-1)}}(\mathbf{x}_i) - y_i) H(\mathbf{x}_i, \mathbf{c}_j) H(\mathbf{c}_j, \cdot) \right) \\ &= \sum_{i=1}^k \beta_i^{(t-1)} \cdot H(\mathbf{c}_i, \cdot) - 2\eta \sum_{i=1}^n \sum_{j=1}^k \left(\sum_{m=1}^k \alpha_m^{(t-1)} \cdot H(\mathbf{x}_i, \mathbf{c}_m) - y_i \right) \cdot H(\mathbf{x}_i, \mathbf{c}_j) H(\mathbf{c}_j, \cdot) \end{aligned}$$

$$\begin{aligned}
&= \sum_{i=1}^k \beta_i^{(t-1)} \cdot H(\mathbf{c}_i, \cdot) - 2\eta \sum_{j=1}^k \sum_{i=1}^n \left(\sum_{m=1}^k \alpha_m^{(t-1)} \cdot H(\mathbf{x}_i, \mathbf{c}_m) - y_i \right) \cdot H(\mathbf{x}_i, \mathbf{c}_j) H(\mathbf{c}_j, \cdot) \\
&= \sum_{i=1}^k \beta_i^{(t-1)} \cdot H(\mathbf{c}_i, \cdot) - 2\eta \sum_{j=1}^k \left(\sum_{i=1}^n \left(\sum_{m=1}^k \alpha_m^{(t-1)} \cdot H(\mathbf{x}_i, \mathbf{c}_m) - y_i \right) \cdot H(\mathbf{x}_i, \mathbf{c}_j) \right) H(\mathbf{c}_j, \cdot) \\
&= \sum_{i=1}^k \beta_i^{(t-1)} \cdot H(\mathbf{c}_i, \cdot) - 2\eta \sum_{j=1}^k \left(\sum_{i=1}^n \left(H(\mathbf{x}_i, \bar{\mathbf{c}}) \alpha^{(t-1)} - y_i \right) \cdot H(\mathbf{x}_i, \mathbf{c}_j) \right) H(\mathbf{c}_j, \cdot) \\
&= \sum_{i=1}^k \beta_i^{(t-1)} \cdot H(\mathbf{c}_i, \cdot) - 2\eta \sum_{j=1}^k \left(\sum_{i=1}^n H(\mathbf{x}_i, \bar{\mathbf{c}}) \alpha^{(t-1)} H(\mathbf{x}_i, \mathbf{c}_j) - H(\bar{\mathbf{x}}, \mathbf{c}_j)^T Y \right) H(\mathbf{c}_j, \cdot) \\
&= \sum_{i=1}^k \beta_i^{(t-1)} \cdot H(\mathbf{c}_i, \cdot) - 2\eta \sum_{j=1}^k \left(H(\bar{\mathbf{x}}, \mathbf{c}_j)^T H(\bar{\mathbf{x}}, \bar{\mathbf{c}}) \alpha^{(t-1)} - H(\bar{\mathbf{x}}, \mathbf{c}_j)^T Y \right) H(\mathbf{c}_j, \cdot) \\
&= \sum_{i=1}^k \beta_i^{(t-1)} \cdot H(\mathbf{c}_i, \cdot) - 2\eta \sum_{j=1}^k H(\bar{\mathbf{x}}, \mathbf{c}_j)^T \left(H(\bar{\mathbf{x}}, \bar{\mathbf{c}}) \alpha^{(t-1)} - Y \right) H(\mathbf{c}_j, \cdot) \\
&= \sum_{i=1}^k \left(\beta_i^{(t-1)} - 2\eta H(\bar{\mathbf{x}}, \mathbf{c}_i)^T \left(H(\bar{\mathbf{x}}, \bar{\mathbf{c}}) \alpha^{(t-1)} - Y \right) \right) H(\mathbf{c}_i, \cdot)
\end{aligned}$$

This implies that $\beta^{(t)} = \beta^{(t-1)} - 2\eta \widehat{H}^T \left(\widehat{H} \alpha^{(t-1)} - Y \right)$. Using Lemma 32, Line 3 follows:

$$\alpha^{(t)} \leftarrow \frac{\text{sgn}(\beta^{(t)}) |\beta^{(t)}|^{q-1}}{\|\beta^{(t)}\|_q^{q-1}}$$

Thus, the full procedure in Algorithm 1 can be written down without retaining evaluations $H(\mathbf{c}_i, \cdot)$'s for linear transformations in the dual space ℓ_q^k (similarly for functions in the primal space ℓ_p^k).

In Algorithm 1, the updates are based on a canonical view of kernel. In the case of adjoint based kernel, the updates for $\beta^{(t)}$ would be slightly different:

$$\beta^{(t)} \leftarrow \beta^{(t-1)} - \eta \cdot (\widehat{H} \alpha^{(t-1)} - Y)$$

Convergence of Algorithm 1: We utilize the square loss functional, which is L -Lipschitz with respect to the p -norm of the RKBS, assuming that α has a bounded norm. According to Kakade et al. (2012) (see Lemma 9), the strong convexity parameter of the mirror map $\frac{1}{2} \|\cdot\|_p^2$ is $(p-1)$ for $p \in (1, 2)$. For $p > 2$, the corresponding squared p -norm is not strongly convex. Consequently, the convergence of Algorithm 1 is theoretically guaranteed at the rate shown in Theorem 25 under p -norm RKBSs for $p \in (1, 2)$. We validate this convergence for various p values within this range for step functions, as demonstrated in Figure 5.

In the following, we discuss our numerical experiments.

Numerical Experiments: We present experimental results that validate our theoretical findings. The dynamics of convergence of Algorithm 1 for different values of p are plotted, for an optimization

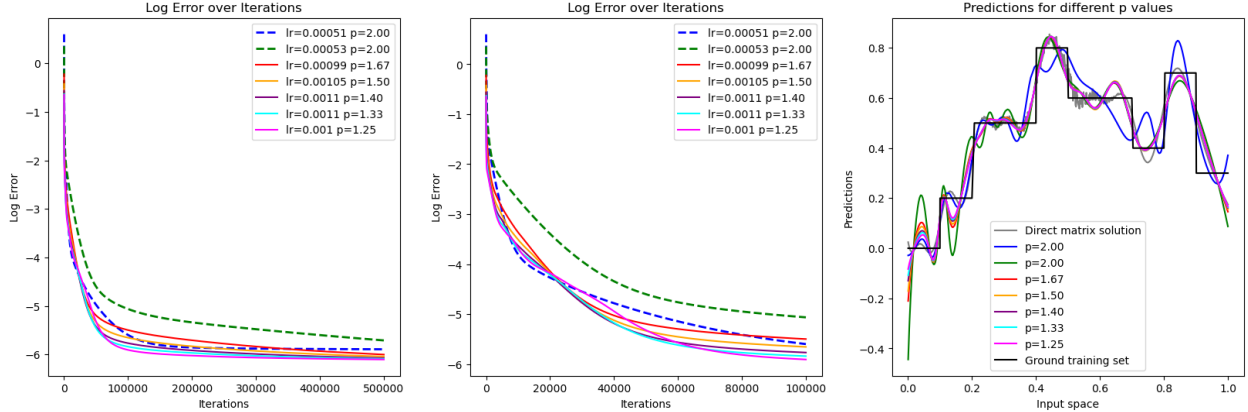


Figure 5: Results from numerical experiments using the mirror descent algorithm (Algorithm 1) for varying p values in a one-dimensional space. We apply squared error loss on 800 training points with 25 centers using the Locally Adaptive-Bandwidths (LAB) RBF kernel. The bivariate function H is defined by $\exp\left(-\frac{\|\theta_i \odot (x - \mathbf{c}_i)\|_2^2}{2}\right)$, with θ_i optimized via gradient descent. **(Leftmost Plot)**: Logarithm of training error versus iterations. **(Middle Plot)**: Zoomed-in log error versus iterations up to 100,000 steps. **(Rightmost Plot)**: Predictions of learned kernel classifiers compared with the direct matrix solution (gray curve) on the training set.

problem in a one-dimensional input space. Specifically, we visualize the impact of different p values on the approximation performance and training error.

To construct the space ℓ_p^k , we employ the Locally Adaptive-Bandwidths (LAB) RBF kernels as introduced in He et al. (2024). Using this kernel, the p -norm RKBS ℓ_p^k is constructed using a bivariate function H , defined as

$$H(x, \mathbf{c}_i) = \exp\left(-\frac{\|\theta_i \odot (x - \mathbf{c}_i)\|_2^2}{2}\right), \quad \forall \mathbf{c}_i \in C,$$

where \odot denotes the Hadamard (element-wise) product, $\|\cdot\|_2$ is the L^2 -norm, and $\theta_i \in \mathbb{R}$ is a center-dependent bandwidth vector. Note that H is asymmetric, which implies that the kernel K for the space ℓ_p^k is also asymmetric. The bandwidths θ_i are computed by performing 10 gradient steps on the loss function $L(\alpha, \boldsymbol{\theta})$, where $\boldsymbol{\theta} = [\theta_1, \theta_2, \dots, \theta_k]^\top$, according to

$$\boldsymbol{\theta}_t \leftarrow \boldsymbol{\theta}_{t-1} - \eta \partial_{\boldsymbol{\theta}_{t-1}} L(\alpha, \boldsymbol{\theta}),$$

where the learning rate η differs from the learning rate used in Algorithm 1.

Experimental Setup: The input space is fixed as $\mathcal{X} = \mathbb{R}$. Using the computed $\boldsymbol{\theta}$ from the gradient update, we run Algorithm 1 to learn a step function in one dimension over the interval $[0, 1]$, as depicted in Figure 5 (black curve in the rightmost plot). The loss functional used is the squared error. Our training set consists of 800 points, with 25 centers randomly selected from this set, yielding $k = 25$ for constructing the space ℓ_p^k . We explore different values of p in the range $(1, 2]$, specifically $\{2, 1.67, 1.5, 1.4, 1.33, 1.25\}$.

We present two types of plots: (1) the logarithm of the training error versus the number of iterations (see the leftmost and middle plots in Figure 5), and (2) the predictions of the learned kernel classifiers $\sum_{i=1}^{25} \alpha_i H(\cdot, \mathbf{c}_i)$ on the training set (see the rightmost plot in Figure 5).

The first plot illustrates the logarithm of the training error over 500,000 iterations for various p values and learning rates. To provide a clearer view of convergence trends, a zoomed-in plot of the log error versus iterations for up to 100,000 steps is also included. A general trend observed across the plots indicates that as p decreases, the approximation quality improves and the training error decreases. This effect is also visible in the predictions on the training points in the rightmost plot of Figure 5. This improvement can be attributed to the expansion of the space as p decreases, as $\|\alpha\|_p$ is an increasing function of p . For a fixed norm $D > 0$, the space ℓ_p^k grows larger with decreasing p , allowing for richer approximations within the space.

Comparison with Direct Matrix Solution: We also compare the predictions obtained using the solution from NumPy (depicted by the gray curve in the rightmost plot of Figure 5). In this comparison, the NumPy solution solves the system of linear equations $K_g \cdot \alpha = Y$, where K_g is the similarity matrix computed over the (training, centers) pairs using a Gaussian kernel, specifically

$$K_g(x_i, c_j) = \exp\left(-\frac{\|x_i - c_j\|_2^2}{\sigma^2}\right),$$

for training point and center pair (x_i, c_j) . Here, K_g is a similarity matrix of dimension 800×25 . The prediction plot uses the best bandwidth σ^2 for the Gaussian kernel, chosen to yield optimal performance across different trials.

Acknowledgments and Disclosure of Funding

We acknowledge support from the National Science Foundation (NSF) and the Simons Foundation for the Collaboration on the Theoretical Foundations of Deep Learning through awards DMS-2031883 and #814639, the TILOS institute (NSF CCF-2112665), and the Office of Naval Research (N8644-NV-ONR). This work was started when P.P. was a postdoctoral fellow at the Halicioğlu Data Science Insititute at UC San Diego.

Appendix A. Technical Proofs

A.1 Mirror maps: Proof of Lemma 17

First, we state a useful lemma on the first-order optimality condition for optimization in functional analysis.

Lemma 33 ((Luenberger, 1997, Chapter 8, Lemma 1)) *Let F be a Fréchet differentiable convex functional on a real normed space \mathcal{B} . Let \mathcal{C}' be a convex cone in \mathcal{B} . A necessary and sufficient condition that $f_0 \in \mathcal{P}$ minimizes F over \mathcal{C}' is that*

$$\begin{aligned}\partial_{f_0} F(f) &\geq 0 \quad \forall f \in \mathcal{C}', \\ \partial_{f_0} F(f_0) &= 0.\end{aligned}$$

We are interested in the case where the convex cone of interest is the space \mathcal{B} itself. The lemma above is stated for any real normed space \mathcal{B} . Given that we are interested in Banach spaces, we could state the following more precise statement on the minimum of the functional F .

Proposition 34 ((Luenberger, 1997, Section 7.4, Theorem 1)) *Let the real-valued functional F have a Gâteaux differential on a vector space \mathcal{B} . A necessary condition for F to have an extremum at $f_0 \in \mathcal{B}$ is that $\partial_{f_0} F(h) = 0$ for all $h \in \mathcal{B}$.*

Proof of Lemma 17: It is straightforward to see that the strict convexity and Gâteaux differentiability of F imply its injectivity. However, we present an alternative proof that introduces notations used to demonstrate surjectivity. To establish the injectivity of the operator $\partial_{(\cdot)} F$, we consider the convex conjugate of F , denoted F^* . Note that $F^* : \mathcal{B}^* \rightarrow \mathbb{R}$, and

$$\forall g \in \mathcal{B}^*, F^*(g) := - \inf_{f \in \mathcal{B}} (-\langle f, g \rangle_{\mathcal{B}} + F(f)) \quad (55)$$

We denote $\hat{F}_g(f) := -\langle f, g \rangle_{\mathcal{B}} + F(f)$. Note that since F is strictly convex, thus \hat{F}_g is strictly convex in f . Thus, the objective $\hat{F}_g(f)$ can only have a *unique* minimizer.

Assume that there exists $\hat{f} \in \mathcal{B}$ such that $g := \partial_{\hat{f}} F$. Using the first-order optimality condition in Lemma 33, \hat{f} is the minimizer of \hat{F}_g , but then it is the unique minimizer. This implies that there doesn't exist $f' \neq \hat{f}$ such that $\partial_{f'} F = \partial_{\hat{f}} F$. This completes the proof of the injectivity of the operator $\partial_{(\cdot)} F$.

Now, we would establish the surjectivity of the operator $\partial_{(\cdot)} F$ over \mathcal{B}^* . Using Theorem 2.3.3 Zălinescu (2002), $F^{**} = F$, i.e.

$$\forall f' \in \mathcal{B}, F(f') := - \inf_{\hat{g} \in \mathcal{B}^*} (-\langle f', \hat{g} \rangle_{\mathcal{B}} + F^*(\hat{g})) \quad (56)$$

Consider a linear transformation $g^* \in \mathcal{B}^*$. Assume that \hat{F}_{g^*} is minimized at $f_0 \in \mathcal{B}$. Thus,

$$F^*(g^*) = \langle f_0, g^* \rangle_{\mathcal{B}} - F(f_0) \quad (57)$$

Using **Young-Fenchel inequality** (Theorem 2.3.1 [Zălinescu \(2002\)](#)), we note that for all $f \in \mathcal{B}$,

$$F(f) + F^*(g^*) \geq \langle f, g^* \rangle_{\mathcal{B}} \quad (58)$$

Subtracting Eq. (57) and Eq. (58), we observe that for all $f \in \mathcal{B}$,

$$F(f) - F(f_0) \geq \langle f - f_0, g^* \rangle_{\mathcal{B}}$$

This implies that $g^* \in \partial_{f_0} F$, i.e g^* is in the subdifferential set of $\partial_{f_0} F$. Since F is Gâteaux differentiable at f_0 it must be the case that $g^* = \partial_{f_0} F$.

Given that g^* was picked arbitrarily, we have shown that for any $g^* \in \mathcal{B}^*$ there exists $f_0 \in \mathcal{B}$ such that $g^* = \partial_{f_0} F$. This completes the proof of the surjectivity of the operator $\partial_{(\cdot)} F$ over \mathcal{B}^* . ■

Using Eq. (57) and Eq. (58), we note that the surjectivity of the mirror map Φ does not necessarily require it is Gâteaux differentiable.

A.2 Existence of a smooth and strongly convex functional: Proof of Lemma 19

In this Appendix, we show that for a functional $F : \mathcal{C} \rightarrow \mathbb{R}$ over a Banach space to be both μ -strongly convex and γ -smooth for $\mu > 0$ and $\gamma < \infty$, it must be isomorphic to a Hilbert space.

The key idea of the proof is to consider the points of second differentiability of a functional F . Essentially, we wish to look at functions $f \in \mathcal{B}$ where a Taylor approximation is possible. If we show that there exists such a function then one can achieve equivalence of the Banach space $\|\cdot\|_{\mathcal{B}}$ to a Hilbert space norm.

We show that indeed such a point exists where a Taylor expansion to the second order is possible if the functional F is continuous and convex. Consider the following lemma:

Lemma 35 *Consider \mathcal{B} be a separable Banach space, and let F be a continuous convex function on \mathcal{B} . Then, there exists an $f \in \mathcal{B}$, $g^* \in \partial_f F$, and a bounded symmetric linear operator $M : \mathcal{B} \rightarrow \mathcal{B}^*$ such that for all $f' \in \mathcal{B}$ and small scalar $\lambda > 0$, we can expand F as follows*

$$F(f + \lambda f') = F(f) + \lambda \langle f', g^* \rangle_{\mathcal{B}} + \frac{\lambda^2}{2} \langle f', M(f') \rangle_{\mathcal{B}} + o(\lambda^2) \quad (\lambda \rightarrow 0). \quad (59)$$

Proof Using Theorem 4.1 [Borwein and Noll \(1994\)](#), it is straight-forward that the set of functions $f \in \mathcal{B}$ where Eq. (59) holds is non-empty implying the statement of the lemma. ■

Now, we state a useful result on the generalization of Parallelogram law, which provides a characterization for showing a Banach space is isomorphic to a Hilbert space. Let $\sum_{\epsilon(n)}$ denote all possible sequences $(\epsilon_1, \dots, \epsilon_n)$ of $\{\pm 1\}$'s.

Proposition 36 (Proposition 3.1, [Kwapień \(1972\)](#)) *A real or complex Banach space \mathcal{B} is isomorphic to a Hilbert space if and only if there exists a constant $C > 0$ such that*

$$C^{-1} \sum_{i=1}^n \|f_i\|_{\mathcal{B}}^2 \leq \frac{1}{2^n} \sum_{\epsilon(n)} \left(\left\| \sum_{i=1}^n \epsilon_i f_i \right\| \right)^2 \leq C \sum_{i=1}^n \|f_i\|_{\mathcal{B}}^2,$$

for any positive integer $n > 1$ and any f_1, f_2, \dots, f_n in \mathcal{B} .

With this we proof the claim of Lemma 19 in the following:

Proof [of Lemma 19] First, assume that \mathcal{B} is separable. The proof of the non-separable case extends from the separable case.

Assume for the sake of contradiction that there exists $F : \mathcal{B} \rightarrow \mathbb{R}$ that is both μ -strongly convex and γ -smooth for some $\mu > 0$ and $\gamma < \infty$. Using Lemma 35, there exists f and $g^* \in \partial_f F$ such that for all $f' \in \mathcal{B}$,

$$F(f + \lambda f') = F(f) + \lambda \langle f', g^* \rangle + \frac{\lambda^2}{2} \langle f', M(f') \rangle + o(\lambda^2) \quad (\lambda \rightarrow 0) \quad (60)$$

for a symmetric operator M . Since F is γ -smooth, we have

$$F(f' + \lambda f') \leq F(f') + \langle \lambda f', g^* \rangle + \frac{\gamma \lambda^2}{2} \|f'\|_{\mathcal{B}}^2 \quad (61)$$

Using μ -strong convexity of F we get

$$F(f' + \lambda f') \geq F(f') + \langle \lambda f', g^* \rangle + \frac{\mu \lambda^2}{2} \|f'\|_{\mathcal{B}}^2 \quad (62)$$

Combining Eq. (60)-(62), we get

$$C_1 \|f'\|_{\mathcal{B}}^2 \leq \langle f', M(f') \rangle_{\mathcal{B}} \leq C_2 \|f'\|_{\mathcal{B}}^2 \quad (63)$$

for some $0 < C_1 \leq C_2$. Now, consider the following inner product $\langle \cdot, \cdot \rangle_H : \mathcal{B} \times \mathcal{B} \rightarrow \mathbb{R}$ on the space \mathcal{B}

$$\forall h_1, h_2 \in \mathcal{B}, \quad \langle h_1, h_2 \rangle_H = \langle h_1, M(h_2) \rangle_{\mathcal{B}}$$

Symmetry of $\langle \cdot, \cdot \rangle_H$ follows by the symmetry of the operator M . Using Eq. (63) we have

$$C_1 \|f'\|_{\mathcal{B}}^2 \leq \langle f', f' \rangle_H \leq C_2 \|f'\|_{\mathcal{B}}^2$$

But alternately we can also write

$$\frac{1}{C_2} \langle f', f' \rangle_H \leq \|f'\|_{\mathcal{B}}^2 \leq \frac{1}{C_1} \langle f', f' \rangle_H \quad (64)$$

Now, we show the condition for Proposition 36 is true for the norm $\|\cdot\|_{\mathcal{B}}$, that yields the isomorphism of $(\mathcal{B}, \|\cdot\|_{\mathcal{B}})$ to $(\mathcal{B}, \|\cdot\|_H)$.

Note that for any n choice of functions $f_1, f_2, \dots, f_n \in \mathcal{B}$

$$\sum_{\epsilon(n)} \left(\left\| \sum_{i=1}^n \epsilon_i f_i \right\|_H^2 \right) = 2^n \sum_i \|f_i\|_H^2$$

But we have

$$\begin{aligned} C_2 \sum_{\epsilon(n)} \left(\left\| \sum_{i=1}^n \epsilon_i f_i \right\|_{\mathcal{B}}^2 \right) &\geq \sum_{\epsilon(n)} \left(\left\| \sum_{i=1}^n \epsilon_i f_i \right\|_H^2 \right) \\ \implies C_2 \sum_{\epsilon(n)} \left(\left\| \sum_{i=1}^n \epsilon_i f_i \right\|_{\mathcal{B}}^2 \right) &\geq C_1 2^n \sum_i \|f_i\|_{\mathcal{B}}^2 \end{aligned}$$

Similarly,

$$\begin{aligned} C_1 \sum_{\epsilon(n)} \left(\left\| \sum_{i=1}^n \epsilon_i f_i \right\|_{\mathcal{B}}^2 \right) &\leq \sum_{\epsilon(n)} \left(\left\| \sum_{i=1}^n \epsilon_i f_i \right\|_H^2 \right) \\ \implies C_1 \sum_{\epsilon(n)} \left(\left\| \sum_{i=1}^n \epsilon_i f_i \right\|_{\mathcal{B}}^2 \right) &\leq C_2 2^n \sum_i \|f_i\|_{\mathcal{B}}^2 \end{aligned}$$

Combining the equations, we have

$$C^{-1} \sum_i \|f_i\|_{\mathcal{B}}^2 \leq \frac{1}{2^n} \sum_{\epsilon(n)} \left(\left\| \sum_{i=1}^n \epsilon_i f_i \right\|_{\mathcal{B}}^2 \right) \leq C \sum_i \|f_i\|_{\mathcal{B}}^2$$

where $C = \frac{C_2}{C_1}$. Thus, we have shown that $(\mathcal{B}, \|\cdot\|_{\mathcal{B}})$ is isomorphic to $(\mathcal{B}, \|\cdot\|_H)$ using Proposition 36.

Now, consider the non-separable case. Using the proof above every separable (open or closed) subspace $\mathcal{S} \subset \mathcal{B}$ is isomorphic to a Hilbert space, which is sufficient to show that the space $(\mathcal{B}, \|\cdot\|_{\mathcal{B}})$ is isomorphic to a Hilbert space using Proposition 36.

For a separable subspace $\mathcal{S} \subset \mathcal{B}$, denote by $C_{\mathcal{S}}$ the infimum of the constants that satisfy Proposition 36. Now, define:

$$C_{\text{sup}} := \sup \{C_{\mathcal{S}} \mid \mathcal{S} \subset \mathcal{B}, \text{ closed, separable}\}$$

Note that if $C_{\text{sup}} < +\infty$, then it clearly satisfies the inequalities above for the whole space \mathcal{B} . Assume on the contrary that $C_{\text{sup}} = +\infty$, so for all $n \in \mathbb{N}$ there exists a separable subspace \mathcal{S}_n such that $C_{\mathcal{S}_n} \geq n$. Now, define by \mathcal{S} the closure of the span of countably many subsets \mathcal{S}_n , i.e., $\mathcal{S} = \overline{\text{span}(\bigcup_n \mathcal{S}_n)}$.

Since closure of countable union separable space is separable we note that \mathcal{S} is separable and $\mathcal{S}_n \subset \mathcal{S}$ for all n . Thus, $C_{\mathcal{S}_n} \leq C_{\mathcal{S}}$, so $C_{\mathcal{S}} \geq n$ for all $n \in \mathbb{N}$. But this gives a contradiction because $C_{\mathcal{S}}$ is finite due to the construction. Thus, $C_{\text{sup}} < +\infty$, and hence the conditions for Proposition 36 are satisfied, which implies that $(\mathcal{B}, \|\cdot\|_{\mathcal{B}})$ is isomorphic to a Hilbert space. \blacksquare

A.3 Proof of Lemma 24

Proof First, note that for any $f' \in \mathcal{C}$, the map $h \mapsto \mathfrak{D}_{\Phi}(h, f')$ is convex. Secondly, note that for any f, f' we can write

$$\partial_f \mathfrak{D}_{\Phi}(\cdot, f') \in \partial_f \Phi - \partial_{f'} \Phi \tag{65}$$

By definition, $\Pi_{\mathcal{B}_0}^{\Phi}(f')$ is the minimizer of $\mathfrak{D}_{\Phi}(f, f')$ in the first component over $\mathcal{B}_0 \cap \mathcal{C}$. This means that for any $g \in \mathcal{B}_0 \cap \mathcal{C}$,

$$\mathfrak{D}_{\Phi}(\Pi_{\mathcal{B}_0}^{\Phi}(f'), f') \leq \mathfrak{D}_{\Phi}(g, f')$$

Using the first-order condition of convexity, we get

$$\left\langle g - \Pi_{\mathcal{B}_0}^\Phi(f'), \partial_{\Pi_{\mathcal{B}_0}^\Phi(f')} \mathfrak{D}_\Phi(\cdot, f') \right\rangle_{\mathcal{B}} \geq 0, \quad \forall g \in \mathcal{B}_0 \cap \mathcal{C}$$

Substituting (65) into the above inequality, we have

$$\left\langle g - \Pi_{\mathcal{B}_0}^\Phi(f'), \partial_{\Pi_{\mathcal{B}_0}^\Phi(f')} \Phi - \partial_{f'} \Phi \right\rangle_{\mathcal{B}} \geq 0$$

Setting $g = f$ in the above inequality, we get

$$\left\langle f - \Pi_{\mathcal{B}_0}^\Phi(f'), \partial_{\Pi_{\mathcal{B}_0}^\Phi(f')} \Phi - \partial_{f'} \Phi \right\rangle_{\mathcal{B}} \geq 0$$

This completes the proof. ■

References

- Andreas Argyriou, Charles A. Micchelli, and Massimiliano Pontil. On spectral learning. *Journal of Machine Learning Research*, 11(31):935–953, 2010. URL <http://jmlr.org/papers/v11/argyriou10a.html>.
- Navid Azizan, Sahin Lale, and Babak Hassibi. Stochastic mirror descent on overparameterized nonlinear models. *IEEE Transactions on Neural Networks and Learning Systems*, 33(12):7717–7727, 2021.
- Francis R. Bach. Breaking the curse of dimensionality with convex neural networks. *ArXiv*, abs/1412.8690, 2014.
- Francesca Bartolucci, Ernesto De Vito, Lorenzo Rosasco, and Stefano Vigogna. Understanding neural networks with reproducing kernel banach spaces. *Applied and Computational Harmonic Analysis*, 62:194–236, 2023. ISSN 1063-5203. doi: <https://doi.org/10.1016/j.acha.2022.08.006>. URL <https://www.sciencedirect.com/science/article/pii/S1063520322000768>.
- Heinz H. Bauschke, Jérôme Bolte, and Marc Teboulle. A descent lemma beyond lipschitz gradient continuity: First-order methods revisited and applications. *Math. Oper. Res.*, 42:330–348, 2017.
- Amir Beck and Marc Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 31(3):167–175, 2003.
- Aharon Ben-Tal, Tamar Margalit, and Arkadi Nemirovski. The ordered subsets mirror descent optimization method with applications to tomography. *SIAM J. Optim.*, 12:79–108, 2001.
- K. P. Bennett and E. J. Bredehsteiner. Duality and geometry in svm classifiers. In *Proceedings of the Seventeenth International Conference on Machine Learning*, pages 57–64, San Francisco, 2000. Morgan Kaufmann.
- Jonathan Michael Borwein and Dominikus Noll. Second order differentiability of convex functions in banach spaces. *Transactions of the American Mathematical Society*, 342:43–81, 1994. URL <https://api.semanticscholar.org/CorpusID:51851560>.
- Léon Bottou and Olivier Bousquet. The tradeoffs of large scale learning. In J. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems*, volume 20. Curran Associates, Inc., 2007. URL https://proceedings.neurips.cc/paper_files/paper/2007/file/0d3180d672e08b4c5312dcdafdf6ef36-Paper.pdf.
- Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, Cambridge, 2004.
- Sébastien Bubeck. *Convex optimization: Algorithms and complexity*, 2015.
- T. Buehler and D.A. Salamon. *Functional Analysis*, volume 191. AMS Graduate Studies in Mathematics, 2018.
- Yair Censor and Stavros A. Zenios. Proximal minimization algorithm withd-functions. *Journal of Optimization Theory and Applications*, 73:451–464, 1992.

- Nicolò Cesa-Bianchi, Pierre Gaillard, Gábor Lugosi, and Gilles Stoltz. Mirror descent meets fixed share (and feels no regret). In *Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 1*, NIPS'12, page 980–988, Red Hook, NY, USA, 2012. Curran Associates Inc.
- Yingyi Chen, Qinghua Tao, Francesco Tonin, and Johan A. K. Suykens. Primal-attention: Self-attention through asymmetric kernel svd in primal representation. *ArXiv*, abs/2305.19798, 2023.
- Seok-Young Chung and Qiyu Sun. Barron space for graph convolution neural networks, 2023.
- Ricky Der and Daniel Lee. Large-margin classification in banach spaces. In Marina Meila and Xiaotong Shen, editors, *Proceedings of the Eleventh International Conference on Artificial Intelligence and Statistics*, volume 2 of *Proceedings of Machine Learning Research*, pages 91–98, San Juan, Puerto Rico, 21–24 Mar 2007. PMLR. URL <https://proceedings.mlr.press/v2/der07a.html>.
- Weinan E, Chao Ma, and Lei Wu. A priori estimates of the population risk for two-layer neural networks. *Communications in Mathematical Sciences*, 2018.
- Jonathan Eckstein. Nonlinear proximal point algorithms using bregman functions, with applications to convex programming. *Math. Oper. Res.*, 18(1):202–226, feb 1993. ISSN 0364-765X.
- Gregory E. Fasshauer, Fred J. Hickernell, and Qi Ye. Solving support vector machines in reproducing kernel banach spaces with positive definite functions. *Applied and Computational Harmonic Analysis*, 38(1):115–139, 2015. ISSN 1063-5203. doi: <https://doi.org/10.1016/j.acha.2014.03.007>. URL <https://www.sciencedirect.com/science/article/pii/S1063520314000475>.
- Kenji Fukumizu, Gert Lanckriet, and Bharath K. Sriperumbudur. Learning in hilbert vs. banach spaces: A measure embedding viewpoint. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc., 2011. URL https://proceedings.neurips.cc/paper_files/paper/2011/file/7b13b2203029ed80337f27127a9f1d28-Paper.pdf.
- Rémi Gribonval, Gitta Kutyniok, Morten Nielsen, and Felix Voigtländer. Approximation spaces of deep neural networks. *Constructive Approximation*, 55:259–367, 2019.
- Fan He, Mingzhen He, Lei Shi, Xiaolin Huang, and Johan A. K. Suykens. Learning analysis of kernel ridgeless regression with asymmetric kernel learning, 2024. URL <https://arxiv.org/abs/2406.01435>.
- Sham M. Kakade, Shai Shalev-Shwartz, and Ambuj Tewari. Regularization techniques for learning with matrices. *Journal of Machine Learning Research*, 13(59):1865–1890, 2012. URL <http://jmlr.org/papers/v13/kakade12a.html>.
- Hamed Karimi, Julie Nutini, and Mark Schmidt. Linear convergence of gradient and proximal-gradient methods under the polyak-łojasiewicz condition. In Paolo Frasconi, Niels Landwehr, Giuseppe Manco, and Jilles Vreeken, editors, *Machine Learning and Knowledge Discovery in Databases*, pages 795–811, Cham, 2016. Springer International Publishing. ISBN 978-3-319-46128-1.
- S. Kwapien. Isomorphic characterizations of inner product spaces by orthogonal series with vector valued coefficients. *Studia Mathematica*, 44(6):583–595, 1972. URL <http://eudml.org/doc/217719>.

- J. Lee, Max Simchowitz, Michael I. Jordan, and Benjamin Recht. Gradient descent only converges to minimizers. In *Annual Conference Computational Learning Theory*, 2016.
- Rong Rong Lin, Hai Zhang Zhang, and Jun Zhang. On reproducing kernel banach spaces: Generic definitions and unified framework of constructions. *Acta Mathematica Sinica, English Series*, 38(8):1459–1483, August 2022. ISSN 1439-7617. doi: 10.1007/s10114-022-1397-7. URL <https://doi.org/10.1007/s10114-022-1397-7>.
- David G. Luenberger. *Optimization by Vector Space Methods*. John Wiley & Sons, Inc., USA, 1st edition, 1997. ISBN 047155359X.
- Charles A. Micchelli and Massimiliano Pontil. A function representation for learning in banach spaces. In John Shawe-Taylor and Yoram Singer, editors, *Learning Theory*, pages 255–269, Berlin, Heidelberg, 2004. Springer Berlin Heidelberg. ISBN 978-3-540-27819-1.
- Charles A. Micchelli and Massimiliano Pontil. Feature space perspectives for learning the kernel. *Mach. Learn.*, 66(2–3):297–319, mar 2007. ISSN 0885-6125. doi: 10.1007/s10994-006-0679-0. URL <https://doi.org/10.1007/s10994-006-0679-0>.
- Arkadii Nemirovski and David Borisovich Yudin. *Problem Complexity and Method Efficiency in Optimization*. 1983.
- Yurii Nesterov. *Introductory Lectures on Convex Optimization: A Basic Course*. Springer Publishing Company, Incorporated, 1 edition, 2014. ISBN 1461346916.
- Yurii Nesterov. *Lectures on Convex Optimization*. Springer Publishing Company, Incorporated, 2nd edition, 2018. ISBN 3319915770.
- Greg Ongie, Rebecca Willett, Daniel Soudry, and Nathan Srebro. A function space view of bounded norm infinite width relu nets: The multivariate case. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=H1lNPxHKDH>.
- Rahul Parhi and Robert D. Nowak. The role of neural network activation functions. *IEEE Signal Processing Letters*, 27:1779–1783, 2019.
- Rahul Parhi and Robert D. Nowak. Banach space representer theorems for neural networks and ridge splines. *J. Mach. Learn. Res.*, 22:43:1–43:40, 2020.
- Rahul Parhi and Robert D. Nowak. What kinds of functions do deep neural networks learn? insights from variational spline theory. *SIAM J. Math. Data Sci.*, 4:464–489, 2021.
- Rahul Parhi and Michael Unser. Function-space optimality of neural architectures with multivariate nonlinearities, 2023.
- B.T. Polyak. Gradient methods for the minimisation of functionals. *USSR Computational Mathematics and Mathematical Physics*, 3(4):864–878, 1963. ISSN 0041-5553. doi: [https://doi.org/10.1016/0041-5553\(63\)90382-3](https://doi.org/10.1016/0041-5553(63)90382-3). URL <https://www.sciencedirect.com/science/article/pii/S0041555363903823>.
- Ming qian He, Fan He, Fanghui Liu, and Xiaolin Huang. Random fourier features for asymmetric kernels. *ArXiv*, abs/2209.08461, 2022.

- Pedro H. P. Savarese, Itay Evron, Daniel Soudry, and Nathan Srebro. How do infinite width bounded norm networks look in function space? In *Annual Conference Computational Learning Theory*, 2019.
- Bernhard Scholkopf and Alex Smola. Learning with kernels: support vector machines, regularization, optimization, and beyond. In *Adaptive computation and machine learning series*, 2001. URL <https://api.semanticscholar.org/CorpusID:52872213>.
- Lei Shi, Yun-Long Feng, and Ding-Xuan Zhou. Concentration estimates for learning with ℓ_1 -regularizer and data dependent hypothesis spaces. *Applied and Computational Harmonic Analysis*, 31(2):286–302, 2011. ISSN 1063-5203. doi: <https://doi.org/10.1016/j.acha.2011.01.001>. URL <https://www.sciencedirect.com/science/article/pii/S1063520311000157>.
- Alistair Shilton, Sunil Gupta, Santu Rana, and Svetha Venkatesh. Gradient descent in neural networks as sequential learning in reproducing kernel banach space. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 31435–31488. PMLR, 23–29 Jul 2023a. URL <https://proceedings.mlr.press/v202/shilton23a.html>.
- Alistair Shilton, Sunil Gupta, Santu Rana, and Svetha Venkatesh. Gradient descent in neural networks as sequential learning in reproducing kernel banach space. In *Proceedings of the 40th International Conference on Machine Learning*, ICML’23. JMLR.org, 2023b.
- Guohui Song, Haizhang Zhang, and Fred J. Hickernell. Reproducing kernel banach spaces with the ℓ^1 norm. *Applied and Computational Harmonic Analysis*, 34(1):96–116, 2013. ISSN 1063-5203. doi: <https://doi.org/10.1016/j.acha.2012.03.009>. URL <https://www.sciencedirect.com/science/article/pii/S1063520312000486>.
- Len Spek, Tjeerd Jan Heeringa, and Christoph Brune. Duality for neural networks through reproducing kernel banach spaces. *ArXiv*, abs/2211.05020, 2022.
- Anand Srinivasan and Jean-Jacques Slotine. Contracting dynamical systems in banach spaces, 2022.
- Haoyuan Sun, Kwangjun Ahn, Christos Thrampoulidis, and Navid Azizan. Mirror descent maximizes generalized margin and can be implemented efficiently. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 31089–31101. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/c9694bf4f9bf3626f7d21158bab74f8e-Paper-Conference.pdf.
- Haoyuan Sun, Khashayar Gatmiry, Kwangjun Ahn, and Navid Azizan. A unified approach to controlling implicit regularization via mirror descent. *ArXiv*, abs/2306.13853, 2023.
- Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, 58(1):267–288, 1996. ISSN 00359246. URL <http://www.jstor.org/stable/2346178>.
- Hongzhi Tong, Di-Rong Chen, and Fenghong Yang. Least square regression with lp-coefficient regularization. *Neural Computation*, 22(12):3221–3235, 2010. doi: 10.1162/NECO_a_00044.

- Michael Unser, Julien Fageot, and Harshit Gupta. Representer theorems for sparsity-promoting ℓ_1 regularization. *IEEE Transactions on Information Theory*, 62(9):5167–5180, 2016. doi: 10.1109/TIT.2016.2590421.
- Daniel Wachsmuth and Gerd Wachsmuth. A simple proof of the baillon-haddad theorem on open subsets of hilbert spaces, 2022. URL <https://arxiv.org/abs/2204.00282>.
- Rui Wang, Yuesheng Xu, and Mingsong Yan. Sparse representer theorems for learning in reproducing kernel banach spaces. *ArXiv*, abs/2305.12584, 2023a.
- Rui Wang, Yuesheng Xu, and Mingsong Yan. Sparse representer theorems for learning in reproducing kernel banach spaces, 2023b.
- Matthew A. Wright and Joseph Gonzalez. Transformers are deep infinite-dimensional non-mercere binary kernel machines. *ArXiv*, abs/2106.01506, 2021.
- Yuesheng Xu. Sparse machine learning in banach spaces. *Applied numerical mathematics : transactions of IMACS*, 187:138–157, 2023.
- Yuesheng Xu and Qi Ye. *Generalized Mercer kernels and reproducing kernel Banach spaces*, volume 258. American Mathematical Society, 2019.
- Qi Ye. Support vector machines in reproducing kernel hilbert spaces versus banach spaces. In Gregory E. Fasshauer and Larry L. Schumaker, editors, *Approximation Theory XIV: San Antonio 2013*, pages 377–395, Cham, 2014. Springer International Publishing. ISBN 978-3-319-06404-8.
- Haizhang Zhang and Jun Zhang. Frames, riesz bases, and sampling expansions in banach spaces via semi-inner products. *Applied and Computational Harmonic Analysis*, 31(1):1–25, 2011. ISSN 1063-5203. doi: <https://doi.org/10.1016/j.acha.2010.09.007>. URL <https://www.sciencedirect.com/science/article/pii/S1063520310001120>.
- Haizhang Zhang and Jun Zhang. Regularized learning in banach spaces as an optimization problem: representer theorems. *Journal of Global Optimization*, 54:235–250, 2012.
- Haizhang Zhang, Yuesheng Xu, and Jun Zhang. Reproducing kernel banach spaces for machine learning. *Journal of Machine Learning Research*, 10(95):2741–2775, 2009. URL <http://jmlr.org/papers/v10/zhang09b.html>.
- J. Zhang and H. Zhang. Categorization based on similarity and features: The reproducing kernel banach space approach. In W. Batchelder, H. Colonius, E. N. Dzhafarov, and J. Myung, editors, *New Handbook of Mathematical Psychology, Volume 2*. Springer, 2018.
- Constantin Zălinescu. On uniformly convex functions. *Journal of Mathematical Analysis and Applications*, 95:344–374, 1983.
- Constantin Zălinescu. Convex analysis in general vector spaces. 2002.
- S. Łojasiewicz. A topological property of real analytic subsets. *Coll. du CNRS, Les equations aux derivees partielles*, pages 87–89, 1963.