

# A General Framework for Consistent Structured Prediction with Implicit Loss Embeddings

**Carlo Ciliberto**

*Department of Electrical and Electronic Engineering,  
Imperial College  
London, UK*

C.CILIBERTO@IMPERIAL.AC.UK

**Lorenzo Rosasco**

*University of Genova, Italy and  
Istituto Italiano di Tecnologia, Genova, Italy  
Massachusetts Institute of Technology, Cambridge, MA, USA*

LROSASCO@MIT.EDU

**Alessandro Rudi**

*INRIA, Paris, France,  
École Normale Supérieure, Paris, France  
PSL Research, France*

ALESSANDRO.RUDI@INRIA.FR

**Editor:** Sebastian Nowozin

## Abstract

We propose and analyze a novel theoretical and algorithmic framework for structured prediction. While so far the term has referred to discrete output spaces, here we consider more general settings, such as manifolds or spaces of probability measures. We define structured prediction as a problem where the output space lacks a vectorial structure. We identify and study a large class of loss functions that implicitly defines a suitable geometry on the problem. The latter is the key to develop an algorithmic framework amenable to a sharp statistical analysis and yielding efficient computations. When dealing with output spaces with infinite cardinality, a suitable implicit formulation of the estimator is shown to be crucial.

**Keywords:** Structured Prediction, Statistical Learning Theory, Kernel Methods.

## 1. Introduction

Statistical learning theory offers a number of methods to deal with supervised problems when the output space is linear (e.g. scalar values or vectors). However, applications involving more general output spaces are becoming increasingly common. Examples include image segmentation (Alahari et al., 2008) or captioning (Karpathy and Fei-Fei, 2015), speech recognition (Bahl et al., 1986; Sutton et al., 2012), manifold regression (Steinke et al., 2010), trajectory planning (Ratliff et al., 2006), protein folding (Joachims et al., 2009), prediction of probability distributions (Frognier et al., 2015), ordinal regression (Pedregosa et al., 2017), information retrieval or ranking (Duchi et al., 2010) to name a few (see Bakir et al., 2007; Nowozin et al., 2011, for more examples). When considering discrete output spaces, these settings are often referred to as *structured prediction* problems, since they require dealing with output spaces that have a specific structure, such as strings, graphs or sequences.

Standard machine learning methods like empirical risk minimization are faced with both modeling and computational challenges in these settings. Therefore, in practice, either one of the following two main strategies are often considered. On the one hand, *surrogate methods* (Bartlett et al., 2006; Mroueh et al., 2012) that design *ad-hoc* algorithms and theory for different learning settings on a case-by-case basis. While this allows to prove strong theoretical guarantees, it makes it difficult to extend previous results to new settings. On the other hand methods such as SVM-struct Tsochantaridis et al. (2005) or Max-Margin Markov Networks Taskar et al. (2004), to which we refer here as *auxiliary function maximization* methods, have broad applicability (Lafferty et al., 2001; Bakir et al., 2007; Nowozin et al., 2011) but typically poor theoretical guarantees (Tewari and Bartlett, 2007).

In this work, we propose a novel structured prediction framework combining the best of both worlds. Our approach extends structured prediction beyond discrete outputs, to include problems such as manifold regression. The lack of linear structure in the output space is the common feature of the different problems we consider. The key observation is that for a very wide range of problems, the associated loss function carries implicitly a natural corresponding geometry. More precisely it admits an *Implicit Loss Embedding (ILE)* into a linear (albeit possibly infinite dimensional) space. Exploiting such a geometry allows us to derive a consistent least squares algorithmic framework. The latter is related to the *Output Kernel Regression (OKR)* framework (Geurts et al., 2006, 2007; Brouard et al., 2011; Kadri et al., 2013; Brouard et al., 2016), which encompasses several methods, including *Kernel Dependency Estimation (KDE)* (Weston et al., 2002; Cortes et al., 2005). OKR approaches lift the output set into a latent Hilbert space via vector-valued kernel embeddings and coincide to learning with a specific loss function corresponding to the distance induced by a positive definite kernel. In this work we consider loss functions that admit a bi-linear representation on a latent Hilbert space, covering a wide range of loss functions used in practice, such as any loss on discrete spaces, divergences, geodesic distances on Riemannian manifolds or smooth functions in  $\mathbb{R}^d$ . Like KDE and OKR for the case of kernel losses, we derive an implicit formulation of the proposed estimator depending only on the loss function. This is crucial when considering structured prediction problems beyond the discrete case.

This paper is the extended version of Ciliberto et al. (2016), which has initiated a recent stream of works where this method has been equivalently referred to as either the *Structured Encoding Loss Function (SELF)* approach or the *Quadratic Surrogate* framework (Osokin et al., 2017; Ciliberto et al., 2017; Korba et al., 2018; Rudi et al., 2018; Struminsky et al., 2018; Luise et al., 2018; Nowak-Vila et al., 2018; Djerrab et al., 2018; Luise et al., 2019; Blondel, 2019; Ciliberto et al., 2019). In this paper we refine the analysis in Ciliberto et al. (2016), providing novel insights and estimators for structured prediction. More precisely, the novel contributions of the current work are: (a) we propose a number of novel estimators for structure prediction (Sec. 3). (b) We study the generalization properties of the proposed estimators, proving consistency and learning rates (Sec. 5). (c) We show that the learning rates and computational costs of the proposed methods are adaptive to the regularity of the learning problem (Sec. 5.4), potentially reducing the complexity of the learning process. (d) We provide a number of sufficient conditions to determine whether a loss admits an ILE, which are easy to verify in practice (Sec. 6). We leverage the latter results to show that most loss functions used in machine learning applications satisfy the ILE definition and therefore that our framework is suited to a large number of settings and applications.

The paper is organized as follows: Sec. 2 introduces structured prediction within the framework of statistical learning theory for supervised problems. In Sec. 3 we present the ILE framework and the novel estimators. In Sec. 4 we draw extensive connections with previous work. Sec. 5 is devoted to study the statistical and computational properties of the proposed estimators. Sec. 6 studies sufficient conditions for a loss function to satisfy the ILE definition. Finally, Sec. 7 concludes the work discussing relevant future directions.

## 2. Problem Setting and Background

We denote by  $\mathcal{X}, \mathcal{Y}$  and  $\mathcal{Z}$  respectively the *input space*, *label space* and *output space* of a learning problem. We let  $\rho$  be a probability measure on  $\mathcal{X} \times \mathcal{Y}$  and  $\Delta : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$  be a loss function measuring prediction errors between a label  $y \in \mathcal{Y}$  and an output  $z \in \mathcal{Z}$ . The distinction between  $\mathcal{Y}$  and  $\mathcal{Z}$  allows to consider applications where labels do not necessarily correspond to the desired outputs (e.g. ranking/information retrieval, see below).

**Supervised Learning.** In supervised learning problems, the goal is to estimate a function  $f^* : \mathcal{X} \rightarrow \mathcal{Z}$  minimizing the *expected risk*

$$\min_{f: \mathcal{X} \rightarrow \mathcal{Z}} \mathcal{E}(f), \quad \text{with} \quad \mathcal{E}(f) = \int \Delta(f(x), y) d\rho(x, y), \quad (1)$$

over the set of measurable functions  $f : \mathcal{X} \rightarrow \mathcal{Z}$ . In practice, the distribution  $\rho$  is given but unknown and only  $n$  examples  $(x_i, y_i)_{i=1}^n$  independently distributed according to  $\rho$  are provided. Given a training set, a learning algorithm needs to find a good approximation  $f_n : \mathcal{X} \rightarrow \mathcal{Z}$  to  $f^*$  such that the corresponding excess risk  $\mathcal{E}(f_n)$  is close to  $\mathcal{E}(f^*)$  and tends to it as the number  $n$  of training points increases.

**Empirical Risk Minimization.** A standard learning approach is *Empirical Risk Minimization (ERM)* (see e.g. Devroye et al., 2013). This method consists in obtaining the estimator  $f_n : \mathcal{X} \rightarrow \mathcal{Z}$  as

$$f_n = \operatorname{argmin}_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \Delta(f(x_i), y_i), \quad (2)$$

by minimizing the empirical version of the expected risk over a suitable space  $\mathcal{H}$  of functions  $f : \mathcal{X} \rightarrow \mathcal{Z}$ . The idea underlying ERM is to use the empirical risk as an approximation of the expected risk, so that the estimation of  $f^*$  via  $f_n$  should become increasingly accurate as the number of training samples grows.

From a statistical perspective, it is sufficient for the loss  $\Delta$  to satisfy very general conditions (e.g. Lipschitz, bounded, etc.) in order for the ERM strategy to enjoy strong statistical guarantees. In particular, a number of results are available proving universal consistency and learning rates (in terms of generalization or excess risk bounds) for the empirical risk estimator  $f_n$ , and a variety of hypotheses spaces  $\mathcal{H}$  (see for instance Shalev-Shwartz and Ben-David, 2014, and references therein).

From a computational perspective, a central question is whether the ERM problem can be solved efficiently. When the output space is linear, such as  $\mathcal{Z} = \mathbb{R}$ , and the loss  $\Delta$  is convex, ERM becomes an efficient strategy for a large family of hypotheses spaces. For

instance, a standard approach is to consider linear parametrization of functions  $f : \mathcal{X} \rightarrow \mathbb{R}$  in a Hilbert space  $\mathcal{H}$  of the form  $f(x) = \langle w, \phi(x) \rangle_{\mathcal{H}}$  with  $w \in \mathcal{H}$  and  $\phi : \mathcal{X} \rightarrow \mathcal{H}$  a feature map. Following this approach, the resulting ERM problem in (2) can be solved efficiently leveraging convex optimization techniques. This same strategy can be naturally extended to the general linear case  $\mathcal{Z} = \mathbb{R}^M$ .

**Structured Prediction and Limitations of ERM.** When the space  $\mathcal{Z}$  does not have a linear structure, applying ERM poses concrete challenges to both modeling and computations:

- **Modeling.** If the output space is non-linear, it is not clear how to design a suitable hypotheses space  $\mathcal{H}$  of candidate estimators. In particular, linear parametrizations of the form  $f(x) = \langle w, \phi(x) \rangle$  introduced above are not possible. For instance, given  $f_1, f_2 : \mathcal{X} \rightarrow \mathcal{Z}$ , it is not guaranteed that  $f_1 + f_2$  takes values in  $\mathcal{Z}$  as well.
- **Computations.** If the hypotheses space is non-linear or the loss is non-convex (e.g. integer-valued), solving ERM can be extremely challenging. Often, approaches based on the regularity of the loss function and the optimization domain, such as gradient methods, cannot be adopted. Hence, it is not clear how to obtain  $f_n$  in practice.

Next we describe a number of problems falling in the above setting.

## 2.1. Examples of Structured Prediction Problems

Below we provide some examples of structured prediction problems according to our definition, that is problems where the output space lacks a linear structure. We refer to (Bakir et al., 2007; Nowozin et al., 2011) for more examples.

- **Classification, Multi-class, Multi-labeling.** In these settings  $\mathcal{Z} = \mathcal{Y} = \{1, \dots, T\}$  is a collection of classes that can be associated to inputs from  $\mathcal{X}$ .
- **Ranking, Ordinal Regression, Information Retrieval.** The goal is to predict an *ordered* list of documents (Bakir et al., 2007; Pedregosa et al., 2017; Duchi et al., 2010). For instance  $x \in \mathcal{X}$  can be a query in a search engine and  $\mathcal{Z}$  is the space of all permutations (ordering) over the documents in a database. The label space  $\mathcal{Y} \neq \mathcal{Z}$  typically contains a set of scalar scores representing the individual relevance of each document to the input query.
- **Sequence Prediction.** The goal is to predict sequences such as time series for financial applications or planning trajectories (Ratliff et al., 2006). Loss functions such as the Dynamic Time Warping (Cuturi and Blondel, 2017) can be used to measure the similarity between two sequences.
- **Predicting Probability Distributions / Histograms.** In these settings, the output  $\mathcal{Z}$  corresponds to a space of probability distributions (Frogner et al., 2015; Luise et al., 2018; Mensch et al., 2019). The loss  $\Delta$  is a metric comparing probabilities, such as the Kullback-Libler divergence or the Hellinger,  $\chi^2$  or Wasserstein distance.

- **Manifold Regression.** Problems where the outputs belong to a smooth Riemannian manifold  $\mathcal{Z}$  (Steinke et al., 2010; Rudi et al., 2018). A natural choice for the loss  $\Delta$  is the squared geodesic distance of the manifold. This setting generalizes the standard regression problem with least-squares loss and  $\mathcal{Z} = \mathbb{R}^M$ , to the manifold scenario.

### 3. A General Framework for Structured Prediction

In this section, we introduce and motivate our structured prediction framework. Our discussion starts from a useful characterization of  $f^* : \mathcal{X} \rightarrow \mathcal{Z}$  the *minimizer of the expected risk* in (1). Assume that  $\rho$  can be factorized as  $\rho(x, y) = \rho(y|x)\rho_{\mathcal{X}}(x)$  with  $\rho(\cdot|x)$  the conditional distribution over  $\mathcal{Y}$  given  $x \in \mathcal{X}$  and  $\rho_{\mathcal{X}}$  the marginal distribution of  $\rho$  over  $\mathcal{X}$ . It can be shown that

$$f^*(x) = \operatorname{argmin}_{z \in \mathcal{Z}} \int_{\mathcal{Y}} \Delta(z, y) d\rho(y|x), \quad (3)$$

that is, the value  $f^*$  at any given  $x \in \mathcal{X}$  corresponds to the minimizer, over the output set  $\mathcal{Z}$ , of the conditional expectation  $\mathbb{E}_{y|x} \Delta(z, y)$ . Indeed, it is possible to show that if  $\Delta$  is measurable, then such a point-wise estimate defines a measurable function. This latter result requires some care and follows from Berge maximum theorem (Aliprantis and Border, 2006) (see also Aumann’s principle (Steinwart and Christmann, 2008)). We refer the reader to Appendix A for a detailed discussion.

In the following, we leverage this characterization of  $f^*$  to develop our approach to structured prediction. First, we consider the case where both output and label spaces are discrete and finite. As noted, this is relevant, since most previous work focused on this setting (Bakir et al., 2007). Additionally, for our presentation, it allows a self-contained introduction of key ideas, deferring the technical details to the general discussion in Sec. 3.2.

#### 3.1. Motivating Analysis: Finite Output Spaces

We begin by discussing how loss functions define a geometry on finite output spaces, and how it can be used to define an estimator. Then, we can consider linear estimators and show how for this class of estimators a useful implicit formulation can be derived.

**Geometry of Structured Loss Functions.** Let  $\mathcal{X} = \mathbb{R}^d$  and assume  $\mathcal{Z} = \mathcal{Y} = \{1, \dots, T\}$  for  $T \in \mathbb{N}$ . In this setting, any loss function  $\Delta : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$  can be characterized in terms of a matrix  $V \in \mathbb{R}^{T \times T}$  such that

$$\Delta(z, y) = e_z^\top V e_y, \quad \forall z \in \mathcal{Z}, y \in \mathcal{Y}, \quad (4)$$

where  $e_y \in \mathbb{R}^T$  denotes the  $y$ -th element of the canonical basis of  $\mathbb{R}^T$ , namely the vector with  $y$ -th entry equal to 1 and the rest equal to 0. Combining this observation with the characterization of  $f^*$  in (3) and using linearity of the integral, we have

$$f^*(x) = \operatorname{argmin}_{z \in \mathcal{Z}} e_z^\top V g^*(x), \quad g^*(x) = \int_{\mathcal{Y}} e_y d\rho(y|x), \quad (5)$$

for any  $x \in \mathcal{X}$ . In particular, note that the function  $g^* : \mathcal{X} \rightarrow \mathbb{R}^T$  is given by

$$g^*(x) = (\rho(1|x), \dots, \rho(T|x))^\top \in \mathbb{R}^T, \quad (6)$$

the vector whose  $y$ -th entry is equal to the probability of observing  $y$ , given  $x$ . This observation is crucial, since it allows to identify  $g^*$  with the *regression function*, that is the minimizer of the expected least squares error (see Thm. 3 for a formal statement)

$$g^* = \operatorname{argmin}_{g: \mathcal{X} \rightarrow \mathbb{R}^T} \int \|e_y - g(x)\|_{\mathbb{R}^T}^2 d\rho(x, y).$$

The above discussion suggests the following approach. Given given  $n$  training points  $(x_i, y_i)_{i=1}^n$ , we could approximate  $g^*$  by a least squares estimate  $g_n$  minimizing

$$g_n = \operatorname{argmin}_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \|e_{y_i} - g(x_i)\|_{\mathbb{R}^T}^2$$

over some function space  $\mathcal{G}$ . Then we obtain the estimator  $f_n : \mathcal{X} \rightarrow \mathcal{Z}$  for any  $x \in \mathcal{X}$  as

$$f_n(x) = \operatorname{argmin}_{z \in \mathcal{Z}} e_z^\top V g_n(x). \quad (7)$$

The advantage of this strategy is that approximating  $g^*$  corresponds to a standard vector-valued regression problem, since the output space is now  $\mathbb{R}^T$  and not the “structured”  $\mathcal{Z}$ . As we discuss next, for linearly parameterized estimators, we can develop a useful implicit formulation. We first add two remarks pointing out connections to related ideas.

**Remark 1 (Conditional mean embedding)** *By construction,  $g^*(x)$  defined in (5) corresponds to the definition of conditional mean embedding of  $\rho(\cdot|x)$  in  $\mathbb{R}^T$  (Song et al., 2009; Grünewälder et al., 2012; Singh et al., 2019). In Sec. 4.4, this connection will provide relevant insights on the structured prediction estimator we propose and its statistical properties.*

**Example 1 (Classification)** *For classification, the estimator  $f_n$  in (7), recovers the least-squares classifier in (Yao et al., 2007; Mroueh et al., 2012). To see this, let  $\mathcal{Z} = \mathcal{Y} = \{1, \dots, T\}$  be a finite set of class labels and let  $\Delta$  be the 0-1 (or mis-classification) loss, namely  $\Delta(z, y) = 1$  if  $z \neq y$  and 0 otherwise. It is easy to see that  $\Delta$  is of the form of (4) with matrix  $V = \mathbb{1}\mathbb{1}^\top - I$ , where  $I$  is the  $T \times T$  identity matrix and  $\mathbb{1}$  is the  $T$ -dimensional vector with all entries equal to 1. For any  $y \in \mathcal{Y}$ , the  $y$ -th entry of  $g_n(x) \in \mathbb{R}^T$  is interpreted as the likelihood of observing a the class  $y$  given the input  $x \in \mathcal{X}$ . Therefore, the classifier  $c_n : \mathcal{X} \rightarrow \mathcal{Y}$  acts by predicting the index of  $g_n(x)$  with higher likelihood. Direct comparison with our approach leads to*

$$c_n(x) = \operatorname{argmax}_{t=1, \dots, T} (g_n(x))_t, \quad f_n(x) = \operatorname{argmin}_{t=1, \dots, T} (V g_n(x))_t. \quad (8)$$

Since  $V = \mathbb{1}\mathbb{1}^\top - I$ , it is straightforward to see that the two methods coincide, namely  $c_n(x) = f_n(x)$  for all  $x \in \mathcal{X}$ .

We next describe a useful representation for linear estimators.

**Implicit formulation for linear estimators.** A possible approach to learn  $g_n$  is by linear ridge regression, namely

$$g_n(x) = W_n x, \quad W_n = \operatorname{argmin}_{W \in \mathbb{R}^{T \times d}} \frac{1}{n} \sum_{i=1}^n \|e_{y_i} - W x_i\|^2 + \lambda \|W\|_F^2, \quad (9)$$

where  $\lambda > 0$  is a hyperparameter and  $\|\cdot\|_F^2$  denotes the squared Frobenius norm of a matrix (sum of the squared entries). The solution of (9) can be obtained in closed form as

$$W_n = Y^\top X (X^\top X + n\lambda I)^{-1} \quad (10)$$

with  $I \in \mathbb{R}^{d \times d}$  the identity matrix and  $X \in \mathbb{R}^{n \times d}$  and  $Y \in \mathbb{R}^{n \times T}$  the matrices with  $i$ -th row corresponding to  $x_i$  and  $e_{y_i}$  respectively. Plugging this solution in (7) leads to an explicit approach to obtain the estimator  $f_n$ . We next discuss a useful observation that will be key to extend our discussion to  $\mathcal{Y}$  and  $\mathcal{Z}$  that are neither finite nor discrete. Specifically, we will show that it is possible to obtain a characterization for  $f_n$  that is equivalent to (7) but in which  $g_n$  *does not appear explicitly*. To see this, first notice that for any  $x \in \mathcal{X}$  we can leverage the closed-form solution for the estimator  $g_n$  to have

$$g_n(x) = W_n x = Y^\top \alpha(x) = \sum_{i=1}^n \alpha_i(x) e_{y_i}, \quad (11)$$

where the weights  $\alpha : \mathcal{X} \rightarrow \mathbb{R}^n$  are such that

$$\alpha(x) = (\alpha_1(x), \dots, \alpha_n(x))^\top = [X(X^\top X + n\lambda I)^{-1}] x \in \mathbb{R}^n. \quad (12)$$

We now plug this characterization of  $g_n$  in the definition of  $f_n$  in (7). Thanks to the linearity of the sum and the matrix-vector product, we have

$$f_n(x) = \operatorname{argmin}_{z \in \mathcal{Z}} \sum_{i=1}^n \alpha_i(x) e_z^\top V e_{y_i} = \operatorname{argmin}_{z \in \mathcal{Z}} \sum_{i=1}^n \alpha_i(x) \Delta(z, y_i), \quad (13)$$

where the last equality follows from the fact that the loss  $\Delta$  is identified by the matrix  $V$  according to (4). Intuitively, for any  $i = 1, \dots, n$ , we can interpret each  $\alpha_i(x)$  as a relevance score encouraging the candidate prediction  $z \in \mathcal{Z}$  to be “similar” to the observed training label  $y_i$ , *according to*  $\Delta(z, y_i)$ .

The key observation in (13) is that the estimator  $f_n$  can be characterized exclusively in terms of the weights  $\alpha$  and the observed labels  $y_i$ . This implies that we do not have to learn the vector-valued estimator  $g_n$  explicitly. In Sec. 3.2 we will leverage this observation to extend the same learning strategy to the case where  $\mathcal{Y}$  and  $\mathcal{Z}$  are not finite or discrete.

**Extension to generic  $\mathcal{X}$ .** We conclude this section by observing that the construction of  $f_n$  can be naturally extended to the case where  $\mathcal{X}$  is a generic set. In particular, consider  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  a positive definite kernel (Aronszajn, 1950). Then, according to standard practice from the kernel methods literature (see e.g. Shawe-Taylor and Cristianini, 2004) we can derive a “dual” formulation for the relevance scores  $\alpha$ . In particular, given the input  $(x_i)_{i=1}^n$  in training, (12) can be equivalently rewritten as

$$\alpha(x) = (\alpha_1(x), \dots, \alpha_n(x))^\top = (K + n\lambda I)^{-1} v(x) \in \mathbb{R}^n, \quad (14)$$

for any  $x \in \mathcal{X}$ , where  $K \in \mathbb{R}^{n \times n}$  is the empirical kernel matrix with entries  $K_{ij} = k(x_i, x_j)$  and  $v(x) \in \mathbb{R}^n$  is the evaluation vector, with entries  $v(x)_i = k(x, x_i)$ , for any  $i, j = 1, \dots, n$ . In the following, we will denote  $\kappa^2 = \sup_{x \in \mathcal{X}} k(x, x)$ . We will always assume  $\kappa < +\infty$  (e.g. by using a normalized kernel or by requiring  $\mathcal{X}$  to be a compact set). It is easy to see that this strategy corresponds indeed to learn the estimator  $g_n$  by solving the empirical risk minimization problem in (9) over the reproducing kernel Hilbert space (RKHS) associated to  $k$  (Aronszajn, 1950). We discuss this in more detail in the following.

### 3.2. General Structured Prediction: Beyond Finite Output Spaces

In this section, we generalize the discussion of Sec. 3.1 to structured prediction problems where  $\mathcal{Y}$  or  $\mathcal{Z}$  are not necessarily finite (or discrete). Also in this case, we show how a relevant geometry can be defined by a corresponding loss function. Further we extend the analysis of linearly parameterized estimators, and show how in this general setting the implicit formulation becomes essential.

**Implicit Loss Embeddings.** The extension to non finite output spaces hinges on a key assumption on the loss that generalizes the observation of (4). We refer to functions satisfying this condition as admitting an *Implicit Loss Embedding*.

**Definition 2 (ILE)** *A continuous map  $\Delta : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$  is said to admit an Implicit Loss Embedding (ILE) if there exists a separable Hilbert space  $\mathcal{H}$  and two measurable bounded maps  $\psi : \mathcal{Z} \rightarrow \mathcal{H}$  and  $\varphi : \mathcal{Y} \rightarrow \mathcal{H}$ , such that for any  $z \in \mathcal{Z}$  and  $y \in \mathcal{Y}$  we have*

$$\Delta(z, y) = \langle \psi(z), \varphi(y) \rangle_{\mathcal{H}}, \tag{15}$$

and  $\|\varphi(y)\|_{\mathcal{H}} \leq 1$ . Additionally, we define  $c_{\Delta} = \sup_{z \in \mathcal{Z}} \|\psi(z)\|_{\mathcal{H}}$ .

The definition of ILE is similar to the characterization of positive definite kernels in terms of feature maps (and indeed it recovers such definition when  $\mathcal{Z} = \mathcal{Y}$  and  $\psi = \varphi$ ), but is significantly more general in that it allows also to consider functions that are not positive definite (for example, distances) or even not symmetric (such as divergences). It is clear that any loss function on finite sets  $\mathcal{Z}$  and  $\mathcal{Y}$  admits an ILE. For instance, in the setting of Sec. 3.1 it is sufficient to choose  $\mathcal{H} = \mathbb{R}^T$ , with maps  $\varphi(y) = e_y$  and  $\psi(z) = V^{\top} e_z$ , to recover the ILE definition. Note that the requirement  $\sup_{y \in \mathcal{Y}} \|\varphi(y)\| \leq 1$  is introduced to simplify the notation but does not limit the generality of the assumption (see Thm. C.2 in the Appendix for more details). We note that in Ciliberto et al. (2016), a variant of the ILE property was introduced (see Asm 1 in such paper). In this work we opted for Thm. 2 since it allows for a more clear notation in the following. However, in Thm. C.1 in the Appendix we provide more details on this point and show that the two notions are actually equivalent.

While the definition of ILE is abstract, it is satisfied by many loss functions often used in structured prediction applications and more generally in machine learning. In Sec. 6 we present a wide range of sufficient conditions to guarantee a function to admit an ILE, which are easier to interpret and verify in practice.

Under the assumption that  $\Delta$  admits an ILE, we can easily retrace the reasoning in Sec. 3.1 to derive the structured prediction estimator. In particular, we have the following result to which we refer to as *Fisher consistency*, a term borrowed from the literature of surrogate methods (see discussion in Sec. 4.1).

**Lemma 3 (Fisher Consistency)** *Let  $\mathcal{Z}$  be compact,  $\Delta : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$  admit an ILE and let  $f^* : \mathcal{X} \rightarrow \mathcal{Z}$  be the solution of (1). Then,*

$$f^*(x) = \operatorname{argmin}_{z \in \mathcal{Z}} \langle \psi(z), g^*(x) \rangle_{\mathcal{H}}, \quad g^*(x) = \int_{\mathcal{Y}} \varphi(y) d\rho(y|x) \tag{16}$$

almost surely with respect to  $\rho_{\mathcal{X}}$ . Moreover,  $g^* : \mathcal{X} \rightarrow \mathcal{H}$  is the minimizer of

$$\mathcal{R}(g) = \int_{\mathcal{Y} \times \mathcal{X}} \|\varphi(y) - g(x)\|_{\mathcal{H}}^2 d\rho(x, y). \tag{17}$$



The result above provides a characterization of  $f^*$  in terms of the conditional expectation of  $\varphi(y)$  with respect to  $x \in \mathcal{X}$ . It generalizes (5) to the case where  $\mathcal{Y}$  and  $\mathcal{Z}$  are not finite. Analogously to (3), the proof of Thm. 3 is reported in Appendix A and leverages Berge’s Maximum theorem. In particular, the compactness of  $\mathcal{Z}$  is a technical requirement to guarantee  $f^*(x)$  to be well-defined. Analogously to the derivation in Sec. 3.1, the result of Thm. 3 motivates us to design a structured prediction estimator by first learning a  $g_n$  to approximate  $g^*$  via least squares over a space  $\mathcal{G}$  of functions  $g : \mathcal{X} \rightarrow \mathcal{H}$

$$g_n = \operatorname{argmin}_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \|\varphi(y_i) - g(x_i)\|_{\mathcal{H}}^2,$$

and then plug  $g_n$  in (16) to obtain an approximation  $f_n$  of  $f^*$  characterized by

$$f_n(x) = \operatorname{argmin}_{z \in \mathcal{Z}} \langle \psi(z), g_n(x) \rangle_{\mathcal{H}},$$

for all  $x \in \mathcal{X}$ . Learning  $g_n$  corresponds to solving a vector-valued regression problem on a (possibly) infinite dimensional output space  $\mathcal{H}$  (Caponnetto and De Vito, 2007). We next discuss in detail the case of linearly parameterized estimators.

**Linearly parameterized estimators and implicit formulation.** For simplicity, instead of directly minimizing the empirical square loss over  $\mathcal{H}$ , we consider again the ridge regression estimator  $g_n : \mathcal{X} \rightarrow \mathcal{H}$  defined as the minimizer of the regularized empirical risk

$$g_n = \operatorname{argmin}_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \|\varphi(y_i) - g(x_i)\|_{\mathcal{H}}^2 + \lambda \|g\|_{\mathcal{G}}^2, \quad (18)$$

over a normed space  $\mathcal{G}$  of vector-valued functions from  $\mathcal{X}$  to  $\mathcal{H}$ . A viable choice for  $\mathcal{G}$  is, given a kernel  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  with associated *reproducing kernel Hilbert space (RKHS)*  $\mathcal{F}$ , to consider  $\mathcal{G} = \mathcal{H} \otimes \mathcal{F}$ , which corresponds to the RKHS of vector-valued functions (see (Micchelli and Pontil, 2004; Alvarez et al., 2012)) with operator-valued kernel  $\Gamma(x, x') = k(x, x')I_{\mathcal{H}}$  and  $I_{\mathcal{H}} : \mathcal{H} \rightarrow \mathcal{H}$  the identity operator on  $\mathcal{H}$ . This approach is a direct generalization of the strategy introduced in the finite setting. Indeed, we have already observed that, when  $\mathcal{Y} = \mathcal{Z}$  is a finite set, we can choose  $\mathcal{H} = \mathbb{R}^T$  to satisfy the ILE definition. Moreover, for the linear kernel  $k(x, x') = \langle x, x' \rangle$  on  $\mathcal{X} = \mathbb{R}^d$ , the associated RKHS  $\mathcal{F}$  is isometric to  $\mathbb{R}^d$ . Therefore, we have  $\mathcal{H} \otimes \mathcal{F} \cong \mathbb{R}^T \otimes \mathbb{R}^d \cong \mathbb{R}^{T \times d}$ . Consequently, any  $g_n \in \mathcal{H} \otimes \mathcal{F}$  can be parametrized by a matrix  $W_n \in \mathbb{R}^{T \times d}$  and the ERM problem in (18) becomes equivalent to the one in (9).

The solution of the ridge regression problem can be obtained in closed form. In particular, it is easy to prove that analogously to (11) in the finite setting, for any  $x \in \mathcal{X}$  we have

$$g_n(x) = \sum_{i=1}^n \alpha_i(x) \varphi(y_i), \quad (19)$$

with the weights  $\alpha(x) = (K + \lambda n I)^{-1} \mathbf{v}(x)$  as in (14). By replacing  $g_n$  to  $g^*$  in Thm. 3 we recover the estimator  $f_n$  of (20), as desired. As we will discuss in more detail in Sec. 4.2, this strategy is related to the *Kernel Dependency Estimator (KDE)* (Weston et al., 2002) and *Input Output Kernel Regression (IOKR)* Brouard et al. (2011, 2016) for the case of loss functions corresponding to distances induced by a reproducing kernel.

The above reasoning can be applied to *any* function  $g_n$  expressed as a linear combination of (embedded) output points  $\varphi(y_i)$ . The following result summarizes this property, which allows to consider a family of novel estimators  $f_n$  parameterized by the weighting function  $\alpha$ .

**Lemma 4** *Let  $\Delta : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$  admit an ILE,  $(y_i)_{i=1}^n$  a set of points in  $\mathcal{Y}$  and  $\alpha : \mathcal{X} \rightarrow \mathbb{R}^n$  a weighting function. Let  $g_n : \mathcal{X} \rightarrow \mathcal{H}$  be such that  $g_n(x) = \sum_{i=1}^n \alpha_i(x) \varphi(y_i)$  for any  $x \in \mathcal{X}$ . Then, the function  $f_n : \mathcal{X} \rightarrow \mathcal{Z}$  such that  $\forall x \in \mathcal{X}$*

$$f_n(x) = \operatorname{argmin}_{z \in \mathcal{Z}} \langle \psi(z), g_n(x) \rangle_{\mathcal{H}} = \operatorname{argmin}_{z \in \mathcal{Z}} \sum_{i=1}^n \alpha_i(x) \Delta(z, y_i). \quad (20)$$

From the observation above, we see that the reasoning in Sec. 3.1 can indeed be generalized to the setting where  $\mathcal{Z}$  and  $\mathcal{Y}$  are not finite and  $\Delta$  admits an ILE. For any  $x \in \mathcal{X}$ , the weights  $\alpha(x)$  are learned from training data according to (14) and only the loss function appears in the form of the estimator. We expand on this in the following remark.

**Remark 5 (“Loss Trick”)** *In practice, learning and evaluating  $f_n$  does not require explicit knowledge of the space  $\mathcal{H}$  or the embeddings  $\psi$  and  $\varphi$  (see (20)), which are implicitly encoded within the definition of ILE and are only required for theoretical purposes (derivation and characterization of generalization properties of  $f_n$  as discussed in Sec. 5). This effect was originally referred to as “loss trick” (Ciliberto et al., 2016) in analogy to the “kernel trick” on the output space, originally introduced in Geurts et al. (2006); Brouard et al. (2011, 2016) for Output Kernel Regression.*

### 3.3. Additional ILE-induced Algorithms and Estimators

In this section we discuss alternative approaches to learn the weighting functions  $\alpha : \mathcal{X} \rightarrow \mathbb{R}^n$  defining the estimator  $f_n$ . These strategies are derived by replacing kernel ridge-regression with a different approximation of  $g^*$  that still satisfies the hypotheses of Thm. 4, namely can be written as the linear combination of training outputs. As already observed in the literature of standard regression settings (see e.g. Rosasco et al., 2005, and references therein), these alternative approaches can offer significant computational advantages over kernel ridge regression while guaranteeing the same generalization performance from the statistical standpoint.

**“Exact” Kernel methods.** A wide family of algorithms that can be used to estimate  $g^*$  are based on *spectral filtering* regularization strategies (Rosasco et al., 2005; Bauer et al., 2007). In particular we consider:

- **$L^2$ -Boosting (L2B).** By considering  $g_n$  the  $t$ -th iterate of gradient descent of the (non-regularized) empirical risk minimization problem in (18), we have

$$\alpha(x) = C_t \mathbf{v}(x) \quad \text{with} \quad C_t = (I - \nu/n K) C_{t-1} + \nu/n I, \quad (21)$$

with  $C_t \in \mathbb{R}^{n \times n}$  defined recursively starting from any  $C_0 \in \mathbb{R}^{n \times n}$  and  $\nu/n$  the gradient descent step size with  $\nu > 0$ . We recall that  $\mathbf{v}(x) \in \mathbb{R}^n$  denotes the evaluation vector in  $x$ , with entries  $\mathbf{v}(x)_i = k(x, x_i)$  for any  $i = 1, \dots, n$ . The number of steps  $t \in \mathbb{N}$  acts as regularization parameter. Accelerated and stochastic versions can be considered. In

the case of loss functions induced by a normalized kernel (e.g. KDE and OKR, see also Sec. 4.2), this estimator corresponds to the one introduced in Geurts et al. (2007).

- **Principal Component Regression (PCR).** Take  $g_n$  to be the estimator obtained by filtering out the eigenvalues of the kernel matrix  $K$  below a threshold  $\lambda > 0$  and inverting the eigenvalues that are above. We have

$$\alpha(x) = U \Sigma_\lambda^\dagger U^\top \mathbf{v}(x), \quad (22)$$

where  $K = U \Sigma U^\top$  is the singular value decomposition of  $K$  and  $\Sigma_\lambda^\dagger$  is the pseudoinverse of the matrix corresponding to  $\Sigma$  with all eigenvalues smaller than  $\lambda$  set to 0.

**Random Projections.** Methods leveraging random projections achieve optimal generalization performance while being significantly more efficient computationally.

- **Random Features (RF).** Let  $(\Omega, \pi)$  be a probability space and  $\zeta : \mathcal{X} \times \Omega \rightarrow \mathbb{R}$  be a map such that  $k(x, x') = \int_\Omega \zeta(x, \omega) \zeta(x', \omega) d\pi(\omega)$  (Rahimi and Recht, 2008)<sup>1</sup>. Let  $M \in \mathbb{N}$  and  $\omega_1, \dots, \omega_M$  be independently sampled from  $\pi$ . Denote by  $\hat{\mathbf{v}}_M : \mathcal{X} \rightarrow \mathbb{R}^M$ , the map

$$\hat{\mathbf{v}}_M(x) = \frac{1}{\sqrt{M}} (\zeta(x, \omega_1), \dots, \zeta(x, \omega_M)). \quad (23)$$

By definition  $\hat{\mathbf{v}}_M(x)^\top \hat{\mathbf{v}}_M(x')$  is a discretization of the integral defining  $k(x, x')$ . The scores  $\alpha$  are learned according to this new feature map

$$\alpha(x) = W \hat{\mathbf{v}}_M(x), \quad W = Q_M (Q_M^\top Q_M + n\lambda I)^{-1}, \quad (24)$$

with  $Q_M \in \mathbb{R}^{n \times M}$ ,  $Q_M = (\hat{\mathbf{v}}_M(x_1), \dots, \hat{\mathbf{v}}_M(x_n))^\top$ . This approach is significantly faster than ridge-regression when  $M \ll n$ .

- **Nystrom Approximation (NY).** Sample  $M \leq n$  points  $\tilde{x}_1, \dots, \tilde{x}_M$  from the input dataset. Denote  $K_{MM} \in \mathbb{R}^{M \times M}$  be the matrix with  $(K_{MM})_{i,j} = k(\tilde{x}_i, \tilde{x}_j)$  and  $K_{nM} \in \mathbb{R}^{n \times M}$  the matrix with elements  $(K_{nM})_{i,j} = k(x_i, \tilde{x}_j)$  (see Smola and Schölkopf, 2000). The scores  $\alpha$  are defined as

$$\alpha(x) = W \tilde{\mathbf{v}}_M(x), \quad W = K_{nM} (K_{nM}^\top K_{nM} + n\lambda K_{MM})^\dagger, \quad (25)$$

with  $\tilde{\mathbf{v}}_M(x) = (k(\tilde{x}_1, x), \dots, k(\tilde{x}_M, x)) \in \mathbb{R}^M$ , for any  $x \in \mathcal{X}$ . These operations are significantly faster than solving ridge-regression when  $M \ll n$ .

**Nadaraya-Watson (NW).**  $g_n$  can be obtained via the Nadaraya-Watson (NW) strategy (Nadaraya, 1964). In this case we have

$$\alpha(x) = \frac{1}{\mathbb{1}^\top \mathbf{v}(x)} \mathbf{v}(x), \quad (26)$$

---

1. E.g. for the Gaussian kernel and  $\mathcal{X} = \mathbb{R}^d$ ,  $d \in \mathbb{N}$ , we have  $\Omega = \mathbb{R}^d \times [0, 1]$ , and for  $(w, b) =: \omega \in \Omega$ ,  $\pi((w, b)) = \mathcal{N}(w)U(b)$ , with  $\mathcal{N}(\cdot)$  standard normal distribution  $U(\cdot)$  uniform distribution, and  $\zeta(x, (w, b)) = \cos(w^\top x + b)$  (see Rahimi and Recht, 2008, for more details).

resulting in the estimator

$$f_n(x) = \operatorname{argmin}_{z \in \mathcal{Z}} \sum_{i=1}^n \frac{k(x, x_i)}{\sum_{j=1}^n k(x_j, x_i)} \Delta(z, y_i) \quad (27)$$

Computation of the NW estimator does not involve the kernel matrix  $K$ . This can be beneficial in large-scale scenarios where the kernel matrix can be large. However, the NW estimator is often less adaptive than ridge-regression to the smoothness properties of the learning problem (Györfi et al., 2006). As a consequence, the learning rates of NW are usually slower than kernel ridge regression in non-worst-case settings.

**Nearest Neighbors (NN).** Given a measure of similarity on the input set (e.g. a kernel), for any test point  $x \in \mathcal{X}$ , the Nearest Neighbor (NN) estimator (Friedman et al., 2001) corresponds to the average, on the space  $\mathcal{H}$  of output training points  $\varphi(y_i)$  whose corresponding inputs  $x_i$  are among the first  $q$  most similar (or closest) to  $x$ . This corresponds to return the *binary* scores  $\alpha(x) \in \{0, 1\}^n$

$$\alpha(x)_i = \begin{cases} 1 & \text{if } x_i \text{ is a } q\text{-nearest-neighbor} \\ 0 & \text{otherwise} \end{cases} \quad (28)$$

NN does not rely on a training phase (except for a possible pre-processing such as kd-trees to allow for a faster search of neighbors at test time). Interestingly, at test time, the cost of the optimization over  $\mathcal{Z}$  in (20) depends also on the number  $q$  of neighbors selected, which is a hyperparameter of the estimator.

So far we have introduced a novel family of estimators for structured prediction and discussed how they can be derived from the notion of Implicit Loss Embedding. We are left with two critical questions that will be addressed in the following: on one hand we need to characterize how the approximation of  $g_n$  can lead to good estimations of  $f^*$ . A second, more concrete question is whether the ILE definition is sufficiently flexible to encompass relevant structured prediction problems or, in other words, which functions admit an ILE. We will address the first question in Sec. 5 and the second one in Sec. 6. Before doing so, in Sec. 4 we draw some connections with previous work. These observations will prove useful to better situate our framework within the literature on structured prediction. We conclude this section with a comment on evaluating the ILE in practice.

### 3.4. Evaluating the ILE Estimator

According to (20), evaluating  $f_n$  on a test point  $x \in \mathcal{X}$  consists in solving an optimization problem over the output space  $\mathcal{Z}$ . This design of the test phase is standard in structured prediction settings and in particular for auxiliary function maximization methods (see Sec. 4.3), where a corresponding optimization protocol is derived on a case-by-case basis depending on the loss and the space  $\mathcal{Z}$ , (see Nowozin et al., 2011, and references therein). The objective functional in our setting can be interpreted as estimating the weighted *barycenter* (or Frechet mean) of the training output points with respect to the “distance” induced by

the loss  $\Delta$ . This perspective can be particularly advantageous when efficient methods for barycenter computation are available for the given loss.

**Example 2 (Distributional Regression with Optimal Transport)** *Luise et al. (2018) addressed the problem of learning to predict probability distributions with respect to optimal transport-based loss. Here  $\mathcal{Z}$  corresponds to the set of probability distributions over a discrete set and  $\Delta$  is the Sinkhorn loss (Cuturi, 2013). The Sinkhorn loss is an entropic regularization of the well-known Wasserstein distance, which enjoys better computational properties and for which it was shown that the ILE property applies (Luise et al., 2018). Efficient methods exists to compute the barycenter with respect to the Sinkhorn loss (Cuturi and Doucet, 2014; Benamou et al., 2015) that can be readily applied to solve (20) in practice.*

When  $\mathcal{Z}$  is locally diffeomorphic to a linear space, the objective functional in our setting allows also to suggest a general stochastic meta-algorithm. In particular, (20) can interpreted as the problem of minimizing an expectation

$$f_n(x) = \operatorname{argmin}_{z \in \mathcal{Z}} \mathbb{E}_{i \sim \alpha(x)} h_i(z, x) \quad \text{with} \quad h_i(z, x) = \frac{\operatorname{sign}(\alpha_i(x))}{\mathbf{a}(x)} \Delta(z, y_i) \quad (29)$$

where  $i \in \{1, \dots, n\}$  is a random variable sampled according to the relevance weights  $\alpha_i(x)$  and  $\mathbf{a}(x) = \sum_{i=1}^n \alpha_i(x)$ . Thus, when  $\Delta$  is (sub)differentiable in the first variable, problems of the form of (29) can be directly addressed addressed by stochastic gradient methods (SGM).

**Example 3 (Manifold Regression)** *If  $\mathcal{Z}$  is a Riemannian manifold, the optimization problem in (29) can be addressed by methods such as stochastic Riemannian descent (Sra and Hosseini, 2016). This strategy was adopted in (Rudi et al., 2018) for the case of  $\mathcal{Z}$  the space of positive semidefinite matrices, the discrete probability simplex and the manifold of orientations of a 2D vector. More recently Marconi et al. (2020) proposed a taxonomy generalization approach for relational data representation, which hinges on solving a manifold structured prediction problem over the hyperbolic space  $\mathcal{Z}$ . In all the above settings, the loss  $\Delta$  was chosen to be the squared geodesic distance on  $\mathcal{Z}$ .*

We conclude this section with a remark on the computational complexity in the discrete setting. In particular we comment on the comparison between the inference problem in (20) related to the overall supervised problem of structured prediction.

**Remark 6 (On the Complexity of Inference)** *In general, given the scores  $\alpha_i(x)$ , solving the inference problem to obtain  $f_n(x)$  requires solving a possibly hard optimization problem. However, in most settings, this approach can be more favorable than ERM. Indeed, consider for simplicity the case where both  $\mathcal{X}$  and  $\mathcal{Z}$  are finite spaces with cardinality  $|\mathcal{X}|$  and  $|\mathcal{Z}|$  respectively. ERM would require solving an optimization problem on the space of all functions from  $\mathcal{X}$  to  $\mathcal{Z}$ , which has cardinality  $|\mathcal{Z}|^{|\mathcal{X}|}$ . On the other hand, the ILE approach acts by first learning the scores  $\alpha_i(x)$ , which is done efficiently by solving a linear system and then finding the best output  $f(x)$  over the space  $\mathcal{Z}$ . This amounts to solving an optimization over a space of cardinality  $|\mathcal{Z}|$ , which is significantly smaller than  $|\mathcal{Z}|^{|\mathcal{X}|}$ .*

## 4. Connections with Previous Work

In this section we highlight some relevant connections between our framework and previous literature. As mentioned in Sec. 2 we show that, although starting from a different perspective, our method can be interpreted as a synthesis of the two main structured prediction strategies considered in the literature, namely *surrogate frameworks* and *auxiliary function maximization methods*. In this sense, our approach represents the best of both worlds, since it is theoretically sound (as we will show in Sec. 5) but it is also applicable to a large family of learning problems. We also draw a connection with the *conditional mean embeddings* literature, which will offer relevant insights on the theoretical analysis of Sec. 5.

### 4.1. Surrogate Frameworks

Surrogate approaches are designed to address specific structured prediction problems such as classification (Bartlett et al., 2006; Mroueh et al., 2012), multi-labeling (Gao and Zhou, 2013), ranking (Duchi et al., 2010) or quantile regression (Steinwart et al., 2011) to name a few. The core idea underlying these methods is to deal with the structure of the problem by: 1) finding an embedding (or *encoding*) of the output variables into a linear space where, 2) a *surrogate* learning problem can be solved efficiently and finally, 3) map back the surrogate solution by means of a suitable *decoding*.

More formally, given a training dataset  $(x_i, y_i)_{i=1}^n$ , a surrogate approach consists in the following three steps:

1. **Encoding.** Choose a coding  $\mathbf{c} : \mathcal{Y} \rightarrow \mathcal{H}$  into a surrogate space  $\mathcal{H}$ .  
Map  $(x_i, y_i)_{i=1}^n$  to the surrogate dataset  $(x_i, \mathbf{c}(y_i))_{i=1}^n$ .
2. **Learning.** Define a surrogate loss  $\mathcal{L} : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$ .  
Learn  $g_n : \mathcal{X} \rightarrow \mathcal{H}$  minimizing  $\mathcal{L}$  on  $(x_i, \mathbf{c}(y_i))_{i=1}^n$ .
3. **Decoding.** Choose a decoding  $\mathbf{d} : \mathcal{H} \rightarrow \mathcal{Z}$  and define  $f_n = \mathbf{d} \circ g_n : \mathcal{X} \rightarrow \mathcal{Y}$ .

A prototypical example of this strategy is represented by binary classification.

**Example 4 (Binary Classification)** *In binary classification the goal is to learn a binary-valued function  $f : \mathcal{X} \rightarrow \mathcal{Z} = \{-1, 1\}$ . The prototypical approach to address this problem is to consider  $\mathbf{c} : \mathcal{Y} = \{-1, 1\} \rightarrow \mathcal{H} = \mathbb{R}$  the identity map and then learn  $g_n : \mathcal{X} \rightarrow \mathbb{R}$  by minimizing a suitable loss  $\mathcal{L} : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$  (e.g. least-squares, hinge, logistic, etc.). The final classifier is then obtained as  $f_n(x) = \text{sign}(g_n(x))$ , with decoding  $\mathbf{d} = \text{sign} : \mathbb{R} \rightarrow \{-1, 1\}$ .*

Surrogate frameworks critically hinge on identifying a suitable candidate for the loss function  $\mathcal{L} : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$ . Indeed, while on one hand  $\mathcal{L}$  should allow to compute the estimator  $g_n : \mathcal{X} \rightarrow \mathcal{H}$  efficiently, on the other hand the surrogate learning process induced by  $\mathcal{L}$  needs to be related to the original structured prediction problem. The requirement for efficiency is typically satisfied by choosing  $\mathcal{L}$  to be a convex loss (e.g. least-squares, hinge or logistic in binary classification, see Example 4). The connection with structured prediction is investigated by studying the ideal learning problem induced by the *surrogate risk*

$$\mathcal{R}(g) = \int \mathcal{L}(g(x), \mathbf{c}(y)) d\rho(x, y), \tag{30}$$

with  $g : \mathcal{X} \rightarrow \mathcal{H}$ . Intuitively, for a “good” surrogate framework, the global minimizer  $g^* : \mathcal{X} \rightarrow \mathcal{H}$  of the risk  $\mathcal{R}$  should allow to recover the original solution  $f^* : \mathcal{X} \rightarrow \mathcal{Z}$  by means of the decoding, for instance  $f^* = \mathbf{d} \circ g^*$ . Moreover, ideally, as the number  $n$  of training points increases and the estimator  $g_n$  provides a better approximation to  $g^*$ , we would like the predictor  $f_n$  to converge to  $f^*$  as well.

Formalizing the observations above, the following two conditions are typically required by Surrogate Frameworks:

- **Fisher Consistency.**  $\mathcal{E}(f^*) = \mathcal{E}(\mathbf{d} \circ g^*)$ ,
- **Comparison Inequality.**  $\mathcal{E}(\mathbf{d} \circ g) - \mathcal{E}(f^*) \leq \gamma(\mathcal{R}(g) - \mathcal{R}(g^*))$  for any  $g : \mathcal{X} \rightarrow \mathcal{H}$ , with  $\gamma : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  continuous, non-decreasing and such that  $\gamma(0) = 0$ .

The Fisher consistency establishes the validity of the surrogate framework in terms of the original problem. It guarantees that we can always recover the ideal  $f^*$  from the surrogate solution  $g^*$ . The comparison inequality allows to automatically extend any result characterizing the learning rates of the surrogate estimator to obtain explicit excess risk bounds for the structured prediction one (possibly accelerated or slowed down by a factor depending on the function  $\gamma$ ). We refer to Mroueh et al. (2012) and references therein for concrete examples of this strategy in a variety of structured prediction settings.

**Surrogate Frameworks and ILE.** While surrogate methods are typically designed on a case-by-case basis for each learning problem, the structured prediction framework proposed in this work can be interpreted as a general form of surrogate approach. In particular, it is natural to choose the encoding map as the ILE feature map on the label space  $\mathcal{Y}$ , namely  $\mathbf{c} = \varphi$ . Moreover, we have observed how any ILE loss function is directly associated to a suitable surrogate output space  $\mathcal{H}$  via the ILE definition itself. In particular, in Thm. 3 we have shown how the corresponding structured prediction problem is naturally associated to the surrogate risk  $\mathcal{R}$  with surrogate loss  $\mathcal{L}(\eta_1, \eta_2) = \|\eta_1 - \eta_2\|_{\mathcal{H}}^2$  the square loss on  $\mathcal{H}$ . It follows that we can choose as decoding the map  $\mathbf{d} : \mathcal{H} \rightarrow \mathcal{Z}$  such that

$$\mathbf{d}(\eta) = \operatorname{argmin}_{z \in \mathcal{Z}} \langle \psi(z), \eta \rangle_{\mathcal{H}} \quad (31)$$

for any  $\eta \in \mathcal{H}$ . Indeed, with the notation of Sec. 3, we have  $f_n = \mathbf{d} \circ g_n$  according to (20) that recovers our structured prediction estimators. Note in particular that Thm. 3 shows that our framework is Fisher consistent, by proving that indeed  $f^* = \mathbf{d} \circ g^*$  as required.

The connection with surrogate methods will be completed by our theoretical analysis of Sec. 5. Indeed, analogously to surrogate approaches, our proof strategy hinges on proving a comparison inequality, which allows to study the generalization properties on the surrogate problem to control the excess risk of the structured prediction estimator. In particular, in Thm. 7 we provide a comparison inequality for our framework with  $\gamma$  corresponding to the square root function.

## 4.2. Output Kernel Regression

Output Kernel Regression (OKR) (Geurts et al., 2006, 2007; Brouard et al., 2011; Kadri et al., 2013; Brouard et al., 2016) is a general framework to address supervised learning

problems where the loss function corresponds to the canonical Euclidean distance induced by a reproducing kernel on the output space. OKR recovers the method originally proposed in Kernel Dependency Estimation (KDE) (Weston et al., 2002; Cortes et al., 2005) as a special case. In these settings we need to assume  $\mathcal{Y} = \mathcal{Z}$ . Consider  $h : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  a reproducing kernel on the output space. Let  $\mathcal{H}$  be the RKHS associated to  $h$ , with feature map  $\varphi : \mathcal{Y} \rightarrow \mathcal{H}$ , namely  $h(y, y') = \langle \varphi(y), \varphi(y') \rangle_{\mathcal{H}}$  for any  $y, y' \in \mathcal{Y}$ . OKR methods address the problem of learning a function  $f : \mathcal{X} \rightarrow \mathcal{Y}$  by minimizing the least-squares loss on  $\mathcal{H}$ ,

$$\Delta(f(x), y) = \|\varphi(f(x)) - \varphi(y)\|_{\mathcal{H}}^2, \quad (32)$$

for  $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$ . Given a dataset  $(x_i, y_i)_{i=1}^n$ , learning  $f_n$  directly might be challenging because of the structure of  $\mathcal{Y}$ . Alternatively, it might be possible to leverage the linear structure of  $\mathcal{H}$  to learn a function  $g_n : \mathcal{X} \rightarrow \mathcal{H}$ . Whenever a test point  $x \in \mathcal{X}$  is provided, the prediction  $f_n(x)$  is then obtained by finding the  $y \in \mathcal{Y}$  for which  $\varphi(y)$  is closest to  $g_n(x)$  according to the canonical distance on  $\mathcal{H}$ . This phase, akin to the decoding of surrogate methods, is referred to as the *preimage problem*. When  $g_n(x) = \sum_{i=1}^n \alpha_i(x) \varphi(y_i)$ , given the score functions  $\alpha$ , this problem can be cast as the optimization

$$f_n(x) = \operatorname{argmin}_{y \in \mathcal{Y}} \|\varphi(y) - g_n(x)\|_{\mathcal{H}}^2 = \operatorname{argmin}_{y \in \mathcal{Y}} h(y, y) - 2 \sum_{i=1}^n \alpha_i(x) h(y, y_i). \quad (33)$$

Indeed, thanks to the reproducing property of the kernel  $h$  we have for every  $y, y' \in \mathcal{Y}$

$$\|\varphi(y) - \varphi(y')\|_{\mathcal{H}}^2 = h(y, y) - 2h(y, y') + h(y', y'). \quad (34)$$

Recently, OKR settings with other loss functions have been considered (hinge (Brouard et al., 2016) or e.g. Huber,  $\epsilon$ -insensitive (Laforgue et al., 2019)).

**Output Kernel Regression and ILE.** There is a clear connection between OKR (and KDE) approaches and the ILE framework considered in this work. Indeed, as we show in Sec. 6 (Thm. 13), the loss function considered by OKR satisfies the ILE definition. Moreover, if the output kernel  $h$  is such that  $h(y, y) = \eta$  for any  $y \in \mathcal{Y}$ , with  $\eta > 0$  a constant, then the ILE estimator in (20) corresponds exactly to OKR (assuming same scores  $\alpha$ ). The latter observation implies that the theoretical analysis reported in Sec. 5 applies also to OKR and KDE. Therefore, a further byproduct of our work is the theoretical analysis of these strategies, which to our knowledge is a novel contribution on its own.

We conclude this section by highlighting two critical differences between our framework and the methods above:

- KDE was designed to address structured prediction problems by substituting the original structured loss with the least-squares induced by a kernel on the output. There is no guarantee in general that the KDE estimator is in any way solving the structured prediction. In this sense KDE could be interpreted as a form of surrogate method for which the surrogate problem does not satisfy neither Fisher consistency nor Comparison Inequality.
- If the condition  $h(y, y) = \eta$  does not hold, the OKR (or KDE) and ILE estimators *do not coincide*. This means that there is no guarantee that the former approaches will enjoy the same generalization properties of ILE methods studied in Sec. 5 below.



### 4.3. Auxiliary Function Maximization Methods

In contrast to surrogate approaches, *auxiliary function maximization* methods (Bakir et al., 2007) have been designed to address a wide range of structured prediction problems within a single, general framework. Given a training dataset  $(x_i, y_i)_{i=1}^n$ , these methods learn a *score function*  $F_n : \mathcal{Z} \times \mathcal{X} \rightarrow \mathbb{R}$  that measures the likelihood of observing a given input-output pair  $(z, x)$ . In these settings, the structured prediction estimator  $f_n : \mathcal{Z} \rightarrow \mathcal{X}$  is defined in terms of an optimization problem, namely as the function selecting the “most likely” output according to the score function  $F_n$ . This amounts to solving the maximization problem

$$f_n(x) = \operatorname{argmax}_{z \in \mathcal{Z}} F_n(z, x), \quad (35)$$

for any input  $x \in \mathcal{X}$  provided at test time.

Auxiliary function maximization methods are identified by the approach used to learn the score function  $F_n$ . A general strategy, adopted for instance by *Structured Support Vector Machines (SVMStruct)* (Tsochantaridis et al., 2005) is to consider a model of the form

$$F_n(z, x) = \langle \Phi(z, x), \mathbf{w}_n \rangle_{\mathcal{G}} \quad (36)$$

for  $x \in \mathcal{X}$  and  $z \in \mathcal{Z}$ , where  $\mathcal{G}$  is a suitable feature space and  $\Phi : \mathcal{Z} \times \mathcal{X} \rightarrow \mathcal{G}$  a joint feature map on the input-output set. The function  $F_n$  is therefore parametrized by the vector  $\mathbf{w}_n \in \mathcal{G}$ , which is learned during the training phase. For instance,  $\mathbf{w}_n$  can be learned by minimizing an upper bound of the empirical risk by extending the strategy used in binary classification settings for the standard SVM approach (Tsochantaridis et al., 2005; Cortes et al., 2016).

Other approaches follow more adherently the interpretation of  $F_n$  as measuring the likelihood of input-output pairs and thus consider models of the form

$$F_n(z, x) = p(z|x) = \frac{e^{-\langle \Phi(z, x), \mathbf{w}_n \rangle}}{\sum_{z' \in \mathcal{Z}} e^{-\langle \Phi(z', x), \mathbf{w}_n \rangle}}, \quad (37)$$

where  $F_n(\cdot, x)$  is a probability distribution over  $\mathcal{Z}$ . These methods consider a similar parametrization of the target function to SVMStruct approaches. However, they differ from the latter methods in that during training the aim to approximate the underlying input-output distribution. This model is often adopted by structured prediction approaches based on *Conditional Random Fields (CRF)* (Vishwanathan et al., 2006; Morency et al., 2007). For an in-depth introduction on auxiliary function maximization methods we refer the reader to Nowozin et al. (2011) and references therein.

We care to point out that, in general, although the approaches above consider models that can be applied to arbitrary output spaces  $\mathcal{Z}$ , these algorithms typically *require  $\mathcal{Z}$  to be finite*. In this sense, a fundamental advantage of the ILE framework considered in this work is to go beyond the finite setting.

**Auxiliary function maximization and ILE.** The structured prediction framework discussed in this work has a natural interpretation as an auxiliary function maximization approach. To see this, consider

$$F_n(z, x) = - \sum_{i=1}^n \alpha_i(x) \Delta(z, y_i) \quad (38)$$

for  $x \in \mathcal{X}$  and  $z \in \mathcal{Z}$ . If  $\Delta$  admits an ILE, our method corresponds to parametrizing  $F_n$  as in (36) above, with  $\mathcal{G} = \mathcal{H} \otimes \mathcal{F}$  and  $\Phi(z, x) = \psi(z) \otimes \phi(x)$ , where  $\mathcal{F}$  is a RKHS on the input set  $\mathcal{X}$  with associated kernel  $\phi : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  and the map  $\phi(x) = k(x, \cdot) \in \mathcal{F}$  can be interpreted as a feature map from  $\mathcal{X}$  to  $\mathcal{F}$ . By leveraging the properties of the tensor product operation, we have

$$F_n(x, z) = \langle \Phi(x, z), \mathbf{w}_n \rangle_{\mathcal{G}} = \langle \psi(z), W_n \phi(x) \rangle_{\mathcal{H}} = \langle \psi(z), g_n(x) \rangle_{\mathcal{H}}, \quad (39)$$

where we have defined  $g_n : \mathcal{X} \rightarrow \mathcal{H}$  as the function such that  $g_n(x) = W_n \phi(x)$  for any  $x \in \mathcal{X}$  and  $W_n : \mathcal{F} \rightarrow \mathcal{H}$  is the operator corresponding to  $\mathbf{w}_n$  given by the canonical isomorphism between  $\mathcal{H} \otimes \mathcal{F}$  and the space  $\text{HS}(\mathcal{F}, \mathcal{H})$  of Hilbert-Schmidt operators from  $\mathcal{F}$  to  $\mathcal{H}$ .

In Sec. 3 we have discussed a number of algorithms to learn  $W_n$  (or, equivalently  $\mathbf{w}_n$ ), whose theoretical properties have then been studied in Sec. 5. This connection opens two relevant questions for future investigation: 1) study approaches to learn the parameters  $W_n$  that do not necessarily converge to the conditional mean (but for which it is still possible to prove the consistency of the resulting structured prediction); 2) While the ILE assumption seem to require a “separable” representation of the form  $\Phi(x, z) = \psi(z) \otimes \phi(x)$ , it would be interesting to consider joint input-output feature maps, which could prove beneficial in settings where structural relations between input and outputs could be leveraged. A potential promising approach to address this question would be to borrow ideas from the literature on vector-valued and multi-task learning with RKHS for vector-valued functions (see for instance Alvarez et al., 2012, and references therein).

#### 4.4. Conditional Mean Embeddings

In this section we highlight the relation between our structured prediction framework and *conditional mean embeddings* (Song et al., 2009). This connection is particularly useful to understand the role played by the surrogate function  $g^*$  within the analysis of Sec. 5.

Let  $\mathcal{H}$  be a RKHS of functions from  $\mathcal{Y}$  to  $\mathbb{R}$  with associated positive definite kernel  $h : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ . We recall that the conditional mean embedding in  $\mathcal{H}$  of  $\rho(\cdot|x)$  given  $x \in \mathcal{X}$  is defined as

$$\mu_{y|x} = \int_{\mathcal{Y}} h(y, \cdot) d\rho(y|x). \quad (40)$$

A key aspect of conditional mean embeddings is that, thanks to the reproducing property of the RKHS, for any  $f \in \mathcal{H}$  we have

$$\langle f, \mu_{y|x} \rangle_{\mathcal{H}} = \mathbb{E}_{y|x} f(y). \quad (41)$$

This allows to evaluate the conditional expectation with respect to  $\rho(\cdot|x)$  of any function  $f \in \mathcal{H}$  by directly performing an inner product with  $\mu_{y|x}$ .

It was observed in Sriperumbudur et al. (2011) that for a wide family of RKHS, called *characteristic*, the kernel mean embedding operator is injective. In other words, two distributions have same embedding in  $\mathcal{H}$  if and only if they coincide. This implies that kernel mean embeddings, and in particular conditional mean embeddings, encode rich information about the associated distribution within a single vector in  $\mathcal{H}$ .

**Conditional Mean Embeddings and ILE.** Let  $\Delta : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$  admit an ILE. Under the same notation of Thm. 2, assume the corresponding surrogate space  $\mathcal{H}$  to be a RKHS and that  $\varphi : \mathcal{Y} \rightarrow \mathcal{H}$  is an associated feature map, namely such that  $h(y, y') = \langle \varphi(y), \varphi(y') \rangle_{\mathcal{H}}$  is the reproducing kernel associated to  $\mathcal{H}$ . Then, according to the characterization of  $g^*$  in (16), for any  $x \in \mathcal{X}$  we have that  $g^*(x) = \mu_{y|x}$  corresponds to the conditional mean embedding of  $\rho(\cdot|x)$  Song et al. (2009). Moreover, by denoting  $\psi(z) = \Delta(z, \cdot)$  and leveraging the reproducing property of the mean embedding, we have

$$\langle \psi(z), g^*(x) \rangle_{\mathcal{H}} = \left\langle \Delta(z, \cdot), \mu_{y|x} \right\rangle_{\mathcal{H}} = \mathbb{E}_{y|x} \Delta(z, y). \quad (42)$$

Interestingly, this observation recovers directly the Fisher consistency result of Thm. 3 *when  $\mathcal{H}$  is an RKHS*. Indeed, we have observed in (3) that the solution  $f^*$  of the structured prediction expected risk corresponds to the minimizer of the conditional expectation  $\mathbb{E}_{y|x} \Delta(z, y)$ . The equation above implies that this is equivalent to have  $f^*(\cdot) = \operatorname{argmin}_{z \in \mathcal{Z}} \langle \psi(z), g^*(\cdot) \rangle$ .

In Sec. 5.3 we will see that in order to prove learning rates for the structured prediction estimator we will need to impose assumptions on  $g^*$ . These could be interpreted as requiring the learning problem to satisfy regularity conditions. Indeed, the connection above between the ILE definition and conditional mean embeddings shows that  $g^*$  is implicitly encoding key properties of the data generating distribution  $\rho$  and, consequently, of the structured prediction problem itself. For more details on the topic, we refer the interested reader to (Muandet et al., 2017) for an in-depth introduction on kernel mean embeddings and to (Song et al., 2009; Grünewälder et al., 2012; Singh et al., 2019) for the special case of conditional mean embeddings.

## 5. Theoretical Analysis

This section is devoted to characterize the statistical properties of the structured prediction estimators introduced in this work. In particular we will prove that under standard hypotheses from the statistical learning literature our approach is universally consistent and enjoys optimal learning rates.

### 5.1. Comparison Inequality

The key result of our analysis, discussed in this section, is to show how the approximation of  $g^*$  via an estimator  $g_n$  (such as those discussed in Sec. 3) allows to characterize the behavior of the corresponding estimator  $f_n = \mathfrak{d} \circ g_n$  with respect to the ideal solution  $f^*$ . The following result provides such characterization to any function  $g : \mathcal{X} \rightarrow \mathcal{H}$ .

**Theorem 7 (Comparison Inequality)** *Let  $\mathcal{Z}$  be a compact set and  $\Delta : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$  admit an ILE. Let  $f^*, g^*$  and the risk  $\mathcal{R}(\cdot)$  be defined as in Thm. 3. Let  $g : \mathcal{X} \rightarrow \mathcal{H}$  be measurable and let  $f : \mathcal{X} \rightarrow \mathcal{Z}$  be such that*

$$f(x) = \operatorname{argmin}_{z \in \mathcal{Z}} \langle \psi(z), g(x) \rangle_{\mathcal{H}}, \quad (43)$$

for any  $x \in \mathcal{X}$ . Then,

$$\mathcal{E}(f) - \mathcal{E}(f^*) \leq 2 c_{\Delta} \sqrt{\mathcal{R}(g) - \mathcal{R}(g^*)} \quad (44)$$

The result in Thm. 7 states that we can control the structured prediction excess risk in terms of the least-squares risk  $\mathcal{R}$  in approximating  $g^*$ . The theorem holds for any function  $g : \mathcal{X} \rightarrow \mathcal{H}$  that is measurable, a technical requirement satisfied in particular by every regression estimator  $g_n$  introduced in Sec. 3.

Thm. 7 shifts the problem of studying the generalization properties of  $f_n$  to that of characterizing the learning rates of the vector-valued estimator  $g_n$ , for which more well-established tools from statistical learning theory can be leveraged. Throughout this work we will refer to (44) as the *comparison inequality* of our structured prediction framework. This notation is borrowed from the literature on surrogate methods, as discussed in more detail in Sec. 4.1. An result analogous to Thm. 7 was shown originally in Ciliberto et al. (2016) for functions satisfying a similar property to ILE. For completeness, in Appendix A we prove it for ILE functions.

## 5.2. Universal Consistency

The comparison inequality in Thm. 7 is instrumental to study the generalization properties of the estimator considered in this work. In particular, the results reported in the rest of this section are obtained by characterizing the statistical properties of the estimator  $g_n$  and then extending them to  $f_n$  by means of the inequality in (44).

We start from the result proving the universal consistency of  $f_n$ . This is a fundamental requirement for a valid learning algorithm, stating that  $\mathcal{E}(f_n)$  converges to the minimum possible risk  $\mathcal{E}(f^*)$  as the number  $n$  of training points grows to infinity. A key assumption in this setting will be that the kernel  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  on the input space, introduced to learn the coefficients  $\alpha_i$  characterizing the solution  $f_n$  in (20), is *universal*. This is a standard assumption in statistical learning theory (Steinwart and Christmann, 2008, see e.g.) and corresponds to requiring the RKHS  $\mathcal{F}$  associated to  $k$  to be dense in the space of continuous function on  $\mathcal{X}$ . Typical examples of universal kernels on  $\mathcal{X} \subseteq \mathbb{R}^d$  are the Gaussian  $k(x, x') = e^{-\|x-x'\|^2/\sigma^2}$  or the Laplacian  $k(x, x') = e^{-\|x-x'\|/\sigma}$  kernels.

**Theorem 8 (Universal Consistency)** *Let  $\mathcal{Z}$  be a compact set and  $\Delta : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$  admit an ILE. Let  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  be a bounded universal reproducing kernel. For any  $n \in \mathbb{N}$  and any distribution  $\rho$  on  $\mathcal{X} \times \mathcal{Y}$  let  $f_n : \mathcal{X} \rightarrow \mathcal{Z}$  be the estimator in (20) trained on  $(x_i, y_i)_{i=1}^n$  points independently sampled from  $\rho$  and with weights  $\alpha$  defined as:*

- (a) (Ridge Regression) in (14) with  $\lambda_n = n^{-1/2}$ , or
- (b) (L2-Boosting) in (21) with step-size  $\nu < 1/\kappa^2$  and  $t_n = n^{1/2}$ , or
- (c) (PCR) in (22) with  $\lambda_n = n^{-1/2}$ .

Then,

$$\lim_{n \rightarrow +\infty} \mathcal{E}(f_n) = \mathcal{E}(f^*) \quad \text{with probability } 1 \quad (45)$$

The proof of Thm. 8 is reported in Appendix B. The main technical step is to show that the estimator  $g_n$  is universally consistent. Then universal consistency of  $f_n$  follows by combining the latter result with the comparison inequality of Thm. 7. We point out that since  $g_n$  is a vector-valued least-squares estimator, the corresponding result in the case where

$\mathcal{H}$  is a finite space has been extensively studied in previous work (see e.g. (Caponnetto and De Vito, 2007)). However, to prove Thm. 8 in the general setting, we extended the work in (Caponnetto and De Vito, 2007) to the case where  $\mathcal{H}$  is infinite dimensional, which was considered an open question (Grünewälder et al., 2012).

### 5.3. Finite Sample Bounds

In order to prove finite sample bounds for structured prediction we need to impose regularity assumptions on the learning problem. This is a standard approach in learning theory (related to the *No-Free-Lunch Theorem* (Devroye et al., 2013)). In particular, we will require the target function  $g^*$  to belong to  $\mathcal{H} \otimes \mathcal{F}$ . This is a standard assumption in learning theory in the context of ridge regression (Caponnetto and De Vito, 2007; Steinwart and Christmann, 2008). In Sec. 4.4 we observed that  $g^*$  is strongly related to the concept of *conditional mean embedding* of the distribution  $\rho(\cdot|x)$  (Song et al., 2009). Therefore, by imposing it to belong to  $\mathcal{H} \otimes \mathcal{F}$  or imposing additional regularity requirements, implicitly corresponds to controlling the regularity of the data generating distribution  $\rho$ .

Below we report the learning rates of the algorithms considered in this work.

**Theorem 9 (Learning Rates)** *Let  $\mathcal{Z}$  be a compact set and  $\Delta : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$  admit an ILE with associated Hilbert space  $\mathcal{H}$ . Let  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  be a continuous reproducing kernel on  $\mathcal{X}$  with associated RKHS  $\mathcal{F}$  such that  $\kappa^2 := \sup_{x \in \mathcal{X}} k(x, x) < +\infty$ . Let  $\rho$  be a distribution on  $\mathcal{X} \times \mathcal{Y}$  and let the corresponding  $g^*$  defined in (16) be such that  $g^* \in \mathcal{H} \otimes \mathcal{F}$ . Let  $\delta \in (0, 1]$  and  $n_0$  sufficiently large such that  $n_0^{-1/2} \geq \frac{9\kappa^2}{n_0} \log \frac{n_0}{\delta}$ . Then, for any  $n \in \mathbb{N}$ , the following estimators  $f_n : \mathcal{X} \rightarrow \mathcal{Z}$  trained on  $n$  points independently sampled from  $\rho$  are such that, with probability at least  $1 - \delta$*

$$\mathcal{E}(f_n) - \mathcal{E}(f^*) \leq c_\Delta \mathfrak{m} \mathfrak{q} \log(4/\delta) n^{-1/4}, \quad (46)$$

with

$$\mathfrak{m} = 16 \left( \kappa(1 + \kappa \|g^*\|) + \kappa \sqrt{1 + \|g^*\|^2 + \|g^*\|} \right), \quad (47)$$

and  $\mathfrak{q}$  defined as follows. This holds for estimators  $f_n$  of the form (20) with corresponding weights  $\alpha$  defined as:

- (a) (Ridge Regression) in (14) with  $\lambda_n = n^{-1/2}$ . With constant  $\mathfrak{q} \leq 3$ .
- (b) (L2-Boosting) in (21) with  $\nu < 1/\kappa^2$  and  $t_n = n^{1/2}$ . With constant  $\mathfrak{q} \leq 2 + 2\gamma + e^{\gamma-1}/\gamma$ .
- (c) (Principal Component Regression) in (22) with  $\lambda_n = n^{-1/2}$ . With constant  $\mathfrak{q} \leq 5$ .

Thm. 9 is obtained as a specialization of Thm. 11 below. This result represents a direct extension of the learning rates known for binary classification (see e.g. (Yao et al., 2007)) to all structured prediction problems with ILE  $\Delta$ . This shows that up to constants, structured prediction problems are in general as challenging as classification, from the statistical perspective. See the next result for more details.

Algorithm	Train time	Train memory	Eval. time	Eval. memory
ILE + RR	$O(n^3 + n^2 c_X)$	$O(n^2)$	$O(nc_X)$	$O(n)$
ILE + L2B	$O(n^2 \sqrt{n} + n^2 c_X)$	$O(n^2)$	$O(nc_X)$	$O(n)$
ILE + PCR	$O(n^2 \sqrt{n} + n^2 c_X)$	$O(n^2)$	$O(nc_X)$	$O(n)$

Table 1: Computational complexity for training and evaluation of  $f_n$  in (20) with weights trained respectively according to (14) (RR), (21) (L2B), (22) (PCR). Hyperparameters chosen according to Thm. 9 to achieve the corresponding learning rate. The term  $c_X$  denotes cost of evaluating the kernel function.

**Remark 10 (Adaptive ILE constants)** *We comment on the constants  $c_\Delta$  and  $\mathbf{m}$  in the bound above (the constant  $\mathbf{q}$  depends exclusively on the chosen algorithm). Note that the ILE characterization of a function  $\Delta$  is not unique in terms of the space  $\mathcal{H}$  and feature maps  $\varphi, \psi$ . Moreover, as observed in Sec. 3, computing the estimator  $f_n$  does not require explicit knowledge of such objects and therefore Thm. 9 holds for any  $(\mathcal{H}, \varphi, \psi)$  such that  $\Delta$  admits an ILE and  $g^* \in \mathcal{H} \otimes \mathcal{F}$ . As a consequence, the bound in (46) implicitly applies for the infimum value of  $c_\Delta \mathbf{m}$  over the set of such triplets.*

*Explicitly estimating this constant is in general an open problem. When  $\mathcal{Z}$  and  $\mathcal{Y}$  have finite cardinality, Nowak-Vila et al. (2018) showed that for a large family of widely used loss functions, such constant is at most polylogarithmic in the cardinality of the sets.*

The result in Thm. 9 provides the suitable hyperparameters for different ILE estimators to achieve same statistical performance. Interestingly, depending on the method, this leads to different computational costs, as reported in Table 1.

#### 5.4. Refined Sample Bounds

Now we refine the analysis above considering additional regularity conditions for the learning problem. In particular we will introduce two standard assumptions in the context of non-parametric regression / conditional mean estimation (Caponnetto and De Vito, 2007). Let  $\mathcal{F}$  be the reproducing kernel Hilbert space associated to the kernel  $k$  on the input space  $\mathcal{X}$ , and  $C : \mathcal{F} \rightarrow \mathcal{F}$  be the linear operator defined as

$$\langle f, Cg \rangle_{\mathcal{F}} = \int f(x)g(x)d\rho_{\mathcal{X}}(x), \quad \forall f, g \in \mathcal{F}. \quad (48)$$

See the notation paragraph in the appendices and Appendix B for more details on the existence and properties of  $\mathcal{F}, k, C$ . Now we can introduce the first condition.

**Assumption 1 (Source condition)** *There exists  $r \geq 0$  and  $h \in \mathcal{H} \otimes \mathcal{F}$  for which*

$$g^* = (C^r \otimes I) h. \quad (49)$$

*The norm of  $\|h\|_{\mathcal{H} \otimes \mathcal{F}}$  will be denoted by  $R := \|h\|_{\mathcal{H} \otimes \mathcal{F}}$ .*

The condition above measures the regularity of  $g^*$  in terms of the eigenspectrum of  $C$ . Note that the assumption is always verified for  $r = 0$  (in that case  $h = g^*$  and  $R = \|g^*\|_{\mathcal{H} \otimes \mathcal{F}}$ ). Moreover, since  $\mathcal{F}$  is separable and  $C$  is trace class (Caponnetto and De Vito, 2007), then  $C$  can be characterized in terms of a non-increasing sequence  $(\sigma_j)_{j \in \mathbb{N}}$  of eigenvalues with associated eigenvectors  $(u_j)_{j \in \mathbb{N}}$ . For simplicity, let  $\mathcal{H} = \mathbb{R}$ . We have  $g^* = \sum_j \beta_j u_j$ , with  $\beta_j = \langle g^*, u_j \rangle_{\mathcal{F}}$ . Then, Assumption 1 is equivalent to require that  $\sum_j \beta_j^2 / \sigma_j^{2r} \leq R$ . Hence, the source condition corresponds to require  $g^*$  to have rapidly decaying coefficients, when expressed in terms of the basis of  $C$ . More generally, when  $\mathcal{H}$  is a separable Hilbert space, we have  $g^* = \sum_j \beta_j \otimes u_j$ , with  $\beta_j \in \mathcal{H}$  defined as  $\beta_j = (u_j \otimes I)g^*$ . Then, Assumption 1 is equivalent to require that  $\sum_j \|\beta_j\|^2 / \sigma_j^{2r} \leq R$ .

The second assumption is expressed with respect to the so called *effective dimension*, defined as

$$d_{\text{eff}}(\lambda) = \text{Tr}(C(C + \lambda I)^{-1}), \quad \forall \lambda > 0, \quad (50)$$

and characterizes the interaction between the measure  $\rho_{\mathcal{X}}$  and the kernel  $k$  on  $\mathcal{X}$ .

**Assumption 2 (Capacity condition)** *There exists  $\gamma \in [0, 1]$  and  $Q > 0$  for which*

$$d_{\text{eff}}(\lambda) \leq Q\lambda^{-\gamma}, \quad \forall \lambda > 0. \quad (51)$$

The condition above is always satisfied with  $\gamma = 1$  when the kernel is bounded. Indeed let  $\kappa^2 := \sup_x k(x, x)$ , then  $d_{\text{eff}}(\lambda) \leq \kappa^2 / \lambda$ . Moreover when the eigenvalues of  $C$  decay as  $\sigma_j(C) \leq Aj^{-\beta}$ , for  $A > 0$ ,  $\beta > 1$  and  $j \in \mathbb{N}$ , then the assumption above is satisfied with  $Q = A$  and  $\gamma = 1/\beta$  (Caponnetto and De Vito, 2007; Rudi et al., 2015; Blanchard and Krämer, 2016). In particular note that: (i) since  $C$  is trace class, the sequence of eigenvalues is summable therefore  $\beta > 1$ ; (ii) the eigenvalue decay is characterized by the choice of the kernel and the marginal probability distribution  $\rho_{\mathcal{X}}$ . For example, when  $\mathcal{X} = [-B, B]^d$ ,  $d \in \mathbb{N}$  for  $B > 0$ ,  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is a Sobolev-kernel of smoothness  $s > d/2$  and  $\rho_{\mathcal{X}}$  is a density bounded from above and away from zero (i.e. there exists  $A \geq a > 0$  such that  $a \leq \rho_{\mathcal{X}}(x) \leq A$  for  $x \in \mathcal{X}$ ), then there exists  $Q$  depending on  $B, s, d$  for which  $\sigma_j(C) \leq Qj^{-2s/d}$  and so  $d_{\text{eff}}(\lambda) \leq Q\lambda^{-d/(2s)}$  (see Wendland, 2004). We can now state the refined version of Thm. 9.

**Theorem 11 (Refined Learning Rates)** *Under the same notation and assumptions of Thm. 9 and under the additional Assumptions 1 and 2, let  $\delta \in (0, 1]$  and  $n_0$  sufficiently large such that  $n_0^{-1/(1+2r+\gamma)} \geq \frac{9\kappa^2}{n_0} \log \frac{n_0}{\delta}$ . For any  $n \geq n_0$ , the following estimators  $f_n : \mathcal{X} \rightarrow \mathcal{Z}$  trained on  $n$  points independently sampled from  $\rho$  are such that, with probability at least  $1 - \delta$*

$$\mathcal{E}(f_n) - \mathcal{E}(f^*) \leq c_{\Delta} \mathfrak{m} \mathfrak{q} \log(4/\delta) n^{-\frac{r+1/2}{2r+\gamma+1}}, \quad (52)$$

with

$$\mathfrak{m} = 16 \left( \kappa(1 + \kappa R) + \kappa \sqrt{Q + R^2 + R} \right), \quad (53)$$

and  $\mathfrak{q}$  defined as follows. This holds for estimators  $f_n$  of the form (20) with corresponding weights  $\alpha$  defined as:

Algorithm	Train time	Train memory	Eval. time	Eval. memory
ILE + RR	$O(n^3 + n^2 c_X)$	$O(n^2)$	$O(nc_X)$	$O(n)$
ILE + L2B	$O(n^{2+\frac{1}{2r+\gamma+1}} + n^2 c_X)$	$O(n^2)$	$O(nc_X)$	$O(n)$
ILE + PCR	$O(n^{2+\frac{\gamma}{2r+\gamma+1}} + n^2 c_X)$	$O(n^2)$	$O(nc_X)$	$O(n)$

Table 2: Computational complexity for training and evaluation of  $f_n$  in (20) with weights trained respectively according to (14) (RR), (21) (L2B), (22) (PCR). Hyperparameters chosen according to Thm. 11 to achieve the corresponding learning rate. The term  $c_X$  denotes cost of evaluating the kernel function.

(a) (Ridge Regression) in (14) with  $\lambda_n = n^{-\frac{1}{2r+\gamma+1}}$ . With  $\mathfrak{q} \leq 3$ .

(b) (L2-Boosting) in (21) with  $\nu < 1/\kappa^2$  and  $t_n = n^{\frac{1}{2r+\gamma+1}}$ . With  $\mathfrak{q} \leq 2 + 2\nu + e^{\nu-1}/\nu$ .

(c) (Principal Component Regression) in (22) with  $\lambda_n = n^{-\frac{1}{2r+\gamma+1}}$ . With  $\mathfrak{q} \leq 5$ .

The theorem above shows that the proposed estimator for structured prediction in (20) has learning rates that are adaptive to the source and capacity condition, when the coefficients are computed according to the algorithms considered in the theorem. Similarly to (45) and Thm. 9, the result is obtained by studying the generalization properties of  $g_n$  and combining such analysis with the comparison inequality. As already pointed out in the commentary of (45), our results in this section generalize those of (Caponnetto and De Vito, 2007) to the case of infinite dimensional output spaces  $\mathcal{H}$ . Interestingly, the result in Thm. 11 refine Thm. 9. In particular, in the worst case  $r = 0, \gamma = 1$  we recover the learning rate  $O(n^{-1/4})$  in Thm. 9, while for stronger regularity assumptions (namely  $r \gg 1$  or  $\gamma \approx 0$ ) the proposed algorithms attain a significantly faster rate close to  $O(n^{-1/2})$ .

When  $\mathcal{Y}$  and  $\mathcal{Z}$  have finite cardinality and the data distribution satisfies additional regularity hypotheses, such as the Tsybakov condition (see Tsybakov et al., 2004; Yao et al., 2007), it is possible to achieve rates of up to  $O(n^{-1})$ , as shown in (Nowak-Vila et al., 2018). A relevant question is whether an analogous notion of the Tsybakov condition could be identified in the case where  $\mathcal{Y}$  and  $\mathcal{Z}$  are not finite. Note that the  $O(n^{-1})$  rate is optimal in the case of binary classification (Bartlett et al., 2006; Tsybakov et al., 2004). This implies that such rate is optimal also for the larger family of structured prediction problems satisfying the ILE assumption. In this sense, a natural question is whether it may be possible to perform a more refined analysis by studying specific structured prediction problems individually.

To conclude, note that the algorithms considered in this work are not only adaptive from the statistical viewpoint but also from a computational perspective. In particular, Table 2 reports the computational costs of running the algorithms described in this work for the choice of hyperparameters reported by Thm. 11 depending on the Assumptions 1 and 2.



## 6. Sufficient Conditions for ILE

In this section we focus our attention to the definition of *Implicit Loss Embedding* (ILE) introduced in Thm. 2. In particular, we provide a number of sufficient conditions that guarantee a loss function to admit an ILE, which are more interpretable and easy to verify than the original definition. We will show that most loss functions used in machine learning and structured prediction settings indeed satisfy the ILE property and therefore that the learning framework proposed in this work applies to a wide family of relevant problems.

**Bounding  $c_\Delta$ .** As a byproduct of our analysis, the results in the following provide also upper bounds for the constant  $c_\Delta$  for a number of loss functions. As observed in Thm. 9 and 11, such constant plays a role in characterizing the learning rates of the ILE estimators. Following the discussion of Thm. 10, it is important to note that the estimates for  $c_\Delta$  reported in this section have been derived for a single parametrization of the ILE definition for  $\Delta$  (namely the space  $\mathcal{H}$  and the feature maps  $\psi$  and  $\varphi$ ). Obtaining sharp bounds for such constants is outside the scope of this work. We refer to (Osokin et al., 2017; Nowak-Vila et al., 2018) for refined analysis in the case where  $\mathcal{Z}$  and  $\mathcal{Y}$  are finite.

### 6.1. ILE on finite Output or Label Spaces

In Sec. 3.1 we provided a preliminary analysis of structured prediction for the case where label and output spaces coincide and are finite, namely  $\mathcal{Y} = \mathcal{Z} = \{1, \dots, T\}$ . This discussion was key in that it motivated the definition of ILE. Indeed, as already mentioned, the ILE definition is satisfied by any loss function acting on finite output and label spaces. The following proposition shows that it is sufficient that only one of the two spaces  $\mathcal{Y}$  or  $\mathcal{Z}$  is finite to guarantee the loss function to admit an ILE.

**Theorem 12 (ILE & finite  $\mathcal{Y}$  or  $\mathcal{Z}$ )** *The function  $\Delta : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$  admits an ILE if one of the following conditions hold:*

- (a)  $\mathcal{Z}$  and  $\mathcal{Y}$  are finite sets. In this case  $c_\Delta \leq \|\Delta\|$  the operator norm of the matrix  $\Delta \in \mathbb{R}^{|\mathcal{Z}| \times |\mathcal{Y}|}$  with entries  $\Delta_{z,y} = \Delta(z, y)$ .
- (b)  $\mathcal{Z}$  is finite,  $\mathcal{Y}$  is compact and  $\Delta(z, \cdot)$  is continuous on  $\mathcal{Y}$  for any  $z \in \mathcal{Z}$ .  
In this case  $c_\Delta \leq \sup_{y \in \mathcal{Y}} \sqrt{\sum_{z \in \mathcal{Z}} |\Delta(z, y)|^2}$ .
- (c)  $\mathcal{Z}$  is compact,  $\mathcal{Y}$  is finite and  $\Delta(\cdot, y)$  is continuous on  $\mathcal{Z}$  for any  $y \in \mathcal{Y}$ .  
In this case  $c_\Delta \leq \sup_{z \in \mathcal{Z}} \sqrt{\sum_{y \in \mathcal{Y}} |\Delta(z, y)|^2}$ .

The result above shows that most loss functions used in typical structured prediction applications admit an ILE. Indeed, previous literature on the topic has been focused on problems where either the output or the label space (or both) are finite, albeit possibly very large (Bakir et al., 2007; Nowozin et al., 2011). In this setting, relevant examples of applications range from computer vision, such as segmentation (Alahari et al., 2008), localization (Blaschko and Lampert, 2008; Lampert et al., 2009), labeling (Karpathy and Fei-Fei, 2015), pixel-wise classification (Szummer et al., 2008), speech recognition (Bahl et al., 1986; Sutton et al., 2012), natural language processing (Tsochantaridis et al., 2005), trajectory planing (Ratliff et al., 2006) or hierarchical classification (Tuia et al., 2011).

The major implication of Thm. 12 is that it justifies the application of the estimator proposed and studied in this paper to address a variety of structured prediction problems previously considered in the literature. Indeed, our analysis in Sec. 5 automatically guarantees that the corresponding estimator has strong theoretical guarantees when applied to these settings.

In the rest of this section we focus on the case where  $\mathcal{Y}$  and  $\mathcal{Z}$  are not necessarily finite, showing that the definition of ILE encompasses a significantly wider family of settings compared to the classic structured prediction literature.

## 6.2. ILE and Reproducing Kernel Hilbert Spaces

We already highlighted the relation between the definition of ILE and the notion of positive definite kernel. The following result provides a more refined characterization of this relation, showing in particular how it is possible to leverage kernels to “build” ILE functions.

**Theorem 13 (ILE & RKHS)** *Let  $\mathcal{Z} = \mathcal{Y}$  be a compact set and  $h : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  a continuous bounded reproducing kernel on  $\mathcal{Y}$  with associated RKHS  $\mathcal{H}$ . Let  $\eta^2 = \sup_{y \in \mathcal{Y}} h(y, y)$ . Then,  $\Delta : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  admits an ILE if one of the following holds:*

- (a) (Kernels).  $\Delta(z, y) = h(z, y)$  for any  $y, z \in \mathcal{Y}$ . In this case  $c_\Delta \leq \eta^2$ .
- (b) (Kernel Dependency Estimation (KDE)).  $\Delta(z, y) = h(z, z) + h(y, y) - 2h(z, y)$  for any  $y, z \in \mathcal{Y}$ . In this case  $c_\Delta = 2(2\eta^4 + 1)$ .
- (c) For every  $y \in \mathcal{Y}$  the functions  $\Delta(\cdot, y) \in \mathcal{H}$  belong to a bounded set of  $\mathcal{H}$ , namely  $\sup_{y \in \mathcal{Y}} \|\Delta(\cdot, y)\|_{\mathcal{H}} = D < +\infty$ . In this case  $c_\Delta \leq \eta D$ . The same holds if the family of functions  $\Delta(z, \cdot)$  parametrized by  $z \in \mathcal{Z}$  belong to a bounded set of  $\mathcal{H}$ .
- (d)  $\Delta$  belongs to  $\mathcal{H} \otimes \mathcal{H}$  the RKHS with associated kernel  $\bar{h} : \mathcal{Y}^2 \times \mathcal{Y}^2 \rightarrow \mathbb{R}$  such that  $\bar{h}((z, y), (z', y')) = h(z, z')h(y, y')$  for any  $z, z', y, y' \in \mathcal{Y}$ . In this case  $c_\Delta \leq \eta^2 \|\Delta\|_{\mathcal{H} \otimes \mathcal{H}}$

Thm. 13 provides four interesting results. First, as already mentioned, the definition of ILE function recovers and is more general than that of positive definite kernel. Second, we see that our framework recovers the *Kernel Dependency Estimation (KDE)* approach (Weston et al., 2002; Cortes et al., 2005), which corresponds to a structured prediction setting with loss function  $\Delta(z, y) = \|h(z, \cdot) - h(y, \cdot)\|_{\mathcal{H}}^2 = h(z, z) + h(y, y) - 2h(z, y)$ .

Point (c) shows that  $\Delta$  admits an ILE if the family of functions  $\{\Delta(\cdot, y)\}_{y \in \mathcal{Y}}$  parametrized by  $y \in \mathcal{Y}$ , is uniformly contained in a ball in  $\mathcal{H}$ . Finally, point (d) of Thm. 13 reports a more general result, showing that *all functions that belong to the RKHS obtained as the tensor product of  $\mathcal{H}$  with itself admit an ILE*. This recovers a large family of loss functions as discussed in the example below.

**Example 5 (Smooth Functions on  $\mathcal{Y} \times \mathcal{Y}$  with  $\mathcal{Y} = [-B, B]^d$  admit an ILE)** *Let  $\Delta \in C^\infty(\mathcal{Y} \times \mathcal{Y})$ , where  $C^\infty(\mathcal{Y})$  denotes the space of smooth functions over  $\mathcal{Y}$ . Let  $\mathcal{H} = W^{d,2}(\mathcal{Y})$  be the Sobolev space of functions over  $\mathcal{Y}$  with up to order  $d$  square integrable weak derivatives (Adams and Fournier, 2003). We have  $C^\infty(\mathcal{Y} \times \mathcal{Y}) = C^\infty(\mathcal{Y}) \otimes C^\infty(\mathcal{Y}) \subset \mathcal{H} \otimes \mathcal{H}$ . Then, Thm. 13 (d) guarantees that  $\Delta$  admits an ILE. For more details see (Lwise et al., 2018).*

### 6.3. ILE and Regularity

The connection between ILE and RKHSs suggest the definition of ILE to be somewhat related to the concept of smoothness or regularity of a function. The following result goes beyond RKHSs and investigates this question in further detail.

**Theorem 14 (ILE & Regularity)** *Let  $\mathcal{Z} = \mathcal{Y} = [-B, B]^d$ ,  $B > 0$ . A function  $\Delta : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  admits an ILE when at least one of the following conditions hold:*

- (a)  $d = 1$  and  $\Delta$  is  $\alpha$ -Hölder continuous with  $\alpha > 1/2$  or it is of bounded variation and  $\alpha$ -Hölder continuous with  $\alpha > 0$ .
- (b)  $\Delta(z, y) = v(z - y)$ , where  $v$  is a function such that  $c_\Delta = \int |\hat{v}(\omega)| d\omega < \infty$  and  $\hat{v}$  is the Fourier transform of  $v$ .
- (c) The mixed partial derivative  $\Delta_{y_1, \dots, y_d} : \mathcal{Y} \rightarrow \mathbb{R}$  of  $\Delta$  exists almost everywhere and  $\Delta_{y_1, \dots, y_d} \in L^p(\mathcal{Y})$  with  $p > 1$ .

Thm. 14 shows that any function that is sufficiently regular admits an ILE. This allows to recover most loss functions used in machine learning and robust estimation as special cases.

**Example 6 (Robust Estimation)** *We have already observed that smooth functions such as the least-squares and logistic loss admit an ILE according to the discussion in Example 5. Here we observe that also the hinge loss, used in binary classification, and most loss functions used for scalar regression on  $\mathcal{Y} = [0, 1]$  admit an ILE, albeit being not smooth. Indeed, most common loss functions used in these contexts are Lipschitz continuous and therefore satisfy Thm. 14 (a). Notable examples are loss functions used for robust estimation such as the absolute value, Huber, Cauchy, German-McLure, “Fair” and  $L_2 - L_1$  (Huber and Ronchetti, 2011). All these functions are differentiable almost everywhere, with uniformly bounded derivatives and thus satisfy Thm. 14 (c).*

### 6.4. Composition Rules for ILE

A natural question is whether some operations over ILE functions preserve the characterization introduced in Thm. 2. Below we provide a set of rules that allow to “build” new ILE functions from known ones.

**Theorem 15** *Let  $\mathcal{Z}$  and  $\mathcal{Y}$  be compact sets. Then  $\Delta : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$  admits an ILE if one of the following holds:*

- (a) (Restriction) *There exist two sets  $\bar{\mathcal{Z}} \supseteq \mathcal{Z}$ ,  $\bar{\mathcal{Y}} \supseteq \mathcal{Y}$  and  $\bar{\Delta} : \bar{\mathcal{Z}} \times \bar{\mathcal{Y}} \rightarrow \mathbb{R}$  such that  $\bar{\Delta}$  admits an ILE and its restriction to  $\mathcal{Z} \times \mathcal{Y}$  corresponds to  $\Delta$ , namely*

$$\Delta = \bar{\Delta}|_{\mathcal{Z} \times \mathcal{Y}}. \quad (54)$$

*In this case  $c_\Delta \leq c_{\bar{\Delta}}$ .*

- (b) (Right Composition) *There exists  $\bar{\mathcal{Z}}, \bar{\mathcal{Y}}$  and a ILE  $\bar{\Delta} : \bar{\mathcal{Z}} \times \bar{\mathcal{Y}} \rightarrow \mathbb{R}$ , such that*

$$\Delta(z, y) = \alpha(z) \bar{\Delta}(A(z), B(y)) \beta(y), \quad (55)$$

*with  $A : \mathcal{Z} \rightarrow \bar{\mathcal{Z}}$ ,  $B : \mathcal{Y} \rightarrow \bar{\mathcal{Y}}$ ,  $\alpha : \mathcal{Z} \rightarrow \mathbb{R}$  and  $\beta : \mathcal{Y} \rightarrow \mathbb{R}$  continuous function, with  $\sup_{z \in \mathcal{Z}} |\alpha(z)| \leq \bar{\alpha}$  and  $\sup_{y \in \mathcal{Y}} |\beta(y)| \leq \bar{\beta}$  with  $\bar{\alpha}, \bar{\beta} \in \mathbb{R}$ . Then  $c_\Delta \leq \bar{\alpha} \bar{\beta} c_{\bar{\Delta}}$ .*

(c) (Left Composition) *There exist  $P \in \mathbb{N}$ , spaces  $(\mathcal{Z}_p)_{p=1}^P, (\mathcal{Y}_p)_{p=1}^P$  and corresponding ILE  $\Delta_p : \mathcal{Z}_p \times \mathcal{Y}_p \rightarrow \mathbb{R}$  such that  $\mathcal{Z} = \mathcal{Z}_1 \times \cdots \times \mathcal{Z}_P, \mathcal{Y} = \mathcal{Y}_1 \times \cdots \times \mathcal{Y}_P$  and*

$$\Delta(z, y) = \Gamma(\Delta_1(z_1, y_1), \dots, \Delta_P(z_P, y_P)), \quad (56)$$

*for any  $z = (z_1, \dots, z_P) \in \mathcal{Z}$  and  $y = (y_1, \dots, y_P) \in \mathcal{Y}$ , where  $\Gamma : \mathbb{R}^P \rightarrow \mathbb{R}$  is an analytic function (e.g. a polynomial).*

The result above provides us several tools to build new ILE functions. In particular, Thm. 15 (a) shows that we can always restrict a ILE function on a smaller pair of output-label sets and still enjoy the same properties of the original loss, also in terms of universal consistency and rates of the resulting structured prediction estimator. Thm. 15 (b) allows to extend a ILE function  $\bar{\Delta}$  to other output-label pairs by means of the embeddings  $A$  and  $B$  and the weighting functions  $\alpha$  and  $\beta$  in (55).

**Example 7 (Restriction of Smooth functions on compact sets admit an ILE)** *Let  $\Delta$  be a smooth function over  $[-B, B]^d$  with  $B > 0$ . Then, by Thm. 15 (a),  $\Delta$  admits an ILE on every compact set  $\mathcal{Y} \subseteq [-B, B]^d$ .*

Finally, Thm. 15 (c) shows that any combination (namely sum and multiplications) of ILE functions is still ILE. To highlight the importance of this result we clarify it in the following.

**Corollary 16** *Let  $\Delta_1 : \mathcal{Z}_1 \times \mathcal{Y}_1 \rightarrow \mathbb{R}$  and  $\Delta_2 : \mathcal{Z}_2 \times \mathcal{Y}_2 \rightarrow \mathbb{R}$  admit an ILE. Then  $\Delta : (\mathcal{Z}_1 \times \mathcal{Z}_2) \times (\mathcal{Y}_1 \times \mathcal{Y}_2) \rightarrow \mathbb{R}$  if, for any  $z_i \in \mathcal{Z}_i, y_i \in \mathcal{Y}_i$  and  $i = 1, 2$ , one of the following conditions hold:*

(a)  $\Delta((z_1, z_2), (y_1, y_2)) = \Delta_1(z_1, y_1) + \Delta_2(z_2, y_2),$

(b)  $\Delta((z_1, z_2), (y_1, y_2)) = \Delta_1(z_1, y_1) \Delta_2(z_2, y_2).$

The result above allows to consider general combinations of loss functions within the framework considered in this work. In particular, the following remark shows how multitask learning problems (possibly with structure on the output) can be recovered in this setting.

**Example 8 (Multitask Learning)** *In multitask learning (MTL) settings the goal is to solve multiple separate supervised problems simultaneously (Evgeniou and Pontil, 2004; Alvarez et al., 2012). The loss functions used in MTL typically consist in the sum of “single task” loss functions over the separate tasks, such as least-squares for regression or logistic/hinge for classification. Since according to Thm. 16 the sum of ILE functions is still ILE, we see that multitask learning is naturally recovered by the framework considered in this work. This fact was observed in (Ciliberto et al., 2017), where the structured prediction perspective on the MTL problem allowed to address the question of how to impose non-linear relations among multiple tasks by introducing the constraint output set  $\mathcal{Z} \subset \mathcal{Y} = \mathbb{R}^T$ .*

## 7. Conclusions

In this work we have presented a general framework for structured prediction. Our work revolves around the key notion of Implicit Loss Embedding (ILE), which allows us to study structured prediction applications where the output space is not finite, differently from most previous work on the topic. This work significantly expanded upon Ciliberto et al. (2016), providing novel insights on the ILE property as well as new algorithms for structured prediction and their corresponding theoretical analysis. Among the main contributions of this work: (a) we showed that the proposed framework can be applied to a wide range of structured prediction problems, providing a systematic approach to derive estimators with strong theoretical guarantees. In particular, we showed that it is possible to leverage existing algorithms from the vector-valued regression literature to obtain novel structured prediction estimators that enjoy equivalent statistical properties of the original method, but with reduced computational requirements. (b) We performed a refined analysis of the excess risk bounds, showing that the statistical rates and computational cost of the considered algorithms are adaptive to standard regularity properties of the learning problem. (c) We provided a number of sufficient conditions to verify whether a given loss admits an ILE. These conditions are significantly easier to verify in practice in comparison to the general definition. Leveraging these conditions we proved that most loss functions used in machine learning indeed admit an ILE and are therefore suited to our framework.

Relevant directions for future work will involve: (a) considering alternative estimators within the ILE framework not necessarily minimizing the square loss in the surrogate space; (b) learning the structure of the output space when it is not fully known a-priori (for instance in manifold regression settings where the output manifold is only accessible via examples). This could be addressed by parametrizing a family of candidate output spaces and finding the optimal parameters while simultaneously fitting the structured prediction model. Finally, (c) an interesting question is to leverage further additional knowledge on the problem structure to improve the overall learning rates of the estimator. This direction has been recently preliminarily investigated in Cortes et al. (2016); Ciliberto et al. (2019), where an explicit factorization of the loss function was used to design problem-specific algorithms and perform a refined analysis of their generalization properties.

## Acknowledgments

We thank Florence D’Alché-Buc and Celine Brouard for their helpful comments and feedback throughout the last revision of this manuscript.

C.C. acknowledges the Royal Society (grant SPREM RGS\R1\201149). L.R. acknowledges the support of the European Research Council (grant SLING 819789), the AFOSR projects FA9550-17-1-0390 and the EU H2020-MSCA-RISE project NoMADS - DLV-777826. A.R. acknowledges the support of the European Research Council (grant SEQUOIA 724063) and the Agence Nationale de la Recherche, French Government, as part of the “Investissements d’avenir” program, reference ANR-19-P3IA-0001 (PRAIRIE 3IA Institute).

## References

- Robert A Adams and John JF Fournier. *Sobolev spaces*, volume 140. Elsevier, 2003.
- Karteek Alahari, Pushmeet Kohli, and Philip HS Torr. Reduce, reuse & recycle: Efficiently solving multi-label mrfs. In *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on*, pages 1–8. IEEE, 2008.
- Charalambos D Aliprantis and Kim Border. *Infinite dimensional analysis: a hitchhiker’s guide*. Springer Science & Business Media, 2006.
- Mauricio A Alvarez, Lorenzo Rosasco, Neil D Lawrence, et al. Kernels for vector-valued functions: A review. *Foundations and Trends in Machine Learning*, 4(3):195–266, 2012.
- Nachman Aronszajn. Theory of reproducing kernels. *Transactions of the American mathematical society*, 68(3):337–404, 1950.
- Lalit Bahl, Peter Brown, Peter De Souza, and Robert Mercer. Maximum mutual information estimation of hidden markov model parameters for speech recognition. In *Acoustics, Speech, and Signal Processing, IEEE International Conference on ICASSP’86.*, volume 11, pages 49–52. IEEE, 1986.
- Gökhan Bakir, Thomas Hofmann, Bernhard Schölkopf, Alexander J. Smola, Ben Taskar, and S.V.N Vishwanathan. *Predicting structured data*. MIT press, 2007.
- Peter L Bartlett, Michael I Jordan, and Jon D McAuliffe. Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101(473):138–156, 2006.
- Frank Bauer, Sergei Pereverzev, and Lorenzo Rosasco. On regularization algorithms in learning theory. *Journal of complexity*, 23(1):52–72, 2007.
- Jean-David Benamou, Guillaume Carlier, Marco Cuturi, Luca Nenna, and Gabriel Peyré. Iterative bregman projections for regularized transportation problems. *SIAM Journal on Scientific Computing*, 37(2):A1111–A1138, 2015.
- Alain Berlinet and Christine Thomas-Agnan. *Reproducing kernel Hilbert spaces in probability and statistics*. Springer Science & Business Media, 2011.
- Gilles Blanchard and Nicole Krämer. Convergence rates of kernel conjugate gradient for random design regression. *Analysis and Applications*, 14(06):763–794, 2016.
- Matthew B Blaschko and Christoph H Lampert. Learning to localize objects with structured output regression. In *European conference on computer vision*, pages 2–15. Springer, 2008.
- Mathieu Blondel. Structured prediction with projection oracles. In *Advances in Neural Information Processing Systems*, pages 12145–12156, 2019.
- Céline Brouard, Florence d’Alché Buc, and Marie Szafranski. Semi-supervised penalized output kernel regression for link prediction. In *International Conference on Machine Learning (ICML)*, 2011.

- Céline Brouard, Marie Szafranski, and Florence d’Alché Buc. Input output kernel regression: Supervised and semi-supervised structured output prediction with operator-valued kernels. *The Journal of Machine Learning Research*, 17(1):6105–6152, 2016.
- Andrea Caponnetto and Ernesto De Vito. Optimal rates for the regularized least-squares algorithm. *Foundations of Computational Mathematics*, 7(3):331–368, 2007.
- Luigi Carratino, Alessandro Rudi, and Lorenzo Rosasco. Learning with sgd and random features. In *Advances in Neural Information Processing Systems*, pages 10192–10203, 2018.
- Charles Castaing and Michel Valadier. *Convex analysis and measurable multifunctions*, volume 580. Springer, 2006.
- Carlo Ciliberto, Lorenzo Rosasco, and Alessandro Rudi. A consistent regularization approach for structured prediction. *Advances in Neural Information Processing Systems 29 (NIPS)*, pages 4412–4420, 2016.
- Carlo Ciliberto, Alessandro Rudi, Lorenzo Rosasco, and Massimiliano Pontil. Consistent multitask learning with nonlinear output relations. In *Advances in Neural Information Processing Systems*, pages 1983–1993, 2017.
- Carlo Ciliberto, Francis Bach, and Alessandro Rudi. Localized structured prediction. *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- Corinna Cortes, Mehryar Mohri, and Jason Weston. A general regression technique for learning transductions. In *Proceedings of the 22nd international conference on Machine learning*, pages 153–160. ACM, 2005.
- Corinna Cortes, Vitaly Kuznetsov, Mehryar Mohri, and Scott Yang. Structured prediction theory based on factor graph complexity. In *Advances in Neural Information Processing Systems*, pages 2514–2522, 2016.
- Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in Neural Information Processing Systems*, pages 2292–2300, 2013.
- Marco Cuturi and Mathieu Blondel. Soft-dtw: a differentiable loss function for time-series. *International Conference on Machine Learning*, 2017.
- Marco Cuturi and Arnaud Doucet. Fast computation of wasserstein barycenters. 2014.
- Luc Devroye, László Györfi, and Gábor Lugosi. *A probabilistic theory of pattern recognition*, volume 31. Springer Science & Business Media, 2013.
- Moussab Djerrab, Alexandre Garcia, Maxime Sangnier, and Florence d’Alché Buc. Output fisher embedding regression. *Machine Learning*, 107(8-10):1229–1256, 2018.
- John C Duchi, Lester W Mackey, and Michael I Jordan. On the consistency of ranking algorithms. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, pages 327–334, 2010.

- Richard M Dudley. *Real analysis and probability*, volume 74. Cambridge University Press, 2002.
- Heinz Werner Engl, Martin Hanke, and Andreas Neubauer. *Regularization of inverse problems*, volume 375. Springer Science & Business Media, 1996.
- Theodoros Evgeniou and Massimiliano Pontil. Regularized multi-task learning. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 109–117. ACM, 2004.
- Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The elements of statistical learning*, volume 10. Springer series in statistics New York, NY, USA:, 2001.
- Charlie Frogner, Chiyuan Zhang, Hossein Mobahi, Mauricio Araya, and Tomaso A Poggio. Learning with a wasserstein loss. In *Advances in Neural Information Processing Systems*, pages 2053–2061, 2015.
- Wei Gao and Zhi-Hua Zhou. On the consistency of multi-label learning. *Artificial Intelligence*, 199:22–44, 2013.
- Pierre Geurts, Louis Wehenkel, and Florence d’Alché Buc. Kernelizing the output of tree-based methods. In *Proceedings of the 23rd international conference on Machine learning*, pages 345–352, 2006.
- Pierre Geurts, Louis Wehenkel, and Florence d’Alché Buc. Gradient boosting for kernelized output spaces. In *Proceedings of the 24th international conference on Machine learning*, pages 289–296, 2007.
- Steffen Grünewälder, Guy Lever, Luca Baldassarre, Sam Patterson, Arthur Gretton, and Massimiliano Pontil. Conditional mean embeddings as regressors. In *International Conference on Machine Learning (ICML)*, volume 5, 2012.
- László Györfi, Michael Kohler, Adam Krzyzak, and Harro Walk. *A distribution-free theory of nonparametric regression*. Springer Science & Business Media, 2006.
- Peter J Huber and Elvezio M Ronchetti. *Robust statistics*. Springer, 2011.
- Thorsten Joachims, Thomas Hofmann, Yisong Yue, and Chun-Nam Yu. Predicting structured objects with support vector machines. *Communications of the ACM*, 52(11):97–104, 2009.
- Hachem Kadri, Mohammad Ghavamzadeh, and Philippe Preux. A generalized kernel approach to structured output learning. In *International Conference on Machine Learning*, pages 471–479, 2013.
- J-P Kahane. *Fourier series and wavelets*. Routledge, 1995.
- Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3128–3137, 2015.



- Anna Korba, Alexandre Garcia, and Florence d’Alché Buc. A structured prediction approach for label ranking. In *Advances in Neural Information Processing Systems*, pages 8994–9004, 2018.
- John Lafferty, Andrew McCallum, and Fernando CN Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of the 18th International Conference on Machine Learning*, 2001.
- Pierre Laforgue, Alex Lambert, Luc Motte, and Florence d’Alché Buc. On the dualization of operator-valued kernel machines. *arXiv preprint arXiv:1910.04621*, 2019.
- Christoph H Lampert, Matthew B Blaschko, and Thomas Hofmann. Efficient subwindow search: A branch and bound framework for object localization. *IEEE transactions on pattern analysis and machine intelligence*, 31(12):2129–2142, 2009.
- Giulia Luise, Alessandro Rudi, Massimiliano Pontil, and Carlo Ciliberto. Differential properties of sinkhorn approximation for learning with wasserstein distance. In *Advances in Neural Information Processing Systems*, pages 5859–5870, 2018.
- Giulia Luise, Dimitris Stamos, Massimiliano Pontil, and Carlo Ciliberto. Leveraging low-rank relations between surrogate tasks in structured prediction. *International Conference on Machine Learning (ICML)*, 2019.
- Gian Maria Marconi, Lorenzo Rosasco, and Carlo Ciliberto. Hyperbolic manifold regression. *Artificial Intelligence and Statistics - AISTATS*, 2020.
- Arthur Mensch, Mathieu Blondel, and Gabriel Peyré. Geometric losses for distributional learning. In *International Conference on Machine Learning (ICML)*, 2019.
- Charles A Micchelli and Massimiliano Pontil. Kernels for multi-task learning. In *Advances in Neural Information Processing Systems*, pages 921–928, 2004.
- Charles A Micchelli, Yuesheng Xu, and Haizhang Zhang. Universal kernels. *The Journal of Machine Learning Research*, 7:2651–2667, 2006.
- Louis-Philippe Morency, Ariadna Quattoni, and Trevor Darrell. Latent-dynamic discriminative models for continuous gesture recognition. In *Computer Vision and Pattern Recognition, 2007. CVPR’07. IEEE Conference on*, pages 1–8. IEEE, 2007.
- Ferenc Móricz and Antal Veres. On the absolute convergence of multiple fourier series. *Acta Mathematica Hungarica*, 117(3):275–292, 2007.
- Youssef Mroueh, Tomaso Poggio, Lorenzo Rosasco, and Jean-Jacques Slotine. Multiclass learning with simplex coding. In *Advances in Neural Information Processing Systems (NIPS) 25*, pages 2798–2806, 2012.
- Krikamol Muandet, Kenji Fukumizu, Bharath Sriperumbudur, Bernhard Schölkopf, et al. Kernel mean embedding of distributions: A review and beyond. *Foundations and Trends® in Machine Learning*, 10(1-2):1–141, 2017.

- Elizbar A Nadaraya. On estimating regression. *Theory of Probability and Its Applications*, 9 (1):141–142, 1964.
- Alex Nowak-Vila, Francis Bach, and Alessandro Rudi. Sharp analysis of learning with discrete losses. *Artificial Intelligence and Statistics - AISTATS*, 2018.
- Sebastian Nowozin, Christoph H Lampert, et al. Structured learning and prediction in computer vision. *Foundations and Trends® in Computer Graphics and Vision*, 6(3–4): 185–365, 2011.
- Anton Osokin, Francis Bach, and Simon Lacoste-Julien. On structured prediction theory with calibrated convex surrogate losses. In *Advances in Neural Information Processing Systems*, pages 302–313, 2017.
- Fabian Pedregosa, Francis Bach, and Alexandre Gramfort. On the consistency of ordinal regression methods. *The Journal of Machine Learning Research*, 18(1):1769–1803, 2017.
- Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances in neural information processing systems*, pages 1177–1184, 2008.
- Nathan D Ratliff, J Andrew Bagnell, and Martin A Zinkevich. Maximum margin planning. In *Proceedings of the 23rd international conference on Machine learning*, pages 729–736. ACM, 2006.
- Lorenzo Rosasco, Ernesto De Vito, and Alessandro Verri. Spectral methods for regularization in learning theory. *DISI, Universita degli Studi di Genova, Italy, Technical Report DISI-TR-05-18*, 2005.
- Alessandro Rudi and Lorenzo Rosasco. Generalization properties of learning with random features. In *Advances in Neural Information Processing Systems*, pages 3215–3225, 2017.
- Alessandro Rudi, Raffaello Camoriano, and Lorenzo Rosasco. Less is more: Nyström computational regularization. In *Advances in Neural Information Processing Systems*, pages 1657–1665, 2015.
- Alessandro Rudi, Carlo Ciliberto, GianMaria Marconi, and Lorenzo Rosasco. Manifold structured prediction. In *Advances in Neural Information Processing Systems*, pages 5610–5621, 2018.
- Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
- John Shawe-Taylor and Nello Cristianini. *Kernel methods for pattern analysis*. Cambridge university press, 2004.
- Rahul Singh, Maneesh Sahani, and Arthur Gretton. Kernel instrumental variable regression. *Advances in Neural Information Processing Systems*, 2019.
- Alex J Smola and Bernhard Schölkopf. Sparse greedy matrix approximation for machine learning. *17th International Conference on Machine Learning*, 2000.

- Le Song, Jonathan Huang, Alex Smola, and Kenji Fukumizu. Hilbert space embeddings of conditional distributions with applications to dynamical systems. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 961–968. ACM, 2009.
- Suvrit Sra and Reshad Hosseini. Geometric optimization in machine learning. In *Algorithmic Advances in Riemannian Geometry and Applications*, pages 73–91. Springer, 2016.
- Bharath K Sriperumbudur, Kenji Fukumizu, and Gert RG Lanckriet. Universality, characteristic kernels and rkhs embedding of measures. *Journal of Machine Learning Research*, 12(Jul):2389–2410, 2011.
- Florian Steinke, Matthias Hein, and Bernhard Schölkopf. Nonparametric regression between general riemannian manifolds. *SIAM Journal on Imaging Sciences*, 3(3):527–563, 2010.
- Ingo Steinwart and Andreas Christmann. *Support Vector Machines*. Information Science and Statistics. Springer New York, 2008.
- Ingo Steinwart, Andreas Christmann, et al. Estimating conditional quantiles with the help of the pinball loss. *Bernoulli*, 17(1):211–225, 2011.
- Kirill Struminsky, Simon Lacoste-Julien, and Anton Osokin. Quantifying learning guarantees for convex but inconsistent surrogates. In *Advances in Neural Information Processing Systems*, pages 669–677, 2018.
- Charles Sutton, Andrew McCallum, et al. An introduction to conditional random fields. *Foundations and Trends® in Machine Learning*, 4(4):267–373, 2012.
- Martin Szummer, Pushmeet Kohli, and Derek Hoiem. Learning crfs using graph cuts. In *European conference on computer vision*, pages 582–595. Springer, 2008.
- Ben Taskar, Carlos Guestrin, and Daphne Koller. Max-margin markov networks. In *Advances in neural information processing systems*, pages 25–32, 2004.
- Ambuj Tewari and Peter L Bartlett. On the consistency of multiclass classification methods. *Journal of Machine Learning Research*, 8(May):1007–1025, 2007.
- Ioannis Tsochantaridis, Thorsten Joachims, Thomas Hofmann, and Yasemin Altun. Large margin methods for structured and interdependent output variables. In *Journal of Machine Learning Research*, pages 1453–1484, 2005.
- Alexander B Tsybakov et al. Optimal aggregation of classifiers in statistical learning. *The Annals of Statistics*, 32(1):135–166, 2004.
- Devis Tuia, Jordi Muñoz-Marí, Mikhail Kanevski, and Gustavo Camps-Valls. Structured output svm for remote sensing image classification. *Journal of signal processing systems*, 65(3):301–310, 2011.
- SVN Vishwanathan, Nicol N Schraudolph, Mark W Schmidt, and Kevin P Murphy. Accelerated training of conditional random fields with stochastic gradient methods. In *Proceedings of the 23rd international conference on Machine learning*, pages 969–976. ACM, 2006.

Holger Wendland. *Scattered data approximation*, volume 17. Cambridge university press, 2004.

Jason Weston, Olivier Chapelle, Vladimir Vapnik, André Elisseeff, and Bernhard Schölkopf. Kernel dependency estimation. In *Advances in neural information processing systems*, pages 873–880, 2002.

Yuan Yao, Lorenzo Rosasco, and Andrea Caponnetto. On early stopping in gradient descent learning. *Constructive Approximation*, 26(2):289–315, 2007.

# Appendix

The appendix are organized in three main parts:

- Appendix A focuses on the general ILE framework, proving the results in Section 3 and the Comparison Inequality of Theorem 7.
- Appendix B covers the details of the theoretical analysis reported in Section 5.
- Appendix C provides the proofs of the results in Sec. 6, offering sufficient conditions to guarantee a loss function to admit an ILE.

**Contributions and connection with previous work.** We recall that this paper is the longer version of (Ciliberto et al., 2016). Therefore, the results reported in Appendix A contain significant overlaps with the original work. We still prove each of the results in detail for the sake of completeness and since in the current work we have extended the framework in (Ciliberto et al., 2016) to the case where  $\Delta : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathcal{H}$  with output space  $\mathcal{Z}$  not necessarily corresponding to the label space  $\mathcal{Y}$ . The results in Appendix B and in particular Appendix C are novel for the most part.

**Setting and Notation.** We assume input, label and output spaces  $\mathcal{X}$ ,  $\mathcal{Y}$  and  $\mathcal{Z}$  to be Polish spaces, namely separable complete metrizable spaces, equipped with the associated Borel sigma-algebra. When referring to the data distribution  $\rho$  on  $\mathcal{X} \times \mathcal{Y}$  we will always assume it to be a Borel probability measure, with  $\rho_{\mathcal{X}}$  the marginal distribution on  $\mathcal{X}$  and  $\rho(\cdot|x)$  the conditional measure on  $\mathcal{Y}$  given  $x \in \mathcal{X}$ . We recall that  $\rho(y|x)$  is a regular conditional distribution (Dudley, 2002). Its domain  $D_{\rho|\mathcal{X}}$  is a measurable set contained in the support of  $\rho_{\mathcal{X}}$  and corresponds to the support of  $\rho_{\mathcal{X}}$  up to a set of measure zero.

For a Hilbert space  $\mathcal{H}$  we denote with  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$  and  $\| \cdot \|_{\mathcal{H}}$  the associated inner product and corresponding norm. Given two Hilbert spaces  $\mathcal{H}_1$  and  $\mathcal{H}_2$  we denote by  $\mathcal{H}_1 \oplus \mathcal{H}_2$  and  $\mathcal{H}_1 \otimes \mathcal{H}_2$  respectively their direct sum and tensor product. In particular, for any  $h_1, h'_1 \in \mathcal{H}_1$  and  $h_2, h'_2 \in \mathcal{H}_2$ , we have

$$\langle h_1 \oplus h_2, h'_1 \oplus h'_2 \rangle_{\mathcal{H}_1 \oplus \mathcal{H}_2} = \langle h_1, h'_1 \rangle_{\mathcal{H}_1} + \langle h_2, h'_2 \rangle_{\mathcal{H}_2} \quad (57)$$

$$\langle h_1 \otimes h_2, h'_1 \otimes h'_2 \rangle_{\mathcal{H}_1 \otimes \mathcal{H}_2} = \langle h_1, h'_1 \rangle_{\mathcal{H}_1} \cdot \langle h_2, h'_2 \rangle_{\mathcal{H}_2}. \quad (58)$$

Given a linear operator  $V : \mathcal{H}_1 \rightarrow \mathcal{H}_2$ , we denote by  $\text{Tr}(V)$  the trace of  $V$  and by  $V^* : \mathcal{H}_2 \rightarrow \mathcal{H}_1$  the adjoint of  $V$ , namely such that  $\langle Vh_1, h_2 \rangle_{\mathcal{H}_2} = \langle h_1, V^*h_2 \rangle_{\mathcal{H}_1}$  for every  $h_1 \in \mathcal{H}_1$ ,  $h_2 \in \mathcal{H}_2$ . Moreover, we denote by  $\|V\| = \sup_{\|h\|_{\mathcal{H}_1} \leq 1} \|Vh\|_{\mathcal{H}_2}$  the operator norm and  $\|V\|_{\text{HS}} = \sqrt{\text{Tr}(V^*V)}$  the Hilbert-Schmidt norm of  $V$ . In particular, we recall that the tensor product  $\mathcal{H}_1 \otimes \mathcal{H}_2$  is isometric to the space of Hilbert-Schmidt operators.

We denote with  $L^2(\mathcal{X}, \rho_{\mathcal{X}}, \mathcal{H})$  the Lebesgue space of square integrable functions on  $\mathcal{X}$  with respect to a measure  $\rho_{\mathcal{X}}$  and with values in a separable Hilbert space  $\mathcal{H}$ . For simplicity we denote with  $L^2(\mathcal{X}, \rho_{\mathcal{X}})$  the space  $L^2(\mathcal{X}, \rho_{\mathcal{X}}, \mathbb{R})$ . We denote with  $\langle f, g \rangle_{\rho_{\mathcal{X}}}$  the inner product  $\int \langle f(x), g(x) \rangle_{\mathcal{H}} d\rho_{\mathcal{X}}(x)$ , for all  $f, g \in L^2(\mathcal{X}, \rho_{\mathcal{X}}, \mathcal{H})$ .

**On the Argmin.** In the main paper we denoted the minimizer of (20) as

$$f_n(x) = \operatorname{argmin}_{y \in \mathcal{Y}} \sum_{i=1}^n \alpha_i(x) \Delta(y, y_i). \quad (59)$$

Clearly, the rigorous notation should be

$$f_n(x) \in \operatorname{argmin}_{y \in \mathcal{Y}} \sum_{i=1}^n \alpha_i(x) \Delta(y, y_i) \quad (60)$$

since it is not guaranteed in general to have one single minimizer for any given  $x \in \mathcal{X}$ . As we will discuss in the following, existence of a measurable function  $f_n$  that satisfies such inclusions requirement for any  $x \in \mathcal{X}$  can be guaranteed under mild assumptions.

## Appendix A. The ILE Framework

This section is devoted to characterize the theoretical properties of the ILE framework introduced in Sec. 3. In particular we prove the results in Theorem 3 (Fischer Consistency) and Theorem 7 (Comparison Inequality), which relate the “surrogate” risk to the original structured prediction one.

We begin by proving that both the structured risk  $\mathcal{E}$  and  $\mathcal{R}$  admit measurable minimizers under very mild conditions.

**Lemma A.1 (Existence of a minimizer for  $\mathcal{E}$ )** *Let  $\Delta : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$  be a continuous function and  $\mathcal{Z}$  a compact set. Then, the expected risk  $\mathcal{E}$  in (1) admits a measurable minimizer  $f^* : \mathcal{X} \rightarrow \mathcal{Z}$  such that*

$$f^*(x) \in \operatorname{argmin}_{z \in \mathcal{Z}} \int_{\mathcal{Y}} \Delta(z, y) d\rho(y|x) \quad (A.1)$$

almost everywhere on  $D_{\rho|\mathcal{X}}$ . Moreover, the function  $m : \mathcal{X} \rightarrow \mathbb{R}$  such that

$$m(x) = \inf_{z \in \mathcal{Z}} r(x, z), \quad \text{with} \quad r(x, z) = \begin{cases} \int_{\mathcal{Y}} \Delta(z, y) d\rho(y|x) & \text{if } x \in D_{\rho|\mathcal{X}} \\ 0 & \text{otherwise} \end{cases} \quad (A.2)$$

for any  $x \in \mathcal{X}$ , is measurable.

**Proof** Since  $\Delta$  is continuous and  $\rho(y|x)$  is a regular conditional distribution, then  $r$  is a Carathéodory function (see Definition 4.50 (pp. 153) in Aliprantis and Border, 2006), namely continuous in  $z$  for each  $x \in \mathcal{X}$  and measurable in  $x$  for each  $z \in \mathcal{Z}$ . Thus, by (Theorem 18.19 pp. 605 in Aliprantis and Border, 2006) (or Aumann’s measurable selection principle (Steinwart and Christmann, 2008; Castaing and Valadier, 2006)), we have that  $m$  is measurable and that there exists a measurable  $f^* : \mathcal{X} \rightarrow \mathcal{Z}$  such that  $r(x, f^*(x)) = m(x)$  for all  $x \in \mathcal{X}$ . Moreover, by definition of  $m$ , given any measurable  $f : \mathcal{X} \rightarrow \mathcal{Z}$ , we have  $m(x) \leq r(x, f(x))$ . Therefore,

$$\mathcal{E}(f^*) = \int_{\mathcal{X}} r(x, f^*(x)) d\rho_{\mathcal{X}}(x) = \int_{\mathcal{X}} m(x) d\rho_{\mathcal{X}}(x) \leq \int_{\mathcal{X}} r(x, f(x)) d\rho_{\mathcal{X}}(x) = \mathcal{E}(f). \quad (A.3)$$

We conclude  $\mathcal{E}(f^*) \leq \inf_{f: \mathcal{X} \rightarrow \mathcal{Z}} \mathcal{E}(f)$  and, since  $f^*$  is measurable,  $\mathcal{E}(f^*) = \min_{f: \mathcal{X} \rightarrow \mathcal{Y}} \mathcal{E}(f)$  and  $f^*$  is a global minimizer.  $\blacksquare$

In the following we will assume  $\Delta : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$  to admit an ILE, with associated Hilbert space  $\mathcal{H}$  and feature maps  $\psi : \mathcal{Z} \rightarrow \mathcal{H}$  and  $\varphi : \mathcal{Y} \rightarrow \mathcal{H}$ . We recall that the surrogate risk associated is defined as

$$\mathcal{R}(g) = \int_{\mathcal{X} \times \mathcal{Y}} \|g(x) - \varphi(y)\|_{\mathcal{H}}^2 d\rho(x, y) \quad (\text{A.4})$$

for any  $g : \mathcal{X} \rightarrow \mathcal{H}$ . Below we show that the global minimizer of  $\mathcal{R}$  corresponds to the conditional expectation of  $\varphi(y)$ .

**Lemma A.2 (Existence of a minimizer for  $\mathcal{R}$ )** *Let  $\mathcal{H}$  a separable Hilbert space and  $\varphi : \mathcal{Y} \rightarrow \mathcal{H}$  measurable and bounded with  $\sup_{y \in \mathcal{Y}} \|\varphi(y)\|_{\mathcal{H}} \leq \Phi$ . Then, the function  $g^* : \mathcal{X} \rightarrow \mathcal{H}$  such that*

$$g^*(x) = \int_{\mathcal{Y}} \varphi(y) d\rho(y|x) \quad \forall x \in D_{\rho|x} \quad (\text{A.5})$$

and  $g^*(x) = 0$  otherwise, belongs to  $L^2(\mathcal{X}, \rho_{\mathcal{X}}, \mathcal{H})$  and is a minimizer of the surrogate risk  $\mathcal{R}$ . Moreover, for any  $g \in L^2(\mathcal{X}, \rho_{\mathcal{X}}, \mathcal{H})$ ,

$$\mathcal{R}(g) - \mathcal{R}(g^*) = \int_{\mathcal{X}} \|g(x) - g^*(x)\|_{\mathcal{H}}^2 d\rho_{\mathcal{X}}(x) \quad (\text{A.6})$$

Hence, any minimizer of  $\mathcal{R}$  is equal to  $g^*$  almost everywhere on the domain of  $\rho_{\mathcal{X}}$ .

**Proof** By hypothesis,  $\|\psi\|_{\mathcal{H}}$  is measurable and bounded. Therefore, since  $\rho(y|x)$  is a regular conditional probability, we have that  $g^*$  is measurable on  $\mathcal{X}$  (see for instance Steinwart and Christmann (2008)). Moreover, the norm of  $g^*$  is dominated by the constant function of value  $\Phi$ , thus  $g^*$  is integrable on  $\mathcal{X}$  with respect to  $\rho_{\mathcal{X}}$  and in particular it is in  $L^2(\mathcal{X}, \rho_{\mathcal{X}}, \mathcal{H})$  since  $\rho_{\mathcal{X}}$  is a finite regular measure. Recall that since  $\rho(y|x)$  is a regular conditional distribution, for any measurable  $g : \mathcal{X} \rightarrow \mathcal{H}$  we have

$$\mathcal{R}(g) = \int_{\mathcal{X} \times \mathcal{Y}} \|g(x) - \psi(y)\|_{\mathcal{H}}^2 d\rho(x, y) = \int_{\mathcal{X}} \int_{\mathcal{Y}} \|g(x) - \psi(y)\|_{\mathcal{H}}^2 d\rho(y|x) d\rho_{\mathcal{X}}(x). \quad (\text{A.7})$$

Notice that  $g^*(x) = \operatorname{argmin}_{h \in \mathcal{H}} \int_{\mathcal{Y}} \|h - \psi(y)\|_{\mathcal{H}}^2 d\rho(y|x)$  almost everywhere on  $D_{\rho|x}$ . Indeed,

$$\int_{\mathcal{Y}} \|h - \psi(y)\|_{\mathcal{H}}^2 d\rho(y|x) = \|h\|_{\mathcal{H}}^2 - 2 \left\langle h, \left( \int_{\mathcal{Y}} \psi(y) d\rho(y|x) \right) \right\rangle + \int_{\mathcal{Y}} \|\psi(y)\|_{\mathcal{H}}^2 d\rho(y|x) \quad (\text{A.8})$$

$$= \|h\|_{\mathcal{H}}^2 - 2 \langle h, g^*(x) \rangle_{\mathcal{H}} + \text{const.} \quad (\text{A.9})$$

for all  $x \in D_{\rho|x}$ , which is minimized by  $h = g^*(x)$  for all  $x \in D_{\rho|x}$ . Therefore, since  $D_{\rho|x}$  is equal to the support of  $\rho_{\mathcal{X}}$  up to a set of measure zero, we conclude that  $\mathcal{R}(g^*) \leq \inf_{g: \mathcal{X} \rightarrow \mathcal{H}} \mathcal{R}(g)$  and, since  $g^*$  is measurable,  $\mathcal{R}(g^*) = \min_{g: \mathcal{X} \rightarrow \mathcal{H}} \mathcal{R}(g)$  and  $g^*$  is a global minimizer as required.

Finally, notice that for any  $g : \mathcal{X} \rightarrow \mathcal{H}$  we have

$$\mathcal{R}(g) - \mathcal{R}(g^*) = \int_{\mathcal{X} \times \mathcal{Y}} \|g(x) - \psi(y)\|_{\mathcal{H}}^2 - \|g^*(x) - \psi(y)\|_{\mathcal{H}}^2 d\rho(x, y) \quad (\text{A.10})$$

$$= \int_{\mathcal{X}} \|g(x)\|_{\mathcal{H}}^2 - 2 \left\langle g(x), \left( \int_{\mathcal{Y}} \psi(y) d\rho(y|x) \right) \right\rangle_{\mathcal{H}} + \|g^*(x)\|_{\mathcal{H}}^2 d\rho_{\mathcal{X}}(x) \quad (\text{A.11})$$

$$= \int_{\mathcal{X}} \|g(x)\|_{\mathcal{H}}^2 - 2 \langle g(x), g^*(x) \rangle_{\mathcal{H}} + \|g^*(x)\|_{\mathcal{H}}^2 d\rho_{\mathcal{X}}(x) \quad (\text{A.12})$$

$$= \int_{\mathcal{X}} \|g(x) - g^*(x)\|_{\mathcal{H}}^2 d\rho_{\mathcal{X}}(x), \quad (\text{A.13})$$

which proves (A.6). Therefore, for any measurable minimizer  $g' : \mathcal{X} \rightarrow \mathcal{H}$  of the surrogate expected risk, we have  $\mathcal{R}(g') - \mathcal{R}(g^*) = 0$  which, by the relation above, implies  $g'(x) = g^*(x)$  a.e. on  $D_{\rho|\mathcal{X}}$ .  $\blacksquare$

Combining the characterizations of the global minimizers of the two risks  $\mathcal{E}$  and  $\mathcal{R}$  we can now prove the following.

**Lemma A.3 (Fisher Consistency)** *Let  $\mathcal{Z}$  be compact,  $\Delta : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$  admit an ILE and let  $f^* : \mathcal{X} \rightarrow \mathcal{Z}$  be the solution of (1). Then,*

$$f^*(x) = \operatorname{argmin}_{z \in \mathcal{Z}} \langle \psi(z), g^*(x) \rangle_{\mathcal{H}}, \quad g^*(x) = \int_{\mathcal{Y}} \varphi(y) d\rho(y|x) \quad (16)$$

almost surely with respect to  $\rho_{\mathcal{X}}$ . Moreover,  $g^* : \mathcal{X} \rightarrow \mathcal{H}$  is the minimizer of

$$\mathcal{R}(g) = \int_{\mathcal{Y} \times \mathcal{X}} \|\varphi(y) - g(x)\|_{\mathcal{H}}^2 d\rho(x, y). \quad (17)$$

**Proof** By Theorem A.2 we know that  $g^*(x) = \int_{\mathcal{Y}} \varphi(y) d\rho(y|x)$  almost everywhere on  $D_{\rho|\mathcal{X}}$  and is the minimizer of  $\mathcal{R}$ . Therefore, for every  $z \in \mathcal{Z}$  we have

$$\langle \psi(z), g^*(x) \rangle_{\mathcal{H}} = \left\langle \psi(z), \int_{\mathcal{Y}} \varphi(y) d\rho(y|x) \right\rangle_{\mathcal{H}} \quad (\text{A.14})$$

$$= \int_{\mathcal{Y}} \langle \psi(z), \varphi(y) \rangle_{\mathcal{H}} d\rho(y|x) = \int_{\mathcal{Y}} \Delta(z, y) d\rho(y|x) \quad (\text{A.15})$$

almost everywhere on  $D_{\rho|\mathcal{X}}$ . Thus, for any measurable function  $f : \mathcal{X} \rightarrow \mathcal{Z}$  we have

$$\mathcal{E}(f) = \int_{\mathcal{X} \times \mathcal{Y}} \Delta(f(x), y) d\rho(x, y) = \int_{\mathcal{X}} \int_{\mathcal{Y}} \Delta(f(x), y) d\rho(y|x) d\rho_{\mathcal{X}}(x) \quad (\text{A.16})$$

$$= \int_{\mathcal{X}} \langle \psi(f(x)), g^*(x) \rangle_{\mathcal{H}} d\rho_{\mathcal{X}}(x). \quad (\text{A.17})$$

We conclude that a minimizer  $f^* : \mathcal{X} \rightarrow \mathcal{Z}$  of  $\mathcal{E}$  can be characterized as a function minimizing pointwise the integral above, namely

$$f^*(x) \in \operatorname{argmin}_{z \in \mathcal{Z}} \langle \psi(f(x)), g^*(x) \rangle_{\mathcal{H}} \quad (\text{A.18})$$



almost everywhere on  $D_{\rho|\mathcal{X}}$ . ■

We now prove Theorem 7, characterizing the relation between the excess risks associated to  $\mathcal{R}$  and  $\mathcal{E}$ .

**Theorem 7 (Comparison Inequality)** *Let  $\mathcal{Z}$  be a compact set and  $\Delta : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$  admit an ILE. Let  $f^*$ ,  $g^*$  and the risk  $\mathcal{R}(\cdot)$  be defined as in Theorem 3. Let  $g : \mathcal{X} \rightarrow \mathcal{H}$  be measurable and let  $f : \mathcal{X} \rightarrow \mathcal{Z}$  be such that*

$$f(x) = \operatorname{argmin}_{z \in \mathcal{Z}} \langle \psi(z), g(x) \rangle_{\mathcal{H}}, \quad (43)$$

for any  $x \in \mathcal{X}$ . Then,

$$\mathcal{E}(f) - \mathcal{E}(f^*) \leq 2 c_{\Delta} \sqrt{\mathcal{R}(g) - \mathcal{R}(g^*)} \quad (44)$$

**Proof** By applying Thm. 3, we have

$$\mathcal{E}(f) - \mathcal{E}(f^*) = \int_{\mathcal{X} \times \mathcal{Y}} \Delta(f(x), y) - \Delta(f^*(x), y) d\rho(x, y) \quad (A.19)$$

$$= \int_{\mathcal{X} \times \mathcal{Y}} \langle \psi(f(x)) - \psi(f^*(x)), \varphi(y) \rangle_{\mathcal{H}} d\rho(x, y) \quad (A.20)$$

$$= \int_{\mathcal{X}} \left\langle \psi(f(x)) - \psi(f^*(x)), \left( \int_{\mathcal{Y}} \varphi(y) d\rho(y|x) \right) \right\rangle_{\mathcal{H}} d\rho_{\mathcal{X}}(x) \quad (A.21)$$

$$= \int_{\mathcal{X}} \langle \psi(f(x)) - \psi(f^*(x)), g^*(x) \rangle_{\mathcal{H}} d\rho_{\mathcal{X}}(x) \quad (A.22)$$

$$= A + B. \quad (A.23)$$

where in the last equation we have removed and added a term  $\int_{\mathcal{X}} \langle \psi(f(x)), g(x) \rangle_{\mathcal{H}} d\rho_{\mathcal{X}}(x)$  leading to

$$A = \int_{\mathcal{X}} \langle \psi(f(x)), (g^*(x) - g(x)) \rangle_{\mathcal{H}} d\rho_{\mathcal{X}}(x) \quad (A.24)$$

$$B = \int_{\mathcal{X}} \langle \psi(f(x)), g(x) \rangle_{\mathcal{H}} - \langle \psi(f^*(x)), g^*(x) \rangle_{\mathcal{H}} d\rho_{\mathcal{X}}(x) \quad (A.25)$$

Now, the term A can be minimized by taking the supremum over  $\mathcal{Z}$  so that

$$A \leq \int_{\mathcal{X}} \sup_{z \in \mathcal{Z}} \left| \langle \psi(z), g^*(x) - g(x) \rangle_{\mathcal{H}} \right| d\rho_{\mathcal{X}}(x). \quad (A.26)$$

For B, we observe that from the characterization of  $f$  in the hypothesis and of  $f^*$  by Theorem 3, we have

$$\langle \psi(f^*(x)), g^*(x) \rangle_{\mathcal{H}} = \inf_{z \in \mathcal{Z}} \langle \psi(z), g^*(x) \rangle_{\mathcal{H}}, \quad (A.27)$$

$$\langle \psi(f(x)), g(x) \rangle_{\mathcal{H}} = \inf_{z \in \mathcal{Z}} \langle \psi(z), g(x) \rangle_{\mathcal{H}}, \quad (A.28)$$

for all  $x \in \mathcal{X}$ . Therefore,

$$B = \int_{\mathcal{X}} \inf_{z \in \mathcal{Z}} \langle \psi(z), g(x) \rangle_{\mathcal{H}} - \inf_{z \in \mathcal{Z}} \langle \psi(z), g^*(x) \rangle_{\mathcal{H}} d\rho_{\mathcal{X}}(x) \quad (\text{A.29})$$

$$\leq \int_{\mathcal{X}} \sup_{z \in \mathcal{Z}} \left| \langle \psi(z), (g(x) - g^*(x)) \rangle_{\mathcal{H}} \right| d\rho_{\mathcal{X}}(x) \quad (\text{A.30})$$

where we have used the fact that for any given two functions  $\eta, \zeta : \mathcal{Z} \rightarrow \mathbb{R}$  we have

$$\left| \inf_{z \in \mathcal{Z}} \eta(z) - \inf_{z \in \mathcal{Z}} \zeta(z) \right| \leq \sup_{z \in \mathcal{Z}} |\eta(z) - \zeta(z)|. \quad (\text{A.31})$$

Therefore, by combining the bounds on  $A$  and  $B$  we have

$$\mathcal{E}(f) - \mathcal{E}(f^*) \leq 2 \int_{\mathcal{X}} \sup_{z \in \mathcal{Z}} \left| \langle \psi(z), g^*(x) - g(x) \rangle_{\mathcal{H}} \right| d\rho_{\mathcal{X}}(x) \quad (\text{A.32})$$

$$\leq 2 \int_{\mathcal{X}} \sup_{z \in \mathcal{Z}} \|\psi(z)\|_{\mathcal{H}} \|g^*(x) - g(x)\|_{\mathcal{H}} d\rho_{\mathcal{X}}(x) \quad (\text{A.33})$$

$$\leq 2c_{\Delta} \int_{\mathcal{X}} \|g^*(x) - g(x)\|_{\mathcal{H}} d\rho_{\mathcal{X}}(x) \quad (\text{A.34})$$

$$\leq 2c_{\Delta} \sqrt{\int_{\mathcal{X}} \|g^*(x) - g(x)\|_{\mathcal{H}}^2 d\rho_{\mathcal{X}}(x)}, \quad (\text{A.35})$$

$$(\text{A.36})$$

where for the last inequality we have used the Jensen's inequality. The proof is concluded recalling that, by Theorem A.2

$$\mathcal{R}(g) - \mathcal{R}(g^*) = \int_{\mathcal{X}} \|g(x) - g^*(x)\|_{\mathcal{H}}^2 d\rho_{\mathcal{X}}(x) \quad (\text{A.37})$$

■

We conclude proving the result in Theorem 4, which is a direct consequence of the linearity induced by the ILE definition.

**Lemma A.4** *Let  $\Delta : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$  admit an ILE,  $(y_i)_{i=1}^n$  a set of points in  $\mathcal{Y}$  and  $\alpha : \mathcal{X} \rightarrow \mathbb{R}^n$  a weighting function. Let  $g_n : \mathcal{X} \rightarrow \mathcal{H}$  be such that  $g_n(x) = \sum_{i=1}^n \alpha_i(x) \varphi(y_i)$  for any  $x \in \mathcal{X}$ . Then, the function  $f_n : \mathcal{X} \rightarrow \mathcal{Z}$  such that  $\forall x \in \mathcal{X}$*

$$f_n(x) = \operatorname{argmin}_{z \in \mathcal{Z}} \langle \psi(z), g_n(x) \rangle_{\mathcal{H}} = \operatorname{argmin}_{z \in \mathcal{Z}} \sum_{i=1}^n \alpha_i(x) \Delta(z, y_i). \quad (20)$$

**Proof** For any  $z \in \mathcal{Z}$  and  $x \in \mathcal{X}$  we have

$$\langle \psi(z), g_n(x) \rangle_{\mathcal{H}} = \left\langle \psi(z), \sum_{i=1}^n \alpha_i(x) \varphi(y_i) \right\rangle_{\mathcal{H}} \quad (\text{A.38})$$

$$= \sum_{i=1}^n \alpha_i(x) \langle \psi(z), \varphi(y_i) \rangle_{\mathcal{H}} \quad (\text{A.39})$$

$$= \sum_{i=1}^n \alpha_i(x) \Delta(z, y_i). \quad (\text{A.40})$$

Therefore, substituting the above equation in the definition of  $f_n$  concludes the proof.  $\blacksquare$

## Appendix B. Universal Consistency and Learning Bounds

**Additional Notation.** Let  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  a positive semidefinite function on  $\mathcal{X}$ . We denote  $\mathcal{F}$  the Hilbert space obtained by the completion

$$\mathcal{F} = \overline{\text{span}\{k(x, \cdot) \mid x \in \mathcal{X}\}} \quad (\text{B.1})$$

according to the norm induced by the inner product  $\langle k(x, \cdot), k(x', \cdot) \rangle_{\mathcal{F}} = k(x, x')$ . Spaces  $\mathcal{F}$  constructed in this way are known as *reproducing kernel Hilbert spaces* and there is a one-to-one relation between a kernel  $k$  and its associated RKHS. For more details on RKHS we refer the reader to Berlinet and Thomas-Agnan (2011). Given a kernel  $k$ , in the following we will denote with  $\phi : \mathcal{X} \rightarrow \mathcal{F}$  the feature map  $\phi(x) = k(x, \cdot) \in \mathcal{F}$  for all  $x \in \mathcal{X}$ . We say that a kernel is bounded if  $\|\phi(x)\|_{\mathcal{F}} \leq \kappa$  with  $\kappa > 0$ . Note that  $k$  is bounded if and only if  $k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{F}} \leq \|\phi(x)\|_{\mathcal{F}} \|\phi(x')\|_{\mathcal{F}} \leq \kappa^2$  for every  $x, x' \in \mathcal{X}$ . In the following we will always assume  $k$  to be continuous and bounded by  $\kappa > 0$ . The continuity of  $k$  with the fact that  $\mathcal{X}$  is Polish implies  $\mathcal{F}$  to be separable Berlinet and Thomas-Agnan (2011).

We introduce here the ideal and empirical operators that we will use in the following to prove the main results of this work.

- $S : \mathcal{F} \rightarrow L^2(\mathcal{X}, \rho_{\mathcal{X}})$  s.t.  $f \in \mathcal{F} \mapsto \langle f, \phi(\cdot) \rangle_{\mathcal{F}} \in L^2(\mathcal{X}, \rho_{\mathcal{X}})$ , with adjoint
- $S^* : L^2(\mathcal{X}, \rho_{\mathcal{X}}) \rightarrow \mathcal{F}$  s.t.  $h \in L^2(\mathcal{X}, \rho_{\mathcal{X}}) \mapsto \int_{\mathcal{X}} h(x) \phi(x) d\rho_{\mathcal{X}}(x) \in \mathcal{F}$ ,
- $Z : \mathcal{H} \rightarrow L^2(\mathcal{X}, \rho_{\mathcal{X}})$  s.t.  $h \in \mathcal{H} \mapsto \langle h, g^*(\cdot) \rangle_{\mathcal{H}} \in L^2(\mathcal{X}, \rho_{\mathcal{X}})$ , with adjoint
- $Z^* : L^2(\mathcal{X}, \rho_{\mathcal{X}}) \rightarrow \mathcal{H}$  s.t.  $h \in L^2(\mathcal{X}, \rho_{\mathcal{X}}) \mapsto \int_{\mathcal{X}} h(x) g^*(x) d\rho_{\mathcal{X}}(x) \in \mathcal{H}$ ,
- $C = S^*S : \mathcal{F} \rightarrow \mathcal{F}$  and  $L = SS^* : L^2(\mathcal{X}, \rho_{\mathcal{X}}) \rightarrow L^2(\mathcal{X}, \rho_{\mathcal{X}})$ ,

with  $g^*(x) = \int_{\mathcal{Y}} \psi(y) d\rho(y|x)$  defined according to (A.5), (see Theorem A.2).

Given a set of input-output pairs  $\{(x_i, y_i)\}_{i=1}^n$  with  $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$  independently sampled according to  $\rho$  on  $\mathcal{X} \times \mathcal{Y}$ , we define the empirical counterparts of the operators just defined as

- $\hat{S} : \mathcal{F} \rightarrow \mathbb{R}^n$  s.t.  $f \in \mathcal{F} \mapsto \frac{1}{\sqrt{n}} (\langle \phi(x_i), f \rangle_{\mathcal{F}})_{i=1}^n \in \mathbb{R}^n$ , with adjoint
- $\hat{S}^* : \mathbb{R}^n \rightarrow \mathcal{F}$  s.t.  $v = (v_i)_{i=1}^n \in \mathbb{R}^n \mapsto \frac{1}{\sqrt{n}} \sum_{i=1}^n v_i \phi(x_i)$ ,
- $\hat{Z} : \mathcal{H} \rightarrow \mathbb{R}^n$  s.t.  $h \in \mathcal{H} \mapsto \frac{1}{\sqrt{n}} (\langle \psi(y_i), h \rangle_{\mathcal{H}})_{i=1}^n \in \mathbb{R}^n$ , with adjoint
- $\hat{Z}^* : \mathbb{R}^n \rightarrow \mathcal{H}$  s.t.  $v = (v_i)_{i=1}^n \in \mathbb{R}^n \mapsto \frac{1}{\sqrt{n}} \sum_{i=1}^n v_i \psi(y_i)$ ,
- $\hat{C} = \hat{S}^* \hat{S} : \mathcal{F} \rightarrow \mathcal{F}$  and  $K = n \hat{S} \hat{S}^* \in \mathbb{R}^{n \times n}$  is the empirical kernel matrix.

In the rest of this section we denote with  $A + \lambda$ , the operator  $A + \lambda I$ , for any symmetric linear operator  $A$ ,  $\lambda \in \mathbb{R}$  and  $I$  the identity operator.

### B.1. Preliminary results

We recall here a basic result characterizing the operators introduced above.

**Proposition B.1** *With the notation introduced above,*

$$C = \int_{\mathcal{X}} \phi(x) \otimes \phi(x) d\rho_{\mathcal{X}}(x) \quad \text{and} \quad Z^*S = \int_{\mathcal{X} \times \mathcal{Y}} \psi(y) \otimes \phi(x) d\rho(x, y) \quad (\text{B.2})$$

where  $\otimes$  denotes the tensor product. Moreover, when  $\phi$  and  $\psi$  are bounded by respectively  $\kappa$  and  $Q$ , we have the following facts

$$(i) \quad \text{Tr}(L) = \text{Tr}(C) = \|S\|_{\text{HS}}^2 = \int_{\mathcal{X}} \|\phi(x)\|_{\mathcal{F}}^2 d\rho_{\mathcal{X}}(x) \leq \kappa^2$$

$$(ii) \quad \|Z\|_{\text{HS}}^2 = \int_{\mathcal{X}} \|g^*(x)\|^2 d\rho_{\mathcal{X}}(x) = \|g^*\|_{\rho_{\mathcal{X}}}^2 < +\infty.$$

**Proof** By definition of  $C = S^*S$ , for each  $h, h' \in \mathcal{F}$  we have

$$\langle h, Ch' \rangle_{\mathcal{F}} = \langle Sh, Sh' \rangle_{\rho_{\mathcal{X}}} = \int_{\mathcal{X}} \langle h, \phi(x) \rangle_{\mathcal{F}} \langle \phi(x), h' \rangle_{\mathcal{F}} d\rho_{\mathcal{X}}(x) \quad (\text{B.3})$$

$$= \int_{\mathcal{X}} \left\langle h, \left( \phi(x) \langle \phi(x), h' \rangle_{\mathcal{F}} \right) \right\rangle_{\mathcal{F}} d\rho_{\mathcal{X}}(x) \quad (\text{B.4})$$

$$= \int_{\mathcal{X}} \left\langle h, \left( \phi(x) \otimes \phi(x) \right) h' \right\rangle d\rho_{\mathcal{X}}(x) \quad (\text{B.5})$$

$$= \left\langle h, \left( \int_{\mathcal{X}} \phi(x) \otimes \phi(x) d\rho_{\mathcal{X}}(x) \right) h' \right\rangle_{\mathcal{F}} \quad (\text{B.6})$$

since  $\phi(x) \otimes \phi(x) : \mathcal{F} \rightarrow \mathcal{F}$  is the operator such that  $h \in \mathcal{F} \mapsto \phi(x) \langle \phi(x), h \rangle_{\mathcal{F}}$ . The characterization for  $Z^*S$  is analogous.

Now, (i). The relation  $\text{Tr}(L) = \text{Tr}(C) = \text{Tr}(S^*S) = \|S\|_{\text{HS}}^2$  holds by definition. Moreover

$$\text{Tr}(C) = \int_{\mathcal{X}} \text{Tr}(\phi(x) \otimes \phi(x)) d\rho_{\mathcal{X}}(x) = \int_{\mathcal{X}} \|\phi(x)\|_{\mathcal{F}}^2 d\rho_{\mathcal{X}}(x) \quad (\text{B.7})$$

by linearity of the trace. (ii) is analogous. Note that  $\|g^*\|_{\rho_{\mathcal{X}}}^2 < +\infty$ . by Theorem A.2 since  $\psi$  is bounded by hypothesis. ■

**Lemma B.2** *Let  $g_n(x) = \widehat{G}^* \phi(x)$  with  $\widehat{G} : \mathcal{H} \rightarrow \mathcal{F}$  a bounded linear operator, then*

$$\mathcal{R}(g_n) - \mathcal{R}(g^*) = \|S\widehat{G} - Z\|_{\text{HS}}^2, \quad (\text{B.8})$$

where  $\|A\|_{\text{HS}}^2 := \text{Tr}(A^*A)$ , for a linear operator  $A$ , is the Hilbert-Schmidt norm.

**Proof** By Theorem A.2, we know that  $g^*(x) = \int_{\mathcal{Y}} \psi(y) d\rho(y|x)$  almost everywhere on the support of  $\rho_{\mathcal{X}}$ , moreover by Theorem B.5  $g_n$ . Therefore, a direct application of Theorem B.1

leads to

$$\mathcal{R}(g_n) - \mathcal{R}(g^*) = \int \|g_n(x) - g^*(x)\|_{\mathcal{H}}^2 d\rho_{\mathcal{X}}(x) = \quad (\text{B.9})$$

$$= \int_{\mathcal{X}} \|\hat{G}^* \phi(x)\|_{\mathcal{H}}^2 - 2 \langle \hat{G}^* \phi(x), g^*(x) \rangle_{\mathcal{H}} + \|g^*(x)\|_{\mathcal{H}}^2 d\rho_{\mathcal{X}}(x) \quad (\text{B.10})$$

$$= \int_{\mathcal{X}} \text{Tr} \left( \hat{G}^* \left( \phi(x) \otimes \phi(x) \right) \hat{G} \right) - 2 \text{Tr} \left( \hat{G}^* \left( \phi(x) \otimes g^*(x) \right) \right) + \text{Tr} \left( g^*(x) \otimes g^*(x) \right) d\rho_{\mathcal{X}}(x) \quad (\text{B.11})$$

$$= \text{Tr}(\hat{G}^* S^* S \hat{G}) - 2 \text{Tr}(\hat{G}^* S^* Z) + \text{Tr}(Z^* Z) \quad (\text{B.12})$$

$$= \|\hat{S} \hat{G} - Z\|_{\text{HS}}^2 \quad (\text{B.13})$$

■

## B.2. Analytic Decomposition for Spectral Filters

To study the various estimators considered in this paper, we need to introduce the notion of *spectral filter*.

**Definition B.3 (spectral filters (Engl et al., 1996))** *Let  $\kappa > 0$ . Then  $\eta_{\lambda} : (0, \kappa^2] \rightarrow \mathbb{R}$  is a spectral filter if there exist  $q_1, q_2 > 0$  s.t. for  $\sigma \in (0, \kappa^2]$  and  $\lambda > 0$*

$$(\sigma + \lambda)\eta_{\lambda}(\sigma) \leq q_1, \quad (1 - \sigma\eta_{\lambda}(\sigma))(\sigma + \lambda) \leq q_2\lambda. \quad (\text{B.14})$$

In this work we have considered a simplified definition of spectral filters with respect to the standard notion. In particular, we do not make a distinction between filters with qualification *larger* than 1 (see e.g. (Engl et al., 1996; Bauer et al., 2007)). The following result gives three concrete examples of spectral filters that will be useful to characterize the estimators  $\hat{g}$  studied in this work.

**Lemma B.4** *The following functions are spectral filters:*

1. (Ridge Regression)  $\eta_{\lambda}(\sigma) = (\sigma + \lambda)^{-1}$ , with  $q_1 = q_2 = 1$
2. (L2-Boosting)  $\eta_{\lambda}(\sigma) = \nu \sum_{j=0}^t (1 - \nu\sigma)^j$ , with step-size  $0 < \nu < 1/\kappa^2$  and  $\lambda = 1/t$ .  
With constants  $q_1 = 1 + 2\nu$  and  $q_2 = e^{\nu-1}/\nu$ .
3. (PCR)  $\eta_{\lambda}(\sigma) = \frac{1}{\sigma} \mathbf{1}_{\sigma > \lambda}$ , where  $\mathbf{1}_{\sigma > \lambda} = 1$  when  $\sigma \geq \lambda$  and 0 otherwise.  
With constants  $q_1 = q_2 = 2$ .

**Proof** (Ridge Regression). It is easy to show that

- $(\sigma + \lambda)\eta_{\lambda}(\sigma) = (\sigma + \lambda)(\sigma + \lambda)^{-1} = 1 = q_1$ ,
- $\frac{1}{\lambda}(1 - \sigma\eta_{\lambda}(\sigma))(\sigma + \lambda) = \frac{1}{\lambda} \left( 1 - \frac{\sigma}{\sigma + \lambda} \right) (\sigma + \lambda) = 1 = q_2$ .

(*L2-Boosting*). Let  $\lambda = 1/t$ . Recall that since  $\nu < 1/\kappa^2$  and  $\sigma \in (0, \kappa^2]$ , we have  $\nu\sigma < 1$ . Therefore,  $\sum_{j=0}^{+\infty} (1 - \nu\sigma)^j = 1/(\nu\sigma)$  and we have

$$\eta_\lambda(\sigma) = \nu \sum_{j=0}^t (1 - \nu\sigma)^j = \frac{1}{\sigma} \left(1 - (1 - \nu\sigma)^{t+1}\right). \quad (\text{B.15})$$

Now  $(\sigma + \lambda)\eta_\lambda(\sigma) = \sigma\eta_\lambda(\sigma) + \lambda\eta_\lambda(\sigma)$ . Then,

$$\sigma\eta_\lambda(\sigma) = \sigma \frac{1}{\sigma} \left(1 - (1 - \nu\sigma)^{t+1}\right) < 1, \quad (\text{B.16})$$

since  $\nu\sigma > 0$ . Moreover, since  $\nu\sigma < 1$ ,

$$\lambda\eta_\lambda(\sigma) = \frac{1}{t}\nu \sum_{j=0}^t (1 - \nu\sigma)^j \leq \frac{(t+1)\nu}{t} < 2\nu. \quad (\text{B.17})$$

Hence we have

$$(\sigma + \lambda)\eta_\lambda(\sigma) \leq 1 + 2\nu = q_1. \quad (\text{B.18})$$

Now, since  $(1 - z) \leq e^{-z}$  and defining  $x = (t+1)\sigma$ , we have

$$\frac{1}{\lambda}(1 - \sigma\eta_\lambda(\sigma))(\sigma + \lambda) = t(1 - \nu\sigma)^{t+1}(\sigma + 1/t) \quad (\text{B.19})$$

$$\leq te^{-\nu\sigma(t+1)}(\sigma + 1/t) \quad (\text{B.20})$$

$$= e^{-\nu\sigma(t+1)} + t\sigma e^{-\nu\sigma(t+1)} \quad (\text{B.21})$$

$$\leq e^{-\nu\sigma(t+1)} + (t+1)\sigma e^{-\nu\sigma(t+1)} \quad (\text{B.22})$$

$$= e^{-\nu x} + x e^{-\nu x} \quad (\text{B.23})$$

$$\leq e^{\nu-1}/\nu = q_2. \quad (\text{B.24})$$

(*PCR*). Let  $\sigma < \lambda$ , then  $\eta_\lambda(\sigma) = 0$  and

- $(\sigma + \lambda)\eta_\lambda(\sigma) = 0 < 2 = q_1,$
- $\frac{1}{\lambda}(1 - \sigma\eta_\lambda(\sigma))(\sigma + \lambda) = \frac{1}{\lambda}(\sigma + \lambda) < \frac{2\lambda}{\lambda} = 2 = q_2.$

If  $\sigma \geq \lambda$ , we have  $\eta_\lambda(\sigma) = 1/\sigma$  and

- $(\sigma + \lambda)\eta_\lambda(\sigma) = \frac{\sigma + \lambda}{\sigma} < \frac{2\sigma}{\sigma} = 2 = q_1,$
- $\frac{1}{\lambda}(1 - \sigma\eta_\lambda(\sigma))(\sigma + \lambda) = 0 < 2 = q_2.$

■

We will be applying filters  $\eta_\lambda$  to the spectrum of an operator as follows. Let  $M : \mathcal{H}_1 \rightarrow \mathcal{H}_2$  be a compact linear operator between two separable Hilbert spaces  $\mathcal{H}_1, \mathcal{H}_2$ . Let  $M =$

$\sum_{i=1}^{+\infty} \sigma_i u_i \otimes v_i$  be the singular value decomposition of  $M$ , with  $(u_i)_{i \in \mathbb{N}}$  and  $(v_i)_{i \in \mathbb{N}}$  a suitable pair of orthonormal bases of  $\mathcal{H}_1$  and  $\mathcal{H}_2$  respectively and  $\sigma_i \geq 0$  for every  $i \in \mathbb{N}$ . We denote the application of  $\eta_\lambda$  to  $M$  as

$$\eta_\lambda(M) = \sum_{i=1}^{+\infty} \eta_\lambda(\sigma_i) u_i \otimes v_i. \quad (\text{B.25})$$

The following results shows that several estimators described in Section 3.3 can be formulated in terms of spectral filters.

**Lemma B.5** *The following algorithms can be represented as*

$$\hat{g}_\lambda(x) = \hat{G}_\lambda^* \phi(x), \quad \hat{G}_\lambda = \eta_\lambda(\hat{C}) \hat{S}^* \hat{Z}, \quad (\text{B.26})$$

where  $\eta_\lambda$  is a spectral filter, in particular

1. (Kernel Ridge Regression)  $\eta_\lambda(\sigma) = (\sigma + \lambda)^{-1}$
2. (Kernel L2-Boosting)  $\eta_\lambda(\sigma) = \nu \sum_{j=0}^t (1 - \nu\sigma)^j$ , with step-size  $\nu$  and  $\lambda = 1/t$ ,
3. (Kernel PCR)  $\eta_\lambda(\sigma) = \frac{1}{\sigma} \mathbf{1}_{\sigma > \lambda}$ , where  $\mathbf{1}_{\sigma > \lambda} = 1$  when  $\sigma \geq \lambda$  and 0 otherwise.

**Proof** Recall that, according to (19) the estimator  $\hat{g}$  is such that, for any  $x \in \mathcal{X}$

$$\hat{g}_\lambda(x) = \sum_{i=1}^n \alpha_i(x) \varphi(y_i). \quad (\text{B.27})$$

It follows by the definition of the three methods considered to learn the vector-valued function  $\alpha$  that, for any  $x \in \mathcal{X}$

$$\alpha(x) = \frac{1}{n} \eta_\lambda \left( \frac{K}{n} \right) \mathbf{v}(x) \quad (\text{B.28})$$

where  $\eta_\lambda$  is the corresponding spectral filter function given in the thesis of this Lemma. Recall that  $\mathbf{v}(x) = (k(x_1, x), \dots, k(x_n, x))$ . By definition of  $\hat{S}$  and  $\hat{Z}$ , we have  $\mathbf{v}(x) = \sqrt{n} \hat{S} \phi(x)$  and  $\sum_{i=1}^n \alpha_i(x) \varphi(y_i) = \frac{1}{\sqrt{n}} \hat{Z}^* \alpha(x)$ . Then

$$\hat{g}_\lambda(x) = \frac{1}{\sqrt{n}} \hat{Z} \alpha(x) = \frac{1}{\sqrt{n}} \hat{Z}^* \eta_\lambda \left( \frac{K}{n} \right) \mathbf{v}(x) = \hat{Z}^* \eta_\lambda \left( \frac{K}{n} \right) \hat{S} \phi(x). \quad (\text{B.29})$$

Since  $K = n \hat{S} \hat{S}^*$ , then

$$\eta_\lambda \left( \frac{K}{n} \right) \hat{S} = \eta_\lambda(\hat{S} \hat{S}^*) \hat{S} = \hat{S} \eta_\lambda(\hat{S}^* \hat{S}) = \hat{S} \eta_\lambda(\hat{C}), \quad (\text{B.30})$$

from which it follows

$$\hat{g}_\lambda(x) = \hat{Z}^* \hat{S} \eta_\lambda(\hat{C}) \phi(x), \quad (\text{B.31})$$

as required. ■

With the characterization provided by the result above, we now proceed in deriving an upper bound for the risk of estimators obtained via spectral filtering methods.

**Theorem B.6** *Let  $\hat{g}$  be characterized as in Theorem B.5 in terms of a spectral filter  $\eta_\lambda$  with constants  $q_1, q_2 > 0$ . Let  $\beta = \|\hat{C}_\lambda^{-1/2} C_\lambda^{1/2}\|^2$ , with  $G_\lambda = S^* L_\lambda^{-1} Z$ . Then,*

$$|\mathcal{R}(g_{n_\lambda}) - \mathcal{R}(g^*)|^{1/2} \leq q_1 \beta \|C_\lambda^{-1/2} (\hat{S}^* \hat{Z} - \hat{C} G_\lambda)\|_{\text{HS}} + 2(1 + q_2) \beta \lambda \|L_\lambda^{-1} Z\|_{\text{HS}}, \quad (\text{B.32})$$

**Proof** From Theorem B.5 we know that  $\hat{g}_\lambda(x) = \hat{G}_\lambda^* \phi(x)$  with  $\hat{G}_\lambda = \eta_\lambda(\hat{C}) \hat{S}^* \hat{Z}$ . From Theorem B.2 we know that  $\mathcal{R}(g_{n_\lambda}) - \mathcal{R}(g^*) = \|\hat{S} \hat{G}_\lambda - Z\|_{\text{HS}}^2$ . We add and remove the term  $S \eta_\lambda(\hat{C}) \hat{S}^* \hat{S} G_\lambda$ , with  $G_\lambda = S^* L_\lambda^{-1} Z$  and  $L_\lambda = L + \lambda I$ , namely

$$S \hat{G}_\lambda - Z = S \eta_\lambda(\hat{C}) \hat{S}^* (\hat{Z} - \hat{S} G_\lambda) + S \eta_\lambda(\hat{C}) \hat{S}^* \hat{S} G_\lambda - Z. \quad (\text{B.33})$$

We decompose the first term in the sum above as

$$S \eta_\lambda(\hat{C}) \hat{S}^* (\hat{Z} - \hat{S} G_\lambda) = (S \hat{C}_\lambda^{-1/2}) (\hat{C}_\lambda^{1/2} \eta_\lambda(\hat{C}) \hat{C}_\lambda^{1/2}) (\hat{C}_\lambda^{-1/2} C_\lambda^{1/2}) [C_\lambda^{-1/2} \hat{S}^* (\hat{Z} - \hat{S} G_\lambda)]. \quad (\text{B.34})$$

Hence, we have

$$\|S \eta_\lambda(\hat{C}) \hat{S}^* (\hat{Z} - \hat{S} G_\lambda)\|_{\text{HS}} \quad (\text{B.35})$$

$$= \|S \hat{C}_\lambda^{-1/2}\| \| \hat{C}_\lambda^{1/2} \eta_\lambda(\hat{C}) \hat{C}_\lambda^{1/2} \| \| \hat{C}_\lambda^{-1/2} C_\lambda^{1/2} \| \| C_\lambda^{-1/2} \hat{S}^* (\hat{Z} - \hat{S} G_\lambda) \|_{\text{HS}}. \quad (\text{B.36})$$

Note that since  $\eta_\lambda$  is a filter and  $\hat{C}$  and  $\hat{C}_\lambda^{1/2}$  have same spectral decomposition,

$$\| \hat{C}_\lambda^{1/2} \eta_\lambda(\hat{C}) \hat{C}_\lambda^{1/2} \| \leq \sup_{\sigma \in (0, \kappa^2]} (\sigma + \lambda) \eta_\lambda(\sigma) \leq q_1. \quad (\text{B.37})$$

Moreover, since  $\hat{C} = \hat{S}^* \hat{S}$ , then

$$C_\lambda^{-1/2} \hat{S}^* (\hat{Z} - \hat{S} G_\lambda) = C_\lambda^{-1/2} (\hat{S}^* \hat{Z} - \hat{C} G_\lambda). \quad (\text{B.38})$$

We now focus on the second term of the sum in (B.33). Let  $r_\lambda(\sigma) = 1 - \sigma \eta_\lambda(\sigma)$ . Since  $\hat{C} = \hat{S}^* \hat{S}$ , we have

$$S \eta_\lambda(\hat{C}) \hat{S}^* \hat{S} G_\lambda - Z = (S G_\lambda - Z) - S r_\lambda(\hat{C}) G_\lambda. \quad (\text{B.39})$$

In particular, since  $L = S S^*$  and by definition of  $G_\lambda$  we have

$$S G_\lambda - Z = -(I - S S^* L_\lambda^{-1}) Z = -(I - L L_\lambda^{-1}) Z = -\lambda L_\lambda^{-1} Z, \quad (\text{B.40})$$

since  $(I - L L_\lambda^{-1}) = (L_\lambda - L) L_\lambda^{-1} = (L + \lambda - L) L_\lambda^{-1} = \lambda L_\lambda^{-1}$ . Moreover

$$S r_\lambda(\hat{C}) G_\lambda = (S \hat{C}_\lambda^{-1/2}) (\hat{C}_\lambda^{1/2} r_\lambda(\hat{C}) \hat{C}_\lambda^{1/2}) (\hat{C}_\lambda^{-1/2} S^*) (L_\lambda^{-1} Z). \quad (\text{B.41})$$

Now note that by definition of  $r_\lambda$  we have

$$\| \hat{C}_\lambda^{1/2} r_\lambda(\hat{C}) \hat{C}_\lambda^{1/2} \| \leq \sup_{\sigma \in (0, \kappa^2]} (1 - \sigma \eta_\lambda(\sigma)) (\sigma + \lambda) \leq q_2 \lambda. \quad (\text{B.42})$$

To conclude, since  $\|S \hat{C}_\lambda^{-1/2}\| \leq \|C_\lambda^{1/2} \hat{C}_\lambda^{-1/2}\|$  (see e.g. Rudi et al., 2015), we have

$$\|S \eta_\lambda(\hat{C}) \hat{S}^* \hat{S} G_\lambda - Z\|_{\text{HS}} \leq 2(1 + q_2) \lambda \| \hat{C}_\lambda^{-1/2} C_\lambda^{1/2} \|^2 \|L_\lambda^{-1} Z\|_{\text{HS}}. \quad (\text{B.43})$$

■



### B.3. Statistical Analysis

In this section we use the decomposition in Theorem B.6 to derive statistical learning rates for the estimators  $\hat{g}$ . To this end, we recall the following result.

**Lemma B.7 (Carratino et al. (2018), Lemma 3)** *Let  $\delta \in (0, 1)$ . When  $\lambda \geq \frac{9\kappa^2}{n} \log \frac{n}{\delta}$  then the following holds with probability at least  $1 - \delta$*

$$\|\widehat{C}_\lambda^{-1/2} C_\lambda^{1/2}\|^2 \leq 2. \quad (\text{B.44})$$

**Proof** Apply Lemma 3 of Carratino et al. (2018) with  $R = n$  and  $\zeta_i = \phi(x_i)$ .  $\blacksquare$

We now prove an intermediate result that will be instrumental in proving the excess risk bounds of the estimators  $\hat{g}$ . We recall the definition of effective dimension given in (50), that will be useful in the following, namely

$$d_{\text{eff}}(\lambda) = \text{Tr}(C(C + \lambda I)^{-1}), \quad \forall \lambda > 0. \quad (\text{B.45})$$

**Proposition B.8** *Let  $\delta \in (0, 1)$  and  $\lambda > 0$ . The following holds with probability at least  $1 - \delta$ ,*

$$\|C_\lambda^{-1/2}(\widehat{S}^* \widehat{Z} - \widehat{C} G_\lambda)\|_{\text{HS}} \leq \lambda \|L_\lambda^{-1} Z\|_{\text{HS}} + \frac{4\kappa \log \frac{2}{\delta}}{\sqrt{\lambda n}} (\kappa \|L_\lambda^{-1/2} Z\|) \quad (\text{B.46})$$

$$+ \sqrt{\frac{16(d_{\text{eff}}(\lambda) + \kappa^2 \lambda \|L_\lambda^{-1} Z\|_{\text{HS}}^2) \log \frac{2}{\delta}}{n}}. \quad (\text{B.47})$$

**Proof** For any  $i = 1, \dots, n$  we consider the random linear operator

$$\zeta_i = C_\lambda^{-1/2} (\phi(x_i) \otimes \varphi(y_i) - (\phi(x_i) \otimes \phi(x_i)) G_\lambda), \quad (\text{B.48})$$

as a vector in the space of Hilbert-Schmidt operators. Hence, taking the expectation with respect to a random sample of training points  $(x_i, y_i)_{i=1}^n$  from  $\rho$ ,

$$\mathbb{E} \zeta_i = C_\lambda (S^* Z - C G_\lambda), \quad (\text{B.49})$$

since  $\mathbb{E}[\varphi(y)|x_i] = g^*(x_i)$ . Moreover, since  $\|C_\lambda^{-1/2} S^*\| = \|C_\lambda^{-1/2} C^{1/2}\| \leq 1$ , following the same reasoning in the proof of Theorem B.6 to obtain (B.40), we have

$$\|\mathbb{E} \zeta_i\|_{\text{HS}} = \|C_\lambda^{-1/2} (S^* Z - C G_\lambda)\|_{\text{HS}} \leq \|C_\lambda^{-1/2} S^*\| \|Z - S G_\lambda\|_{\text{HS}} \leq \lambda \|L_\lambda^{-1} Z\|_{\text{HS}}. \quad (\text{B.50})$$

Now we need to study the moments of  $Z_i$  to obtain the final result. First note that

$$\|G_\lambda\| \leq \|S^* L_\lambda^{-1/2}\| \|L_\lambda^{-1/2} Z\| \leq \|L_\lambda^{-1/2} Z\|, \quad (\text{B.51})$$

since  $\|S^* L_\lambda^{-1/2}\| = \|L^{-1/2} L_\lambda^{-1/2}\| \leq 1$ . Recall that  $\sup_{y \in \mathcal{Y}} \|\varphi(y)\| \leq 1$ . Then, we have

$$\|\zeta_i\|_{\text{HS}} \leq \|C_\lambda^{-1/2}\| \|\phi(x_i)\| (\|\varphi(y_i)\| + \|G_\lambda\| \|\phi(x_i)\|) \leq \lambda^{-1/2} \kappa (1 + \kappa \|L_\lambda^{-1/2} Z\|). \quad (\text{B.52})$$

Hence, for any  $p \geq 2$

$$\mathbb{E}\|\zeta_i - \mathbb{E}\zeta_i\|_{\text{HS}}^p \leq \mathbb{E}\|\zeta_i - \zeta'_i\|_{\text{HS}}^p \leq 2^p \mathbb{E}\|\zeta_i\|_{\text{HS}}^p \quad (\text{B.53})$$

$$\leq 4 \left( 2\lambda^{-1/2} \kappa (1 + \kappa \|L_\lambda^{-1/2} Z\|) \right)^{p-2} \mathbb{E}\|\zeta_i\|_{\text{HS}}^2. \quad (\text{B.54})$$

Moreover, denote by  $\sigma^2(x)$  the conditional variance  $\sigma^2(x) = \mathbb{E}_{y_i}[\|\varphi(y_i) - g^*(x_i)\|^2 | x_i]$ . Since  $g^*(x) = \mathbb{E}[\varphi(y) | x]$  for any  $x$  in the domain of  $\rho_{\mathcal{X}}$ , we have  $\sigma(x) \leq 1$ . Hence,

$$\mathbb{E}\|\zeta_i\|_{\text{HS}}^2 = \mathbb{E}\|\varphi(y_i) - g_\lambda(x_i)\|^2 \|C_\lambda^{-1/2} \phi(x_i)\|^2 \quad (\text{B.55})$$

$$= \mathbb{E}_{x_i} \|C_\lambda^{-1/2} \phi(x_i)\|^2 \mathbb{E}_{y_i}[\|\varphi(y_i) - g_\lambda(x_i)\|^2 | x_i] \quad (\text{B.56})$$

$$= \mathbb{E}_{x_i} \|C_\lambda^{-1/2} \phi(x_i)\|^2 (\sigma(x)^2 + \|g^*(x) - g_\lambda(x)\|^2) \quad (\text{B.57})$$

$$\leq d_{\text{eff}}(\lambda) + \kappa^2/\lambda \mathbb{E}\|g^*(x) - g_\lambda(x)\|^2, \quad (\text{B.58})$$

where the last inequality follows by observing that

$$\mathbb{E}\|C_\lambda^{-1/2} \phi(x)\|^2 = \mathbb{E} \text{Tr}(C_\lambda^{-1} \phi(x) \otimes \phi(x)) = \text{Tr}(C_\lambda^{-1} C) = d_{\text{eff}}(\lambda), \quad (\text{B.59})$$

and also

$$d_{\text{eff}}(\lambda) = \mathbb{E} \|C_\lambda^{-1/2} \phi(x)\|^2 \leq \mathbb{E} \|C_\lambda^{-1/2}\|^2 \|\phi(x)\|^2 \leq \lambda^{-1} \kappa^2. \quad (\text{B.60})$$

Recall from Theorem A.2 that  $\mathcal{R}(g) - \mathcal{R}(g^*) = \mathbb{E}\|g(x) - g^*(x)\|^2$ . Then, by applying Theorem B.2, we have

$$\mathbb{E}\|g^*(x) - g_\lambda(x)\|^2 = \|g_\lambda - g^*\|_{L^2(X, \rho)}^2 \quad (\text{B.61})$$

$$= \|SG_\lambda - Z\|_{\text{HS}}^2 \quad (\text{B.62})$$

$$= \|(SS^* L_\lambda^{-1} - I)Z\|_{\text{HS}}^2 \quad (\text{B.63})$$

$$= \lambda^2 \|L_\lambda^{-1} Z\|_{\text{HS}}^2. \quad (\text{B.64})$$

Finally, we have for any  $p \geq 2$

$$\mathbb{E}\|\zeta_i - \mathbb{E}\zeta_i\|_{\text{HS}}^p \leq \frac{1}{2} p! \left( 8d_{\text{eff}}(\lambda) + 8\kappa^2 \lambda \|L_\lambda^{-1} Z\|_{\text{HS}}^2 \right) \left( \frac{2\kappa}{\sqrt{\lambda}} (1 + \kappa \|L_\lambda^{-1/2} Z\|) \right)^{p-2}. \quad (\text{B.65})$$

We can now apply Bernstein inequality as in (Rudi and Rosasco (2017) Proposition 2). We have

$$\left\| \frac{1}{n} \sum_{i=1}^n \zeta_i - \mathbb{E}\zeta_i \right\|_{\text{HS}} \leq \frac{4\kappa \log \frac{2}{\delta}}{\sqrt{\lambda n}} (1 + \kappa \|L_\lambda^{-1/2} Z\|) + \sqrt{\frac{16(d_{\text{eff}}(\lambda) + \kappa^2 \lambda \|L_\lambda^{-1} Z\|_{\text{HS}}^2) \log \frac{2}{\delta}}{n}} \quad (\text{B.66})$$

holds with probability  $1 - \delta$ .

The proof is concluded by observing that

$$\|C_\lambda^{-1/2} (\widehat{S}^* \widehat{Z} - \widehat{C} G_\lambda)\|_{\text{HS}} = \left\| \frac{1}{n} \sum_{i=1}^n \zeta_i \right\|_{\text{HS}} \leq \left\| \frac{1}{n} \sum_{i=1}^n \zeta_i - \mathbb{E}\zeta_i \right\|_{\text{HS}} + \|\mathbb{E}\zeta\|_{\text{HS}}. \quad (\text{B.67})$$

■

**Theorem B.9** *Under the assumptions of Theorem B.6, let  $\delta \in (0, 1)$ . Then, for  $\lambda \geq \frac{9\kappa^2}{n} \log \frac{n}{\delta}$  the following holds with probability at least  $1 - \delta$*

$$|\mathcal{R}(g_{n\lambda}) - \mathcal{R}(g^*)|^{1/2} \leq \frac{8q_1\kappa \log \frac{2}{\delta}}{\sqrt{\lambda n}} (1 + \kappa \|L_\lambda^{-1/2} Z\|) \quad (\text{B.68})$$

$$+ \sqrt{\frac{64q_1^2(d_{\text{eff}}(\lambda) + \kappa^2\lambda \|L_\lambda^{-1} Z\|_{\text{HS}}^2) \log \frac{4}{\delta}}{n}} \quad (\text{B.69})$$

$$+ 2(2 + q_1 + 2q_2) \lambda \|L_\lambda^{-1} Z\|_{\text{HS}} \quad (\text{B.70})$$

**Proof** The result is obtained by decomposing the risk with Theorem B.6, and controlling in high probability both terms  $\beta = \|\widehat{C}_\lambda^{-1/2} C_\lambda^{1/2}\|^2$  and  $\|C_\lambda^{-1/2}(\widehat{S}^* \widehat{Z} - \widehat{C}G_\lambda)\|_{\text{HS}}$  and then taking the intersection bound of the two events.  $\blacksquare$

#### B.4. Universal Consistency

Now we are ready to give the universal consistency result.

**Theorem 8 (Universal Consistency)** *Let  $\mathcal{Z}$  be a compact set and  $\Delta : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$  admit an ILE. Let  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  be a bounded universal reproducing kernel. For any  $n \in \mathbb{N}$  and any distribution  $\rho$  on  $\mathcal{X} \times \mathcal{Y}$  let  $f_n : \mathcal{X} \rightarrow \mathcal{Z}$  be the estimator in (20) trained on  $(x_i, y_i)_{i=1}^n$  points independently sampled from  $\rho$  and with weights  $\alpha$  defined as:*

(a) (Ridge Regression) in (14) with  $\lambda_n = n^{-1/2}$ , or

(b) (L2-Boosting) in (21) with step-size  $\nu < 1/\kappa^2$  and  $t_n = n^{1/2}$ , or

(c) (PCR) in (22) with  $\lambda_n = n^{-1/2}$ .

Then,

$$\lim_{n \rightarrow +\infty} \mathcal{E}(f_n) = \mathcal{E}(f^*) \quad \text{with probability } 1 \quad (45)$$

**Proof** Recall that by Theorem B.8, for any  $\lambda \geq \frac{9\kappa^2}{n} \log \frac{n}{\delta}$  the following holds with probability at least  $1 - \delta$

$$|\mathcal{R}(g_{n\lambda}) - \mathcal{R}(g^*)|^{1/2} \leq \frac{8q_1\kappa \log \frac{4}{\delta}}{\sqrt{\lambda n}} (1 + \kappa \|L_\lambda^{-1/2} Z\|) \quad (\text{B.71})$$

$$+ \sqrt{\frac{64q_1^2(d_{\text{eff}}(\lambda) + \kappa^2\lambda \|L_\lambda^{-1} Z\|_{\text{HS}}^2) \log \frac{4}{\delta}}{n}} \quad (\text{B.72})$$

$$+ 2(2 + q_1 + 2q_2) \lambda \|L_\lambda^{-1} Z\|_{\text{HS}}. \quad (\text{B.73})$$

In particular, let  $n_0$  be sufficiently large, such that  $n_0^{-1/2} \geq \frac{9\kappa^2}{n_0} \log \frac{n_0}{\delta}$ . For any  $n \in \mathbb{N}$  with  $n \geq n_0$ , let  $\lambda_n = n^{-1/2}$ . Recall that  $d_{\text{eff}}(\lambda) \leq \kappa^2/\lambda$  (see (B.60)) and note that

$\|L_\lambda^{-1/2}Z\| \leq \lambda^{-1/2}\|Z\|$ ,  $\|L_\lambda^{-1}Z\|_{\text{HS}} \leq \lambda^{-1}\|Z\|_{\text{HS}}$  and  $\|Z\| \leq \|Z\|_{\text{HS}}$ . Applying Theorem B.8 guarantees that the inequality

$$\begin{aligned} |\mathcal{R}(g_{n\lambda}) - \mathcal{R}(g^*)|^{1/2} &\leq [8q_1\kappa(1 + \kappa\|Z\|)] n^{-1/2} \log(4/\delta) \\ &+ [64q_1^2\kappa^2(1 + \|Z\|_{\text{HS}}^2)]^{1/2} n^{-1/4} (\log(4/\delta))^{1/2} \\ &+ 2(2 + q_1 + 2q_2) \lambda \|L_\lambda^{-1}Z\|_{\text{HS}}, \end{aligned} \quad (\text{B.74})$$

holds with probability at least  $1 - \delta$ . We denote this event by  $E_{n,\delta}$ .

Recall that  $L$  is a compact operator (actually Hilbert-Schmidt), hence it admits an eigendecomposition  $L = \sum_{i \in \mathbb{N}} \sigma_i u_i \otimes u_i$ , with  $\sigma_i \geq \sigma_j > 0$  for  $1 \leq i \leq j \in \mathbb{N}$  and  $(u_i)_{i \in \mathbb{N}}$  is a set of orthonormal functions in  $L^2(\mathcal{X}, \rho_{\mathcal{X}})$ . Note that  $(u_i)_{i \in \mathbb{N}}$  is an orthonormal basis of  $L^2(\mathcal{X}, \rho_{\mathcal{X}})$ . To show this, consider  $W \subseteq \mathcal{X}$  the support of  $\rho_{\mathcal{X}}$ . Note that  $W$  is compact and Polish since it is a closed subset of the compact Polish space  $\mathcal{X}$ . Let  $\mathcal{L}$  be the RKHS  $\mathcal{L} = \overline{\text{span}\{k(x, \cdot) \mid x \in W\}}$ , with same inner product of  $\mathcal{F}$ . Note that  $\mathcal{L}$  is separable since it is the image of a compact space via a continuous function. By definition of universality for the kernel  $k$ , the set  $\mathcal{L}$  is dense in  $C(W)$ . Additionally, by Corollary 5 in (Micchelli et al., 2006) we have  $C(W) = \overline{\text{span}\{u_i \mid i \in \mathbb{N}\}}$ . Thus, since  $C(W)$  is dense in  $L^2(\mathcal{X}, \rho_{\mathcal{X}})$ , we can conclude that  $(u_i)_{i \in \mathbb{N}}$  is a basis of  $L^2(\mathcal{X}, \rho_{\mathcal{X}})$ .

We now focus on  $\lambda \|L_\lambda^{-1}Z\|_{\text{HS}}$ . In particular, we want to express  $\|L_\lambda^{-1}Z\|_{\text{HS}}^2$  in terms of the basis  $(u_i)_{i \in \mathbb{N}}$  associated to  $L$ . In particular, since  $(u_i)_{i \in \mathbb{N}}$  is a basis for  $L^2(\mathcal{X}, \rho_{\mathcal{X}})$  and  $(L + \lambda I)^{-1}u_i = (\sigma_i + \lambda)^{-1}u_i$  for any  $i \in \mathbb{N}$ , we have

$$\lambda_n^2 \|(L + \lambda_n)^{-1}Z\|_{\text{HS}}^2 = \lambda_n^2 \text{Tr}((L + \lambda_n)^{-1}ZZ^*(L + \lambda_n)^{-1}) = \lambda_n^2 \sum_{i \in \mathbb{N}} \frac{\langle u_i, ZZ^*u_i \rangle_{L^2}}{(\sigma_i + \lambda)^2}. \quad (\text{B.75})$$

Now let  $t_n = n^{-1/4}$ , and  $T_n = \{i \in \mathbb{N} \mid \sigma_i \geq t_n\} \subset \mathbb{N}$ . Denote by  $w_i^2 := \langle u_i, ZZ^*u_i \rangle_{L^2}$  and note that  $(w_i)_{i \in \mathbb{N}}$  is square summable and  $\sum_{i \in \mathbb{N}} w_i^2 = \|Z\|_{\text{HS}}^2 < \infty$ . For any  $n \in \mathbb{N}$ , we have

$$\lambda_n^2 \|(L + \lambda_n)^{-1}Z\|_{\text{HS}}^2 = \sum_{i \in T_n} \frac{\lambda_n^2 w_i^2}{(\sigma_i + \lambda_n)^2} + \sum_{i \in \mathbb{N} \setminus T_n} \frac{\lambda_n^2 w_i^2}{(\sigma_i + \lambda_n)^2} \quad (\text{B.76})$$

$$\leq \frac{\lambda_n^2}{t_n^2} \sum_{i \in T_n} w_i^2 + \sum_{i \in \mathbb{N} \setminus T_n} w_i^2 \leq \|Z\|_{\text{HS}}^2 n^{-1/4} + \sum_{i \in \mathbb{N} \setminus T_n} w_i^2 \quad (\text{B.77})$$

since  $\lambda_n/t_n = n^{-1/2}/n^{-1/4} = n^{-1/4}$ . Since the series  $\sum_{i \in \mathbb{N}} w_i^2$  is convergent, we have  $\sum_{i \in \mathbb{N} \setminus T_n} w_i^2 \rightarrow 0$  as  $n \rightarrow +\infty$ . We conclude that

$$0 \leq \lim_{n \rightarrow \infty} \lambda_n^2 \|(L + \lambda_n)^{-1}Z\|_{\text{HS}}^2 \leq \lim_{n \rightarrow \infty} \|Z\|_{\text{HS}}^2 n^{-1/4} + \sum_{i \in \mathbb{N} \setminus T_n} w_i^2 = 0. \quad (\text{B.78})$$

Now, let  $\delta_n = n^{-2}$  and  $A_n = E_{n,\delta_n}^c$  be the complementary event to  $E_{n,\delta_n}$  characterized by (B.74). For any  $n \geq n_0 := 200(1 + \kappa^2)$  we have  $\lambda_n \geq \frac{9\kappa^2}{n} \log n^3$  and the event  $E_{n,\delta_n}$  holds with probability at least  $1 - \delta_n$ . Equivalently, the probability of  $A_n$  is upper bounded by  $\delta_n$ . Since  $\sum_{n=n_0+1}^{+\infty} \delta_n < +\infty$ , we can apply the Borel-Cantelli lemma (Theorem 8.3.4. pag 263 of Dudley (2002)) on the sequence  $(E_{n,\delta_n})_{n \in \mathbb{N}}$  and conclude that the statement

$$\lim_{n \rightarrow \infty} \mathcal{R}(g_{n\lambda_n}) - \mathcal{R}(g^*) > 0, \quad (\text{B.79})$$

holds with probability 0. Thus, the converse statement

$$\lim_{n \rightarrow \infty} \mathcal{R}(g_{n\lambda_n}) - \mathcal{R}(g^*) = 0. \quad (\text{B.80})$$

holds with probability 1. The final result is obtained by applying the comparison inequality between the surrogate problem and the original excess risk from Theorem 7.  $\blacksquare$

## B.5. Learning Rates

In this section we study the generalization properties of the proposed estimators. We address this question by considering the special case where the solution  $g^*$  of the expected surrogate risk belongs to the same hypotheses space  $\mathcal{H} \otimes \mathcal{F}$  where our estimator belongs to. We start this analysis by recalling that in this case,  $g^*$  admits a closed form solution.

**Lemma B.10** *Let  $\Delta : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$  satisfy Theorem 2 for suitable  $\psi, \varphi$  and  $\mathcal{H}$ . Assume that the surrogate expected risk minimization of  $\mathcal{R}$  at (17) attains a minimum on  $\mathcal{H} \otimes \mathcal{F}$ . Then the minimizer  $g^* \in \mathcal{H} \otimes \mathcal{F}$  of  $\mathcal{R}$  with minimal norm  $\|\cdot\|_{\mathcal{H} \otimes \mathcal{F}}$  is of the form*

$$g^*(x) = G_* \phi(x), \quad \forall x \in \mathcal{X} \quad \text{with} \quad G_* = C^\dagger S^* Z : \mathcal{F} \rightarrow \mathcal{H}. \quad (\text{B.81})$$

**Proof** Let  $g \in \mathcal{H} \otimes \mathcal{F}$  such that  $g(x) = G\phi(x), \forall x \in \mathcal{X}$  for some linear operator  $G \in \mathcal{H} \otimes \mathcal{F}$ . We have

$$\mathcal{R}(g) = \int_{\mathcal{X} \times \mathcal{Y}} \|G\phi(x) - \psi(y)\|_{\mathcal{H}}^2 d\rho(x, y) \quad (\text{B.82})$$

$$= \int_{\mathcal{X} \times \mathcal{Y}} \text{Tr}(G(\phi(x) \otimes \phi(x))G^*) - 2\text{Tr}(G(\phi(x) \otimes \psi(y))) + \|\psi(y)\|_{\mathcal{H}}^2 d\rho_{\mathcal{X}}(x, y) \quad (\text{B.83})$$

$$= \text{Tr}(GC G^*) - 2\text{Tr}(GS^* Z) + \text{const}, \quad (\text{B.84})$$

where we have used Theorem B.1 and the linearity of the trace. The derivation above implies that  $\mathcal{R}$  is a convex quadratic functional since  $C$  is positive semidefinite. Hence,  $\mathcal{R}$  attains a minimum on  $\mathcal{H} \otimes \mathcal{F}$  if and only if the range of  $S^* Z$  is contained in the range of  $C$ , namely  $\text{Ran}(S^* Z) \subseteq \text{Ran}(C) \subset \mathcal{F}$  (see Engl et al. (1996) Chap. 2). In this case  $G = C^\dagger S^* Z : \mathcal{F} \rightarrow \mathcal{H}$  exists and is the minimum norm minimizer for  $\mathcal{R}$ , as desired.  $\blacksquare$

We recall here the two main assumptions required in the following.

**Assumption 1 (Source condition)** *There exists  $r \geq 0$  and  $h \in \mathcal{H} \otimes \mathcal{F}$  for which*

$$g^* = (C^r \otimes I) h. \quad (49)$$

*The norm of  $\|h\|_{\mathcal{H} \otimes \mathcal{F}}$  will be denoted by  $R := \|h\|_{\mathcal{H} \otimes \mathcal{F}}$ .*

**Assumption 2 (Capacity condition)** *There exists  $\gamma \in [0, 1]$  and  $Q > 0$  for which*

$$d_{\text{eff}}(\lambda) \leq Q\lambda^{-\gamma}, \quad \forall \lambda > 0. \quad (51)$$

We are ready to prove the main result characterizing the learning rates of the proposed estimators.

**Theorem B.11** *Let  $\mathcal{H}$  be a Hilbert space and let  $\varphi : \mathcal{Y} \rightarrow \mathcal{H}$  a continuous map from  $\mathcal{Y}$  to  $\mathcal{H}$ . Let  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  be a continuous reproducing kernel on  $\mathcal{X}$  with associated RKHS  $\mathcal{F}$  such that  $\kappa^2 := \sup_{x \in \mathcal{X}} k(x, x) < +\infty$ . Let  $\rho$  be a distribution on  $\mathcal{X} \times \mathcal{Y}$  and let the corresponding  $g^*$  defined in (16) satisfy Assumption 1. Let also Assumption 2 hold. Let  $\delta \in (0, 1]$  and  $n_0$  sufficiently large such that  $n_0^{-1/(1+2r+\gamma)} \geq \frac{9\kappa^2}{n_0} \log \frac{n_0}{\delta}$ . Let  $n \in \mathbb{N}$ ,  $n \geq n_0$  and  $\lambda \geq \frac{9\kappa^2}{n} \log \frac{n}{\delta}$ . Let  $g_n : \mathcal{X} \rightarrow \mathcal{H}$  be a spectral filtering estimator of the form introduced in Theorem B.5 trained on  $n$  points randomly sampled from  $\rho$ . Then, the following holds with probability at least  $1 - \delta$ ,*

$$|\mathcal{R}(g_n) - \mathcal{R}(g^*)|^{1/2} \leq M n^{-\frac{r+1/2}{2r+1+\gamma}} \quad (\text{B.85})$$

Where the constant  $M$  is

$$M = M(Q, \|H\|, \delta, q_1, q_2) = 8q_1 \left[ \kappa(1 + \kappa\|H\|_{\text{HS}}) + \sqrt{(Q + \kappa^2\|H\|_{\text{HS}}^2)} \right] \log \frac{4}{\delta} \quad (\text{B.86})$$

$$+ 2(2 + q_1 + 2q_2)\|H\|_{\text{HS}}. \quad (\text{B.87})$$

**Proof** According to the excess risk bound in Theorem B.9 we have that the following holds with probability at least  $1 - \delta$

$$|\mathcal{R}(g_{n\lambda}) - \mathcal{R}(g^*)|^{1/2} \leq \frac{8q_1\kappa \log \frac{4}{\delta}}{\sqrt{\lambda n}} (1 + \kappa\|L_\lambda^{-1/2}Z\|) \quad (\text{B.88})$$

$$+ \sqrt{\frac{64 q_1^2 (d_{\text{eff}}(\lambda) + \kappa^2 \lambda \|L_\lambda^{-1}Z\|_{\text{HS}}^2) \log \frac{4}{\delta}}{n}} \quad (\text{B.89})$$

$$+ 2(2 + q_1 + 2q_2) \lambda \|L_\lambda^{-1}Z\|_{\text{HS}}. \quad (\text{B.90})$$

From Assumption 1, we have  $g^*(x) = (C^r \otimes I) h \phi(x) = H^* C^r \phi(x)$  for any  $x \in \mathcal{X}$ , where  $H : \mathcal{F} \rightarrow \mathcal{H}$  is the Hilbert-Schmidt operator corresponding to  $h$  under the canonical isomorphism between  $\text{HS}(\mathcal{F}, \mathcal{H})$  and  $\mathcal{H} \otimes \mathcal{F}$ . In particular,  $\|H\|_{\text{HS}} = \|h\|_{\mathcal{H} \otimes \mathcal{F}}$ . By Assumption 1  $\|h\|_{\mathcal{H} \otimes \mathcal{F}} = R$ . Therefore, we can characterize  $Z = S C^r H$ . Indeed, recall that for any  $w \in \mathcal{H}$ , by definition of  $Z$ , we have

$$(Zw)(\cdot) = \langle w, g^*(\cdot) \rangle_{\mathcal{H}} = \langle w, H^* C^r \phi(\cdot) \rangle_{\mathcal{H}} = \langle C^r H w, \phi(\cdot) \rangle_{\mathcal{H}} = (S C^r H w)(\cdot) \quad (\text{B.91})$$

Then, denote by  $(1/2 - r)_+ = \max(0, 1/2 - r)$ , we have

$$\|L_\lambda^{-1}Z\|_{\text{HS}} = \|L_\lambda^{-1}S C^r H\|_{\text{HS}} \leq \lambda^{-(1/2-r)_+} R, \quad (\text{B.92})$$

and, analogously,

$$\|L_\lambda^{-1/2}Z\|_{\text{HS}} \leq R. \quad (\text{B.93})$$

Then, let  $\lambda = n^{-1/(1+2r+\gamma)}$  and  $n \geq n_0$  with  $n_0$  such that  $n_0^{-1/(1+2r+\gamma)} \geq \frac{9\kappa^2}{n_0} \log \frac{n_0}{\delta}$ . From Assumption 2, we have  $d_{\text{eff}}(\lambda) \leq Q\lambda^{-\gamma}$ . Therefore,

$$|\mathcal{R}(g_{n\lambda}) - \mathcal{R}(g^*)| \leq 8q_1 \left[ \kappa(1 + \kappa R) + \sqrt{(Q + \kappa^2 R^2)} \right] n^{-\frac{r+1/2}{2r+1+\gamma}} \log \frac{4}{\delta} \quad (\text{B.94})$$

$$+ 2(2 + q_1 + 2q_2)R n^{-\frac{r+1/2}{2r+1+\gamma}}, \quad (\text{B.95})$$

with probability at least  $1 - \delta$ . ■

Theorem 11 from the main paper is a direct consequence the result above when considering specific spectral filters.

**Theorem 11 (Refined Learning Rates)** *Under the same notation and assumptions of Thm. 9 and under the additional Assumptions 1 and 2, let  $\delta \in (0, 1]$  and  $n_0$  sufficiently large such that  $n_0^{-1/(1+2r+\gamma)} \geq \frac{9\kappa^2}{n_0} \log \frac{n_0}{\delta}$ . For any  $n \geq n_0$ , the following estimators  $f_n : \mathcal{X} \rightarrow \mathcal{Z}$  trained on  $n$  points independently sampled from  $\rho$  are such that, with probability at least  $1 - \delta$*

$$\mathcal{E}(f_n) - \mathcal{E}(f^*) \leq c_\Delta \mathbf{m} \mathbf{q} \log(4/\delta) n^{-\frac{r+1/2}{2r+\gamma+1}}, \quad (52)$$

with

$$\mathbf{m} = 16 \left( \kappa(1 + \kappa R) + \kappa \sqrt{Q + R^2 + R} \right), \quad (53)$$

and  $\mathbf{q}$  defined as follows. This holds for estimators  $f_n$  of the form (20) with corresponding weights  $\alpha$  defined as:

(a) (Ridge Regression) in (14) with  $\lambda_n = n^{-\frac{1}{2r+\gamma+1}}$ . With  $\mathbf{q} \leq 3$ .

(b) (L2-Boosting) in (21) with  $\nu < 1/\kappa^2$  and  $t_n = n^{\frac{1}{2r+\gamma+1}}$ . With  $\mathbf{q} \leq 2 + 2\nu + e^{\nu-1}/\nu$ .

(c) (Principal Component Regression) in (22) with  $\lambda_n = n^{-\frac{1}{2r+\gamma+1}}$ . With  $\mathbf{q} \leq 5$ .

**Proof** Let  $g_n : \mathcal{X} \rightarrow \mathcal{H}$  be an estimator satisfying the hypotheses of Theorem B.11 and let  $f_n : \mathcal{X} \rightarrow \mathcal{Z}$  be such that for any  $x \in \mathcal{X}$ ,  $f_n(x) = \text{argmin}_{z \in \mathcal{Z}} \langle \psi(z), g_n(x) \rangle_{\mathcal{H}}$ . By applying the comparison inequality from Theorem 7, we have that

$$\mathcal{E}(f_n) - \mathcal{E}(f^*) \leq 2c_\Delta |\mathcal{R}(g_n) - \mathcal{R}(g^*)|^{1/2} \leq 2c_\Delta \mathbf{M} n^{-\frac{r+1/2}{2r+1+\gamma}}, \quad (\text{B.96})$$

holds with probability at least  $1 - \delta$ . Now, note that  $\log(4/\delta) > 1$  since  $\delta \leq 1$ , we have that the constant  $\mathbf{M}$  is upper bounded by

$$\mathbf{M} \leq \left( 8q_1 \left[ \kappa(1 + \kappa \|H\|_{\text{HS}}) + \sqrt{(Q + \kappa^2 \|H\|_{\text{HS}}^2)} \right] + 2(2 + q_1 + 2q_2) \|H\|_{\text{HS}} \right) \log \frac{4}{\delta}. \quad (\text{B.97})$$

Replacing the quantities  $q_1$  and  $q_2$  from Theorem B.4 associated to the corresponding estimators we obtain the required upper bounds for  $\mathbf{m}$  stated in the thesis of the theorem. ■

We note that Theorem 9 is a corollary of Theorem 11 as shown below.

**Theorem 9 (Learning Rates)** *Let  $\mathcal{Z}$  be a compact set and  $\Delta : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$  admit an ILE with associated Hilbert space  $\mathcal{H}$ . Let  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  be a continuous reproducing kernel on  $\mathcal{X}$  with associated RKHS  $\mathcal{F}$  such that  $\kappa^2 := \sup_{x \in \mathcal{X}} k(x, x) < +\infty$ . Let  $\rho$  be a distribution on  $\mathcal{X} \times \mathcal{Y}$  and let the corresponding  $g^*$  defined in (16) be such that  $g^* \in \mathcal{H} \otimes \mathcal{F}$ . Let  $\delta \in (0, 1]$  and  $n_0$  sufficiently large such that  $n_0^{-1/2} \geq \frac{9\kappa^2}{n_0} \log \frac{n_0}{\delta}$ . Then, for any  $n \in \mathbb{N}$ , the following estimators  $f_n : \mathcal{X} \rightarrow \mathcal{Z}$  trained on  $n$  points independently sampled from  $\rho$  are such that, with probability at least  $1 - \delta$*

$$\mathcal{E}(f_n) - \mathcal{E}(f^*) \leq c_\Delta \mathfrak{m} \mathfrak{q} \log(4/\delta) n^{-1/4}, \quad (46)$$

with

$$\mathfrak{m} = 16 \left( \kappa(1 + \kappa \|g^*\|) + \kappa \sqrt{1 + \|g^*\|^2} + \|g^*\| \right), \quad (47)$$

and  $\mathfrak{q}$  defined as follows. This holds for estimators  $f_n$  of the form (20) with corresponding weights  $\alpha$  defined as:

- (a) (Ridge Regression) in (14) with  $\lambda_n = n^{-1/2}$ . With constant  $\mathfrak{q} \leq 3$ .
- (b) (L2-Boosting) in (21) with  $\nu < 1/\kappa^2$  and  $t_n = n^{1/2}$ . With constant  $\mathfrak{q} \leq 2 + 2\gamma + e^{\gamma-1}/\gamma$ .
- (c) (Principal Component Regression) in (22) with  $\lambda_n = n^{-1/2}$ . With constant  $\mathfrak{q} \leq 5$ .

**Proof** The result is a corollary of the theorem above, by considering that Assumption 2 is always satisfied with  $Q = \kappa^2$  and  $\gamma = 1$  and that when  $g^*$  is in  $\mathcal{G}$ , then Assumption 1 is satisfied with  $r = 0$  and  $h = g$ . ■

## Appendix C. Sufficient Conditions for ILE

In this section we provide more details related to the ILE definition introduced in Theorem 2. In particular we discuss the connection with the original framework considered in (Ciliberto et al., 2016) and prove the results reported in Section 6 providing sufficient conditions to determine whether a function admits an ILE.

### C.1. Relations with the “ILE” definition in (Ciliberto et al., 2016)

In Ciliberto et al. (2016), the ILE definition was introduced as the following assumption in the case  $\mathcal{Z} = \mathcal{Y}$  (see Assumption 1 in Ciliberto et al., 2016).

**Assumption 3** *There exists a separable Hilbert space  $\mathcal{H}$  with inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ , a continuous embedding  $\zeta : \mathcal{Y} \rightarrow \mathcal{H}$  and a bounded linear operator such that*

$$\Delta(z, y) = \langle \zeta(z), V\zeta(y) \rangle_{\mathcal{H}} \quad \forall z, y \in \mathcal{Y} \quad (C.1)$$

It can be noticed that the two definitions are quite similar one to the other. Both require the existence of a separable Hilbert space  $\mathcal{H}$  where the function  $\Delta$  assumes a “bilinear”



structure. However the definition above requires  $\mathcal{Z} = \mathcal{Y}$  and the existence of a linear operator combining a single feature map  $\zeta$ .

Despite these differences, the following result shows that the above assumption is equivalent to the ILE definition (Theorem 2) in the main paper.

**Proposition C.1 (Equivalence of ILE Definitions)** *A loss  $\Delta : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  admits an ILE if and only if it satisfies Assumption 3.*

**Proof** ( $\Rightarrow$ ). Let  $\Delta$  satisfy the ILE definition with Hilbert space  $\mathcal{H}$  and feature maps  $\psi, \varphi : \mathcal{Y} \rightarrow \mathcal{H}$ . We define  $\overline{\mathcal{H}} = \mathcal{H} \oplus \mathcal{H}$  and consider the map  $\zeta : \mathcal{Y} \rightarrow \overline{\mathcal{H}}$  such that  $\zeta(y) = (\psi(y), \varphi(y))^\top \in \overline{\mathcal{H}}$  for any  $y \in \mathcal{Y}$ . Moreover, we define  $V : \overline{\mathcal{H}} \rightarrow \overline{\mathcal{H}}$  the linear operator such that  $V(h_1, h_2) = (h_2, 0)$  for any  $h = (h_1, h_2) \in \overline{\mathcal{H}}$ . It is easy to see that  $V$  is bounded, and actually it has operator norm  $\|V\| = 1$ . Therefore, we have

$$\langle \zeta(z), V\zeta(y) \rangle_{\overline{\mathcal{H}}} = \langle (\psi(z), \varphi(z)), (\varphi(y), 0) \rangle_{\overline{\mathcal{H}}} = \langle \psi(z), \varphi(y) \rangle_{\mathcal{H}} + \langle \varphi(z), 0 \rangle_{\mathcal{H}} = \Delta(z, y) \quad (\text{C.2})$$

for any  $y, z \in \mathcal{Y}$  as desired. Hence  $\Delta$  satisfies Assumption 3 with associated Hilbert space  $\overline{\mathcal{H}}$ . Note that  $\zeta$  is continuous since both  $\psi$  and  $\varphi$  are continuous by the ILE definition and  $V$  is linear and bounded by construction.

( $\Leftarrow$ ). Let  $\Delta$  satisfy Assumption 3 with Hilbert space  $\mathcal{H}$ , feature map  $\zeta : \mathcal{Y} \rightarrow \mathcal{H}$  and linear operator  $V : \mathcal{H} \rightarrow \mathcal{H}$ . Let  $\Phi = \sup_{y \in \mathcal{Y}} \|\zeta(y)\|_{\mathcal{H}}$ . Then we take  $\psi, \varphi : \mathcal{Y} \rightarrow \mathcal{H}$  the functions such that  $\psi(z) = \Phi V\varphi(z)$  and  $\varphi(y) = \zeta(y)/\Phi$  for any  $z, y \in \mathcal{Y}$ . Clearly, we have

$$\langle \psi(z), \varphi(y) \rangle_{\mathcal{H}} = \langle \zeta(z), V\zeta(y) \rangle_{\mathcal{H}} = \Delta(z, y). \quad (\text{C.3})$$

By construction, we have  $\sup_{y \in \mathcal{Y}} \|\varphi(y)\| \leq 1$  and  $c_\Delta = \sup_{z \in \mathcal{Z}} \|\psi(z)\|_{\mathcal{H}} \leq \Phi^2 \|V\|$ .  $\blacksquare$

## C.2. ILE definition without “normalization” of $\varphi$

We point out that the requirement for  $\sup_{y \in \mathcal{Y}} \|\varphi(y)\|_{\mathcal{H}} \leq 1$  is not necessary but was introduced for the sake of exposition. We formalize this in the following.

**Lemma C.2 (“Unnormalized” ILE)** *Let  $\Delta : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$  be such that there exists a separable Hilbert space  $\mathcal{H}$  and two continuous bounded maps  $\bar{\psi} : \mathcal{Z} \rightarrow \mathcal{H}$  and  $\bar{\varphi} : \mathcal{Y} \rightarrow \mathcal{H}$ , such that  $\sup_{z \in \mathcal{Z}} \|\bar{\psi}(z)\|_{\mathcal{H}} \leq \psi_\Delta$  and  $\sup_{y \in \mathcal{Y}} \|\bar{\varphi}(y)\|_{\mathcal{H}} \leq \Phi_\Delta$ , with  $\psi_\Delta > 0$ ,  $\Phi_\Delta > 0$  and*

$$\Delta(z, y) = \langle \bar{\psi}(z), \bar{\varphi}(y) \rangle_{\mathcal{H}}, \quad (\text{C.4})$$

for every  $z \in \mathcal{Z}$  and  $y \in \mathcal{Y}$ . Then  $\Delta$  admits an ILE with  $c_\Delta \leq \psi_\Delta \Phi_\Delta$ .

**Proof** The result is easy to prove by taking  $\psi : \mathcal{Z} \rightarrow \mathcal{H}$  and  $\varphi : \mathcal{Y} \rightarrow \mathcal{H}$  such that  $\psi(z) = \sup_{y' \in \mathcal{Y}} \Phi_\Delta \bar{\psi}(z)$  for any  $z \in \mathcal{Z}$  and  $\varphi(y) = \bar{\varphi}(y)/\Phi_\Delta$  for any  $y \in \mathcal{Y}$ . Indeed it is straightforward to see that the characterization of  $\Delta$  in terms of the inner product between  $\psi$  and  $\varphi$  still holds and, by construction,  $\sup_{y \in \mathcal{Y}} \|\varphi(y)\|_{\mathcal{H}} \leq 1$  and  $c_\Delta = \sup_{z \in \mathcal{Z}} \|\psi(z)\|_{\mathcal{H}} \leq \psi_\Delta \Phi_\Delta$ .  $\blacksquare$

### C.3. Finite $\mathcal{Y}$ or $\mathcal{Z}$

We now focus on proving the sufficient conditions to guarantee  $\Delta$  to admit an ILE. We begin from the case where either the label set  $\mathcal{Y}$  or the output set  $\mathcal{Z}$  are finite.

**Theorem 12 (ILE & finite  $\mathcal{Y}$  or  $\mathcal{Z}$ )** *The function  $\Delta : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$  admits an ILE if one of the following conditions hold:*

(a)  $\mathcal{Z}$  and  $\mathcal{Y}$  are finite sets. In this case  $c_\Delta \leq \|\Delta\|$  the operator norm of the matrix  $\Delta \in \mathbb{R}^{|\mathcal{Z}| \times |\mathcal{Y}|}$  with entries  $\Delta_{z,y} = \Delta(z, y)$ .

(b)  $\mathcal{Z}$  is finite,  $\mathcal{Y}$  is compact and  $\Delta(z, \cdot)$  is continuous on  $\mathcal{Y}$  for any  $z \in \mathcal{Z}$ . In this case  $c_\Delta \leq \sup_{y \in \mathcal{Y}} \sqrt{\sum_{z \in \mathcal{Z}} |\Delta(z, y)|^2}$ .

(c)  $\mathcal{Z}$  is compact,  $\mathcal{Y}$  is finite and  $\Delta(\cdot, y)$  is continuous on  $\mathcal{Z}$  for any  $y \in \mathcal{Y}$ . In this case  $c_\Delta \leq \sup_{z \in \mathcal{Z}} \sqrt{\sum_{y \in \mathcal{Y}} |\Delta(z, y)|^2}$ .

**Proof** (a). The proof of point (a) has been already given in the discussion after Theorem 2. We recall it here for completeness. By hypothesis we have  $\mathcal{Z} = \{z_1, \dots, z_p\}$  and  $\mathcal{Y} = \{y_1, \dots, y_q\}$  for some  $p, q \in \mathbb{N}$ . Let  $V \in \mathbb{R}^{p \times q}$  be the matrix whose entries correspond to the values of  $\Delta$  on pairs of points in  $\mathcal{Z} \times \mathcal{Y}$ . More precisely

$$V_{ij} = \Delta(z_i, y_j) \quad \forall i = 1, \dots, p, \quad j = 1, \dots, q. \quad (\text{C.5})$$

It is easy to prove that the ILE definition holds for  $\mathcal{H} = \mathbb{R}^q$  and feature maps  $\psi : \mathcal{Z} \rightarrow \mathcal{H}$  and  $\varphi : \mathcal{Y} \rightarrow \mathcal{H}$  corresponding to

$$\psi(z_i) = V^\top e_i^{(p)} \quad \varphi(y_j) = e_j^{(q)} \quad (\text{C.6})$$

for any  $i = 1, \dots, p$  and any  $j = 1, \dots, q$ , with  $e_i^{(p)} \in \mathbb{R}^p$  denoting the  $i$ -th element of the canonical basis of  $\mathbb{R}^p$ , namely the  $p$ -dimensional vector with  $i$ -th entry equal to 1 and all others equal to 0. Indeed, by construction we have

$$\langle \psi(z_i), \varphi(y_j) \rangle_{\mathcal{H}} = \langle e_i^{(p)}, V e_j^{(q)} \rangle = V_{ij} = \Delta(z_i, y_j). \quad (\text{C.7})$$

Finally, we have

$$c_\Delta = \sup_{z \in \mathcal{Z}} \|\psi(z)\|_{\mathcal{H}} = \sup_{i=1, \dots, p} \|V e_i^{(p)}\| \leq \|V\|, \quad (\text{C.8})$$

as required.

(b).  $\mathcal{Z} = \{z_1, \dots, z_p\}$  with  $p \in \mathbb{N}$ . We choose  $\mathcal{H} = \mathbb{R}^p$  and the feature maps  $\bar{\psi} : \mathcal{Z} \rightarrow \mathcal{H}$  and  $\bar{\varphi} : \mathcal{Y} \rightarrow \mathcal{H}$  such that  $\psi(z_i) = e_i^{(p)}$  for every  $i = 1, \dots, p$  and

$$\varphi(y) = (\Delta(z_1, y), \dots, \Delta(z_p, y))^\top \in \mathbb{R}^p, \quad (\text{C.9})$$

for any  $y \in \mathcal{Y}$ . Now, let

$$r = \sup_{y \in \mathcal{Y}} \|\varphi(y)\|_{\mathcal{H}} = \sup_{y \in \mathcal{Y}} \sqrt{\sum_{z \in \mathcal{Z}} \Delta(z, y)^2} \quad (\text{C.10})$$

we can define  $\psi = r\bar{\psi}$  and  $\varphi/r$ . We have that the ILE definition is satisfied, since

$$\langle \psi(z_i), \varphi(y) \rangle_{\mathcal{H}} = \left\langle e_i^{(p)}, \varphi(y) \right\rangle = \Delta(z_i, y), \quad (\text{C.11})$$

for every  $i = 1, \dots, p$ . Moreover, since  $\|\psi(z)\| = 1$  for every  $z \in \mathcal{Z}$ , we conclude that  $c_{\Delta} = r$  as required.

(c). The proof of point (c) is analogous to (b) with the difference that for  $\mathcal{Y} = \{y_1, \dots, y_q\}$  with  $q \in \mathbb{N}$  we choose  $\mathcal{H} = \mathbb{R}^q$ , and feature maps  $\varphi(y_j) = e_j^{(q)}$  for any  $j = 1, \dots, q$  and

$$\psi(z) = (\Delta(z, y_1), \dots, \Delta(z, y_q))^{\top} \in \mathbb{R}^q, \quad (\text{C.12})$$

for any  $z \in \mathcal{Z}$ . The proof follows identically to (b).  $\blacksquare$

#### C.4. ILE and Reproducing Kernel Hilbert Spaces

We now focus on the relation between ILE and reproducing kernel Hilbert spaces.

**Theorem 13 (ILE & RKHS)** *Let  $\mathcal{Z} = \mathcal{Y}$  be a compact set and  $h : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  a continuous bounded reproducing kernel on  $\mathcal{Y}$  with associated RKHS  $\mathcal{H}$ . Let  $\eta^2 = \sup_{y \in \mathcal{Y}} h(y, y)$ . Then,  $\Delta : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  admits an ILE if one of the following holds:*

- (a) (Kernels).  $\Delta(z, y) = h(z, y)$  for any  $y, z \in \mathcal{Y}$ . In this case  $c_{\Delta} \leq \eta^2$ .
- (b) (Kernel Dependency Estimation (KDE)).  $\Delta(z, y) = h(z, z) + h(y, y) - 2h(z, y)$  for any  $y, z \in \mathcal{Y}$ . In this case  $c_{\Delta} = 2(2\eta^4 + 1)$ .
- (c) For every  $y \in \mathcal{Y}$  the functions  $\Delta(\cdot, y) \in \mathcal{H}$  belong to a bounded set of  $\mathcal{H}$ , namely  $\sup_{y \in \mathcal{Y}} \|\Delta(\cdot, y)\|_{\mathcal{H}} = D < +\infty$ . In this case  $c_{\Delta} \leq \eta D$ . The same holds if the family of functions  $\Delta(z, \cdot)$  parametrized by  $z \in \mathcal{Z}$  belong to a bounded set of  $\mathcal{H}$ .
- (d)  $\Delta$  belongs to  $\mathcal{H} \otimes \mathcal{H}$  the RKHS with associated kernel  $\bar{h} : \mathcal{Y}^2 \times \mathcal{Y}^2 \rightarrow \mathbb{R}$  such that  $\bar{h}((z, y), (z', y')) = h(z, z')h(y, y')$  for any  $z, z', y, y' \in \mathcal{Y}$ . In this case  $c_{\Delta} \leq \eta^2 \|\Delta\|_{\mathcal{H} \otimes \mathcal{H}}$

**Proof** (a). Follows directly from point (c). Indeed  $h(z, \cdot) \in \mathcal{H}$  and  $\sup_{z \in \mathcal{Y}} \|h(z, \cdot)\| \leq \eta$  by hypothesis.

(b). We note that point (b) follows directly from Theorem 16, which guarantees the finite sums and products of ILE functions to be ILE as well. Here we give a more direct proof for completeness.

Let  $\bar{\mathcal{H}} = \mathbb{R} \oplus \mathcal{H} \oplus \mathbb{R}$  equipped with the canonical inner product of the direct sum and let  $V : \bar{\mathcal{H}} \rightarrow \bar{\mathcal{H}}$  the linear operator such that

$$V(\alpha, h, \beta) = (\beta, -2h, \alpha) \quad (\text{C.13})$$

for any  $h \in \mathcal{H}$  and  $\alpha, \beta \in \mathbb{R}$ . Let  $\zeta : \mathcal{Y} \rightarrow \bar{\mathcal{H}}$  be such that

$$\zeta(y) = (h(y, y), h(y, \cdot), 1)^{\top}, \quad (\text{C.14})$$

for any  $y \in \mathcal{Y}$ . Let

$$r = \sup_{y \in \mathcal{Y}} \|\zeta(y)\|_{\overline{\mathcal{H}}} = \sup_{y \in \mathcal{Y}} \sqrt{h(y, y)^2 + \|h(y, \cdot)\|^2 + 1} \leq \sqrt{2\eta^4 + 1}. \quad (\text{C.15})$$

We choose  $\psi, \varphi : \mathcal{Y} \rightarrow \overline{\mathcal{H}}$  as

$$\psi(z) = r V \zeta(z) \quad \varphi(y) = \frac{1}{r} \zeta(y), \quad (\text{C.16})$$

for any  $z, y \in \mathcal{Y}$ . Then, by construction

$$\langle \psi(z), \varphi(y) \rangle_{\overline{\mathcal{H}}} = \left\langle (1, -2h(z, \cdot), h(z, z))^\top, (h(y, y), h(y, \cdot), 1)^\top \right\rangle_{\overline{\mathcal{H}}} \quad (\text{C.17})$$

$$= h(y, y) - 2h(z, y) + h(z, z) \quad (\text{C.18})$$

$$= \Delta(z, y), \quad (\text{C.19})$$

as desired. Moreover,

$$c_\Delta = \sup_{z \in \mathcal{Z}} \psi(z) = r \sup_{z \in \mathcal{Z}} \|V \zeta(z)\| \leq r^2 \|V\| \leq 2r^2 = 2(2\eta^4 + 1), \quad (\text{C.20})$$

as desired, with  $\|V\| = 2$  denoting the operator norm of  $V$ .

(c). We prove the statement for the case  $\sup_{z \in \mathcal{Z}} \|\Delta(z, \cdot)\| = D < +\infty$ . We consider the feature maps  $\psi, \varphi : \mathcal{Y} \rightarrow \mathcal{H}$  such that

$$\psi(z) = \eta \Delta(z, \cdot) \quad \text{and} \quad \varphi(y) = \frac{1}{\eta} h(y, \cdot), \quad (\text{C.21})$$

for any  $z, y \in \mathcal{Y}$ . Then, by construction we have

$$\langle \psi(z), \varphi(y) \rangle_{\mathcal{H}} = \langle \Delta(z, \cdot), h(y, \cdot) \rangle_{\mathcal{H}} = \Delta(z, y), \quad (\text{C.22})$$

where the last inequality follows from the reproducing property of the kernel  $h$  and the fact that  $\Delta(z, \cdot) \in \mathcal{H}$ . Moreover we have  $c_\Delta = \sup_{z \in \mathcal{Z}} \|\psi(z)\|_{\mathcal{H}} = \eta \sup_{z \in \mathcal{Z}} \|\Delta(z, \cdot)\| = \eta D$ . The case  $\sup_{y \in \mathcal{Y}} \|\Delta(\cdot, y)\| \leq D$  follows from an analogous reasoning.

(d). Note that the kernel  $\bar{h}$  has feature map  $(z, y) \mapsto h(z, \cdot) \otimes h(y, \cdot)$  for any  $z, y \in \mathcal{Y}$ . Since by hypothesis  $\Delta \in \mathcal{H} \otimes \mathcal{H}$ , the reproducing property for  $\bar{h}$  implies

$$\Delta(z, y) = \langle \Delta, h(z, \cdot) \otimes h(y, \cdot) \rangle. \quad (\text{C.23})$$

Since  $\mathcal{H} \otimes \mathcal{H}$  is isometric to the space of Hilbert-Schmidt operators from  $\mathcal{H}$  to  $\mathcal{H}$ , there exists an operator  $V : \mathcal{H} \rightarrow \mathcal{H}$  such that  $\|V\|_{\text{HS}} = \|\Delta\|_{\mathcal{H} \otimes \mathcal{H}}$  and

$$\langle \Delta, h(z, \cdot) \otimes h(y, \cdot) \rangle = \langle V, h(z, \cdot) \otimes h(y, \cdot) \rangle_{\text{HS}} = \langle V h(z, \cdot), h(y, \cdot) \rangle_{\mathcal{H}}, \quad (\text{C.24})$$

where the last inequality follows from the standard properties of tensor products. We can therefore choose  $\psi, \varphi : \mathcal{Y} \rightarrow \mathcal{H}$  such that

$$\psi(z) = \eta V h(z, \cdot) \quad \text{and} \quad \varphi(y) = \frac{1}{\eta} h(y, \cdot), \quad (\text{C.25})$$

for any  $y \in \mathcal{Y}$  to guarantee the ILE definition to hold. Moreover, by construction  $c_\Delta = \sup_{z \in \mathcal{Z}} \|\psi(z)\|_{\mathcal{H}} \leq \eta \|V\|_{\text{HS}} \sup_{z \in \mathcal{Z}} \|h(z, \cdot)\|_{\mathcal{H}} \leq \eta^2 \|V\|_{\text{HS}}$  which concludes the proof since  $\|V\|_{\text{HS}} = \|\Delta\|_{\mathcal{H} \otimes \mathcal{H}}$  by construction.  $\blacksquare$

We now report the result relating general notions of regularity for the loss and the ILE definition. Before proving the main result in Theorem 14 we need the following two Lemmas.

**Lemma C.3 (Multiple Fourier Series)** *Let  $\mathcal{Y} = [-B, B]^d$  with  $d \in \mathbb{N}$  and  $B > 0$ . Let  $(\hat{f}_h)_{h \in \mathbb{Z}^d} \in \mathbb{C}$  and  $f : \mathcal{Y} \rightarrow \mathbb{C}$  defined as*

$$f(y) = \sum_{h \in \mathbb{Z}^d} \hat{f}_h e^{2\pi i h^\top y}, \quad \forall y \in \mathcal{Y}, \quad \text{with} \quad \sum_{h \in \mathbb{Z}^d} |\hat{f}_h| \leq M, \quad (\text{C.26})$$

for  $0 < M < \infty$  and  $i = \sqrt{-1}$ . Then the function  $f$  is continuous and

$$\sup_{y \in \mathcal{Y}} |f(y)| \leq M. \quad (\text{C.27})$$

**Proof** The continuity of  $f$  follows from (Kahane, 1995), pag. 129 and Example 2. To show that  $f$  is uniformly bounded on  $\mathcal{Y}$  it is sufficient to see that

$$\sup_{y \in \mathcal{Y}} |f(y)| \leq \sup_{y \in \mathcal{Y}} \sum_{h \in \mathbb{Z}^d} |\hat{f}_h| |e^{2\pi i h^\top y}| \leq \sum_{h \in \mathbb{Z}^d} |\hat{f}_h| \leq M. \quad (\text{C.28})$$

$\blacksquare$

**Lemma C.4** *Let  $\mathcal{Y} = [-B, B]^d$  with  $d \in \mathbb{N}$  and  $B > 0$ . Let  $\Delta : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  be such that*

$$\Delta(y, z) = \sum_{h, k \in \mathbb{Z}^d} \hat{\Delta}_{h, k} e_h(y) e_k(z), \quad \forall y, z \in \mathcal{Y}, \quad (\text{C.29})$$

with  $e_h(y) = e^{2\pi i h^\top y}$  for any  $y \in \mathcal{Y}$ ,  $i = \sqrt{-1}$  and  $\hat{\Delta}_{h, k} \in \mathbb{C}$  for any  $h, k \in \mathbb{Z}^d$ . If

$$c_\Delta = \sum_{h, k \in \mathbb{Z}^d} |\hat{\Delta}_{h, k}| < \infty, \quad (\text{C.30})$$

then  $\Delta$  admits an ILE.

**Proof** We start by applying Theorem C.3 for the input domain  $\mathcal{Y} \times \mathcal{Y}$ , which guarantees that the function  $\Delta$  is bounded continuous. We introduce the following sequences

$$\alpha_h = \sum_{k \in \mathbb{Z}^d} |\hat{\Delta}_{h, k}|, \quad f_h(z) = \frac{1}{\alpha_h} \sum_{k \in \mathbb{Z}^d} \hat{\Delta}_{h, k} e_k(z) \quad \forall h \in \mathbb{Z}^d, z \in \mathcal{Y}. \quad (\text{C.31})$$

For any  $p > 0$ , let  $\ell_p(\mathbb{Z}^d)$  denote the set of sequences  $(a_k)_{k \in \mathbb{Z}^d}$  such that  $\sum_{h \in \mathbb{Z}^d} |a_k|^p < +\infty$ . Note that by hypothesis  $(\alpha_h)_{h \in \mathbb{Z}^d} \in \ell_1(\mathbb{Z}^d)$ . Moreover, by applying again Theorem C.3, we have that the functions  $f_h$  are continuous and bounded by 1 for any  $h \in \mathbb{Z}^d$ .

Denote  $A^2 = \|(\alpha_h)_{h \in \mathbb{Z}^d}\|_{\ell_1(\mathbb{Z}^d)} = \sum_{h,k \in \mathbb{Z}^d} |\widehat{\Delta}_{h,k}|$ . Let  $\mathcal{H} = \ell_2(\mathbb{Z}^d)$  and let  $\psi, \varphi : \mathcal{Y} \rightarrow \mathcal{H}$  be such that

$$\psi(z) = A \left( \sqrt{\alpha_h} f_h(z) \right)_{h \in \mathbb{Z}^d} \quad \text{and} \quad \varphi(y) = \frac{1}{A} \left( \sqrt{\alpha_h} e_h(y) \right)_{h \in \mathbb{Z}^d}, \quad (\text{C.32})$$

for all  $z, y \in \mathcal{Y}$ . By construction, we have

$$\langle \psi(z), \varphi(y) \rangle_{\mathcal{H}} = \sum_{h \in \mathbb{Z}^d} \alpha_h e_h(y) f_h(z) = \sum_{h,k \in \mathbb{Z}^d} \widehat{\Delta}_{h,k} e_h(y) e_h(z) = \Delta(y, z). \quad (\text{C.33})$$

Moreover, by construction we also have  $\|\varphi(y)\|_{\mathcal{H}} \leq 1$  and  $\mathbf{c}_{\Delta} = \sup_{z \in \mathcal{Z}} \|\psi(z)\| \leq A^2$ .

To conclude the proof, we need to show that the two feature maps are continuous. Define  $\zeta_1(y, z) = \langle \varphi(y), \varphi(z) \rangle_{\mathcal{H}}$  and  $\zeta_2(y, z) = \langle \psi(y), \psi(z) \rangle_{\mathcal{H}_0}$  for all  $y, z \in \mathcal{Y}$ . We have

$$\zeta_1(y, z) = \sum_{h \in \mathbb{Z}^d} \alpha_h \overline{e_h(y)} e_h(z), \quad (\text{C.34})$$

$$\zeta_2(y, z) = \sum_{h \in \mathbb{Z}^d} \alpha_h \overline{f_h(y)} f_h(z) = \sum_{k,l \in \mathbb{Z}^d} \beta_{k,l} \overline{e_k(y)} e_l(z) \quad (\text{C.35})$$

with  $\beta_{k,l} = \sum_{h \in \mathbb{Z}^d} \frac{\widehat{\Delta}_{h,k} \widehat{\Delta}_{h,l}}{\alpha_h}$ , for  $k, l \in \mathbb{Z}^d$ , therefore  $\zeta_1, \zeta_2$  are bounded and continuous by Theorem C.3, since  $\sum_{h \in \mathbb{Z}^d} \alpha_h < \infty$  and  $\sum_{k,l \in \mathbb{Z}^d} |\beta_{k,l}| < \infty$ . Note that  $\psi$  and  $\varphi$  are bounded, since  $\zeta_1$  and  $\zeta_2$  are. Moreover for any  $y, z \in \mathcal{Y}$ , we have

$$\|\varphi(y) - \varphi(z)\|_{\mathcal{H}}^2 = \langle \varphi(z), \varphi(z) \rangle_{\mathcal{H}} + \langle \varphi(y), \varphi(y) \rangle_{\mathcal{H}} - 2 \langle \varphi(z), \varphi(y) \rangle_{\mathcal{H}} \quad (\text{C.36})$$

$$= \zeta_1(z, z) + \zeta_1(y, y) - 2\zeta_1(z, y) \quad (\text{C.37})$$

$$\leq |\zeta_1(z, z) - \zeta_1(z, y)| + |\zeta_1(z, y) - \zeta_1(y, y)|, \quad (\text{C.38})$$

and the same holds for  $\psi$  with respect to  $\zeta_2$ . Thus the continuity of  $\varphi$  is ensured by the continuity of  $\zeta_1$  and the same for  $\psi$  with respect to  $\zeta_2$ .  $\blacksquare$

We are ready to prove the following result.

**Theorem 14 (ILE & Regularity)** *Let  $\mathcal{Z} = \mathcal{Y} = [-B, B]^d$ ,  $B > 0$ . A function  $\Delta : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  admits an ILE when at least one of the following conditions hold:*

- (a)  *$d = 1$  and  $\Delta$  is  $\alpha$ -Hölder continuous with  $\alpha > 1/2$  or it is of bounded variation and  $\alpha$ -Hölder continuous with  $\alpha > 0$ .*
- (b)  *$\Delta(z, y) = v(z - y)$ , where  $v$  is a function such that  $\mathbf{c}_{\Delta} = \int |\widehat{v}(\omega)| d\omega < \infty$  and  $\widehat{v}$  is the Fourier transform of  $v$ .*
- (c) *The mixed partial derivative  $\Delta_{y_1, \dots, y_d} : \mathcal{Y} \rightarrow \mathbb{R}$  of  $\Delta$  exists almost everywhere and  $\Delta_{y_1, \dots, y_d} \in L^p(\mathcal{Y})$  with  $p > 1$ .*

**Proof** (a-b). Either hypotheses in (a) or (b) are sufficient to guarantee that the Fourier expansion of  $\Delta$  is absolutely summable (see Theorem 5' and Theorem 6' pag. 291 of Móricz and Veres, 2007). By Theorem C.4 we can conclude that  $\Delta$  admits an ILE.

(c). Let  $\gamma^2 = \int_{-\infty}^{+\infty} |\widehat{v}(\omega)| < +\infty$  we have that for any  $z, y \in \mathcal{Y}$ , the anti-Fouier transform of  $\widehat{v}(\omega)$  in  $z - y$  is

$$v(z - y) = \int_{-\infty}^{+\infty} \widehat{v}(\omega) e^{i\langle \omega, z-y \rangle} d\omega = \int_{-\infty}^{+\infty} \widehat{v}(\omega) e^{i\langle \omega, z \rangle} e^{-i\langle \omega, y \rangle} d\omega. \quad (\text{C.39})$$

Let now  $\mathcal{H} = L^2(\mathbb{R}, \mathbb{C})$  the space of square integrable functions from  $\mathbb{R}$  to  $\mathbb{C}$  with respect to the Lebesgue measure and let  $\psi, \varphi : \mathcal{Y} \rightarrow \mathcal{H}$  be such that

$$\psi(z) = \gamma \sqrt{\widehat{v}(\cdot)} e^{i\langle \cdot, z \rangle} \quad \text{and} \quad \varphi(y) = \frac{1}{\gamma} \sqrt{\widehat{v}(\cdot)} e^{-i\langle \cdot, y \rangle}. \quad (\text{C.40})$$

By the anti-Fourier transform we have

$$\langle \psi(z), \varphi(y) \rangle_{\mathcal{H}} = v(z - y) = \Delta(z, y) \quad (\text{C.41})$$

for every  $z, y \in \mathcal{Y}$ . Moreover, by construction  $\|\varphi(y)\|_{\mathcal{H}} = 1$  and  $\mathbf{c}_{\Delta} = \sup_{z \in \mathcal{Z}} \|\psi(z)\|_{\mathcal{H}} = \gamma^2$  as required.  $\blacksquare$

### C.5. Composition Rules for ILE

We conclude with the result characterizing composition rules for the ILE property.

**Theorem 15** *Let  $\mathcal{Z}$  and  $\mathcal{Y}$  be compact sets. Then  $\Delta : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$  admits an ILE if one of the following holds:*

(a) (Restriction) *There exist two sets  $\bar{\mathcal{Z}} \supseteq \mathcal{Z}$ ,  $\bar{\mathcal{Y}} \supseteq \mathcal{Y}$  and  $\bar{\Delta} : \bar{\mathcal{Z}} \times \bar{\mathcal{Y}} \rightarrow \mathbb{R}$  such that  $\bar{\Delta}$  admits an ILE and its restriction to  $\mathcal{Z} \times \mathcal{Y}$  corresponds to  $\Delta$ , namely*

$$\Delta = \bar{\Delta}|_{\mathcal{Z} \times \mathcal{Y}}. \quad (54)$$

*In this case  $\mathbf{c}_{\Delta} \leq \mathbf{c}_{\bar{\Delta}}$ .*

(b) (Right Composition) *There exists  $\bar{\mathcal{Z}}, \bar{\mathcal{Y}}$  and a ILE  $\bar{\Delta} : \bar{\mathcal{Z}} \times \bar{\mathcal{Y}} \rightarrow \mathbb{R}$ , such that*

$$\Delta(z, y) = \alpha(z) \bar{\Delta}(A(z), B(y)) \beta(y), \quad (55)$$

*with  $A : \mathcal{Z} \rightarrow \bar{\mathcal{Z}}$ ,  $B : \mathcal{Y} \rightarrow \bar{\mathcal{Y}}$ ,  $\alpha : \mathcal{Z} \rightarrow \mathbb{R}$  and  $\beta : \mathcal{Y} \rightarrow \mathbb{R}$  continuous function, with  $\sup_{z \in \mathcal{Z}} |\alpha(z)| \leq \bar{\alpha}$  and  $\sup_{y \in \mathcal{Y}} |\beta(y)| \leq \bar{\beta}$  with  $\bar{\alpha}, \bar{\beta} \in \mathbb{R}$ . Then  $\mathbf{c}_{\Delta} \leq \bar{\alpha} \bar{\beta} \mathbf{c}_{\bar{\Delta}}$ .*

(c) (Left Composition) *There exist  $P \in \mathbb{N}$ , spaces  $(\mathcal{Z}_p)_{p=1}^P, (\mathcal{Y}_p)_{p=1}^P$  and corresponding ILE  $\Delta_p : \mathcal{Z}_p \times \mathcal{Y}_p \rightarrow \mathbb{R}$  such that  $\mathcal{Z} = \mathcal{Z}_1 \times \cdots \times \mathcal{Z}_P$ ,  $\mathcal{Y} = \mathcal{Y}_1 \times \cdots \times \mathcal{Y}_P$  and*

$$\Delta(z, y) = \Gamma(\Delta_1(z_1, y_1), \dots, \Delta_P(z_P, y_P)), \quad (56)$$

*for any  $z = (z_1, \dots, z_P) \in \mathcal{Z}$  and  $y = (y_1, \dots, y_P) \in \mathcal{Y}$ , where  $\Gamma : \mathbb{R}^P \rightarrow \mathbb{R}$  is an analytic function (e.g. a polynomial).*

**Proof** (a). Let  $\bar{\Delta}$  admit an ILE with  $\mathcal{H}$  separable Hilbert space and feature maps  $\bar{\psi} : \bar{\mathcal{Z}} \rightarrow \mathcal{H}$  and  $\bar{\varphi} : \bar{\mathcal{Y}} \rightarrow \mathcal{H}$ . Clearly, for any  $\mathcal{Y} \subseteq \bar{\mathcal{Y}}$ ,  $\mathcal{Z} \subseteq \bar{\mathcal{Z}}$ , we have that the restriction of the feature maps  $\psi = \bar{\psi}|_{\mathcal{Z}}$  and  $\varphi = \bar{\varphi}|_{\mathcal{Y}}$  are such that

$$\langle \psi(z), \varphi(y) \rangle_{\mathcal{H}} = \langle \bar{\psi}(z), \bar{\varphi}(y) \rangle_{\mathcal{H}} = \bar{\Delta}(z, y) = \bar{\Delta}|_{\bar{\mathcal{Z}} \times \bar{\mathcal{Y}}} = \Delta(z, y), \quad (\text{C.42})$$

for any  $z \in \mathcal{Z}$  and  $y \in \mathcal{Y}$ . The proof is concluded by observing that  $\sup_{y \in \mathcal{Y}} \|\bar{\varphi}(y)\|_{\mathcal{H}} \leq \sup_{y \in \bar{\mathcal{Y}}} \|\bar{\varphi}(y)\|_{\mathcal{H}} \leq 1$  and  $\mathbf{c}_{\Delta} = \sup_{z \in \mathcal{Z}} \|\bar{\psi}(z)\|_{\mathcal{H}} \leq \sup_{z \in \bar{\mathcal{Z}}} \|\bar{\psi}(z)\|_{\mathcal{H}} = \mathbf{c}_{\bar{\Delta}}$ .

(b). Let  $\bar{\Delta}$  admit an ILE with  $\mathcal{H}$  separable Hilbert space and feature maps  $\bar{\psi} : \bar{\mathcal{Z}} \rightarrow \mathcal{H}$  and  $\bar{\varphi} : \bar{\mathcal{Y}} \rightarrow \mathcal{H}$ . Let  $\bar{\beta} = \sup_{y \in \mathcal{Y}} |\beta(y)|$ . We consider the feature maps  $\psi : \mathcal{Z} \rightarrow \mathcal{H}$  and  $\varphi : \mathcal{Y} \rightarrow \mathcal{H}$  such that

$$\psi(z) = \bar{\beta} \alpha(z) \bar{\psi}(A(z)) \quad \text{and} \quad \varphi(y) = \frac{\beta(z)}{\bar{\beta}} \bar{\varphi}(B(z)), \quad (\text{C.43})$$

for all  $z \in \mathcal{Z}$  and  $y \in \mathcal{Y}$ . By construction we have

$$\langle \psi(z), \varphi(y) \rangle = \alpha(z) \langle \bar{\psi}(A(z)), \bar{\varphi}(B(y)) \rangle \beta(y) = \alpha(z) \bar{\Delta}(A(z), B(y)) \beta(y). \quad (\text{C.44})$$

Moreover, we have  $\sup_{y \in \mathcal{Y}} \|\varphi(y)\|_{\mathcal{H}} = \frac{1}{\bar{\beta}} \sup_{y \in \mathcal{Y}} \beta(y) \|\bar{\varphi}(y)\|_{\mathcal{H}} \leq 1$  and  $\mathbf{c}_{\Delta} = \sup_{z \in \mathcal{Z}} \|\psi(z)\|_{\mathcal{H}} = \bar{\beta} \sup_{z \in \mathcal{Z}} \alpha(z) \|\bar{\psi}(z)\|_{\mathcal{H}} \leq \bar{\alpha} \bar{\beta} \mathbf{c}_{\bar{\Delta}}$  as required.

(c). By definition of analytic functions, we have that  $\Gamma$  has form

$$\Gamma(r_1, \dots, r_P) = \sum_{t \in \mathbb{N}^P} \alpha_t \prod_{p=1}^P r_p^{t_p} \quad \forall r = (r_1, \dots, r_P)^{\top} \in \mathbb{R}^P, \quad (\text{C.45})$$

for some scalar weights  $\alpha_t$  with  $t \in \mathbb{N}^P$ . Therefore, for any  $z = (z_1, \dots, z_P) \in \mathcal{Z}$  and  $y = (y_1, \dots, y_P) \in \mathcal{Y}$ , we have

$$\Delta(z, y) = \Gamma(\Delta_1(z_1, y_1), \dots, \Delta_P(z_P, y_P)) = \sum_{t \in \mathbb{N}^P} \alpha_t \prod_{p=1}^P \Delta_p(z_p, y_p)^{t_p}. \quad (\text{C.46})$$

Recall that for any two Hilbert spaces  $\mathcal{H}$  and  $\mathcal{H}'$ , by definition of direct sum  $\mathcal{H} \oplus \mathcal{H}'$  and tensor product  $\mathcal{H} \otimes \mathcal{H}'$ , we have

$$\langle h_1 \oplus h'_1, h_2 \oplus h'_2 \rangle_{\mathcal{H} \oplus \mathcal{H}'} = \langle h_1, h_2 \rangle_{\mathcal{H}} + \langle h'_1, h'_2 \rangle_{\mathcal{H}'}, \quad (\text{C.47})$$

$$\langle h_1 \otimes h'_1, h_2 \otimes h'_2 \rangle_{\mathcal{H} \otimes \mathcal{H}'} = \langle h_1, h_2 \rangle_{\mathcal{H}} \cdot \langle h'_1, h'_2 \rangle_{\mathcal{H}'}, \quad (\text{C.48})$$

for any  $h_1, h_2 \in \mathcal{H}$  and  $h'_1, h'_2 \in \mathcal{H}'$ . Moreover, for any  $p \in \mathbb{N}$ , we denote  $\mathcal{H}^{\otimes p}$  the tensor product of  $\mathcal{H}$  with itself  $p$  times (with  $\mathcal{H}^{\otimes 0} = \mathbb{R}$ ) and denote  $h^{\otimes p} \in \mathcal{H}^{\otimes p}$  the tensor product of  $h$  with itself  $p$  times, for any  $h \in \mathcal{H}$  (with  $h^{\otimes 0} = 1$ ). This implies in particular that for any  $t \in \mathbb{N}^P$  we have

$$\prod_{p=1}^P \Delta_p(z_p, y_p)^{t_p} = \left\langle \bigotimes_{p=1}^P \psi_p(z_p)^{\otimes t_p}, \bigotimes_{p=1}^P \varphi_p(y_p)^{\otimes t_p} \right\rangle_{\bar{\mathcal{H}}_t}, \quad (\text{C.49})$$



$\mathcal{H}_t = \bigotimes_{p=1}^P \mathcal{H}_p^{\otimes t_p}$  and we have denoted with  $\psi_p(z_p)^{\otimes t_p}$  the tensor product of  $\psi_p(z_p)$  with itself  $t_p$  times.

For any  $t \in \mathbb{N}^P$ , let  $\beta_t = \text{sign}(\alpha_t) \sqrt{|\alpha_t|}$  and  $\gamma_t = \sqrt{|\alpha_t|}$ . Then, we have

$$\sum_{t \in \mathbb{N}^P} \alpha_t \prod_{p=1}^P \Delta_p(z_p, y_p)^{t_p} = \sum_{t \in \mathbb{N}^P} \alpha_t \left\langle \bigotimes_{p=1}^P \psi_p(z_p)^{\otimes t_p}, \bigotimes_{p=1}^P \varphi_p(y_p)^{\otimes t_p} \right\rangle_{\bar{\mathcal{H}}_t} \quad (\text{C.50})$$

$$= \left\langle \bigoplus_{t \in \mathbb{N}^P} \beta_t \left[ \bigotimes_{p=1}^P \psi_p(z_p)^{\otimes t_p} \right], \bigoplus_{t \in \mathbb{N}^P} \gamma_t \left[ \bigotimes_{p=1}^P \varphi_p(y_p)^{\otimes t_p} \right] \right\rangle_{\mathcal{H}} \quad (\text{C.51})$$

$$= \langle \psi(z), \varphi(y) \rangle_{\mathcal{H}}. \quad (\text{C.52})$$

where  $\mathcal{H} = \bigoplus_{t \in \mathbb{N}^P} \bar{\mathcal{H}}_t$  and we  $\psi : \mathcal{Z} \rightarrow \mathcal{H}$  and  $\varphi : \mathcal{Y} \rightarrow \mathcal{H}$  are feature maps such that

$$\psi(z) = \bigoplus_{t \in \mathbb{N}^P} \beta_t \left[ \bigotimes_{p=1}^P \psi_p(z_p)^{\otimes t_p} \right] \quad \text{and} \quad \varphi(y) = \bigoplus_{t \in \mathbb{N}^P} \gamma_t \left[ \bigotimes_{p=1}^P \varphi_p(y_p)^{\otimes t_p} \right], \quad (\text{C.53})$$

for any  $z = (z_1, \dots, z_P) \in \mathcal{Z}$  and  $y = (y_1, \dots, y_P) \in \mathcal{Y}$ . First note that such feature maps are well defined, namely that they indeed take value in  $\mathcal{H}$ . In particular, we have

$$\|\varphi(y)\|_{\mathcal{H}} = \left\| \bigoplus_{t \in \mathbb{N}^P} \gamma_t \left[ \bigotimes_{p=1}^P \varphi_p(y_p)^{\otimes t_p} \right] \right\|_{\bar{\mathcal{H}}}^2 \quad (\text{C.54})$$

$$= \sum_{t \in \mathbb{N}^P} |\alpha_t| \left\| \bigotimes_{p=1}^P \varphi_p(y_p)^{\otimes t_p} \right\|_{\bar{\mathcal{H}}_t}^2 \quad (\text{C.55})$$

$$= \sum_{t \in \mathbb{N}^P} |\alpha_t| \prod_{p=1}^P \left( \|\varphi_p(y_p)\|_{\mathcal{H}_p}^2 \right)^{t_p}. \quad (\text{C.56})$$

To show that the above series is finite, recall that the power series defining  $\Gamma$  is absolutely convergent for any  $r = (r_1, \dots, r_P) \in \mathbb{R}^P$ . Indeed, let  $\bar{r} = (\bar{r}_1, \dots, \bar{r}_P) \in \mathbb{R}^P$  such that  $|r_p| < |\bar{r}_p|$  for any  $p = 1, \dots, P$ . Since  $\Gamma$  is analytic also in  $\bar{r}$ , the associated power series is convergent and therefore, for  $\|t\| \rightarrow +\infty$  we have  $\alpha_t \prod_{p=1}^P \bar{r}_p^{t_p} \rightarrow 0$ . This implies that there exists  $Q > 0$  such that, for any  $t \in \mathbb{N}^P$  with  $\|t\| > Q$ ,

$$|\alpha_t| \prod_{p=1}^P |\bar{r}_p|^{t_p} \leq 1. \quad (\text{C.57})$$

By multiplying both sides of the inequality above by  $\prod_{p=1}^P |r_p/\bar{r}_p|^{t_p}$ , we have

$$|\alpha_t| \prod_{p=1}^P |r_p|^{t_p} < \prod_{p=1}^P \left| \frac{r_p}{\bar{r}_p} \right|^{t_p}. \quad (\text{C.58})$$

Since  $|r_p/\bar{r}_p| < 1$  for  $p = 1, \dots, P$  by construction, we can conclude that

$$\sum_{t \in \mathbb{N}^P} \left| \alpha_t \prod_{p=1}^P r_p^{t_p} \right| < +\infty. \quad (\text{C.59})$$

In particular, we have that the domain of the function  $\bar{\Gamma} : \mathbb{R}^P \rightarrow \mathbb{R}$ , such that

$$\bar{\Gamma}(r_1, \dots, r_P) = \sum_{t \in \mathbb{N}^P} \left| \alpha_t \prod_{p=1}^P r_p^{t_p} \right| \quad (\text{C.60})$$

corresponds to  $\mathbb{R}^P$ , namely  $\bar{\Gamma}(r_1, \dots, r_P) < +\infty$  for any  $r = (r_1, \dots, r_P)^\top \in \mathbb{R}^P$ .

Therefore, using the fact that  $\sup_{y_p \in \mathcal{Y}_p} \|\varphi_p(y_p)\| \leq 1$  for any  $p = 1, \dots, P$ , we have

$$\|\varphi(y)\|_{\mathcal{H}} \leq \sqrt{\bar{\Gamma}(1, \dots, 1)} < +\infty, \quad (\text{C.61})$$

for any  $y = (y_1, \dots, y_P) \in \mathcal{Y}$ . By following an analogous reasoning for  $\psi$ , we have

$$\|\psi(z)\|_{\mathcal{H}} \leq \sqrt{\bar{\Gamma}(c_{\Delta_1}^2, \dots, c_{\Delta_P}^2)} < +\infty, \quad (\text{C.62})$$

for any  $z = (z_1, \dots, z_P) \in \mathcal{Z}$ .

We need to show that the maps  $\psi$  and  $\varphi$  are continuous. To see this it is sufficient to prove that they are the uniform limit of continuous functions. In particular for any  $Q \in \mathbb{R}$ , let  $\varphi^{(Q)} : \mathcal{Y} \rightarrow \mathcal{H}$  be such that

$$\varphi^{(Q)}(y) = \left( \bigoplus_{t \in \mathbb{N}^P, \|t\| \leq Q} \gamma_t \left[ \bigotimes_{p=1}^P \varphi_p(y_p)^{\otimes t_p} \right] \right) \oplus \left( \bigoplus_{t \in \mathbb{N}^P, \|t\| > Q} 0 \right), \quad (\text{C.63})$$

for any  $y = (y_1, \dots, y_P) \in \mathcal{Y}$ . Note that  $\varphi^{(Q)}(y) \in \mathcal{H}$ , since by construction  $\mathcal{H} = \left( \bigoplus_{\|t\| \leq Q} \mathcal{H}_t \right) \oplus \left( \bigoplus_{\|t\| > Q} \mathcal{H}_t \right)$ . Moreover,  $\varphi^{(Q)} : \mathcal{Y} \rightarrow \mathcal{H}$  is continuous for any  $Q \in \mathbb{R}$  since it is the direct sum of a finite number of continuous functions.

Therefore, for any  $y \in \mathcal{Y}$ , we have

$$\|\varphi(y) - \varphi^{(Q)}(y)\|_{\mathcal{H}}^2 = \left\| \bigoplus_{t \in \mathbb{N}^P, \|t\| > Q} \gamma_t \left[ \bigotimes_{p=1}^P \varphi_p(y_p)^{\otimes t_p} \right] \right\|_{\mathcal{H}}^2 \quad (\text{C.64})$$

$$= \sum_{t \in \mathbb{N}^P, \|t\| > Q} |\alpha_t| \prod_{p=1}^P \left( \|\varphi_p(y_p)\|_{\mathcal{H}_p}^2 \right)^{t_p} \quad (\text{C.65})$$

$$\leq \sum_{t \in \mathbb{N}^P, \|t\| > Q} |\alpha_t|, \quad (\text{C.66})$$

where we have made use of the fact that  $\sup_{y_p \in \mathcal{Y}_p} \|\varphi_p(y_p)\|_{\mathcal{H}_p} \leq 1$  for any  $p = 1, \dots, P$ . Since  $\sum_{t \in \mathbb{N}^P, \|t\| \in \mathbb{N}^P} |\alpha_t| < +\infty$ , we have that for  $Q \rightarrow +\infty$ , the residual  $\sum_{t \in \mathbb{N}^P, \|t\| > Q} |\alpha_t|$  will tend to zero. We conclude that

$$\lim_{Q \rightarrow +\infty} \|\varphi(y) - \varphi^{(Q)}(y)\|_{\mathcal{H}} \rightarrow 0, \quad (\text{C.67})$$

showing that  $\varphi$  is uniform limit of continuous functions and hence continuous itself, as desired. The exact same argument holds for  $\psi : \mathcal{Z} \rightarrow \mathcal{H}$ .

Clearly,  $\varphi$  is not necessarily taking values in the ball of radius 1 in  $\mathcal{H}$ . To this end we can invoke Theorem C.2 by replacing  $\varphi(y)$  with  $\bar{\varphi}(y) = \varphi(y) / \sqrt{\bar{\Gamma}(1, \dots, 1)}$  and  $\psi(z)$  with

$\bar{\psi}(z) = \sqrt{\bar{\Gamma}(1, \dots, 1)}\psi(z)$ . In this way the ILE definition in Theorem 2 is satisfied, with  $\mathbf{c}_\Delta \leq \sqrt{\bar{\Gamma}(1, \dots, 1)\bar{\Gamma}(\mathbf{c}_{\Delta_1}^2, \dots, \mathbf{c}_{\Delta_P}^2)}$ . ■